

Chapter 4

Bayesian adaptation

In the previous chapter the usage of the MCE procedure in sparse data scenarios was investigated. In certain scenarios, however, one has a reasonable amount of non-task-specific training data available which can be used for training purposes. This chapter therefore investigates the usage of adaptation techniques that can use non-task-specific data as well as the limited task-specific data to create a better recognition system than could be attained though only using the task-specific data.

The standard *maximum a-posteriori* (MAP) adaptation technique is introduced and discussed. The MAP algorithm, however, makes assumptions about the form of the prior probability distribution used. This can be a problem when the prior is relatively complex. Two new adaptation algorithms are therefore proposed, a *gradient-based* MAP algorithm and an MCE-based adaptation algorithm. These adaptation algorithms make no assumptions about the form of the prior distribution used.

Section 4.1 introduces the basic Bayesian theory necessary for the algorithms in this chapter. The choice of prior distribution is discussed in Section 4.2. Section 4.3 reviews the maximum *a posteriori* estimation (MAP) algorithm as proposed by Gauvain and Lee [45]. An alternative *gradient-based* method of obtaining the MAP estimate

is described in Section 4.4. A Bayesian inspired *modification* to the MCE training procedure is then proposed in Section 4.5. Finally, the different methods discussed in this chapter are experimentally compared in Section 4.6.

4.1 Introduction

We can define a probabilistic model $P(\mathbf{X}|\theta)$ for any random process \mathbf{X} , in which a set of parameters θ determine its probability distribution. $P(\mathbf{X}|\theta)$ is called the *likelihood function*. We however want to infer θ in our probabilistic model, using data which we have obtained.

The result of Bayesian learning is a probability distribution $P(\theta|\mathbf{X})$ which expresses our beliefs of how likely individual parameters θ are, given the training data \mathbf{X} . $P(\theta|\mathbf{X})$ is called the *posterior distribution*. When classifying an unknown observation, the probability that the unknown observation was generated by the same process as that which generated the observations for a given class must be calculated. This is done by integrating the likelihood function with respect to the posterior distribution. Section 4.1.3 will present a more detailed discussion of this process.

Bayesian methods can be used for the inference of parameter values in a model, given the data. We can, however, also use Bayesian methods for the purpose of model comparison, where preferences are assigned to alternative models of differing complexity. David Mackay [72] focused primarily on the usage of Bayesian methods for the training and comparison of neural network models. Mackay used Bayesian methods to compare models of differing complexity and topology. Most people would include the above two uses of Bayesian methods in the data modeling process.

A further application of Bayesian methods is in the adaptation of existing models. An example of this in speech recognition is maximum a-posteriori (MAP) parameter estimation. MAP is a point estimate and one does not integrate with respect to the

posterior distribution. Many would therefore claim that MAP is not a true Bayesian method. If, however, one assumes that the posterior distribution is sufficiently peaked about the most probable θ (maximum *a-posteriori* probability), then the MAP procedure is a reasonable approximation of Bayesian learning. MAP can be used for several purposes, including parameter smoothing and adaptation. Parameter smoothing applies extra constraints to the model parameters so as to reduce the effect of insufficient training data. This can be achieved using MAP by incorporating vague heuristic information in the prior distribution. In adaptation, however, non-task-specific information is available and is used to determine the prior distribution used in MAP estimation.

The remainder of this section will summarize the pertinent Bayesian theory used in this chapter. For a more complete introductory text on Bayesian statistics, the reader is referred to Box and Tiao [13] and DeGroot [27]. The theory and discussions in this section will be biased towards speech recognition applications of Bayesian adaptation.

4.1.1 Bayes' theorem

Given a vector $\mathbf{y} = (y_1, \dots, y_n)$ of n observations, with probability distribution $P(\mathbf{y}|\theta)$, which depends on the k parameters $\theta^T = (\theta_1, \dots, \theta_k)$ with probability distribution $P(\theta)$, then given the observed data \mathbf{y} , the conditional distribution of θ is

$$P(\theta|\mathbf{y}) = \frac{P(\mathbf{y}|\theta)P(\theta)}{P(\mathbf{y})}. \quad (4.1)$$

The denominator in Eq. (4.1), $P(\mathbf{y})$, is a normalizing factor, which ensures that the integral of $P(\theta|\mathbf{y})$ is equal to one. It can be written as follows:

$$P(\mathbf{y}) = \int P(\mathbf{y}|\theta)P(\theta)d\theta. \quad (4.2)$$

Equation (4.1) is referred to as Bayes' theorem. The distribution $P(\theta)$, is called the

prior distribution and expresses what is known about the model parameters before any data is observed. The *posterior* distribution $P(\theta|\mathbf{y})$, tells us what is known about the model parameters, given that data has been observed. In what follows, the prior distribution and posterior distribution will sometimes simply be referred to as the “prior” and “posterior” respectively.

The distribution $P(\mathbf{y}|\theta)$ is often referred to as the data *likelihood* and can be written $L(\theta|\mathbf{y})$. Then in effect, $P(\mathbf{y}|\theta)$ is regarded as a function of \mathbf{y} and not of θ .

In many Bayesian methods, the normalizing constant is not necessary and Eq. (4.1) is written as

$$P(\theta|\mathbf{y}) \propto L(\theta|\mathbf{y})P(\theta). \quad (4.3)$$

4.1.2 Sequential nature of Bayes’ theorem

If we assume that the observations are independent, then we can write Bayes’ theorem as follows

$$P(\theta|\mathbf{y}) \propto P(\theta) \prod_{i=1}^n L(\theta|\mathbf{y}^{(i)}) = \left[P(\theta) \prod_{i=1}^k L(\theta|\mathbf{y}^{(i)}) \right] \prod_{i=k+1}^n L(\theta|\mathbf{y}^{(i)}) \quad (4.4)$$

$$\propto P(\theta|\mathbf{y}_1, \dots, \mathbf{y}_k) \prod_{i=k+1}^n L(\theta|\mathbf{y}^{(i)}) \quad (k < n). \quad (4.5)$$

Equation (4.4) is exactly the same as Eq. (4.5), except that $P(\theta|\mathbf{y}_1, \dots, \mathbf{y}_k)$, the posterior distribution of θ given $\mathbf{y}_1, \dots, \mathbf{y}_k$, now acts as the prior distribution for $\mathbf{y}_{k+1}, \dots, \mathbf{y}_n$.

Bayes’ theorem, therefore describes the process of learning as data becomes available. We can therefore, as Eq. (4.5) suggests, compute the posterior for a given set of data

and then use that posterior as a “prior” when more data becomes available. This result is of utmost importance to the methods proposed later in this chapter.

4.1.3 Bayesian learning and prediction

The result of Bayesian learning is a probability distribution (posterior) which expresses our beliefs of how likely individual parameters values are. This is crucial, as it allows learning to be performed using probability theory.

In a Bayesian approach to HMM parameter estimation and recognition, the objective is to find a predictive distribution ($P(\mathbf{x}|\mathbf{X}) = \int P(\mathbf{x}|\theta)P(\theta|\mathbf{X})d\theta$) for an unknown utterance, given the utterance observations, as well as the training observations. Let the observation sequence for the i th example be written as \mathbf{O}_i . For n training examples $\mathbf{O} = (\mathbf{O}_1, \dots, \mathbf{O}_n)$, Bayes’ theorem (Eq. (4.1)) can be written as

$$\begin{aligned}
 P(\theta|\mathbf{O}) &= \frac{P(\mathbf{O}|\theta)P(\theta)}{P(\mathbf{O})} \\
 &\propto P(\mathbf{O}|\theta)P(\theta).
 \end{aligned}
 \tag{4.6}$$

Assuming independence of the observations we can write the likelihood as follows:

$$P(\mathbf{O}_i|\theta) = \prod_{j=1}^{n_u} P(\mathbf{o}_{ij}|\theta),
 \tag{4.7}$$

where n_u is the number of observations in the training utterance \mathbf{O}_i . In a Bayesian framework, when we wish to classify an unknown input, we need to calculate the following probability,

$$P(\mathbf{O}_{unknown}|\mathbf{O}_1^{(i)}, \dots, \mathbf{O}_n^{(i)}) = \int P(\mathbf{O}_{unknown}|\theta)P(\theta|\mathbf{O}_1^{(i)}, \dots, \mathbf{O}_n^{(i)})d\theta,
 \tag{4.8}$$

where i is the class and $\mathbf{O}_{unknown}$ is the unknown observation sequence. The classifier decision is the class resulting in the highest value of Eq. (4.8), i.e.

$$C(\mathbf{O}_{unknown}) = i \quad \text{where} \quad i = \underset{j}{\operatorname{argmax}} P(\mathbf{O}_{unknown} | \mathbf{O}_1^{(j)}, \dots, \mathbf{O}_n^{(j)}), \quad (4.9)$$

where $C(\mathbf{O}_{unknown})$ is the classifier's decision for the unknown observation.

4.1.4 Maximum *a-posteriori* probability estimate

Assuming the posterior is sufficiently peaked around the most probable point (θ_{MAP}), we can approximate Eq. (4.8) as

$$P(\mathbf{O}_{unknown} | \mathbf{O}_1, \dots, \mathbf{O}_n) \approx P(\mathbf{O}_{unknown} | \theta_{MAP}), \quad (4.10)$$

where θ_{MAP} is the Maximum *a-posteriori* (MAP) estimate of the parameter vector θ .

The MAP point is the set of parameters that maximize Eq. (4.6), or

$$\begin{aligned} \theta_{MAP} &= \underset{\theta}{\operatorname{argmax}} P(\theta | \mathbf{O}_1, \dots, \mathbf{O}_n) \\ &= \underset{\theta}{\operatorname{argmax}} \left[P(\mathbf{O}_1, \dots, \mathbf{O}_n | \theta) P(\theta) \right]. \end{aligned} \quad (4.11)$$

If we had no prior knowledge about θ , then we would choose a non-informative (improper) prior to be used in Eq. (4.11), i.e., $P(\theta) = \text{constant}$. Equation (4.11) then reduces to the normal maximum likelihood (ML) formulation.

4.1.5 MAP adaptation in speech recognition

This section presents a concise literature survey of the usage of the MAP procedure for HMM parameter estimation in speech recognition.

Lee *et al.* [67] introduced a MAP algorithm, where the parameters of multivariate Gaussian state observation densities of HMM models were adapted for speaker adaptation. They showed that for an alpha digit task, with only a small amount of speaker specific data, their MAP estimated HMM gave better results than an ML estimated HMM.

Gauvain and Lee [43] extended the MAP formulation for HMMs to handle parameters of *mixtures* of Gaussian densities. It was shown that MAP estimation can be used for parameter smoothing, speaker adaptation, speaker clustering and corrective training. Gauvain and Lee [45] later presented a theoretical framework for MAP estimation for HMMs with Gaussian mixture state densities, where they proposed using an expectation-maximization (EM) approach to finding the MAP estimate. In this work, the MAP formulation was also extended to include the estimation of the transition probabilities and initial state probabilities. The application of their MAP algorithm to corrective training and parameter smoothing [44] and speaker adaptation [66] were also reported.

Huo *et al.* [53, 52] studied the usage of MAP estimation for semi-continuous (or tied mixture) HMMs. Zavaliagkos *et al.* [120] also investigated using various degrees of parameter tying, so as to force MAP to adapt parameters for which adaptation data is not available.

The MAP estimation algorithm can be used for various purposes, including parameter smoothing [44], speaker adaptation [67, 66, 32], dialect adaptation [41] and cross-language adaptation [85, 84]. In this chapter, we are primarily interested in the usage of the MAP estimation algorithm for adaptation purposes as opposed to applications

such as parameter smoothing and corrective training. Model adaptation is a process for adjusting seed models to create more specialized models using a small amount of adaptation data. Figure 4.1 illustrates an abstraction of the MAP adaptation algorithm as it would typically be used in speaker adaptation. The process is very similar for other adaptation applications of the algorithm, such as for cross-language adaptation.

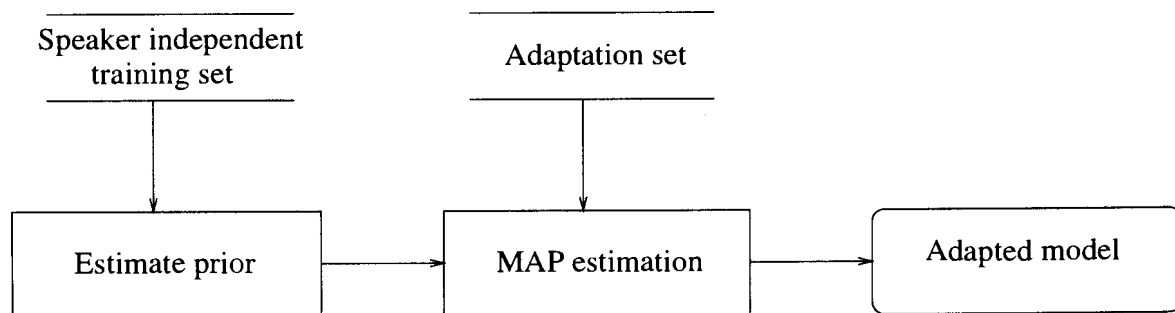


Figure 4.1: Abstraction of the typical MAP speaker adaptation process

4.2 Prior distributions

In this section, the choice of the prior density family is discussed. The prior distribution is an important part of any Bayesian method, as it expresses our knowledge about the distributions prior to any data being observed. It is especially important when there is little data available.

Prior distributions must therefore be chosen for the following HMM parameters:

- transition probabilities a_{ij} (Eq. (2.3)),
- mixture weights C_{jk} (Eq. (2.5)),
- Gaussian means μ_{jk} (Eq. (2.5)), and
- Gaussian variances Σ_{jk} (Eq. (2.5)).

Conjugate prior (see Appendix A) distributions have been chosen in this work, not only because of the convenient relationship between prior and posterior, but also due to their usage in the literature. If the prior distribution of θ belongs to a conjugate family of distributions, then for a likelihood of a specific form, the posterior distribution of θ must also belong to the same family as that of the prior distribution. Note that with most of the approaches developed in this and later chapters it is not necessary nor even convenient to have conjugate prior distributions. It is, however, essential to use a conjugate prior for the application of the EM algorithm to the MAP estimation problem (Section 4.3).

The prior distribution for all the parameters of the Gaussian mixtures of the HMMs is chosen to be a normal-Wishart distribution, where the conditional of μ (mean) and R (precision matrix, with $R = \Sigma^{-1}$) are given as follows: The conditional distribution of μ (when $R = r$), is a multivariate normal distribution with mean vector m and precision matrix νr for $m \in \mathbb{R}^D$ and $\nu > 0$. The marginal distribution of R is a Wishart distribution with α degrees of freedom and precision matrix τ . The priors for the transition probabilities and Gaussian mixture weights are chosen to be Dirichlet distributions. Table 4.1 summarizes the HMM parameters and their prior distributions.

Table 4.1: Summary of the HMM parameters and chosen prior distributions. A Wishart distribution is represented with \mathcal{W} and a Dirichlet with \mathcal{D}

Parameter	Prior distribution	Prior parameters
$\mu_{jk} R_{jk} = r_{jk}$	$\mathcal{N}(m_{jk}, (\nu_{jk} r_{jk})^{-1/2})$	<i>mean</i> : m_{jk} , <i>precision</i> : $\nu_{jk} r_{jk}$
$R_{jk} (\Sigma_{jk}^{-1})$	$\mathcal{W}(n_{jk}, \tau_{jk})$	<i>degrees of freedom</i> : n_{jk} , <i>precision</i> : τ_{jk}
c_j	$\mathcal{D}(\delta_j)$	δ_j
a_i	$\mathcal{D}(\alpha_i)$	α_i

The prior of the parameters (μ_{jk} and r_{jk}) of the Gaussian mixture component k of state j can therefore be written (from Eqs. (A.5) and (A.1)) as [27]:

$$g_{\text{gaussian}}(r_{jk}, \mu_{jk} | n_{jk}, \nu_{jk}, m_{jk}, \tau_{jk}) = (2\pi)^{-\frac{D}{2}} |\nu_{jk} r_{jk}|^{\frac{1}{2}} e^{-\frac{1}{2} \nu_{jk} (\mu_{jk} - m_{jk})^T r_{jk} (\mu_{jk} - m_{jk})} \\ c |\tau_{jk}|^{\frac{n_{jk}}{2}} |r_{jk}|^{\frac{(n-D-1)}{2}} e^{-\frac{1}{2} \text{tr}(r_{jk} \tau_{jk})}, \quad (4.12)$$

where $(n_{jk}, \nu_{jk}, m_{jk}, r_{jk})$ are the prior distribution parameters. The value c is a normalizing constant which ensures that the integral of the prior is equal to one. Assuming a diagonal covariance (or precision) matrix, the log of the Gaussian distribution prior (Eq. (4.12)) can be written as

$$\log g_{\text{gaussian}}(\Sigma_{jk}, \mu_{jk} | n_{jk}, \nu_{jk}, m_{jk}, \tau_{jk}) = -\frac{D}{2} \log(2\pi \nu_{jk}) + \log(c) \\ - \frac{n_{jk} - D}{2} \log\left(\prod_{l=1}^D \sigma_{jkl}^2\right) - \frac{1}{2} \nu_{jk} \sum_{l=1}^D \sigma_{jkl}^{-2} (\mu_{jkl} - m_{jkl})^2 \\ + \frac{n_{jk}}{2} \cdot \log\left(\prod_{l=1}^D \tau_{jkl}\right) - \frac{1}{2} \sum_{l=1}^D \frac{\tau_{jkl}}{\sigma_{jkl}^2}. \quad (4.13)$$

The prior for the mixture weights can be written as:

$$g_{\text{weight}}(c_j | \delta_j) = \frac{\Gamma(\delta_{j1} + \delta_{j2} + \dots + \delta_{jM})}{\Gamma(\delta_{j1})\Gamma(\delta_{j2})\dots\Gamma(\delta_{jM})} c_{j1}^{\delta_{j1}-1} c_{j2}^{\delta_{j2}-1} \dots c_{jM}^{\delta_{jM}-1}, \quad (4.14)$$

where δ_j is the parameter vector associated with the mixture weight priors for state j .

The prior for the transition probabilities can be written as:

$$g_{\text{trans}}(a_i | \alpha_i) = \frac{\Gamma(\alpha_{i1} + \alpha_{i2} + \dots + \alpha_{iN})}{\Gamma(\alpha_{i1})\Gamma(\alpha_{i2})\dots\Gamma(\alpha_{iN})} a_{i1}^{\alpha_{i1}-1} a_{i2}^{\alpha_{i2}-1} \dots a_{ik}^{\alpha_{ik}-1}, \quad (4.15)$$

where α_i is the prior parameter vector associated with the transition probabilities a_i from state i and N being the number of states in the HMM.

If we assume independence for parameters of the Gaussians, mixture weights and transition probabilities, the joint prior density $g(\theta)$ for all the parameters of the HMM can be written as the product of the prior p.d.f.'s defined in Eqs. (4.12), (4.14) and (4.15), i.e.

$$g(\theta) = \prod_{i=1}^N \left(g_{trans}(a_i|\alpha_i) g_{weight}(c_i|\delta_i) \prod_{k=1}^M g_{gaussian}(\Sigma_{ik}, \mu_{ik}|n_{ik}, \nu_{ik}, m_{ik}, r_{ik}) \right), \quad (4.16)$$

where M is the number of Gaussian mixture components per state. The joint prior density for all HMMs can, assuming independence, be written as the product of the joint prior distributions (Eq. (4.16)) of the individual HMMs.

4.3 Expectation-Maximization MAP

Gauvain and Lee [45] introduced a method of estimating the MAP point for all the HMM parameters as defined in Table 4.1 using an expectation maximization (EM) approach. Their method and related theory will be briefly reviewed in this section. This algorithm will, for the rest of this thesis, often be referred to as simply the MAP algorithm.

MAP estimation is relatively simple if the family of p.d.f.'s $P(\cdot|\theta), \theta \in \Theta$ possesses a sufficient statistic of fixed dimension for the parameter θ we wish to estimate. A sufficient statistic T only exists if $P(\mathbf{O}|\theta)$ can be factored [27] as follows for all values of \mathbf{O} and θ :

$$P(\mathbf{O}|\theta) = u(\mathbf{O})v(T(\mathbf{O}), \theta). \quad (4.17)$$

Here, the function u is positive and does not depend on θ , while the function v is non-negative and depends on the data (observations) \mathbf{O} only through $T(\mathbf{O})$. For hidden

Markov models, due to the underlying hidden process, a sufficient statistic of fixed dimension does not exist. This is known as an “incomplete data” problem. When using hidden Markov models, the true state sequence $\mathbf{q} = \{q_1, \dots, q_T\}$ and the true sequence of associated mixture components $\mathbf{l} = \{l_1, \dots, l_t\}$ are not observed or known. As a result, the observed data \mathbf{O} is not sufficient to be able to directly estimate the HMM parameters (*incomplete data*). For HMMs, the *complete data* \mathbf{x} is the combination of the observations \mathbf{O} , state sequence \mathbf{q} and mixture component sequence \mathbf{l} , i.e. $\mathbf{x} = (\mathbf{O}, \mathbf{q}, \mathbf{l})$.

The term *incomplete data* (\mathbf{y}) implies two sample spaces \mathcal{X} and \mathcal{Y} , and a many-to-one mapping from \mathcal{X} to \mathcal{Y} . Given $\mathbf{x} \in \mathcal{X}$ and $\mathbf{y} \in \mathcal{Y}$, where \mathbf{x} is not observed directly but indirectly through the observed data \mathbf{y} , a mapping $\mathbf{x} \rightarrow \mathbf{y}(\mathbf{x})$ exists from \mathcal{X} to \mathcal{Y} , with \mathbf{x} known only to be in $\mathcal{X}(\mathbf{x})$, the subset of \mathcal{X} determined by the equation $\mathbf{y} = \mathbf{y}(\mathbf{x})$. We refer to \mathbf{x} as the “complete data”. Note that there are many *complete-data* specifications $f(\mathbf{x}|\theta)$ that will generate a given *incomplete-data* specification $g(\mathbf{y}|\theta)$.

The expectation-maximization(EM) algorithm [5, 28, 71] discussed next is an iterative procedure for approximating ML estimates in cases involving *incomplete data*.

4.3.1 Expectation maximization

Expectation maximization, as described in this section, is typically used for maximum likelihood (ML) estimation. Expectation maximization re-estimation is based on an auxiliary function $Q(\theta, \bar{\theta})$ defined in terms of the current parameter set θ and the new parameter set $\bar{\theta}$. The auxiliary function is defined as the expectation of the *complete-data* log-likelihood $\log[f(\mathbf{x}|\bar{\theta})]$ given the observed incomplete data \mathbf{y} and the current parameter set θ , i.e.

$$Q(\theta, \bar{\theta}) = E[\log f(\mathbf{x}|\bar{\theta})|\mathbf{y}, \theta] \quad (4.18)$$

which exists for all pairs $(\theta, \bar{\theta})$. The EM iteration is defined as follows:

Expectation (E-step): Compute $Q(\theta, \bar{\theta})$.

Maximization (M-step): Choose $\bar{\theta}$ which maximizes $Q(\theta, \bar{\theta})$. We simply differentiate $Q(\theta, \bar{\theta})$ with respect to $\bar{\theta}$ and find a maximum (i.e. solve $\frac{\partial Q(\theta, \bar{\theta})}{\partial \bar{\theta}} = 0$).

For an HMM, the *complete-data* likelihood is the joint likelihood of \mathbf{O} , $\mathbf{q} = \{q_1, \dots, q_T\}$ the unobserved state sequence and $\mathbf{l} = \{l_1, \dots, l_t\}$ the unobserved sequence of associated mixture components ($\mathbf{y} = (\mathbf{O}, \mathbf{q}, \mathbf{l})$).

For an HMM, the auxiliary function $Q(\theta, \bar{\theta})$ in Eq. (4.18) takes the following form [71]:

$$Q(\theta, \bar{\theta}) = \sum_{\mathbf{q}} E[f(\mathbf{O}, \mathbf{q}|\theta) \log f(\mathbf{O}, \mathbf{q}|\bar{\theta})] \quad (4.19)$$

where $f(\mathbf{O}, \mathbf{q}|\theta)$ is the probability of observing the data \mathbf{O} for state sequence \mathbf{q} given the parameter set θ . The utility of $Q(\theta, \bar{\theta})$ stems in part from the fact that if $Q(\theta, \bar{\theta}) > Q(\theta, \theta)$, then $f(\mathbf{O}|\bar{\theta}) > f(\mathbf{O}|\theta)$. This property is shown in Eqs. (4.20) and (4.21).

$$Q(\theta, \bar{\theta}) - Q(\theta, \theta) = \sum_{\mathbf{q}} E \left[f(\mathbf{O}, \mathbf{q}|\theta) \log \left(\frac{f(\mathbf{O}, \mathbf{q}|\bar{\theta})}{f(\mathbf{O}, \mathbf{q}|\theta)} \right) \right] \quad (4.20)$$

Using the inequality $\log(x) \leq x - 1$, we can write Eq. (4.20) as:

$$\begin{aligned} Q(\theta, \bar{\theta}) - Q(\theta, \theta) &\leq \sum_{\mathbf{q}} E \left[f(\mathbf{O}, \mathbf{q}|\theta) \left(\frac{f(\mathbf{O}, \mathbf{q}|\bar{\theta})}{f(\mathbf{O}, \mathbf{q}|\theta)} - 1 \right) \right] \\ &\leq f(\mathbf{O}|\bar{\theta}) - f(\mathbf{O}|\theta). \end{aligned} \quad (4.21)$$

Hence, $Q(\theta, \bar{\theta}) > Q(\theta, \theta)$ implies that $f(\mathbf{O}|\bar{\theta}) > f(\mathbf{O}|\theta)$. Therefore, when $Q(\theta, \bar{\theta})$ as defined in Eq. (4.19) is maximized in the maximization step, then $f(\mathbf{O}|\bar{\theta})$ will also

be maximized. The EM algorithm described here will be used in the next section to iteratively obtain the MAP estimate.

4.3.2 MAP estimate for HMMs

Although the expectation maximization procedure (Eq. (4.19)) is typically used to find the ML estimate, we wish to use it to obtain the MAP estimate. Dempster *et al.* [28] pointed out that the posterior mode (MAP point) can also be estimated using the EM algorithm by maximizing $Q(\theta, \bar{\theta}) + \log[g(\theta)]$ at the M-step of each EM iteration. It was also shown that $\log[f(\mathbf{x}|\theta)] + \log[g(\theta)]$ increases at each EM iteration and an expression for the rate of convergence was provided.

It is relatively straightforward to show [45] that the auxiliary function of the EM algorithm applied to the MAP problem can be decomposed into the sum of three auxiliary functions, $Q_\pi(\bar{\pi}, \theta)$, $Q_A(\bar{A}, \theta)$ and $Q_\Upsilon(\Upsilon, \theta)$, with $\Upsilon = (\bar{c}, \bar{\mu}, \bar{\Sigma})$. The three auxiliary functions are maximized independently, resulting in the following re-estimation formulas:

$$\bar{\pi}_i = \frac{(\alpha_{0i} - 1) + \gamma_0(i)}{\sum_{j=1}^N (\alpha_{0j} - 1) + \sum_{j=1}^N \gamma_0(j)} \quad (4.22)$$

$$\bar{a}_{ij} = \frac{(\alpha_{ij} - 1) + \sum_{t=1}^T \xi_t(i, j)}{\sum_{j=1}^N (\alpha_{ij} - 1) + \sum_{j=1}^N \sum_{t=1}^T \xi_t(i, j)} \quad (4.23)$$

$$\bar{c}_{ij} = \frac{(\delta_{ik} - 1) + \sum_{t=1}^T w_{ikt}}{\sum_{k=1}^M (\delta_{ik} - 1) + \sum_{k=1}^M \sum_{t=1}^T w_{ikt}} \quad (4.24)$$

$$\bar{\mu}_{ik} = \frac{\nu_{ik} m_{ik} + \sum_{t=1}^T w_{ikt} \mathbf{o}_t}{\nu_{ik} + \sum_{t=1}^T w_{ikt}} \quad (4.25)$$

$$\bar{\Sigma}_{ik} = \frac{\tau_{ik} + \sum_{t=1}^T w_{ikt} (\mathbf{o}_t - \mu_{ik})(\mathbf{o}_t - \mu_{ik})^T + \nu_{ik} (m_{ik} - \mu_{ik})(m_{ik} - \mu_{ik})^T}{(n_{ik} - D) + \sum_{t=1}^T w_{ikt}}, \quad (4.26)$$

where w_{ikt} is defined as

$$w_{ikt} = \gamma_t(i) \frac{c_{ik} \mathcal{N}(\mathbf{o}_t | \mu_{ik}, \Sigma_{ik})}{\sum_{l=1}^M c_{il} \mathcal{N}(\mathbf{o}_t | \mu_{il}, \Sigma_{il})}. \quad (4.27)$$

The values $\gamma_t(i)$, $\gamma_t(i, j)$ and $\xi_t(i, j)$ are obtained using the forward-backward algorithm [90, p. 334].

Segmental MAP estimate

Instead of maximizing $f(\theta | \mathbf{O})$ (using the forward-backward algorithm), we can, using Viterbi alignment, maximize $f(\theta, \mathbf{q} | \mathbf{O})$, the joint posterior density of the parameter set θ and the state sequence \mathbf{q} . It can easily be shown, as for the segmental k -means algorithm [61], that alternate maximization over \mathbf{q} and θ results in a non-decreasing $f(\theta^{(n)}, \mathbf{q}^{(n)} | \mathbf{O})$, with

$$\mathbf{q}^{(n+1)} = \underset{\mathbf{q}}{\operatorname{argmax}} f(\mathbf{O}, \mathbf{q} | \theta^{(n)}) \quad (4.28)$$

$$\theta^{(n+1)} = \underset{\theta}{\operatorname{argmax}} f(\mathbf{O}, \mathbf{q}^{(n+1)} | \theta). \quad (4.29)$$

The most likely state sequence in Eq. (4.28) is obtained using the Viterbi algorithm. The EM algorithm can once again be used to perform the maximization in Eq. (4.29). The re-estimation formulae (4.22) to (4.26) used in the forward-backward MAP estimate are used here, however, the probabilities $\xi_t(i, j)$, $\gamma_t(i)$ and $\gamma_t(j, k)$ are obtained from the best state sequence (and not using the forward-backward algorithm) as follows:

$$\xi_t(i, j) = \delta(\mathbf{q}_t - i)\delta(\mathbf{q}_{t+1} - j), \quad (4.30)$$

$$\gamma_t(i) = \delta(\mathbf{q}_t - i), \quad (4.31)$$

$$\text{and } \gamma_t(j, k) = \delta(\mathbf{q}_t - j) \frac{c_{jk} \mathcal{N}(\mathbf{o}_t, \mu_{jk}, \Sigma_{jk})}{b_j(\mathbf{o}_t)}, \quad (4.32)$$

with δ being the Kronecker delta function.

The segmental MAP approach described above has been implemented and used in this work. An embedded version of the MAP algorithm has also been implemented and used, which uses the trellis search algorithm (Section 2.1.5) to obtain the best state and HMM sequence for a given utterance. The embedded MAP algorithm uses a search to automatically segment and label the acoustic units, and therefore does not rely on a labelled database being available. This is essential for databases, such as TIDIGITS, which are not manually labelled either at phoneme or word levels.

4.3.3 Prior density estimation

In Section 4.2 the form of the prior distributions and the distribution parameters were described. In this section, the estimation of the prior parameters will be discussed. The prior distribution should, in a true Bayesian approach, incorporate *a-priori* knowledge of the parameters we are attempting to estimate.

Two methods of estimating the prior parameters will be discussed here, namely empirical Bayes and a simpler method typically used in MAP adaptation. The empirical Bayes method is not used, but is included so that the two methods can be compared. When estimating the prior parameters we will only be using the prior data. The adaptation data cannot be used in any way when determining the prior parameters.

Empirical Bayes methods

In Empirical Bayesian methods [73, 97] the prior distribution $g(\theta)$ is estimated by finding a distribution function g that satisfies the relationship

$$h(\mathbf{O}) = \int f(\mathbf{O}|\theta)g(\theta)d\theta, \quad (4.33)$$

where $h(\mathbf{O})$ is the probability distribution of the data ($P(\mathbf{y})$ in Eq. (4.2)).

It is, however, necessary to ensure that the distribution $g(\theta)$ obtained is unique. The search for $g(\theta)$ could be an almost impossible task if we don't choose $g(\theta)$ to be part of a given parametric family $g(\theta; \alpha, \beta, \dots)$ of distributions, where α, β, \dots are the unknown prior parameters.

In this case, we do not know $h(\mathbf{O})$ exactly, but estimate an empirical distribution function $h_n(\mathbf{O})$ obtained from a sample of n observations on the random variable whose distribution function is $h(\mathbf{O})$. We therefore have the approximate relationship

$$h_n(\mathbf{O}) \approx \int f(\mathbf{O}|\theta)g(\theta)d\theta. \quad (4.34)$$

There are several methods for solving Eq. (4.34), i.e. finding $g(\theta)$, including maximum-likelihood (ML) and the method of moments. The solution of Eq. (4.34) is not a trivial task.

Though the empirical Bayes approach is not implemented or used for the MAP estimation in this work, it is important to understand the empirical Bayes approach and we therefore digress at this point to study a simple example taken from [73].

Given observed data \mathbf{x} , a model $f(\mathbf{O}|\theta) = f(\mathbf{x}|M)$, being a simple Normal distribution with mean M and standard deviation of 1 and a prior $g(M)$ that is a Normal

distribution with mean μ and variance σ^2 , i.e.

$$f(\mathbf{x}|M) = \mathcal{N}(M, 1),$$

$$g(M) = \mathcal{N}(\mu, \sigma^2),$$

then $h_n(\mathbf{O}) = f_G(\mathbf{x}) = \int f(x|M)g(M; \mu, \sigma)dM = \mathcal{N}(\mu, 1 + \sigma^2)$. The maximum likelihood estimates of the prior parameters $\tilde{\mu}$ and $\tilde{\sigma}^2$ are [73, p. 53]

$$\tilde{\mu} = \bar{\mathbf{x}}$$

$$\tilde{\sigma}^2 = \max(0, (\frac{n-1}{n})s^2 - 1),$$

where $\bar{\mathbf{x}}$ is the sample mean and s^2 is the sample variance of past observations. Note that the estimate of the prior is not directly affected by the number of samples used, as would be the case with the posterior distribution.

Mode of posterior

A significantly simpler method [45] of estimating the prior distribution from a given set of data would be to maximize the joint likelihood of \mathbf{O} and θ , i.e. $f(\mathbf{O}, \theta|\phi)$, over θ and ϕ . Here ϕ is the parameter vector of the prior distribution.

Starting with an initial estimate of ϕ , and iteratively using alternate maximization over θ and ϕ , i.e.

$$\theta^{(n)} = \underset{\theta}{\operatorname{argmax}}[f(\mathbf{O}, \theta|\phi^{(n)})] \quad (4.35)$$

$$\phi^{(n)} = \underset{\phi}{\operatorname{argmax}}[g(\theta^{(n)}|\phi)], \quad (4.36)$$

we can estimate the prior parameters ϕ and model parameters θ which maximize $f(\mathbf{O}, \theta | \phi)$. The solution of Eq. (4.35) is the mode of the posterior (MAP estimate) for the current prior parameter set. The solution of Eq. (4.36) is the ML estimate of the prior parameters based on the current HMM parameters.

Solving Eq. (4.36) poses two problems:

- ML estimation thereof is not simple as a result of the chosen prior distribution of Section 4.2.
- More prior density parameters must be estimated than for the HMM itself. This is called overparameterization and is a problem.

The overparameterization problem can be overcome by adding certain constraints to the prior parameters, so as to reduce the number of prior parameters to be estimated. The prior family is limited to the posterior density family of the complete data model when no prior information is available. It is then easily shown [45] that the following constraints can be imposed:

$$\delta_{ik} = \nu_{ik} + 1 \quad (4.37)$$

$$n_{ik} = \nu_{ik} + D, \quad (4.38)$$

solving the overparameterization problem to some extent. The parameter δ_{ik} is the mixture weight prior parameter in Eq. (4.14), with n_{ik} and ν_{jk} being prior parameters in Eq. (4.12) for state i and mixture component k .

Note that using Eqs. (4.35) and (4.36) when no prior information is available will result in θ being the mode of the likelihood function (ML estimate) and we therefore set the mode of the prior to be equal to the parameters of a given HMM. Given that the prior

family has been chosen to be the same as that of the complete-data likelihood it makes sense that the mode of the prior will be the ML point estimate.

The prior parameters m_{jk} , τ_{jk} and α_{ij} are therefore directly estimated from the ML HMM models, while δ_{ik} and n_{ik} are obtained using Eqs. (4.37) and (4.38). The parameters n_{ik} and ν_{ik} do not determine the mode of the prior distribution, but determine the degree to which the prior is peaked about its mode. The parameter ν_{ik} is therefore a parameter chosen by the user and is typically chosen as a global parameter.

The value of the parameter ν should therefore incorporate *a priori* knowledge about the suitability of the ML model for the task at hand. If the data used to obtain the prior was from the same task as that used to obtain the final MAP point, then we would expect to use a relatively large value for ν . However, if there was a large mismatch between the prior data and data used to obtain the MAP point, then we are not that sure of the prior and a smaller ν would therefore be chosen. An example of where this might occur is in cross-language adaptation, where data from a given language is used to create a prior and the target language's data, along with the prior is used to obtain the MAP point estimate. The influence of ν on the performance of the MAP algorithm will be experimentally be determined in Section 4.6.

Let us investigate the usage of the “mode of posterior” method to determine the prior parameters for the example discussed earlier. The forms of the prior and likelihood functions are identical and therefore there is no problem of overparameterization. If we were to choose the prior mode directly from the ML point, then one would choose $\mu = \bar{x}$ which is exactly the same as that obtained using the empirical Bayes method. The value of σ^2 is, however, still unknown and it determines the degree to which the prior is peaked about its mode. We could arbitrarily fix it, as we have done with ν above, but we could also choose it to be proportional to the sample variance s^2 .

Discussion on prior estimation

If the data samples used were taken from the same source, it would have been better to assume a non-informative (constant) prior distribution and used all the examples to determine the posterior. Note that from Section 4.1.2 we would deduce that if one needed a prior, then the sequential nature of Bayes' theorem tells us that the posterior for the observed data would be the appropriate prior for any subsequently observed data.

If we once again look at the example used in the previous two sections, it is relatively simple to show that the posterior $f(M|\mathbf{x})$ for the given data \mathbf{x} assuming a non-informative prior would be

$$f(M|\mathbf{x}) = \mathcal{N}(\mu, \sigma^2), \quad (4.39)$$

where $\mu = \bar{\mathbf{x}}$ and $\sigma^2 = \frac{1}{N}s^2$. As expected, the posterior becomes more peaked around its mode μ as more data is observed. Following Eq. (4.5), we should use this posterior, as the prior for any new data that is observed.

Neither of the two methods described previously result in a prior which becomes more peaked as the amount of observed data increases. This is a potential problem, as it does not account for the fact that some HMMs (or states or mixtures) will have been observed more often than others. However, the two methods could potentially help the algorithm when there is a reasonable mismatch between prior data and the task specific data.

The sequential nature of Bayes can be used to determine the MAP estimate by using the posterior of the prior data as the prior distribution (Eq. (4.5)). However, there will typically be more prior data than adaptation data, and the prior distribution will therefore tend to dominate in the calculation of the posterior using Eq. (4.5). Any

reasonable mismatch between the prior data and task-specific data will also tend to be a problem. These problems can, however, be addressed by simply weighing the prior distribution (posterior of previously observed data) with a value which is a function of the mismatch (*a-priori* knowledge/belief) and the amount of prior and adaptation data.

4.4 Gradient based MAP

The MAP estimation method proposed by Gauvain and Lee [45], assumes that the prior used is of a specific form. This is potentially a limiting factor in the performance of that MAP algorithm. In this section a gradient based MAP estimation algorithm is developed which does not make assumptions about the form of the prior distribution. This algorithm will, so as to prevent confusion, be referred to as the GMAP algorithm.

The above statement is not entirely true, as a prior of fixed form is used. It is, however, a non-informative prior which is used in the calculation of the new prior, which in turn is then used in the adaptation process. Though, if true prior knowledge about the model or system is available, it can be expressed through this parametric prior.

In the proposed algorithm, the prior will not be estimated at all, but will be implicitly included in the update procedure. It will, however, be far more computationally expensive than the regular MAP algorithm. The improved performance will, however, offset the extra computational difficulties for certain tasks. This adaptation algorithm will probably not find a place in rapid adaptation needed in some speaker adaptation tasks. It should, however, be more than useful in tasks such as cross-language adaptation and some speaker adaptation tasks where training time is not critical.

In the discussion in section 4.3.3, I explained under which circumstances one can use the posterior (or weighted version thereof) of the given “prior” data as the prior distribution. There is no problem with this, except when the same data is used to estimate

the prior and the MAP point of the posterior distribution and should not be done. In Section 4.1.2, it was pointed out that the posterior of a given set of data can be used as the prior for another independent set of data when calculating the posterior of the union of the two sets.

We can rewrite Eq. (4.5) as

$$P(\theta|\mathbf{O}) \propto \left[P(\theta) \prod_{i=1}^k L(\theta|\mathbf{O}_i) \right] \prod_{i=k+1}^n L(\theta|\mathbf{O}_i) \quad (k < n). \quad (4.40)$$

The first part in square brackets is the posterior of a set of data $(1 \dots k)$ and is used as the prior for the remainder of the data. The posterior of this reference or “prior” set, used as the prior in the adaptation framework will tend to dominate Eq. (4.40) when there are more training examples than adaptation examples (i.e. when $k \gg n - k$). MAP adaptation is typically used in situations where this will occur. It therefore makes sense to weigh the “prior” in some way so as to ensure that it does not dominate. In our implementation, the weighting is done as follows

$$P(\theta|\mathbf{O}) \propto \left[P(\theta) \prod_{i=1}^k L(\theta|\mathbf{O}_i) \right]^\varphi \prod_{i=k+1}^n L(\theta|\mathbf{O}_i) \quad (k < n). \quad (4.41)$$

The value φ has the effect of flattening and widening the prior when $0 \leq \varphi < 1$ and making it more peaked around the mode when $\varphi > 1$. Figure 4.2 presents an example of a Normal distribution with a mean of zero and standard deviation of two, which has been raised to the value of $\varphi = 0.2$ and $\varphi = 2$.

It is convenient at this point to express the posterior in terms of an energy function, which will be optimized to find the MAP estimate. This will become increasingly relevant in Chapter 5.

We can write the posterior probability function as follows

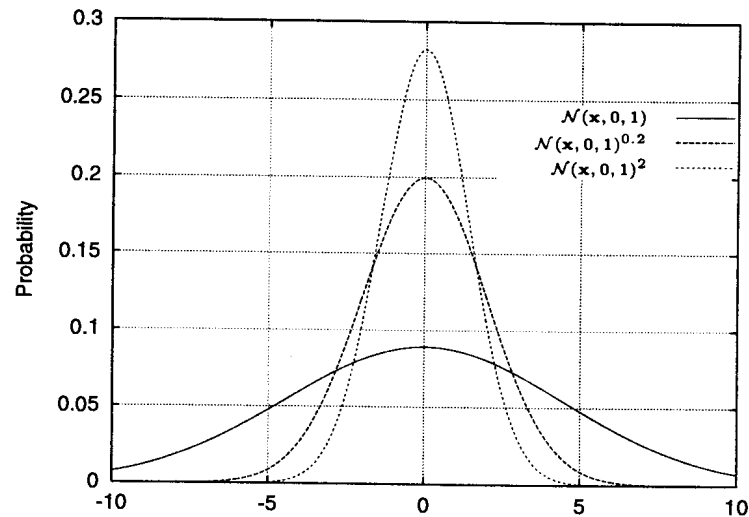


Figure 4.2: Example of the effects of raising a distribution to a power (assuming it is normalized)

$$P(\theta|\mathbf{O}) \propto e^{-E(\theta;\mathbf{O})} \quad (4.42)$$

where $E(\theta; \mathbf{O})$ is the “potential energy” function. Any probability function can be written in this way by defining $E(\theta; \mathbf{O}) = -\log P(\theta|\mathbf{O}) - \log(Q)$ for a given constant Q ($Q = 1$).

Writing the posterior (Eq. (4.6)) in the above form, and assuming independence of the observations, we get

$$E(\theta; \mathbf{O}) = -\log[P(\theta)] - \sum_{i=1}^n \log[PO_i|\theta]. \quad (4.43)$$

The above equation assumes that an informative prior $P(\theta)$ is available. We could again use a parametric prior estimated from the training data as done in the normal MAP approach (Section 4.3). Instead of estimating a prior of fixed form, the prior has been directly included into the posterior calculation (Eq. (4.41)), which can be written as

$$E(\theta; \mathbf{O}) = \varphi \left[-\log[P(\theta)] - \sum_{i=1}^k \log[P(\mathbf{O}_i|\theta)] \right] - \sum_{i=k+1}^n \log[P(\mathbf{O}_i|\theta)], \quad (4.44)$$

where examples $1, \dots, k$ are the prior (reference) set and examples $(k+1), \dots, n$ are the adaptation set.

The steepest descent algorithm [12] can now be used to iteratively estimate the MAP point on the posterior defined in Eq. (4.44), i.e.

$$\theta^{(i)}(n+1) = \theta^{(i)}(n) - \epsilon \frac{\partial E(\theta)}{\partial \theta^{(i)}}, \quad (4.45)$$

where ϵ is the learning rate or step size of the update. The steepest descent algorithm (Eq. (4.45)) is an unconstrained optimization technique and given that certain constraints must be maintained for HMMs, some modifications are required. The next section will investigate the usage of transformations to ensure that the constraints are maintained.

4.4.1 Parameter transformations

Instead of using a complicated constrained steepest descent algorithm, we can, as with minimum classification error (MCE) training (Section 3.2), use transformations to maintain the above constraints during parameter adaptation.

It is necessary to maintain the original HMM constraints, namely

1. $\sum_j a_{ij} = 1$ and $a_{ij} \geq 0$,
2. $\sum_k c_{jk} = 1$ and $c_{jk} \geq 0$, and
3. $\sigma_{jkl} \geq 0$.

The standard parameter transformations given in Table 3.1 are repeated here in Table 4.2 for convenience. These transformations ensure that the unconstrained steepest descent algorithm can be used in the transformed parameter space.

Table 4.2: Parameter transformations used in MCE

Parameter	Transformed parameter	Forward transform	Reverse transform	
a_{ij}	\tilde{a}_{ij}	$\tilde{a}_{ij} = \ln(a_{ij})$	$a_{ij} = \frac{e^{\tilde{a}_{ij}}}{\sum_k e^{\tilde{a}_{ij}}}$	Transition probabilities
c_{jk}	\tilde{c}_{jk}	$\tilde{c}_{jk} = \ln(c_{jk})$	$c_{jk} = \frac{e^{\tilde{c}_{jk}}}{\sum_k e^{\tilde{c}_{jk}}}$	Mixture weights
μ_{jkl}	$\tilde{\mu}_{jkl}$	$\tilde{\mu}_{jkl} = \frac{\mu_{jkl}}{\sigma_{jkl}}$	$\mu_{jkl} = \sigma_{jkl} \tilde{\mu}_{jkl}$	Gaussian mean
σ_{jkl}	$\tilde{\sigma}_{jkl}$	$\tilde{\sigma}_{jkl} = \ln(\sigma_{jkl})$	$\sigma_{jkl} = e^{\tilde{\sigma}_{jkl}}$	Gaussian std. dev.

These transformations should be kept in mind when calculating both the prior and likelihood derivatives.

4.4.2 Parameter adaptation

The updates for the individual parameter types will now be derived using Eqs. (4.44) and (4.45). Note that the adaptation is done in the transformed parameter space and the parameters are then transformed back to the original parameter space. The derivation of the parameter updates are given in this section. Some of the derivatives will also be used later in Chapter 5.

Gaussian mixture means

The parameter update for the Gaussian mixture mean $\tilde{\mu}_{jkl}$ using the steepest descent algorithm of state j , mixture k and element l is

$$\tilde{\mu}_{jkl}^{(i)}(n+1) = \tilde{\mu}_{jkl}^{(i)}(n) - \epsilon \frac{\partial E(\tilde{\theta}_n)}{\partial \tilde{\mu}_{jkl}^{(i)}}, \quad (4.46)$$

where

$$\begin{aligned} \frac{\partial E(\tilde{\theta}_n)}{\partial \tilde{\mu}_{jkl}^{(i)}} &= -\varphi \frac{\partial \ln[P(\tilde{\theta})]}{\partial \tilde{\mu}_{jkl}^{(i)}} - \varphi \sum_{i=1}^k \frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{\mu}_{jkl}^{(i)}} \\ &\quad - \sum_{i=k+1}^n \frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{\mu}_{jkl}^{(i)}}. \end{aligned} \quad (4.47)$$

From Eq. (2.9) we have

$$\frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{\mu}_{jkl}^{(i)}} = \sum_{t=1}^T \delta(\bar{q}_t - j) \frac{\partial \ln(b_j^{(i)}(\mathbf{o}_t))}{\partial \tilde{\mu}_{jkl}^{(i)}}, \quad (4.48)$$

where $\delta()$ is the Kronecker delta function and $b_j^{(i)}(\mathbf{o}_t)$ is the observation probability (Eq. (2.5)) of state j of HMM i at time t . Assuming a diagonal covariance matrix, we get

$$\begin{aligned} \frac{\partial \ln(b_j^{(i)}(\mathbf{o}_t))}{\partial \tilde{\mu}_{jkl}^{(i)}} &= c_{jk}^{(i)} (2\pi)^{-d/2} |\Sigma_{jk}^{(i)}|^{-1/2} (b_j^{(i)}(\mathbf{o}_t))^{-1} \left(\frac{\mathbf{o}_{tl}}{\sigma_{jkl}^{(i)}} - \tilde{\mu}_{jkl}^{(i)} \right) \\ &\quad e^{-\frac{1}{2} \sum_{l=1}^D \left(\frac{\mathbf{o}_{tl}}{\sigma_{jkl}^{(i)}} - \tilde{\mu}_{jkl}^{(i)} \right)^2} \\ &= \frac{c_{jk} \mathcal{N}(\mathbf{o}_t; \mu_{jk}, \Sigma_{jk})}{b_j^{(i)}(\mathbf{o}_t)} \cdot \left(\frac{\mathbf{o}_{tl} - \mu_{jkl}}{\sigma_{jkl}} \right). \end{aligned} \quad (4.49)$$

Note that the derivatives are sometimes written in terms of the original parameter and not in terms of the transformed parameter. The partial derivative of $\ln[P(\tilde{\theta})]$ with respect to $\tilde{\mu}_{jkl}^{(i)}$ can be obtained from Eq. (4.13),

$$\begin{aligned} \frac{\partial g(\hat{\theta})}{\partial \tilde{\mu}_{jkl}^{(i)}} &= -\frac{1}{2} \nu \sigma_{jkl}^{-2} \cdot 2(\tilde{\mu}_{jkl} \sigma_{jkl} - M_{jkl}) \cdot \sigma_{jkl} \\ &= -\nu \sigma_{jkl}^{-1} (\mu_{jkl} - M_{jkl}). \end{aligned} \quad (4.50)$$

Having updated the transformed mean using Eq. (4.46), the correct mean can be found using the inverse transformation $\mu_{jkl}^{(i)}(n+1) = \tilde{\mu}_{jkl}^{(i)}(n)\sigma_{jkl}^{(i)}(n+1)$ from Table 4.2.

Gaussian mixture variance

The steepest descent update for the Gaussian mixture variance $\tilde{\sigma}_{jkl}^2$ for state j , mixture k and element l is

$$\tilde{\sigma}_{jkl}^{(i)}(n+1) = \tilde{\sigma}_{jkl}^{(i)}(n) - \epsilon \frac{\partial E(\tilde{\theta}_n)}{\partial \tilde{\sigma}_{jkl}^{(i)}} \quad (4.51)$$

where

$$\begin{aligned} \frac{\partial E(\tilde{\theta}_n)}{\partial \tilde{\sigma}_{jkl}^{(i)}} = & -\varphi \frac{\partial \ln[P(\tilde{\theta})]}{\partial \tilde{\sigma}_{jkl}^{(i)}} - \varphi \sum_{i=1}^k \frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{\sigma}_{jkl}^{(i)}} \\ & - \sum_{i=k+1}^n \frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{\sigma}_{jkl}^{(i)}}. \end{aligned} \quad (4.52)$$

From Eq. (2.9) we can again write

$$\frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{\sigma}_{jkl}^{(i)}} = \sum_{t=1}^T \delta(\bar{q}_t - j) \frac{\partial \ln(b_j^{(i)}(\mathbf{o}_t))}{\partial \tilde{\sigma}_{jkl}^{(i)}}, \quad (4.53)$$

where (from Eq. (2.5))

$$\begin{aligned}
\frac{\partial \ln(b_j^{(i)}(\mathbf{o}_t))}{\partial \tilde{\sigma}_{jkl}^{(i)}} &= c_{jk}^{(i)} (2\pi)^{-d/2} |\Sigma_{jk}^{(i)}|^{-1/2} (b_j^{(i)}(\mathbf{o}_t))^{-1} \\
&\quad \left[\left(\frac{\mathbf{o}_{tl} - \mu_{jkl}^{(i)}}{\sigma_{jkl}^{(i)}} \right)^2 - 1 \right] e^{-\frac{1}{2} \sum_{l=1}^D \left(\frac{\mathbf{o}_{tl} - \mu_{jkl}^{(i)}}{\sigma_{jkl}^{(i)}} \right)^2} \\
&= \frac{c_{jk} \mathcal{N}(\mathbf{o}_t; \mu_{jk}, \Sigma_{jk})}{b_j^{(i)}(\mathbf{o}_t)} \cdot \left[\left(\frac{\mathbf{o}_{tl} - \mu_{jkl}^{(i)}}{\sigma_{jkl}^{(i)}} \right)^2 - 1 \right].
\end{aligned} \tag{4.54}$$

The partial derivative of $\ln[P(\tilde{\theta})]$ with respect to $\tilde{\sigma}_{jkl}^{(i)}$ can again be obtained from Eq. (4.13),

$$\begin{aligned}
\frac{\partial g(\hat{\theta})}{\partial \tilde{\sigma}_{jkl}^{(i)}} &= \frac{\partial g(\hat{\theta})}{\partial \sigma_{jkl}^{(i)}} \frac{\partial \sigma_{jkl}^{(i)}}{\partial \tilde{\sigma}_{jkl}^{(i)}} = \left[\nu \frac{(\mu_{jkl} - M_{jkl})^2}{\sigma_{jkl}^3} - \frac{(n-D)}{\sigma_{jkl}} + \frac{\tau_{jkl}}{\sigma_{jkl}^3} \right] \cdot (e^{\tilde{\sigma}_{jkl}}) \\
&= \nu \frac{(\mu_{jkl} - M_{jkl})^2}{\sigma_{jkl}^2} - (n-D) + \frac{\tau_{jkl}}{\sigma_{jkl}^2}.
\end{aligned} \tag{4.55}$$

After the update in the transformed space the new standard deviation can be found using $\sigma_{jkl}^{(i)}(n+1) = e^{\tilde{\sigma}_{jkl}(n+1)}$.

Gaussian mixture weights

The steepest descent update for the Gaussian mixture weights is

$$\tilde{c}_{jk}^{(i)}(n+1) = \tilde{c}_{jk}^{(i)}(n) - \epsilon \frac{\partial E(\tilde{\theta}_n)}{\partial \tilde{c}_{jk}^{(i)}} \tag{4.56}$$

where

$$\begin{aligned} \frac{\partial E(\tilde{\theta}_n)}{\partial \tilde{c}_{jk}^{(i)}} &= -\varphi \frac{\partial \ln[P(\tilde{\theta})]}{\partial \tilde{c}_{jk}^{(i)}} - \varphi \sum_{i=1}^k \frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{c}_{jk}^{(i)}} \\ &\quad - \sum_{i=k+1}^n \frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{c}_{jk}^{(i)}}. \end{aligned} \quad (4.57)$$

The partial derivative of Eq. (2.9) with respect to $\tilde{c}_{jk}^{(i)}$ is

$$\frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{c}_{jk}^{(i)}} = \sum_{t=1}^T \delta(\bar{q}_t - j) \frac{\partial \ln(b_j^{(i)}(\mathbf{o}_t))}{\partial \tilde{c}_{jk}^{(i)}}, \quad (4.58)$$

where

$$\begin{aligned} \frac{\partial \ln(b_j^{(i)}(\mathbf{o}_t))}{\partial \tilde{c}_{jk}^{(i)}} &= c_{jk}^{(i)} \left[(2\pi)^{-d/2} |\sigma_{jk}^{(i)}|^{-1/2} (b_j^{(i)}(\mathbf{o}_t))^{-1} \right. \\ &\quad \left. e^{-\frac{1}{2} \sum_{l=1}^D \left(\frac{\mathbf{o}_{tl} - \mu_{jkl}^{(i)}}{\sigma_{jkl}^{(i)}} \right)^2} - 1 \right] \\ &= c_{jk}^{(i)} \left[\frac{\mathcal{N}(\mathbf{o}_t; \mu_{jk}, \Sigma_{jk})}{b_j^{(i)}(\mathbf{o}_t)} - 1 \right]. \end{aligned} \quad (4.59)$$

The prior distribution for the mixture weights is a Dirichlet distribution and we therefore obtain the partial derivative of $\ln[P(\tilde{\theta})]$ with respect to $\tilde{c}_{jk}^{(i)}$ from Eq. (4.14).

$$\begin{aligned} \frac{\partial \ln[P(\tilde{\theta})]}{\partial \tilde{c}_{jk}^{(i)}} &= \frac{\partial g(\hat{\theta})}{\partial c_{jk}^{(i)}} \frac{\partial c_{jk}^{(i)}}{\partial \tilde{c}_{jk}^{(i)}} \\ &= \sum_m \left[(\delta_{jm} - 1) \frac{\partial \ln(c_{jm})}{\partial c_{jm}} \cdot \frac{\partial c_{jm}^{(i)}}{\partial \tilde{c}_{jk}^{(i)}} \right] \\ &= \sum_m \left[(\delta_{jm} - 1) \cdot c_{jm}^{-1} \cdot (\delta(k - m)c_{jm} - c_{jk}c_{jm}) \right] \\ &= (\delta_{jk} - 1) - \sum_m (\delta_{jm} - 1)c_{jm}. \end{aligned} \quad (4.60)$$

The new weight can be found by transforming back to the true weight space, using

$$c_{jk} = \frac{e^{\tilde{c}_{jk}}}{\sum_k e^{\tilde{c}_{jk}}}.$$

Transition probabilities

The steepest descent update for the transition probabilities is

$$\tilde{a}_{ij}^{(i)}(n+1) = \tilde{a}_{ij}^{(i)}(n) - \epsilon \frac{\partial E(\tilde{\theta}_n)}{\partial \tilde{a}_{ij}^{(i)}} \quad (4.61)$$

where

$$\begin{aligned} \frac{\partial E(\tilde{\theta}_n)}{\partial \tilde{a}_{ij}^{(i)}} = & -\varphi \frac{\partial \ln[P(\tilde{\theta})]}{\partial \tilde{a}_{ij}^{(i)}} - \varphi \sum_{i=1}^k \frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{a}_{ij}^{(i)}} \\ & - \sum_{i=k+1}^n \frac{\partial \ln[P(\mathbf{O}_i|\theta)]}{\partial \tilde{a}_{ij}^{(i)}}. \end{aligned} \quad (4.62)$$

The partial derivative of Eq. (2.9) with respect to $\tilde{a}_{ij}^{(i)}$ is

$$\frac{\partial \ln[P(\mathbf{O}_i)]}{\partial \tilde{a}_{ij}^{(i)}} = \sum_{t=1}^T \sum_s \delta(\bar{q}_{t-1} - i) \delta(\bar{q}_t - s) \left[\delta(j - s) - a_{ij} \right]. \quad (4.63)$$

The prior distribution for the transition probabilities is a Dirichlet distribution and we therefore obtain the partial derivative of $\ln[P(\tilde{\theta})]$ with respect to $\tilde{c}_{jk}^{(i)}$ from Eq. (4.15).

$$\begin{aligned}
\frac{\partial \ln[P(\tilde{\theta})]}{\partial \tilde{a}_{ij}^{(i)}} &= \frac{\partial g(\hat{\theta})}{\partial a_{ij}^{(i)}} \frac{\partial a_{ij}^{(i)}}{\partial \tilde{a}_{ij}^{(i)}} \\
&= \sum_k \left[(\alpha_{ik} - 1) \frac{\partial \ln(a_{ik})}{\partial a_{ij}} \cdot \frac{\partial a_{ij}^{(i)}}{\partial \tilde{a}_{ij}^{(i)}} \right] \\
&= \sum_k \left[(\alpha_{ik} - 1) \cdot a_{ik}^{-1} \cdot (\delta(j - k)a_{ik} - a_{ij}a_{ik}) \right] \\
&= (\alpha_{ij} - 1) - \sum_k (\alpha_{ik} - 1)a_{ik}.
\end{aligned} \tag{4.64}$$

Finally, the updated transition probability is found by transforming back using $a_{ij} = \frac{e^{\tilde{a}_{ij}}}{\sum_k e^{\tilde{a}_{ik}}}$.

4.4.3 Initial prior

Although the GMAP algorithm has been developed such that an initial parametric prior can be used (to incorporate true *a-priori* information), it is unlikely that any such information will exist. This algorithm will therefore often be used with a non-informative initial prior. This results in the prior gradient being zero, with no effect on the resultant update.

4.5 MAPMCE

The term MAP-MCE is sometimes used to refer to the usage of the MCE algorithm, but starting at the MAP point. This section describes an entirely different technique, which has been named MAPMCE in this work for the reason that it estimates the MAP point of the posterior of the probability of choosing the correct class (MCE), i.e. it optimizes the classification error (as criterion function) with respect to the parameters while incorporating a prior in the formulation.

In Section 3.2 it was mentioned that the MCE loss function is a reasonable estimate of the error probability. Given that this is the case, we can estimate the MAP point for the posterior distribution of the probability of choosing the correct class, namely

$$P(\theta, C_i | \mathbf{O}) \propto P(C_i | \theta, \mathbf{O}) P(\theta), \quad (4.65)$$

where C_i is the correct class and $P(C_i | \theta, \mathbf{O})$ is the probability of choosing the correct class C_i given the current model parameters θ and the observations \mathbf{O} .

If the MCE loss function is a good estimate of the error probability, then we can write

$$P(C_i | \theta, \mathbf{O}) \approx 1 - l(\mathbf{O}; \theta), \quad (4.66)$$

where $l(\mathbf{O}; \theta)$ is the loss function defined in Eq. (3.15).

Rewriting Eq. (4.41) for the posterior in Eq. (4.65) and using the MCE loss function, $P(C_i | \theta, \mathbf{O})$ in Eq. (4.66) instead of the likelihood function $L(\theta | \mathbf{O})$ (or $P(\mathbf{O} | \theta)$), we get

$$P(\theta | \mathbf{O}) \propto \left[P(\theta) \prod_{i=1}^k (1 - l(\mathbf{O}_i; \theta)) \right]^\varphi \prod_{i=k+1}^n (1 - l(\mathbf{O}_i; \theta)) \quad (k < n). \quad (4.67)$$

As mentioned in Section 4.4.3, there is little chance of a parametric prior being available for the gradient based MAP algorithm and this is also the case for the MAPMCE algorithm. The rest of this procedure will therefore be developed without the parametric prior $P(\theta)$.

Following the same reasoning as in Section 4.4, we can express Eq. (4.67) in terms of an energy function,

$$E(\theta; \mathbf{O}) = -\varphi \sum_{i=1}^k \ln[(1 - l(\mathbf{O}_i; \theta))] - \sum_{i=k+1}^n \ln[(1 - l(\mathbf{O}_i; \theta))]. \quad (4.68)$$

Unfortunately the gradient ($\frac{\partial E(\theta)}{\partial \theta}$) for Eq. (4.68) is not finite, with $\frac{\partial E(\theta)}{\partial \theta} = \infty$ for any $l(\mathbf{O}_i; \theta) = 1$.

A different estimate of the probability $P(C_i|\theta, \mathbf{O})$ has therefore been chosen, namely

$$P(C_i|\theta, \mathbf{O}) \approx e^{-l(\mathbf{O}_i; \theta)}. \quad (4.69)$$

Although this is not an entirely accurate estimate of the probability, it results in the following, more convenient, energy function

$$E(\theta; \mathbf{O}) = \varphi \sum_{i=1}^k l(\mathbf{O}_i; \theta) + \sum_{i=k+1}^n l(\mathbf{O}_i; \theta), \quad (4.70)$$

where $l(\mathbf{O}_i; \theta)$ is the MCE loss function (Eq. (3.15)) for observation i . Equation (4.70) is intuitively pleasing as it is simply a weighted version of the standard MCE algorithm in the sense that the “prior” terms ($i = 1 \dots k$) are weighted by φ . An error in the “prior” or reference set will therefore not be penalized as heavily as an error in the adaptation set ($0 \leq \varphi \leq 1$). The algorithm will therefore try to minimize errors in both sets, but will place more emphasis on errors that occur in the adaptation set.

The implementation of this algorithm, requires the following simple modification to the standard MCE algorithm which was described in Chapter 3: if an observation is in the “prior” dataset, then the gradient with respect to the model parameters ($\nabla l(\mathbf{O}; \theta)$) used in the GPD update (Eq. (3.18)) must be weighed with φ .

4.6 Experiments

The goal of this section is to experimentally compare the three algorithms discussed earlier in this chapter (MAP, GMAP and MAPMCE) in conditions where limited training data is available, with a reasonable amount of non-task-specific data available for adaptation purposes.

Before continuing, it is necessary to explain the convention I have used to describe (or label) the algorithms used:

- MAP - A MAP algorithm labelled as “MAP $T_X (T_Y)$ ” uses the dataset T_X as its adaptation set and the ML model created using T_Y to determine the prior distribution. Here, Eqs. (4.22) through (4.26) are used.
- GMAP - A GMAP algorithm labelled as “GMAP $T_X (T_Y)$ ” uses the dataset T_X as its adaptation set and the ML model created using T_Y as a starting point. Note that the prior dataset is not included in the description as it is a constant for each experiment. Here, Eqs. (4.46) through (4.64) are used.
- MAPMCE - A MAPMCE algorithm labelled as “MAPMCE $T_X (T_Y)$ ” uses the dataset T_X as its adaptation set and the ML model created using T_Y as a starting point. As with GMAP, the prior dataset is not included in the description thereof. Here, Eq. (4.70) is implemented using Eq. (3.25) through Eq. (3.46).

The number of iterations used for each procedure differed. When using MAP, 10 iteration typically proved to be sufficient with the testing set performance converging at or before 10 iterations. GMAP required between 10 and 30 iterations for the testing set performance to converge, while the MAPMCE algorithm also required between 10 and 30 iterations to attain peak testing set performance. The number of iterations required for the gradient-based algorithms (GMAP and MAPMCE) is dependent on the step size ϵ , the size of the datasets and the weighting factor φ .

In situations where HMMs of differing complexity are used, the number of states and mixtures will be included in the description. For example, the algorithm description “MAP 3,5 $T_X (T_Y)$ ” refers to the MAP algorithm using a 3 state, 5 mixture HMM.

4.6.1 SUNSpeech

This section compares the three adaptation algorithms within a language adaptation framework. The SUNSpeech dataset described in Section 2.4 is used for this purpose. The subsets of the dataset described in Section 2.4.3 are used for this experiment.

Table 4.3 presents the Afrikaans test set accuracy of the base system (described in Section 2.1) for the different training sets. As expected, the performance of the system trained using the English subset (31.9% and 33.9% for 5 and 10 mixtures respectively) is worse than that trained using either the full Afrikaans training set (48.6% and 51.5%) or the reduced Afrikaans training subset (42.5% and 41.2%). Training using combined English and Afrikaans sets produces relatively poor results as compared to only using the associated Afrikaans set.

Table 4.3: Base system results for SUNSpeech Afrikaans test set

Training set	Mixtures	
	5	10
English (E)	31.9%	33.9%
Afrikaans train (A)	48.6%	51.5%
Afrikaans train + English ($A + E$)	37.9%	41.4%
Afrikaans adapt (A_S)	42.5%	41.2%
Afrikaans adapt + English ($A_S + E$)	34.2%	37.6%

Using 10 mixture components as opposed to 5 improves results, except when using the reduced Afrikaans training set where a reduction in performance occurs. This phenomenon can be ascribed to the bias/variance problem discussed in Section 2.2.1.

The more complex 10 mixture HMM has a lower bias, and will therefore work better in situations where there is sufficient data available. In situations where there is little data available, the variance term becomes dominant and so the less complex model (5 mixture HMM) works best.

The recognition results appear relatively low, as compared to that obtained for continuous recognition on other continuous phoneme recognition tasks, such as TIMIT. The recognition performance of the same system using the TIMIT database is 57.0% and 60.0% for 3 state models using 5 and 10 mixtures respectively. There are several reasons that could account for the disparity in recognition performance, including:

- There are a total of 59 phonetic categories (including *silence* and *unknown*) used in SUNSpeech. This is as compared to, for example, TIMIT where there are 39 phonetic categories. The task is therefore somewhat more complex.
- The SUNSpeech dataset labels include an *unknown* class, which accounts for 1.9% of the labels in the test set. The *unknown* class is not modeled, this will result in at best a 1.9% poorer error rate. The *unknown* class includes any sounds or speech that cannot be included in any of the other 58 phonetic categories.

MAP

This section presents the results obtained for the MAP algorithm described in Section 4.3 using the SUNSpeech dataset. The adaptation set and prior are an integral part of the MAP algorithm and will therefore be included in any description thereof.

Figure 4.3 presents the MAP adaptation results for a 3 state, 5 mixture system when only the reduced Afrikaans training set (A_S) is available. The line labelled “ML A_S ” gives the base maximum likelihood performance (42.5%) of a 3 state, 5 mixture model trained on the Afrikaans adaptation set only. The MAP algorithms “MAP A_S (E)”

and “ $MAP_{A_S}(A_S + E)$ ” use a prior created using the English training set (E) and a prior made using the union of the English and Afrikaans adaptation sets ($A_S + E$).

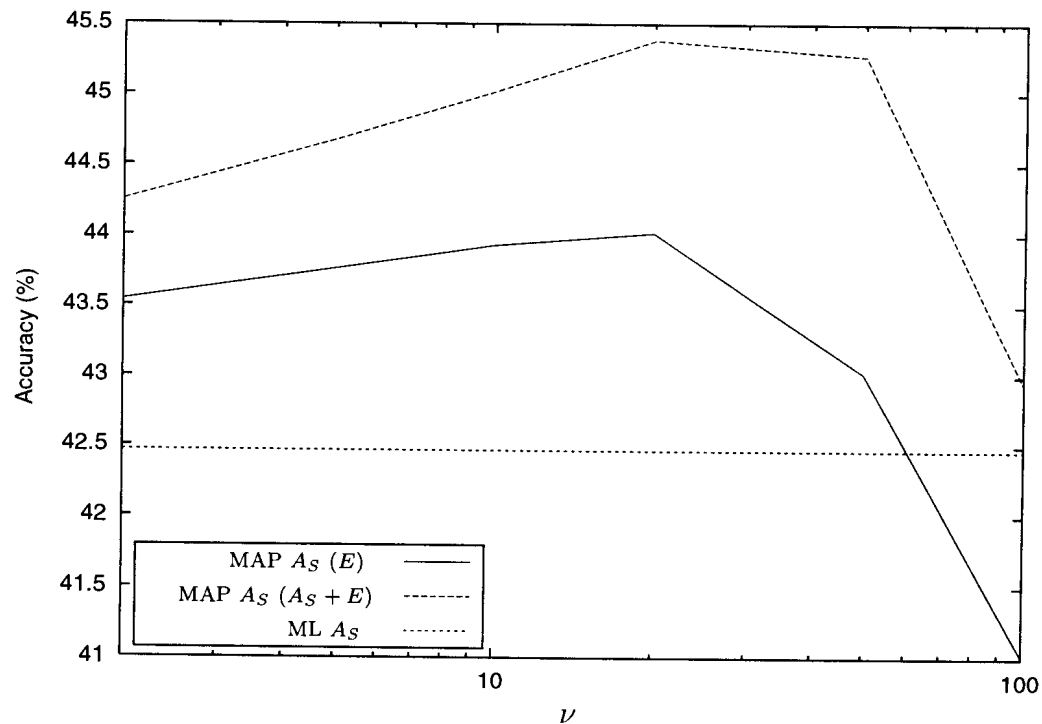


Figure 4.3: MAP adaptation results plotted versus the prior parameter ν (Eq. (4.12)) for adaptation set A_S using an English prior (E) and a prior created using pooled English and Afrikaans adaptation set ($A_S + E$) for a 3 state, 5 mixture model.

Peak performance of 45.4% phoneme accuracy is attained using the pooled set ($A_S + E$) to create the prior (5% relative improvement in error rate). The MAP algorithm using a prior created using the English training set only performed considerably worse, reaching a peak performance of 44.0% (2.6% relative improvement in error rate).

Both implementations attain peak performance with $\nu = 20$. As expected, the performance of the MAP algorithm is highly dependent on the value chosen for the parameter ν . If ν is too small, the prior would not be informative enough and little prior information is therefore obtained from the prior training set. A large value of ν , on the other hand, would result in a prior that is too restrictive, thereby keeping the resultant model very close to the model trained on the prior training set. The recognition accuracy is therefore expected to tend towards that of an ML model trained using the only

the adaptation set (ML A_S) for very small values of ν . Conversely, the recognition accuracy is expected to tend towards that of the ML model trained using the prior training set (E) for large values of ν .

There is, however, an optimal value for ν between these two extremes. This occurs when information from both sets is being used optimally.

Figure 4.4 shows the the MAP adaptation results for a 3 state, 5 mixture HMM system when the full Afrikaans training set (A) is used as the adaptation set. The results for the MAP algorithm using both the English set prior and a pooled prior created using both the English set and full Afrikaans training set ($A + E$) are compared.

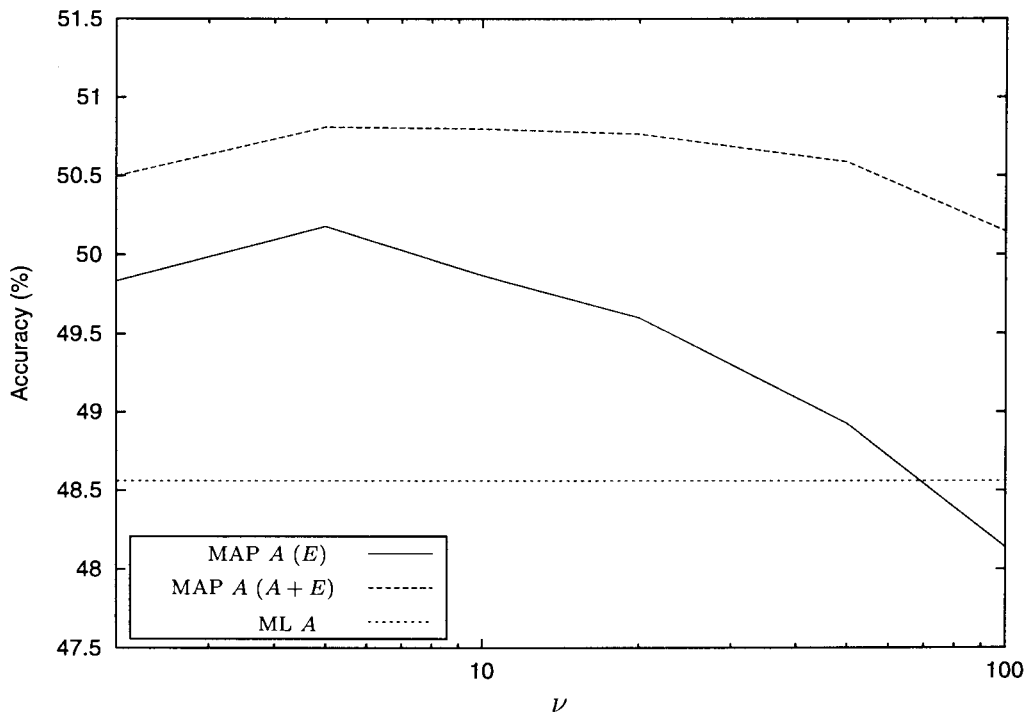


Figure 4.4: MAP adaptation results using full Afrikaans set

The recognition accuracy for the base ML system (ML A) is 48.6%. The MAP algorithm using the pooled prior ($A + E$) with peak performance of 50.8% (4.3% relative improvement of error rate), once again performed better than the MAP algorithm using only the English prior (E) with a peak performance of 50.2% (3.1% relative

improvement of error rate).

Here, peak performance was attained for $\nu = 5$ for both algorithms, which is less than that required when the smaller Afrikaans training set is used. As the amount of task-specific training data increases, one would expect the optimal value of ν to decrease, as has been demonstrated in this experiment.

One would also expect the MAP algorithm to become less useful as the amount of task-specific training data increases. This is evident here, with the relative improvement in performance (4.3%) being less than that obtained for the small Afrikaans training set (5%).

The results for a 3 state, 10 mixture HMM model experiment using only the reduced Afrikaans set (A_S) are shown in Figure 4.5 (the priors are either E or $A_S + E$). Maximum recognition accuracy of 44.8% is attained using the pooled prior, while the MAP algorithm using only the English prior again does not work as well managing a maximum of 43.8%.

The results are very similar to that obtained for the less complex 3 state, 5 mixture HMM. It is, however, apparent that the absolute performance is worse than that of the the 3 state, 5 mixture model (44.8% versus 45.4%).

Table 4.4 summarizes the best accuracies attained using the MAP algorithm for the different configurations. Note that the 3 state, 10 mixture model is consistently worse than the 5 mixture variant when the small Afrikaans training set is used. It is also noticeable that the effect of the MAP algorithm is reduced as the amount of Afrikaans data (adaptation data) is increased.

Importantly, the MAP algorithm was not able to improve on the ML results for the 3 state 10 mixture model, when using the English prior with the full Afrikaans training set.

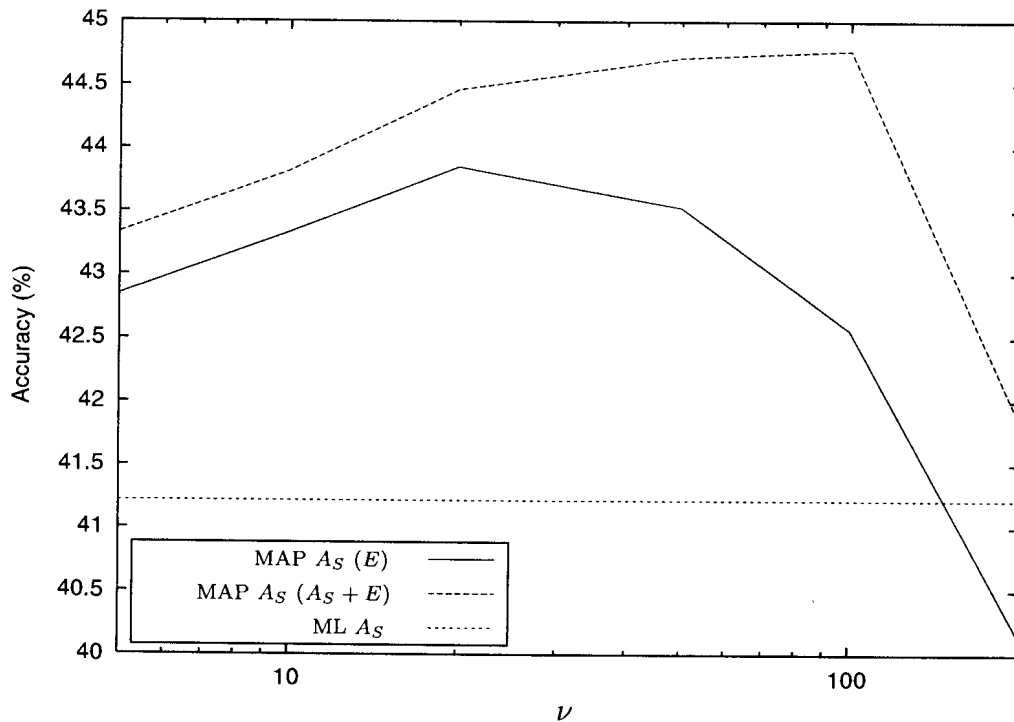


Figure 4.5: MAP adaptation results for a 3 state, 10 mixture HMM using the small Afrikaans training set

GMAP

The gradient based MAP algorithm presented in Section 4.4 is tested using the SUN-Speech dataset here. Incorporating the adaptation data in both the prior and adaptation datasets results in the adaptation data (Eq. (4.44)) having an effective weight of $(1 + \varphi)$ with the prior data having a weight of φ . The same result can be accomplished, without incorporating the adaptation data in the prior data, by choosing a new weighting factor $\varphi' = \frac{\varphi}{(1+\varphi)}$. There is therefore no point in pooling the English and Afrikaans datasets and using the resultant dataset as a prior dataset, as was done for the standard MAP algorithm. As a result, the English dataset is used as the prior dataset for all experiments conducted in this section.

As a result of using the steepest descent algorithm, one would expect this algorithm to easily fall into local minima. The starting point, should therefore influence the

Table 4.4: Summary of the MAP accuracy results obtained for the SUNSpeech database. Relative improvement in error rate over baseline performance is given in brackets.

Description	Mixtures	
	5	10
Baseline ML (A_S)	42.5% (0.0%)	41.2% (0.0%)
MAP A_S (E)	44.0% (2.6%)	43.9% (4.6%)
MAP A_S ($A_S + E$)	45.4% (5.0%)	44.8% (6.1%)
Baseline ML (A)	48.6% (0.0%)	51.5% (0.0%)
MAP A (E)	50.2% (3.1%)	51.25% (-0.5%)
MAP A ($A + E$)	51.1% (4.9%)	52.4% (1.9%)

performance of the algorithm to some degree. Figure 4.6 presents the results using the gradient based MAP algorithm for a 3 state, 5 mixture HMM, starting at (a) the English ML model, (b) the ML model for the small Afrikaans training set and (c) the best MAP estimate from the previous section. The baseline (ML) accuracies for the large dataset (ML A) and small training set (ML A_S) are also plotted.

Starting the algorithm using the English ML model, results in relatively poor accuracy (peak accuracy of 35%), considerably worse than that of the ML model of the small Afrikaans training set. Using the ML model of the small Afrikaans training set and the best MAP estimate (which produced an accuracy of 45.4%) as starting points, produces good improvements in accuracy relative to that of the ML model (peak accuracy of 46.0% and 47.1% respectively).

Figure 4.7 shows the performance of the gradient based MAP algorithm for both 5 and 10 mixture HMMs. The starting point was chosen as the best MAP estimate for all algorithms presented in this graphic.

It is important to note that here, the performance of the 10 mixture HMM is better than that of the 5 mixture HMM. This is, as opposed to that of standard MAP and the ML algorithms, where the 3 state 10 mixture HMM performed slightly worse.

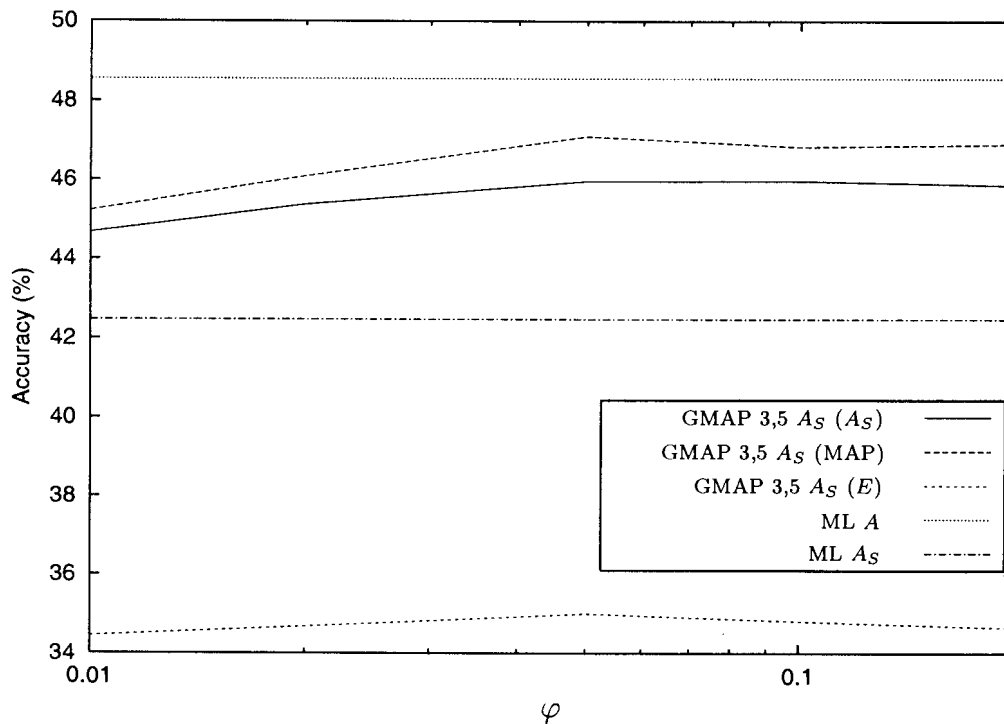


Figure 4.6: Gradient MAP adaptation results for the small Afrikaans training set (A_S)

The weighting factor φ (Eq. (4.44)) has much the same effect as the parameter ν (Eq. (4.12)) has in the standard MAP algorithm. Although it is not easily observable from the results presented here, this algorithm exhibits asymptotic behaviour in the extremes of φ . For $\varphi = 0$ the training examples in the prior set would be ignored and the algorithm would become an ML algorithm. A value of 1 would weight prior and adaptation examples equally and the algorithm becomes an ML algorithm trained on the pooled dataset.

Table 4.5 summarizes the results obtained when using the GMAP algorithm. The GMAP algorithm produces significant improvements on the baseline performance when the small Afrikaans training set is available. Note that this is not the case when the English set is used as starting point, and the starting point should therefore be chosen carefully. Smaller, but significant improvements are realized when the full Afrikaans training set is used.

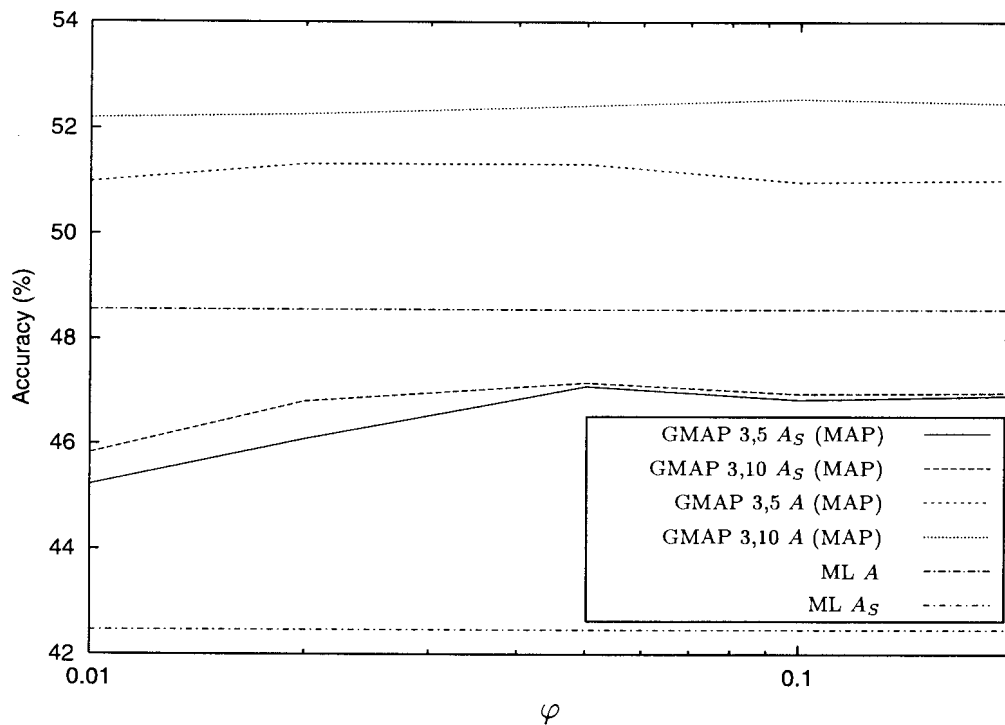


Figure 4.7: Gradient MAP adaptation results for 5 and 10 mixture HMMs

MAPMCE

The MAPMCE algorithm presented in Section 4.5 is tested using the SUNSpeech dataset here. The English dataset is used as the prior dataset for all experiments conducted in this section.

In Chapter 3 it was argued that it was not necessary to use a sigmoid (or smoothed zero-one) loss function and results indicated that there was little to be gained from the use thereof. It is therefore necessary to again investigate the effect of a smoothed zero-one loss function within the MAPMCE framework.

Figure 4.8 presents the results using the MAPMCE algorithm for a 3 state 5 mixture model using the small Afrikaans training set, starting from the English ML model. The results for both sigmoid and no-sigmoid loss functions are shown. The baseline performance, namely that of the ML model trained using the small Afrikaans set, is

Table 4.5: Summary of the accuracy results obtained using the GMAP algorithm

Description	Mixtures	
	5	10
Baseline ML (A_S)	42.5%	41.2%
GMAP A_S (E)	35.0%	34.9%
GMAP A_S (A_S)	46.0%	46.2%
GMAP A_S (MAP)	47.1%	47.2%
Baseline ML (A)	48.6%	51.5%
GMAP A (MAP)	51.3%	52.5%

shown.

The MAPMCE algorithm using a sigmoid loss function attains a peak accuracy of 32.4% for small values of φ . This is not sufficient to improve on the baseline performance of the ML model. The MAPMCE algorithm which does not make use of a sigmoid loss function, however, performs considerably better and manages to improve on the baseline performance. A peak accuracy of 43.5% is attained, resulting in a relative improvement in error rate of 1.8% as compared to that of the baseline ML system.

These results confirm the discussion in Chapter 3, showing that the sigmoid function merely serves to introduce additional local minima, and therefore hinders the algorithm from reaching a meaningful local minimum. All subsequent experiments using the MAPMCE algorithm will therefore not use a smoothed zero-one loss function.

Figure 4.9 shows the results of the MAPMCE algorithm for various configurations. The ML results for a 3 state, 5 mixture HMM are shown for both the full (ML 3,5 A) and reduced Afrikaans (ML 3,5 A_S) training sets for reference purposes. The starting point for all the results presented here is the best associated MAP estimate.

It is evident from Figure 4.9 that although the algorithm does exhibit a degree of

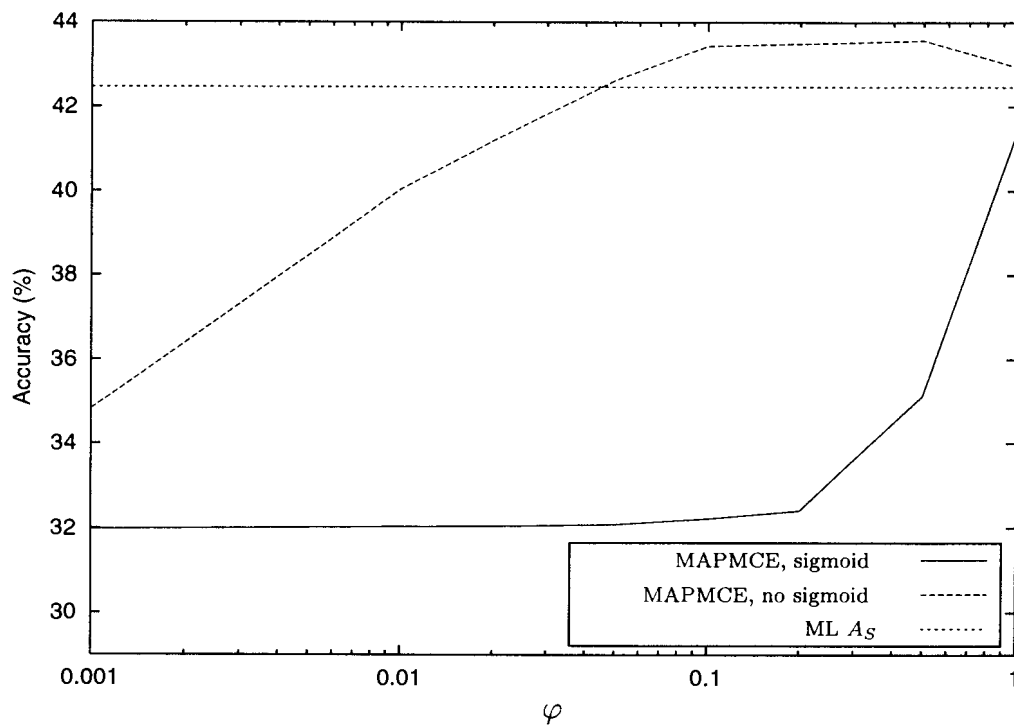


Figure 4.8: MAPMCE adaptation results for the small Afrikaans training set when using the English ML model as starting point

dependence on the weighting parameter φ , it is not overly sensitive to it.

We would prefer an algorithm to perform better with more complex models, with performance increasing as the model complexity increased. Unfortunately, here the algorithm performs best for the reduced dataset, when using the less complex 5 mixture per state HMM. As expected, the opposite is true when the full training set is used, resulting in the more complex model performing best.

Table 4.6 summarizes the results obtained for the MAPMCE algorithm when using the SUNSpeech dataset. Here, unlike the GMAP algorithm an improvement in performance is realized when using the English ML model as starting point, although this is not the case when using a sigmoid loss function. The MAP point estimate was once again found to be the best starting point for the MAPMCE algorithm. Reasonable improvements are attained for both the small and full Afrikaans training sets.

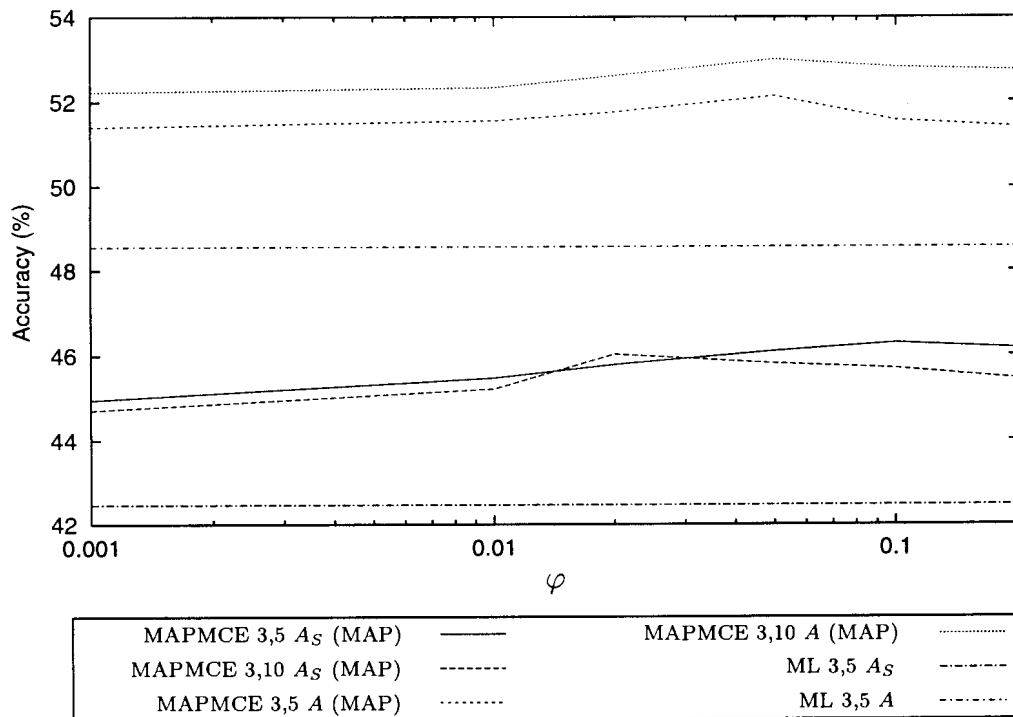


Figure 4.9: MAPMCE adaptation results for both Afrikaans training sets using the best MAP point as starting point

Comparison

Table 4.7 summarizes the best results obtained using the different algorithms for the SUNSpeech dataset. Considerable improvements are realized from the usage of the three algorithms when the small Afrikaans training set is used, with the GMAP algorithm performing best under these conditions (8.0% and 10.2% relative improvement in error rate for the 5 and 10 mixture models respectively).

Smaller, but finite, improvements are attained when the full Afrikaans training set is available. Here, the MAPMCE algorithm produces the largest improvements (6.8% and 3.1% relative improvement in error rate for the 5 and 10 mixture HMMs respectively). The MAP algorithm, although it manages to improve on the baseline performance, produces the worst results of the three algorithms for the configurations and scenarios tested. Note, however, that the MAP algorithm was integral to the optimal perfor-

Table 4.6: Summary of the results obtained using the MAPMCE algorithm

Description	Mixtures	
	5	10
MAPMCE A_S (E)	43.5%	43.3%
MAPMCE A_S (MAP)	46.3%	46.0%
MAPMCE A (MAP)	52.1%	53.0%

mance of the GMAP and MAPMCE algorithms, as the best MAP estimate was used as starting point.

Table 4.7: Summary of the best results obtained for the SUNSpeech dataset. Relative improvement in error rate over baseline is given in braces.

Description	Mixtures	
	5	10
Baseline ML A_S	42.5 (0.0%)	41.2 (0.0%)
MAP A_S ($A_S + E$)	45.4 (5.0%)	44.8 (6.1%)
GMAP A_S (MAP)	47.1 (8.0%)	47.2 (10.2%)
MAPMCE A_S (MAP)	46.3 (6.6%)	46.0 (8.2%)
Baseline ML A	48.6 (0.0%)	51.5 (0.0%)
MAP A ($A + E$)	51.1 (4.9%)	52.4 (1.9%)
GMAP A (MAP)	51.3 (5.3%)	52.5 (2.1%)
MAPMCE A (MAP)	52.1 (6.8%)	53.0 (3.1%)

4.6.2 TIMIT

The collection of large speech databases is not a trivial task (if done properly). It is not always possible to collect, segment and annotate large databases for every task or language. It is also often the case that there are imbalances in the databases. An example of one such imbalance is the fact that there is often more male speakers than

female speakers (or *vice-versa*). If there is, for example, far fewer female speakers than male speakers, then the recognizers will tend to work poorly for female speakers (as compared to performance for male speakers).

This experimental section attempts to recreate such a scenario, where relatively little female speaker training data is available, so as to test the three algorithms discussed earlier in this chapter under these conditions. Table 4.8 details the TIMIT training and testing sets used in the experiments presented in this section. Note that the default TIMIT training set has an imbalance with respect to male and female speakers, with 326 male speakers and 136 female speakers. The standard TIMIT training set is therefore used as one such example.

Table 4.8: Description of TIMIT training and testing sets used

Description	Label	Number of speakers			Sentences	Duration (minutes)
		Male	Female	Total		
Training sets:						
Full (standard)	T	326	136	462	4620	236.5
Male	T_M	326	0	326	3260	165.2
Female	T_F	0	136	136	1360	71.3
Female-small	T_{FS}	0	16	16	160	8.2
Testing set (standard)						
Male test set		112	0	112	1120	56.8
Female test set		0	56	56	560	29.7

Although there is an imbalance in the standard training set, with the female speaker accounting for only 29.5% of the training set, it is necessary to investigate the effect of an even smaller number of female speakers. It is for this purpose, that a small female training set has been created, consisting of 16 speakers (2 from each of the 8 dialect regions in TIMIT). The small female training set (T_{FS}) is therefore approximately 5% the size of the full male training set (T_M) and 12% of the full female training set (T_F). A more detailed description of the TIMIT database can be found in Section 2.4.

Two different scenarios therefore exist for the purpose of this section, namely

- The full TIMIT training set (T) is available for training purposes.
- The full male training set (T_M) and the small or reduced female training set (T_{FS}) are available for training.

The algorithms will be tested using two criteria, namely:

1. Adaptation performance - the ability of the algorithm to improve the accuracy of the test set associated with the adaptation set,
2. Training performance - the ability of the algorithm to improve the accuracy of the combined testing set (“standard” testing set in Table 4.8).

The rationale behind the above two criteria, is that although the MAP algorithm is a good adaptation technique, it will possibly not perform as well in situations where we want good performance for both adaptation and reference (prior) sets. There are numerous situations where one would prefer better combined performance. Often, for example, when a severe imbalance in the number of female and male speakers occurs and a gender independent recognizer is used, the performance for the one grouping will tend to be relatively poor (and so too the combined performance).

The combined performance, however, can be somewhat misleading. A reasonable combined performance can be produced if a system works well for one group and poorly for another. For example, a system that attains 90% correct recognition for male speakers but only 30% for female speakers would, assuming equal numbers of male and female speakers, therefore give 60% correct recognition for the combined set. The minimum performance for the two sets will also be presented, i.e. in the above example the 30% recognition rate will be reported. When summarizing results, only those algorithms

and their configurations which maximize either the female, combined and minimum testing set performance will be listed.

Table 4.9 presents the baseline ML results obtained for the TIMIT database. It is evident, as expected, that if one were to create a gender specific recognizer, that the models trained using the gender dependent subsets are best. The small female training set produces results which are worse than that obtained using the full female training set for both male and female testing sets (and consequently the combined testing set).

Table 4.9: Base system (ML) results for TIMIT

Training set	Test set accuracy			
	Male	Female	Combined	Minimum
Full training set T	57.3%	56.5%	57.0%	56.5%
Male training subset T_M	59.1%	43.0%	53.7%	43.0%
Female training subset T_F	40.1%	61.0%	47.2%	40.1%
Small female training subset T_{FS}	39.9%	56.8%	45.6%	39.9%
$T_{FS} + T_M$	59.1%	46.7%	54.9%	46.7%

The ML model resulting from the combination (pooling) of the male and the small female training subsets ($T_{FS} + T_M$) results in better combined performance than that attained by any ML model other than that of the full training dataset. The performance for the female testing set is, however, poor and it is only better than that of the ML model trained using the male training subset only.

A 3 state, 5 mixture HMM is used for all experiments using the TIMIT database in this section. The general configuration of the HMM recognizer, is as described in Section 2.1.

MAP

The usage of the MAP algorithm for both adaptation and training purposes will now be tested for both scenarios where limited training data is available for female speakers.

Figure 4.10 presents the MAP results when the small female training set is used. A pooled prior created from the pooled dataset ($T_M + T_{FS}$) is used. The performance for the female testing set peaks at $\nu = 100$, giving a peak accuracy of 57.4%. The accuracy for the male testing set improves as ν gets larger, i.e. as the prior becomes more restrictive and the final result is closer to the ML estimate obtained from the prior.

Combined performance peaks at $\nu = 2000$ giving a peaked combined performance of 56%. The solution resulting in the best combined performance would probably not be acceptable due to the disparity in results for the male and female testing sets (57.8% and 52.4% respectively). A better solution is that which maximizes the minimum performance of the two sets, this occurs for $\nu = 700$ which gives 55.1% accuracy for both male and female sets.

Figure 4.11 shows the MAP results when using the full female training set (T_F). A pooled prior ($T_M + T_F$) is used. Peak performance of 61.0% is attained for the female training set for $\nu = 50$. The performance for the male testing set, as with the previous experiment, increases as ν becomes larger.

Although the graphic does not show the point at which the male and female performance is equal or at which point maximum combined performance is produced, it will undoubtedly be very close to that at which $\nu = \infty$ or the performance of the ML model (of the pooled set).

The MAP results for the TIMIT database are summarized in Table 4.10. Using MAP for gender adaptation (improving female test set performance) when using the small training set results in a relatively small improvement in error rate (1.4% for $\nu = 100$)

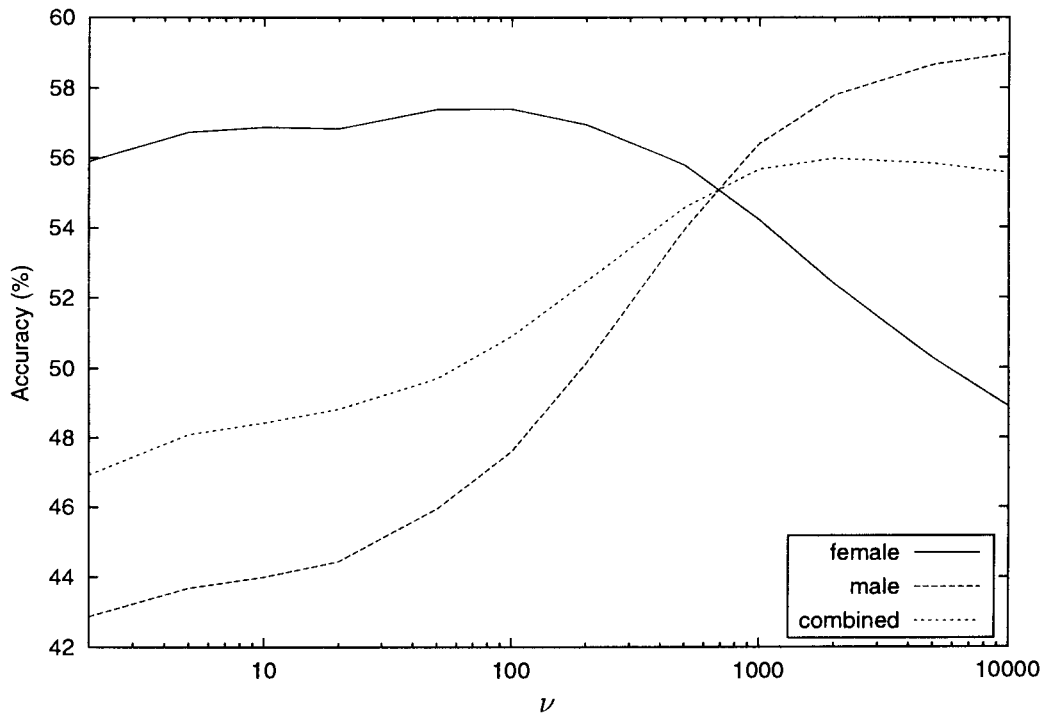


Figure 4.10: MAP adaptation results for the small female training set, using a pooled prior ($T_M + T_{FS}$)

when using the small female subset. No improvement is observed when the full female training set is used, where 61.0% accuracy is attained, equal to that realized using the ML model for the full female testing set.

A considerable improvement in minimum accuracy is attained when using reduced female training set, where a 15.8% relative reduction in error is realized (46.7% \rightarrow 55.1%). Note that the best minimum accuracy (55.1%) using MAP is not much worse as that attained when using the full training set (56.5%). The MAP algorithm does not improve on the minimum accuracy attained of the ML model when using the full training set.

The results for the MAP algorithm using a prior created using the pooled dataset have been presented. The results for the MAP algorithm using a prior created from the male training set have not been included as they are similar for adaptation (female testing set accuracy) and worse in terms of the combined and minimum test set accuracies.

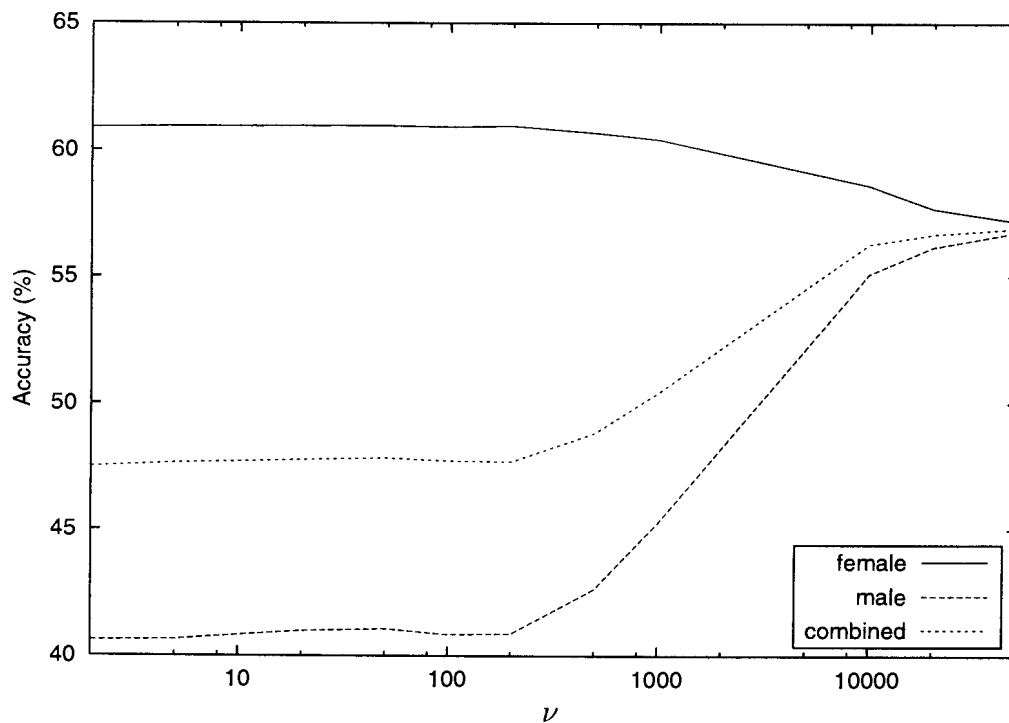


Figure 4.11: MAP adaptation results for the full female training set, using a pooled prior ($T_M + T_F$)

GMAP

The GMAP algorithm is experimentally evaluated here, for the two situations described earlier. The full male training set (T_M) is available in both situations and is used as the prior dataset throughout.

The results obtained using the GMAP algorithm for the small female training set are shown in Figure 4.12. The initial starting point used here is the ML model created from the pooled data ($T_M + T_{FS}$). The performance for the female testing set increases as φ decreases. It is therefore evident that for this scenario, the GMAP algorithm will give best performance for the female test set when $\varphi = 0$ (the ML estimate using the small female testing set).

Combined performance is maximized when $\varphi = 0.1$, producing a peak combined accu-

Table 4.10: MAP results for gender adaptation experiments using TIMIT

Algorithm	Test set accuracy (%)			
	Male	Female	Combined	Minimum
MAP T_{FS} ($T_{FS} + T_M$), $\nu = 100$	47.6	57.4	50.9	47.6
MAP T_{FS} ($T_{FS} + T_M$), $\nu = 700$	55.1	55.1	55.1	55.1
MAP T_{FS} ($T_{FS} + T_M$), $\nu = 2000$	57.8	52.4	56.0	52.4
MAP T_F ($T_F + T_M$), $\nu = 50$	41.1	61.0	47.8	41.1
MAP T_F ($T_F + T_M$), $\nu = \infty$	57.3	56.5	57.0	56.5

racy of 56.4%, while the minimum accuracy is maximized when $\varphi = 0.02$ resulting in a 55.3% peak minimum accuracy. Although it was not an objective of this experiment, a slight improvement in male testing set performance is attained for $\varphi = 0.5$, where a 59.2% accuracy is realized.

Figure 4.13 shows the GMAP results the full female training set. The best female testing set performance for the range of φ shown occurs for $\varphi = 0.02$, where an accuracy of 60.5% is attained. The female testing accuracy for $\varphi = 0$ should, however, be the same as that of the ML model estimated using the full female training set (61.0%). The peak combined performance (57.0%) is attained when $\varphi = 1$, while the peak minimum performance (56.7%) results when $\varphi = 0.8$.

Note that the value of φ resulting in optimal minimum testing set accuracy is larger when the entire female training set is used. If there is no imbalance in the numbers of male and female speakers, then one would expect the optimal value of φ to be one (assuming that recognition accuracy would be the same for equivalent amounts of data). As more and more training data becomes available for female speakers, the contribution thereof to the posterior distribution becomes larger. The optimal value of φ therefore increases, so that the adaptation data does not dominate.

The best GMAP results using the TIMIT dataset are summarized in Table 4.11. Rea-

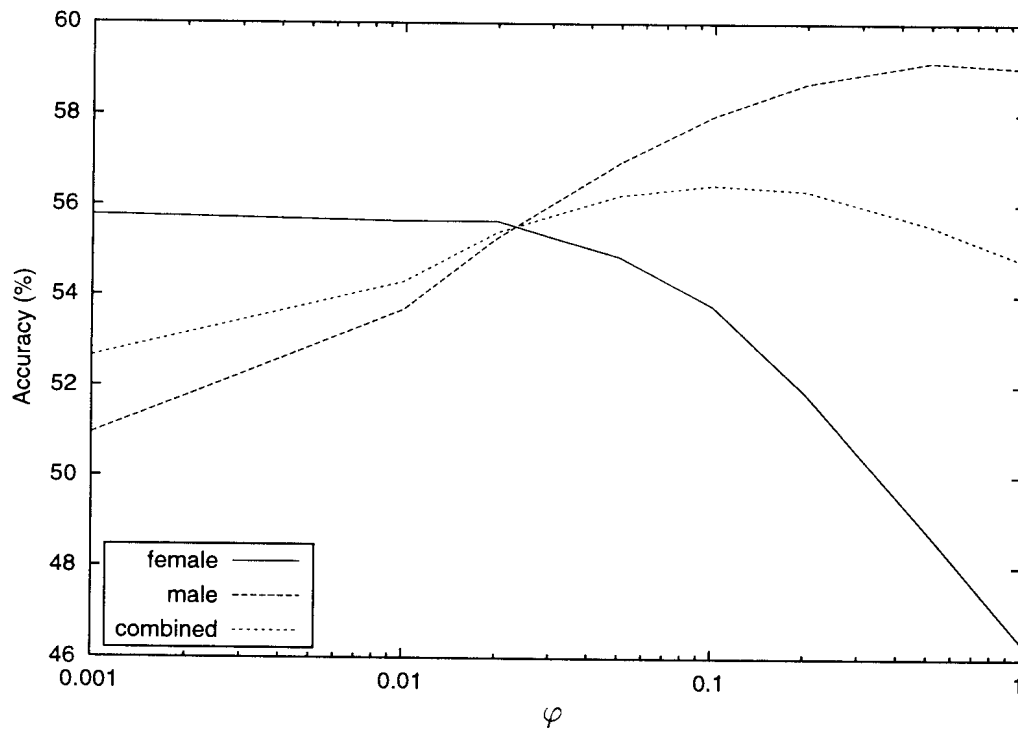


Figure 4.12: GMAP adaptation results for the small female training set, using the pooled data ($T_M + T_{FS}$) ML model as a starting point

sonable improvements in minimum and combined testing set results are realized using the GMAP algorithm with the small female training set, where a 16.1% and 3.3% relative improvement in error rates are achieved respectively.

The GMAP algorithm, however, failed to improve on the female testing set performance attained using the ML trained model using the corresponding female training set, for both scenarios where the small and full female training sets are used.

MAPMCE

The MAPMCE algorithm presented in Section 4.5 is tested using the TIMIT dataset given the conditions described earlier. The male speaker training set is used as the prior dataset for all experiments conducted in this section.

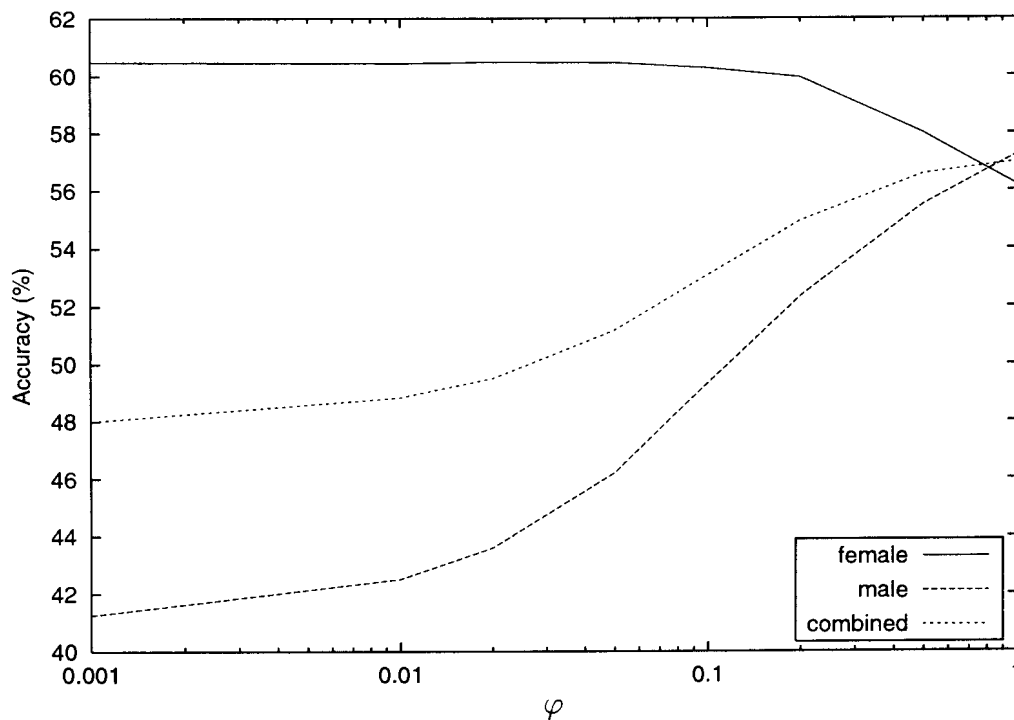


Figure 4.13: GMAP adaptation results for the full female training set, using a pooled data ($T_M + T_F$) ML model as starting point

Figure 4.14 presents the MAPMCE results for the small female training set using the pooled ML model as a starting point. Note that here the female testing set performance does not improve to the extent that it becomes better than that of the male testing set. The peak minimum and female testing set performances are therefore obtained for the same value of $\varphi = 0.1$, which produces a peak accuracy of 56.3% for both criteria. The best combined performance of 62.2% is attained when $\varphi = 0.5$, though this is only slightly better than that attained at the peak minimum accuracy point, which results in a combined performance of 62.1%.

Figure 4.15 shows the MAPMCE results for the small female training set, but with the best MAP point used as starting point. Here, the female testing set performance is considerably better, with a peak performance of 60.3% being attained for $\varphi = 0.2$. The accuracies attained for the male testing set are, however, somewhat worse with a peak male testing set accuracy of 63.0% resulting for $\varphi = 1.0$. The best combined

Table 4.11: GMAP results for gender adaptation experiments using TIMIT

Algorithm	Test set accuracy (%)			
	Male	Female	Combined	Minimum
GMAP $T_{FS} (T_{FS} + T_M)$, $\varphi = 0.001$	50.5	55.8	52.6	50.5
GMAP $T_{FS} (T_{FS} + T_M)$, $\varphi = 0.02$	55.3	55.6	55.4	55.3
GMAP $T_{FS} (T_{FS} + T_M)$, $\varphi = 0.1$	58.0	53.8	56.4	53.8
GMAP $T_{FS} (MAP)$, $\varphi = 0.01$	46.2	57.0	50.0	46.2
GMAP $T_{FS} (MAP)$, $\varphi = 0.08$	55.0	55.1	55.1	55.0
GMAP $T_{FS} (MAP)$, $\varphi = 0.2$	57.2	52.0	55.4	52.0
GMAP $T_F (T_F + T_M)$, $\varphi = 0.02$	43.6	60.5	49.5	43.6
GMAP $T_F (T_F + T_M)$, $\varphi = 0.8$	56.7	56.7	56.7	56.7
GMAP $T_F (T_F + T_M)$, $\varphi = 1.0$	57.2	56.2	57.0	56.2
GMAP $T_F (MAP)$, $\varphi = 0.01$	41.2	60.9	47.8	41.2
GMAP $T_F (MAP)$, $\varphi = 0.7$	57.0	57.0	57.0	57.0
GMAP $T_F (MAP)$, $\varphi = 1.0$	57.7	55.9	57.1	55.9

accuracy is 61.4% for $\varphi = 1.0$, which is worse than that produced by the MAPMCE algorithm using the ML model trained using the pooled dataset as starting point.

The peak minimum test set accuracy is 60.3% ($\varphi = 0.2$), which is considerably better than that attained (56.3%) using the pool dataset ML model as starting point. Although the combined testing set performance is better using the previous variant of the MAPMCE algorithm, one would rather choose to use the later version due to the considerably improved minimum and female testing set performance.

Table 4.12 summarizes the MAPMCE results using both the small and full female training sets. Large improvements in the minimum, combined and female testing set performances are attained for both the small and full female training sets. Using the MAP point estimate as starting point produced the best female and minimum testing set performances. Maximum combined testing set performance was, however, attained

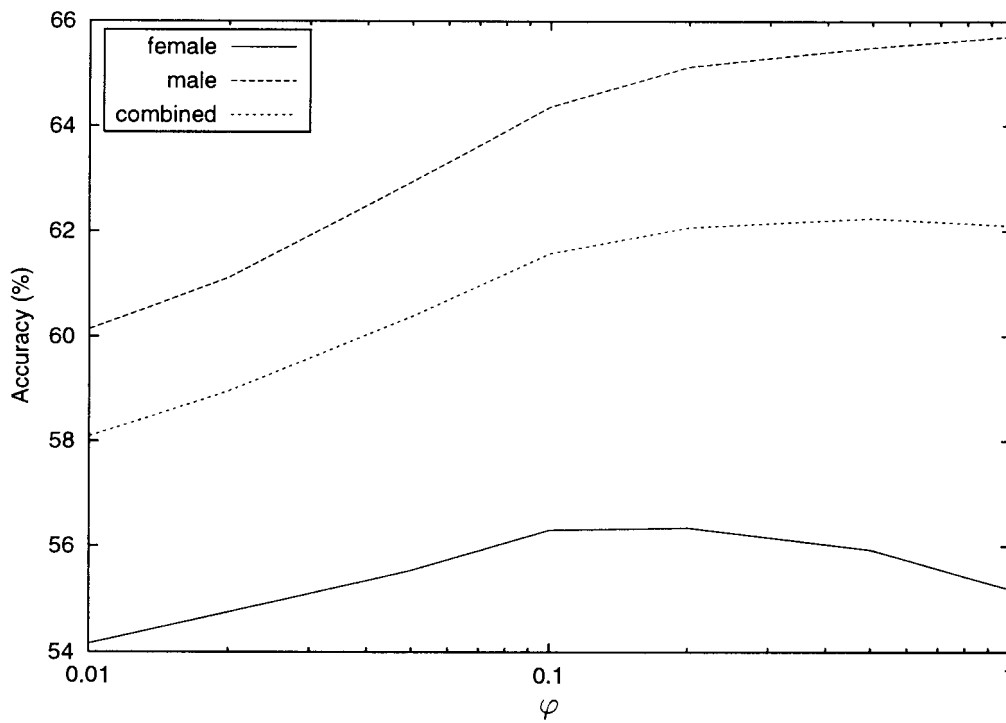


Figure 4.14: MAPMCE adaptation results for the small female training set, using the full male training set as prior (T_M) with the pooled ML model as starting point

when the ML model trained using the pooled data was used as starting point.

Comparison

Table 4.13 summarizes the best adaptation results (female testing set accuracy) obtained using the different algorithms for the TIMIT dataset. The MAP and GMAP algorithms did not improve female testing set performance when the full female training set was used. MAP improved the female test set accuracy to a limited degree when the small female training set was available. The MAPMCE algorithm performed best, with relative improvements in error rate of 8.1% and 13.1% resulting for the small T_{FS} and full T_F female training sets respectively.

Table 4.14 summarizes the best training results (minimum testing set accuracy) obtained using the different algorithms for the TIMIT dataset. The MAP and GMAP

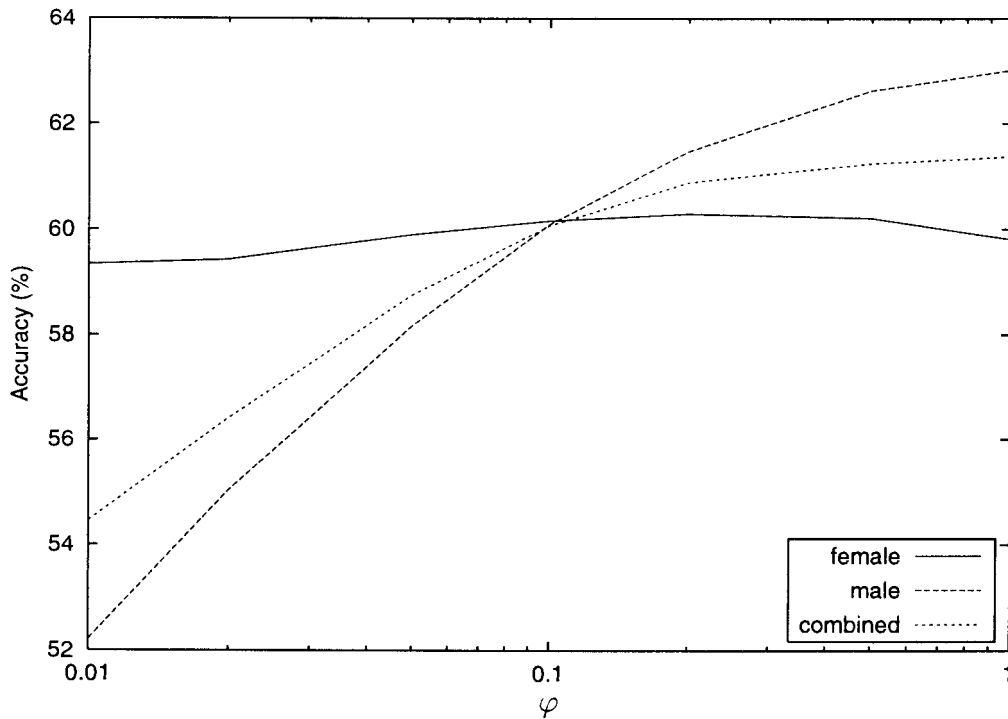


Figure 4.15: MAPMCE adaptation results for the small female training set, using the full male training set as prior (T_M) with the best MAP estimate as starting point

algorithm produced reasonable improvements in error rate when the small female training set was used (15.8% and 16.1% respectively). Little or no improvement was attained using these two algorithms when the full female training set was used. MAPMCE, once again produced marked improvements in the minimum testing set accuracy, with 25.5% and 14.5% relative improvement in error rate for the small and full female training sets respectively.

4.6.3 TIDIGITS

In this section, the three algorithms (MAP, GMAP and MAPMCE) will be compared within the framework of a continuous digit recognition task. As in the previous section, a situation is created where the number of female speakers is limited. The adaptation algorithms are therefore required to improve female test set performance, as well as for

Table 4.12: MAPMCE results for gender adaptation experiments using TIMIT

Description	Test set accuracy (%)			
	Male	Female	Combined	Minimum
MAPMCE T_{FS} ($T_{FS} + T_M$), $\varphi = 0.2$	65.1	56.3	62.1	56.3
MAPMCE T_{FS} ($T_{FS} + T_M$), $\varphi = 0.5$	65.5	55.9	62.2	55.9
MAPMCE T_{FS} ($T_{FS} + T_M$), $\varphi = 1.0$	65.7	55.2	62.1	55.2
MAPMCE T_{FS} (MAP) ($\varphi = 0.2$)	61.5	60.3	60.9	60.3
MAPMCE T_{FS} (MAP) ($\varphi = 1.0$)	63.0	59.8	61.4	59.8
MAPMCE T_F (MAP) ($\varphi = 0.2$)	60.0	66.1	61.8	60.0
MAPMCE T_F (MAP) ($\varphi = 1.0$)	62.8	65.1	63.4	62.8

the combined and minimum test set performance.

Table 2.6 presents the training and testing sets used in the experiments conducted in this section. The standard TIDIGITS training set comprises 77 digit sequences from each of the 57 female speakers in the dataset and 55 male speakers. There is therefore not an imbalance in the numbers of male and female speakers. Two subsets of the female training set are therefore created, so as to simulate the situation where relatively little data from female speakers is available.

The first subset (T_{WS}) consists of five female speakers selected at random from the complete female training set and all of the associated 77 digit sequences. The total duration of the resultant training set (T_{WS}) is 10.1 minutes long, which is approximately 7.7% the size of the full female training set. The second, smaller female training set (T_{WVS}) is a subset of the first, where the digit sequences of each of the speakers has been reduced to one each of the 1,2,3,4,5 and 7 digit sequences per speaker (randomly selected). The number of speakers is also reduced to 4, resulting in a female training set which is 55 seconds in duration and approximately one tenth the size of the first female training subset.

Table 4.13: Summary of the best adaptation results obtained for the TIMIT dataset. Relative improvement in error rate is given in braces.

Description	Female test set accuracy (%)
Baseline ML T_{FS}	56.8 (0.0%)
MAP T_{FS} ($T_{FS} + T_M$), $\nu = 100$	57.4 (1.4%)
GMAP T_{FS} ($T_{FS} + T_M$), $\varphi = 0$	56.8 (0.0%)
MAPMCE T_{FS} (MAP), $\varphi = 0.2$	60.3 (8.1%)
Baseline ML T_F	61.0 (0.0%)
MAP T_F ($T_F + T_M$), $\nu = 50$	61.0 (0.0%)
GMAP T_F ($T_F + T_M$), $\varphi = 0$	61.0 (0.0%)
MAPMCE T_F (MAP), $\varphi = 0.2$	66.1 (13.1%)

The second subset (T_{WVS}) is extremely small, but it will be shown that reasonable performance can be realized even for such a small training set. Note that although the female training sets are small the full male training set is assumed to be available for training purposes.

As in the last section, the algorithms will be compared by their ability to adapt models (female testing set performance) and performance of gender independent training (combined and minimum test set accuracy).

Table 4.16 details the performance of an 8 state, 5 mixture ML trained HMM trained using the datasets described above. The results using the full training sets are presented as a reference for the experimental results using the reduced datasets.

Reasonable accuracies are realized when only the small female training set (T_{WS}) is available. Peak female testing set accuracy of 95.6% is attained when using a model trained using only the small female training set. Pooling the small female training set and the full male training set ($T_M + T_{WS}$) decreases the performance for the female testing set slightly, but improves male, combined and minimum testing set performances markedly.

Table 4.14: Summary of the best minimum test set accuracy obtained for the TIMIT dataset. Relative improvement in error rate is given in braces.

Description	Minimum test set accuracy (%)
Baseline ML $T_{FS} + T_M$	46.7 (0.0%)
MAP T_{FS} ($T_{FS} + T_M$), $\nu = 700$	55.1 (15.8%)
GMAP T_{FS} ($T_{FS} + T_M$), $\varphi = 0.02$	55.3 (16.1%)
MAPMCE T_{FS} (MAP), $\varphi = 0.2$	60.3 (25.5%)
Baseline ML T	56.5 (0.0%)
MAP T_F ($T_F + T_M$), $\nu = 50$	56.5 (0.0%)
GMAP T_F ($T_F + T_M$), $\varphi = 0.8$	56.7 (0.5%)
MAPMCE T_F (MAP), $\varphi = 1.0$	62.8 (14.5%)

Using the very small female training set (T_{WVS}) results in poor results (25.2% accuracy for the female testing set). Pooling with the male training set, once again greatly improves results. It is, however, interesting to note that the ML trained model using only the male training set results in better results.

When computing the relative improvement obtained using the algorithms, the best results obtained when using only the available training data will be used. Table 4.17 shows the best ML results that can be achieved for the two scenarios, where the amount of female training data is limited. Note that, as mentioned, the best results obtained when the very small female training set is available is that of the ML model trained using only the male training set and the associated results are therefore used.

As discussed, the combined performance is not a good measure of the suitability of a particular model or algorithm. The ML model trained using only the male (or female) training set is a good example of this problem, a combined performance of 90.6% seems reasonable, but there is a large difference between female and male testing set performance (97.4% vs. 84.0%), which is not desired.

Table 4.15: Training and testing sets used with the TIDIGIT database

Description	Label	Number of speakers			Digit sequences	Duration (minutes)
		Male	Female	Total		
Training sets:						
Full (standard)	T	55	57	112	8623	253.4
Man	T_M	55	0	55	4235	121.9
Woman	T_W	0	57	57	4389	131.5
Woman-small	T_{WS}	0	5	5	385	10.1
Woman-very-small	T_{WVS}	0	4	4	28	55 s
Testing set		55	57	113	8623	254.4

MAP

The standard MAP algorithm presented in Section 4.3 is now tested within connected digit recognition task using the TIDIGITS database. As with the previous databases used to test the MAP algorithm, a pooled prior was found to perform better than that created from only the reference (male) training set. The prior used in all experiments presented in this section is therefore created using the pooled dataset containing the male training set and the relevant female training set.

Figure 4.16 shows the MAP adaptation results for the scenario where the small female training set (T_{WS}) is available. The pooled data set ($T_{WS} + T_M$) is used to create the prior. Peak female testing set performance of 97.4% is attained for $\nu = 200$, which also results in a peak combined performance of 97.2%. The peak minimum accuracy of 97.2% occurs at $\nu = 300$, relatively close to the combined performance peak.

Figure 4.17 presents the MAP adaptation results for the scenario where the smallest female training set (T_{WVS}) is available. The amount of female adaptation data is extremely small, and as a result the female testing set performance cannot be improved to the extent that it is better than that of the male testing set (for the specific config-

Table 4.16: Base system results for an 8 state, 5 mixture HMM using TIDIGITS

Training set	Test set accuracy (%)			
	Male	Female	Combined	Minimum
Full training set (T)	97.5	98.4	97.9	97.5
Man, training set (T_M)	97.4	84.0	90.6	84.0
Woman, training set (T_W)	90.3	98.8	94.6	90.3
Woman, small training set (T_{WS})	78.0	95.6	87.0	78.0
$T_M + T_{WS}$	97.6	94.4	95.9	94.4
Woman, very small set (T_{WVS})	15.6	25.2	20.5	15.6
$T_M + T_{WVS}$	91.3	80.1	85.6	80.1

Table 4.17: Best base system results, given the available data

Available data	Test set accuracy (%)			
	Male	Female	Combined	Minimum
T_M, T_{WS}	97.6	95.6	95.9	94.4
T_M, T_{WVS}	97.4	84.0	90.6	84.0

uration). Peak female, combined and minimum testing set accuracies of 84.4%, 87.2% and 84.4% respectively are attained for $\nu = 200$. The female and minimum testing set performances are slightly better than that of the ML model trained using the male training set only (a relative improvement in error rate of 2.5%). The combined test set accuracy is, however, considerably worse than that obtained using the best ML model (-36.2%).

Table 4.18 summarizes the results obtained using the MAP adaptation algorithm for the two training datasets. Reasonable improvements in error rates are attained when the larger female training set (T_{WS}) is available. However, the MAP algorithm fails to significantly improve the minimum and female testing set accuracy when the smallest female training (T_{WVS}) set is used. Importantly, a decrease in combined performance

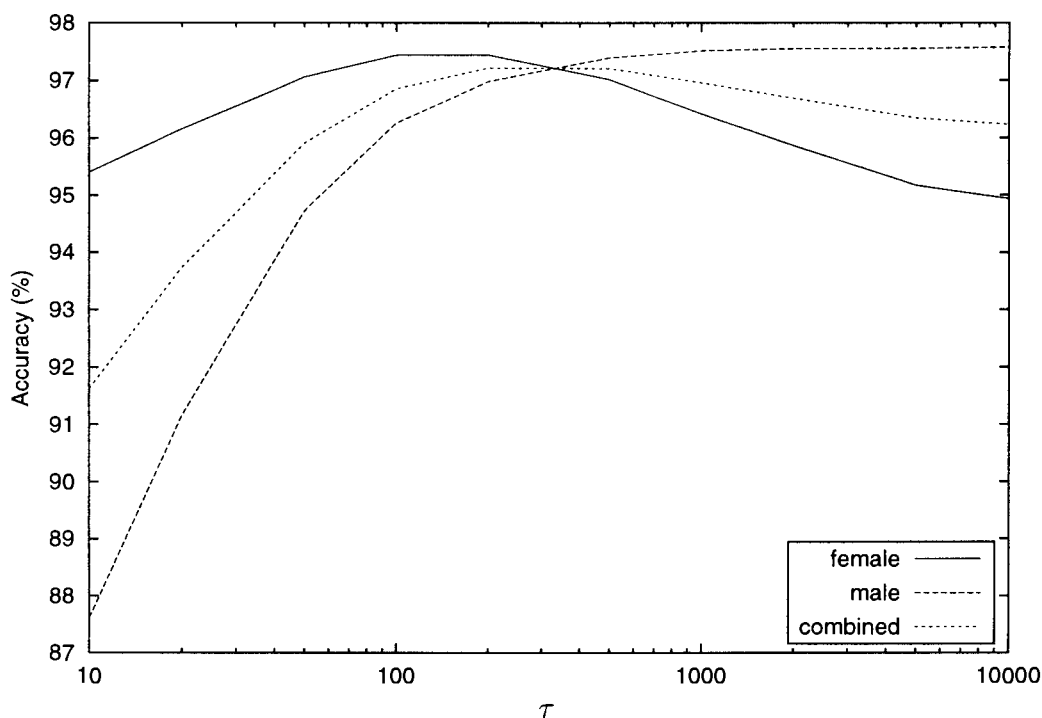


Figure 4.16: MAP adaptation results for an 8 state, 5 mixture HMM using the small female training set (T_{WS}) as the adaptation set and the pooled training set ($T_{WS} + T_M$) to create the prior.

results when using the MAP algorithm in the extreme situation, where the smallest female training set is used.

GMAP

The GMAP algorithm presented in Section 4.4 is tested using the TIDIGIT dataset given the conditions described earlier. The male speaker training set is used as the prior dataset for all experiments conducted in this section.

Figure 4.18 details the GMAP adaptation results for an 8 state, 5 mixture HMM when the small female training set (T_{WS}) is available. The algorithm attains peak performance of 97.6% for the female testing set when $\varphi = 0.05$, which is a 45.5% relative improvement in error rate. A slight, though insignificant improvement in accuracy for

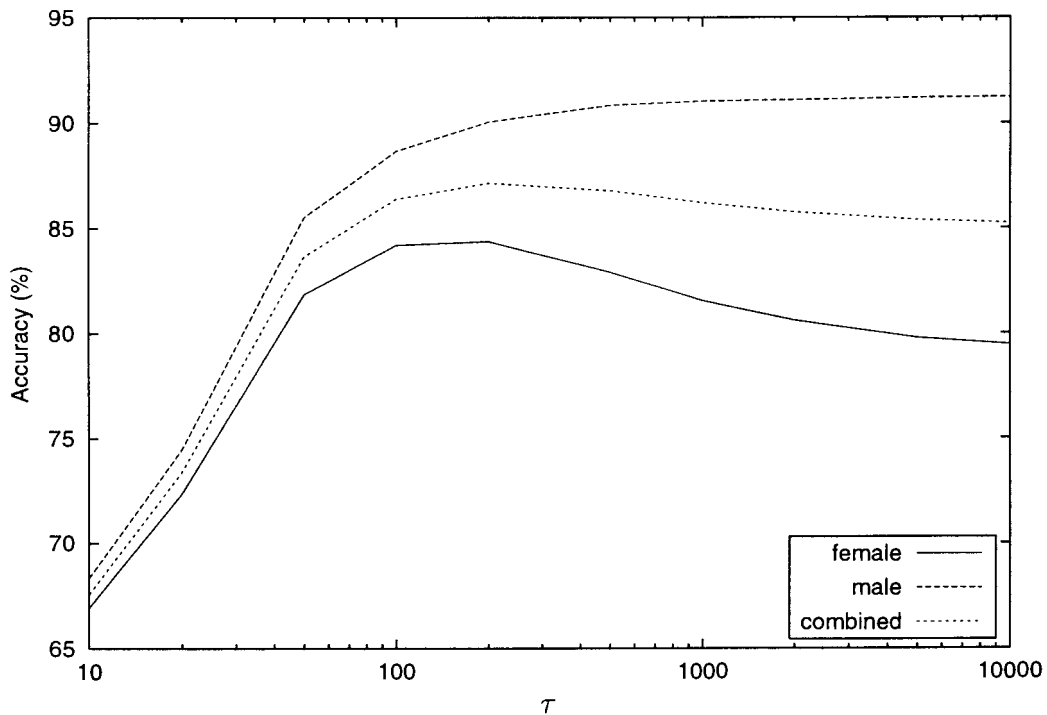


Figure 4.17: MAP adaptation results when the smallest female training set (T_{WVS}) is available

the male testing set can be attained if we set $\varphi = 0.5$. Peak combined and minimum testing set performances both occur when $\varphi = 0.2$, where a 97.5% accuracy is realized (for both). This equates to a 39% and 55.4% relative improvement in error rate for the combined and minimum testing set accuracies respectively.

The GMAP results for scenario where the very small female training set is available are presented in Figure 4.19. Here, as with MAP, the female testing set performance can not be improved to the degree that it is higher than that for the male testing set (for the range of φ presented). One would, however, expect the male training set performance to drop below that of the female testing set as $\varphi \rightarrow 0$, which is the ML estimate using only the female training data. A peak female testing set accuracy of 88.8% results for $\varphi = 0.01$, which is a 30.0% relative improvement in error rate (compared to the male training set (T_M) ML model).

Consequently, the minimum testing set performance peaks at the same point and accu-

Table 4.18: MAP results for an 8 state, 5 mixture HMM for TIDIGITS

Description	Test set accuracy (%)			
	Male	Female	Combined	Minimum
MAP 8,5 $T_{WS}(T_{WS} + T_M), \nu = 200$	97.0	97.4	97.2	97.0
MAP 8,5 $T_{WS}(T_{WS} + T_M), \nu = 300$	97.2	97.2	97.2	97.2
MAP 8,5 $T_{WS}(T_{WS} + T_M), \nu = 10^5$	97.6	94.9	96.2	94.9
MAP 8,5 $T_{WVS}(T_{WVS} + T_M), \nu = 200$	90.0	84.4	87.2	84.4
MAP 8,5 $T_{WVS}(T_{WVS} + T_M), \nu = 10^5$	91.2	79.5	85.2	79.5

racy as that of the female testing set (88.8% at $\varphi = 0.01$), which is also a 30.0% in error rate. The combined performance is a maximum at either $\varphi = 0.02$ or $\varphi = 0.01$ where a combined accuracy of 89.9% results (a relative increase in error rate of 7.4%). Here, the combined error rate is still worse than that attained using the ML model trained using the male dataset, though as discussed, this is not a good measure of model or algorithm performance.

Table 4.19 summarizes the results obtained using the GMAP algorithms for the various scenarios created using the TIDIGIT dataset. The results for the GMAP algorithm using the ML model of the pooled dataset ($T_{WVS} + T_M$) as a starting point are also presented (for the small female dataset). It is noticeable that here, the usage of the best MAP point estimate results in a small, but significant improvement in error rates (for the three criteria used to evaluate the algorithms).

MAPMCE

The MAPMCE algorithm presented in Section 4.5 is tested using the TIDIGIT dataset given the conditions described earlier in this section. The male speaker training set is used as the prior dataset for all experiments conducted in this section.

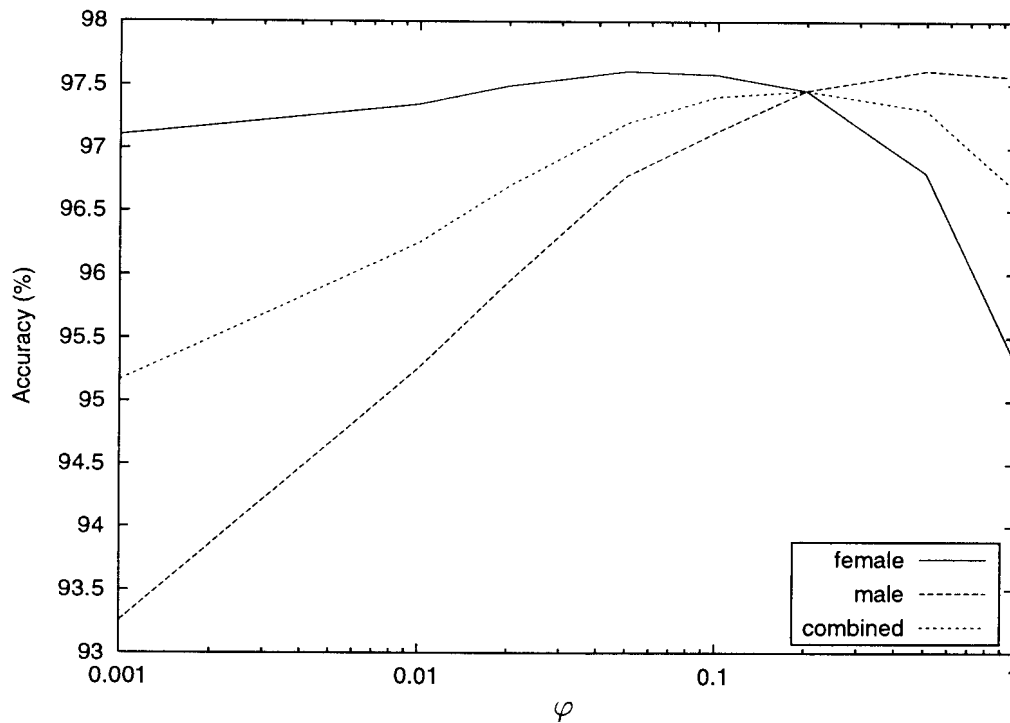


Figure 4.18: GMAP adaptation results for an 8 state, 5 mixture HMM using the small female training set (T_{WS}) as the adaptation set and the male training set (T_M) as prior. The best MAP point is used as starting point.

Figure 4.20 presents the MAPMCE results for an 8 state, 5 mixture HMM when the small female training set is available. The male training set is used as the prior dataset and the best MAP estimate is used as a starting point. It is interesting to note that the female testing set accuracy is above that of the male testing set for the range of φ presented. The male, female, combined and minimum test set accuracies peak at $\varphi = 0.5$, where their respective accuracies are 98.1%, 98.5%, 98.3% and 98.1%. The associated relative improvement in error rate is therefore 65.9% for the female testing set, 58.5% for the combined testing set performance and 66.1% for the minimum testing set accuracy.

The female testing set results are, however, considerably worse than that of the male testing set when using the MAPMCE algorithm for the scenario where the very small female training set is available, as shown in Figure 4.21. The peak female, combined and minimum testing set accuracies all occur at $\varphi = 0.05$. A peak female testing set

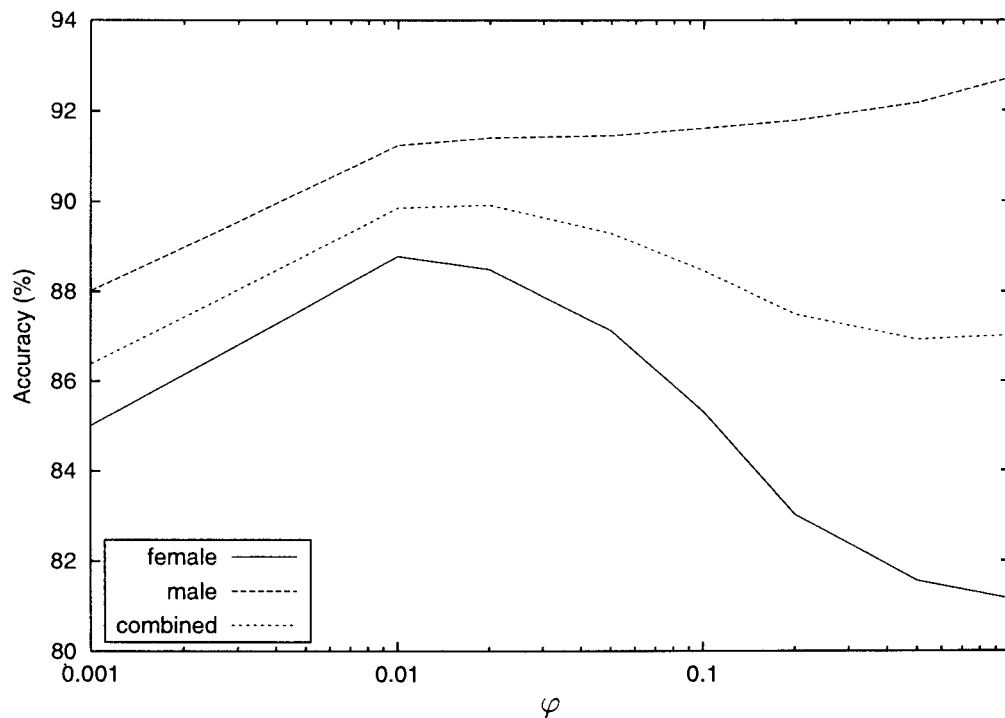


Figure 4.19: GMAP adaptation results for an 8 state, 5 mixture HMM using the very small female training set (T_{WVS}) as the adaptation set and the male training set (T_M) as prior. The ML estimate trained using the pooled training set is used as starting point.

accuracy of 93.0% is attained, which equates to a 56.3% relative improvement in error rate. Large improvements are also realized for the combined and minimum testing set performance measures, where a 95.1% combined testing set accuracy (47.9% relative improvement) and 93.0% minimum testing set accuracy (56.3% relative improvement) are attained.

Table 4.20 summarizes the results obtained using the MAPMCE algorithm under the two scenarios created using the TIDIGITS dataset. The results for the MAPMCE algorithm using the ML model of the pooled dataset ($T_{WS} + T_M$) as starting point are also presented (for the small female training set). Once again, as with the GMAP algorithm, using the best MAP point as the starting point has proved to be a better choice (compared to using the ML estimate of the pooled dataset). The results for the MAPMCE algorithm for the smallest female training set using the best MAP point as starting point have not been included as they are very similar to that obtained using

Table 4.19: GMAP results for an 8 state, 5 mixture HMM for TIDIGITS

Description	Test set accuracy (%)			
	Male	Female	Combined	Minimum
GMAP T_{WS} (MAP), $\varphi = 0.05$	96.8	97.6	97.2	96.8
GMAP T_{WS} (MAP), $\varphi = 0.2$	97.5	97.5	97.5	97.5
GMAP T_{WS} (MAP), $\varphi = 0.5$	97.6	96.8	97.3	96.8
GMAP T_{WS} ($T_{WS} + T_M$), $\varphi = 0.05$	97.4	97.4	97.4	97.4
GMAP T_{WS} ($T_{WS} + T_M$), $\varphi = 1.0$	97.6	95.4	95.9	95.4
GMAP T_{WVS} ($T_{WVS} + T_M$), $\varphi = 0.01$	91.2	88.8	89.9	88.8
GMAP T_{WVS} ($T_{WVS} + T_M$), $\varphi = 1.0$	92.7	81.2	87.0	81.2

the ML model of the pooled dataset.

Table 4.20: MAPMCE results for an 8 state, 5 mixture HMM for TIDIGITS

Description	Test set accuracy (%)			
	Male	Female	Combined	Minimum
MAPMCE $T_{WS}(T_{WS} + T_M)$, $\varphi = 0.05$	98.7	97.7	98.2	97.7
MAPMCE $T_{WS}(T_{WS} + T_M)$, $\varphi = 0.5$	98.7	97.2	97.9	97.2
MAPMCE T_{WS} (MAP), $\varphi = 0.5$	98.1	98.5	98.3	98.1
MAPMCE $T_{WVS}(T_{WVS} + T_M)$, $\varphi = 0.05$	97.9	93.0	95.1	93.0
MAPMCE $T_{WVS}(T_{WVS} + T_M)$, $\varphi = 1.0$	97.9	91.1	93.8	91.1

Comparison

Table 4.21 summarizes the female testing set results obtained for the TIDIGITS dataset using the three algorithms which have been evaluated in this section. The best minimum testing set accuracies are summarized in Table 4.22. Looking at the normal results, it is evident that the MAPMCE algorithm performs best, with the GMAP algorithm performing well.

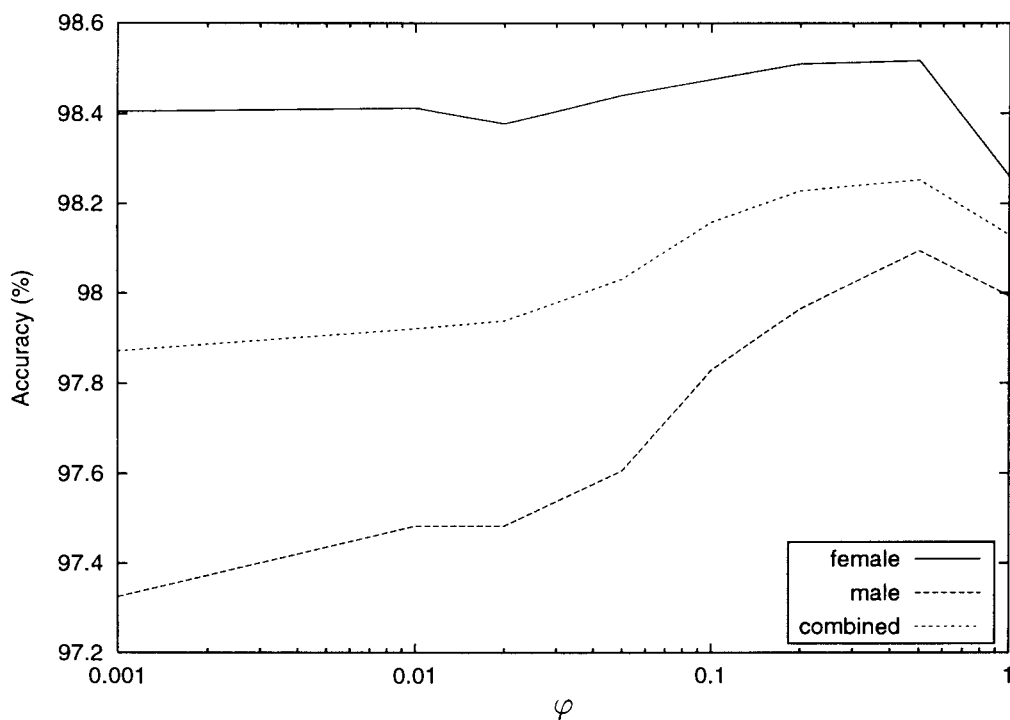


Figure 4.20: MAPMCE adaptation results for an 8 state, 5 mixture HMM using the small female training set (T_{WS}) as the adaptation set and the male training set (T_M) as prior. The best MAP point is used as starting point.

The MAP algorithm, although it does not perform as well as the other algorithms, does manage reasonable improvements in testing set performances for the larger, small female training set. It does not, however, significantly improve results when the very small female training set is used. The GMAP and MAPMCE algorithms produce far better results under these extreme conditions.

Duration modeling Duration modeling is a technique, which is often used to improve recognition accuracy in continuous digit recognition applications. It is therefore important to determine the effect of duration modeling on the performance of the algorithms tested. Note that the MAPMCE (and GMAP) algorithm can be used to estimate the duration modeling parameters. So as to ensure a fair comparison, this potential improvement has not been used and the duration parameters are therefore estimated with the ML algorithm using the relevant model and available training data.

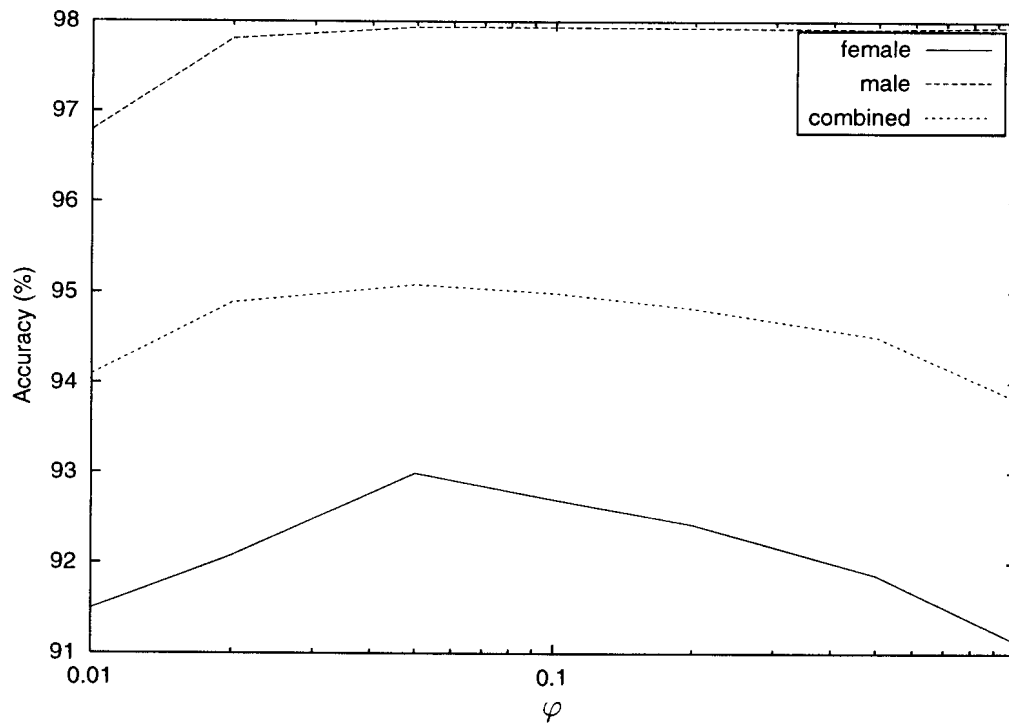


Figure 4.21: MAPMCE adaptation results for an 8 state, 5 mixture HMM using the very small female training set (T_{WVS}) as the adaptation set and the male training set (T_M) as prior. The best MAP point is used as starting point.

It is evident, from tables 4.21 and 4.22 that MAPMCE does not perform as well when used with duration modeling, and even results in a decrease in accuracy for the minimum test set accuracy measure when used with the small female training set (compared to the baseline system with duration modeling).

The MAP algorithm works better, both in absolute and relative terms, when using duration modeling. The GMAP algorithm, however, works best of all when used with duration modeling, resulting in recognition performances which are considerably better than that realized by the MAPMCE and MAP algorithms.

The MCE (MAPMCE) algorithm, due to its discriminative nature, incorporates duration information into the HMM without us explicitly modeling it. This is as a result of the MCE algorithm trying to reduce insertions and deletions. Much of the improvements realized when using the MCE algorithm therefore overlap with those attained

Table 4.21: Comparison of best results for the female test set, with and without duration modeling. The relative improvement in error rate compared to the baseline ML accuracy is given in brackets.

Description	Female test set accuracy (%)	
	Normal	Duration
ML $T_M + T_{WS}$	95.6 (0.0%)	97.6 (0.0%)
MAP $T_{WS} (T_{WS} + T_M)$, $\nu = 200$	97.4 (40.9%)	98.5 (37.5%)
GMAP $T_{WS} (MAP)$, $\varphi = 0.05$	97.6 (45.5%)	98.5 (37.5%)
MAPMCE $T_{WS} (MAP)$, $\varphi = 0.5$	98.5 (65.9%)	98.5 (37.5%)
ML T_M	84.0 (0.0%)	91.8 (0.0%)
MAP $T_{WVS} (T_{WVS} + T_M)$, $\nu = 200$	84.4 (2.5%)	94.6 (34.1%)
GMAP $T_{WVS} (T_{WVS} + T_M)$, $\varphi = 0.01$	88.8 (30.0%)	96.3 (54.9%)
MAPMCE $T_{WVS} (T_{WVS} + T_M)$, $\varphi = 0.05$	93.0 (56.3%)	94.4 (31.7%)

when using duration modeling. It is for this reason that the MAPMCE algorithm does not work as well when duration modeling is used.

4.6.4 Summary of results

In Section 4.6.1 the three algorithms, MAP, GMAP and MAPMCE were used within a language adaptation framework. The algorithms were used to adapt English seed data or models to create a language specific recognizer for Afrikaans. Either the full Afrikaans training set or a reduced subset thereof was used as the adaptation set.

The MAPMCE algorithm performed best when the full Afrikaans dataset was used, resulting in relative improvements in error rate of up to 6.8% and 3.1% for a 5 and 10 mixture HMM respectively (3 states). The GMAP algorithm resulted in reasonable performance increases, with a 5.3% and 2.1% relative improvement in error rate being attained for a 5 and 10 mixture HMM respectively. The standard MAP algorithm, resulted in similar, but slightly worse, improvements as that produced by the GMAP

Table 4.22: Comparison of best results for the minimum accuracy criterion, with and without duration modeling. The improvement in error rate due to the usage of duration modeling is also given.

Description	Minimum accuracy (%)	
	Normal	Duration
ML $T_M + T_{WS}$	94.4 (0.0%)	97.6 (0.0%)
MAP $T_{WS} (T_{WS} + T_M)$, $\nu = 200$	97.0 (46.4%)	97.9 (12.5%)
GMAP $T_{WS} (MAP)$, $\varphi = 0.2$	97.5 (55.3%)	98.4 (33.3%)
MAPMCE $T_{WS} (MAP)$, $\varphi = 0.5$	98.1 (66.0%)	97.5 (-4.2%)
ML T_M	84.0 (0.0%)	91.8 (0.0%)
MAP $T_{WVS} (T_{WVS} + T_M)$, $\nu = 200$	84.4 (2.5%)	94.6 (34.1%)
GMAP $T_{WVS} (T_{WVS} + T_M)$, $\varphi = 0.01$	88.8 (30.0%)	96.3 (54.9%)
MAPMCE $T_{WVS} (T_{WVS} + T_M)$, $\varphi = 0.05$	93.0 (56.3%)	94.4 (31.7%)

algorithm.

When the small Afrikaans training set was used, the MAPMCE algorithm did not perform as well as the GMAP algorithm, which produced relative improvements in error rate of 8.0% and 10.2% for a 5 and 10 mixture HMMs respectively. The MAP algorithm resulted in relatively poor improvements in accuracy (compared to the GMAP algorithm).

In Section 4.6.2 the three algorithms were evaluated within a gender adaptation framework, using the TIMIT dataset. The TIMIT dataset has an imbalance in the number of male and female speakers, and the male training set was therefore used as the prior dataset and the female dataset as the adaptation set. A smaller female dataset was also created, so as to determine the effect of the algorithms when an even larger imbalance exists. Two criteria were used to evaluate the algorithms, namely adaptation performance and training performance. The female testing set accuracy and the minimum of the female and male testing set accuracies were used to determine the performance for these two criteria.

Here, the MAPMCE algorithm performed best for both datasets and both testing criteria, resulting in marked improvements in accuracy. The MAP and GMAP algorithms resulted in little or no improvement in accuracy for the female testing set (adaptation). Reasonable improvements in the minimum testing set accuracy were attained by these two algorithms when the small female training set is used, with the GMAP algorithm performing slightly better than the standard MAP algorithm.

In Section 4.6.3 the algorithms were tested within a gender adaptation framework for a connected digit task. The TIDIGIT database was used for this purpose. Here, once again the MAPMCE algorithm worked best, followed by the GMAP algorithm and the MAP algorithm performing worst. However, when duration modeling is used, the GMAP algorithm results in the best performance, with relative improvements in digit accuracy of up to 54.9%.

When comparing MCE or MAPMCE to non-discriminative algorithms such as MAP and GMAP, one must be careful as they are often not directly comparable. MCE as mentioned, due to its discriminative nature, automatically tends to incorporate duration and language modeling information into the HMMs. This means that when duration modeling or language modeling is used, the effect thereof will be smaller when using MCE. McDermott [75] noted this, when he used a bigram language model with MCE and compared it to using a bigram model with ML trained HMMs.

Table 4.23 gives the execution times for the three algorithms (MAP, GMAP and MAPMCE) on a Pentium III 600MHz computer. The times given are for one iteration only; total execution time is equal to execution time given in Table 4.23 times the number of iterations used. Given that the GMAP and MAPMCE algorithms required more iterations (around 30) than that required by MAP (10 iterations), it is evident that the GMAP and MAPMCE algorithms are considerably more computationally expensive than the MAP algorithm. The improved performance, however, warrants the additional computational expense.

Table 4.23: Execution times (in minutes) for one iteration (update) of the three algorithms (MAP, GMAP and MAPMCE) on a Pentium III 600MHz computer. The prior data set (square brackets) and adaptation set (round brackets) are included in the dataset description.

Dataset	MAP	GMAP	MAPMCE
SUNSpeech (A) [E]	1.3	13.4	37.7
TIMIT (T_F) [T_M]	4.2	16.9	54.1
TIDIGITS (T_W) [T_M]	0.3	4.8	18.0

4.7 Summary

This chapter introduced Bayesian adaptation and its usage within a continuous speech recognition framework. The MAP algorithm of Gauvain and Lee [45] was described. A gradient-based MAP algorithm which makes no assumption about the form of the prior was proposed and the implementation thereof was discussed. A Bayesian inspired MCE-based adaptation algorithm (MAPMCE) was also proposed. The MAPMCE algorithm is an extension of the MCE algorithm and the implementation thereof is relatively simple.

The three algorithms were experimentally evaluated in Section 4.6, using the SUN-Speech database for language adaptation, and the TIMIT and TIDIGIT databases for gender adaptation experiments. On the whole, the MAPMCE algorithm proved to work best, with the GMAP algorithm performing reasonably well. The MAP algorithm, in general, resulted in considerably worse performance than either the GMAP or MAPMCE algorithms.

Chapter 5

Bayesian learning

This chapter develops a Bayesian approach to learning for HMMs in speech recognition. Markov chain Monte Carlo methods which can be used to numerically integrate the posterior distribution as required by the Bayesian learning approach form part of this discussion. The implementation of Bayesian learning for HMMs in speech recognition is presented, including the requirement of maintaining the original HMM constraints, choice of prior and utterance recognition. This work shows that the Bayesian learning approach can be successfully applied to complex models when the amount of training data is small. This is contrary to the notion that one must limit the complexity of the model when training data is limited, as was discussed in Section 2.2.1.

This work was inspired by the work of Neal [82], who proposed a Markov chain Monte Carlo based Bayesian learning procedure for neural networks. In this chapter, the Bayesian learning procedure used by Neal will be implemented and adapted for usage with hidden Markov models and speech recognition. Previous applications of Bayesian learning in speech recognition have concentrated on using approximations. One such approximation is that of Huo *et al.* [54, 58] who used a Gaussian distribution to approximate the posterior distribution. The MAP estimation algorithm (Chapter 4) can also be regarded as an approximate Bayesian approach. However, my formulation

is the first Markov chain Monte Carlo based Bayesian learning approach used for hidden Markov model speech recognition systems.

The Bayesian learning framework is introduced here and past work both within the field of speech recognition and in the more general field of neural networks is discussed.

5.1 Introduction

Bayesian methods can be used for the inference of parameter values in a model given the data. Bayesian methods have also been used for the purpose of model comparison. David Mackay [72] focused primarily on the usage of Bayesian methods for the comparison and training of neural networks. Most people would include the above two uses of Bayesian methods in the data modeling process.

The remainder of this section will summarize the relevant Bayesian theory used in this chapter. Certain sections from Chapter 4 have been reproduced for readability. For a more complete introductory text on Bayesian statistics, the reader is referred to Box and Tiao [13], DeGroot [27] and Bishop [12]. The theory and discussions in this section will be biased towards speech recognition applications of Bayesian learning.

5.1.1 Bayes' theorem

The fundamental concept of Bayesian analysis is that the plausibilities of alternative hypotheses are represented by probabilities, with inference being performed by evaluating these probabilities.

Given a vector $\mathbf{y} = (y_1, \dots, y_n)$ of n observations, with probability distribution $P(\mathbf{y}|\theta)$, which depends on the k parameters $\theta^T = (\theta_1, \dots, \theta_k)$. The parameter vector θ has the probability distribution $P(\theta)$. Given the observed data, the conditional distribution of

θ is

$$P(\theta|\mathbf{y}) = \frac{P(\mathbf{y}|\theta)P(\theta)}{P(\mathbf{y})}. \quad (5.1)$$

The denominator in Eq. (5.1), $P(\mathbf{y})$, is a normalizing factor, which ensures that the integral of $P(\theta|\mathbf{y})$ is equal to one. It can be written as follows:

$$P(\mathbf{y}) = \int P(\mathbf{y}|\theta)P(\theta)d\theta. \quad (5.2)$$

Equation (5.1) is referred to as Bayes' theorem. The distribution $P(\theta)$, is called the *prior* distribution and expresses what is known about the model parameters before any data is observed. The *posterior* distribution $P(\theta|\mathbf{y})$, tells us what is known about the model parameters given that data has been observed. In what follows, the prior distribution and posterior distribution will again sometimes simply be referred to as the “prior” and “posterior” respectively.

The distribution $P(\mathbf{y}|\theta)$ is often referred to as the data *likelihood* and can be written $L(\theta|\mathbf{y})$. This is valid when $P(\mathbf{y}|\theta)$ is regarded as a function of \mathbf{y} and not of θ .

In many Bayesian methods, the normalizing constant is not necessary and Eq. (5.1) is written as

$$P(\theta|\mathbf{y}) \propto L(\theta|\mathbf{y})P(\theta). \quad (5.3)$$

5.1.2 Bayesian learning and prediction

The result of Bayesian learning is a probability distribution (posterior) which expresses our beliefs of how likely individual parameters values are. This is the basis for Bayesian

learning, as it allows learning to be performed using probability theory.

In a Bayesian approach to HMM parameter estimation and recognition, the objective is to find a predictive distribution for an unknown utterance, given the observations of the utterance, as well as the training observations. Let the observations for the i th utterance be written as \mathbf{O}_i . For n training utterance examples $\mathbf{O} = (\mathbf{O}_1, \dots, \mathbf{O}_n)$, Bayes' theorem (Eq. (4.1)) can be written as

$$P(\theta|\mathbf{O}) = \frac{P(\mathbf{O}|\theta)P(\theta)}{P(\mathbf{O})} \quad (5.4)$$

$$\propto P(\mathbf{O}|\theta)P(\theta).$$

Assuming independence of the observations we can write the likelihood as follows:

$$P(\mathbf{O}|\theta) = \prod_{i=1}^n P(\mathbf{O}_i|\theta). \quad (5.5)$$

In a Bayesian framework, when we wish to classify an unknown input, we need to calculate the following probability,

$$P(\mathbf{O}_{unknown}|\mathbf{O}_1^{(i)}, \dots, \mathbf{O}_n^{(i)}) = \int P(\mathbf{O}_{unknown}|\theta)P(\theta|\mathbf{O}_1^{(i)}, \dots, \mathbf{O}_n^{(i)})d\theta, \quad (5.6)$$

where i is the class and $\mathbf{O}_{unknown}$ is the unknown observation. The classifier decision C is the class resulting in the highest value of Eq. (5.6), i.e.

$$C(\mathbf{O}_{unknown}) = i \quad \text{where} \quad i = \underset{j}{\operatorname{argmax}} P(\mathbf{O}_{unknown}|\mathbf{O}_1^{(j)}, \dots, \mathbf{O}_n^{(j)}), \quad (5.7)$$

where $C(\mathbf{O}_{unknown})$ is the classifier's decision for the unknown observation.

Unfortunately, due to the nature of the incomplete data problem caused by the underlying hidden processes of an HMM, the evaluation of Eq. (5.6) is non-trivial. If, however, the posterior ($P(\theta|\mathbf{O})$) is well approximated by a Gaussian, then Eq. (5.6) can be approximated as follows [72]:

$$P(\mathbf{O}_{unknown}|\mathbf{O}_1^{(i)}, \dots, \mathbf{O}_n^{(i)}) \approx P(\mathbf{O}_{unknown}|\theta_{MAP})P(\theta_{MAP})(2\pi)^{D/2}|A|^{1/2} \quad (5.8)$$

where θ is D -dimensional, θ_{MAP} is the maximum *a-posteriori* point (mode of the posterior $P(\theta|\mathbf{O})$) and A is the modal dispersion matrix, i.e., $A = -V^{-1}$, where V is the Hessian matrix of second derivatives of the log of the posterior evaluated at the mode of the posterior.

This approximation has been used extensively in Bayesian approaches. MacKay [72] used this approximation in his work on model selection for neural networks. Huo *et al.* [54, 58] proposed a *quasi-Bayesian predictive classification* approach for continuous density HMMs, in which they used this approximation. Another approximation was also proposed by Huo *et al.* [57], which used the following Viterbi-based approximation,

$$P(\mathbf{O}_{unknown}|\mathbf{O}_1, \dots, \mathbf{O}_n) \approx \max_{\mathbf{q}, \mathbf{l}} \int P(\mathbf{O}_{unknown}|\theta)P(\theta|\mathbf{O}_1, \dots, \mathbf{O}_n)d\theta, \quad (5.9)$$

where \mathbf{q} is a state sequence and \mathbf{l} is the sequence of associated mixture components. A modified Viterbi algorithm was used to compute the above approximation.

We can, however, use Monte Carlo (MC) methods to obtain a better approximation of Eq. (5.6) than Eq. (5.8) or Eq. (5.9). MC methods make no assumption concerning the form of the distribution, as done in the above approximations. In theory, MC methods can approximate Eq. (5.6) for complex distributions with multiple modes, as well as distributions for which the dominant contribution of the integral results from areas in parameter space which are not near a mode. Markov chain Monte Carlo methods [12, 82] will therefore be used in this implementation and will be described in Section

5.2.

In the field of neural networks, Neal [82] used the “Hybrid Monte Carlo” method (described later in Section 5.2.3) for Bayesian learning. The following are some of the applications using Monte Carlo algorithms that have been reported for speech recognition or speech processing. Vermaak and Niranjan [114] used a Markov chain Monte Carlo algorithm for speech enhancement. Robert *et al.* [98] presented a Markov chain Monte Carlo strategy to obtain a marginal MAP estimate. Godsill and Andrieu [48] used Markov chain Monte Carlo methods for the separation and recovery of convolatively mixed autoregressive processes. We will, however, use Markov chain MC methods to implement Bayesian learning for HMM speech recognizers.

5.1.3 Bias/variance problem

Let us once again look at the bias/variance problem discussed in Section 2.2.1 as it is central to the Bayesian learning approach. Integrating over the posterior, as in Eq. (5.6), results in the variance term of Eq. (2.25) being greatly reduced. Adjusting the complexity of models based on the amount of training data in a Bayesian framework therefore makes little sense as the variance term effectively disappears (except for extreme sparse training data). More complex models, which perform worse when using a single point estimate (ML), will therefore perform better than less complex models in a Bayesian implementation. We will, however, prefer simpler models due to other reasons, such as computational complexity.

Numerical integration with respect to the posterior $P(\theta|\mathbf{O})$ using a fixed number of samples N will, however, increase the variance of the solution. The effect of the number of samples used will be investigated in the experimental section later in this chapter.

5.1.4 Hierarchical models

Hidden Markov models are relatively complex and have many parameters to estimate. It is, therefore useful to specify the joint distribution of some of these parameters in terms of a common hyperparameter γ which has a prior distribution of its own. This is known as a *hierarchical model*.

The prior distribution $P(\theta)$ can then be written in terms of the hyperparameters as follows (assuming independence),

$$P(\theta) = \int P(\gamma) \prod_{i=1}^D P(\theta_i|\gamma) d\gamma \quad (5.10)$$

where $P(\gamma)$ is the prior distribution of the hyperparameter γ , $P(\theta_i|\gamma)$ is the prior distribution of the parameter θ_i given the hyperparameter.

A hierarchical model, if well formulated, can be considerably more intelligible than using a direct prior distribution. We can also in this way, incorporate vague heuristic information into the prior, as will be done in Section 5.3.

5.2 Monte Carlo methods

As mentioned in Section 5.1.2, we want to evaluate Eq. (5.6), which is the expectation of the function $P(\mathbf{O}_{unknown}|\theta)$ with respect to the posterior distribution $P(\theta|\mathbf{O}_1, \dots, \mathbf{O}_n)$. Such expectations can be estimated using Monte Carlo methods, by summing $P(\mathbf{O}_{unknown}|\theta^{(j)})$ using N samples of θ (the i th sample denoted by $\theta^{(i)}$) generated from the posterior distribution $P(\theta|\mathbf{O})$ for $j = \{1, \dots, N\}$, i.e.

$$P(\mathbf{O}_{unknown}|\mathbf{O}_1, \dots, \mathbf{O}_n) \approx \sum_{j=1}^N P(\mathbf{O}_{unknown}|\theta^{(j)}), \quad (5.11)$$

where the samples $\theta_1, \dots, \theta_N$ are generated by a process such that the distribution thereof is that of the posterior $P(\theta|\mathbf{O})$.

The sampling methods described in this section were developed for situations where probability distribution cannot be directly sampled. Sampling from a one-dimensional Gaussian distribution can, for example, be done directly.

Suppose we wish to generate a sample from a distribution $P(\theta)$ for $\theta \in \Theta$, but cannot do so directly. This can be done by constructing a Markov chain with state space Θ , which is easy to simulate, and whose equilibrium distribution is $P(\theta)$.

The following are sufficient conditions for such an algorithm to approach the desired distribution [111]:

- Invariance with respect to the distribution P . If for all pairs of configurations θ and θ' ,

$$\frac{P(\theta \rightarrow \theta')}{P(\theta' \rightarrow \theta)} = \frac{P_\infty(\theta')}{P_\infty(\theta)} \quad (5.12)$$

and at step n we have $P_n(\theta) = P_\infty(\theta)$, then at step $n + 1$ we will have $P_{n+1}(\theta) = P_\infty(\theta)$. The desired distribution is therefore an equilibrium distribution. This condition is called the *detailed balance* condition, and any chain which satisfies it is said to be reversible. The resulting distribution $P(\theta)$ persists once established and is therefore invariant (or stationary).

- Ergodicity. This condition specifies that the probability distribution at step $n + 1$ should be closer to $P_\infty(\theta)$ than at step n . An algorithm which complies with condition will converge to its equilibrium condition from any initial configuration.

In the following sections, three Markov chain Monte Carlo methods will be described. The above conditions (detailed balance and ergodicity) will be used to determine the suitability of each method for the implementation of Bayesian learning. We will in

particular find that *Gibbs sampling* (Section 5.2.1) is not ergodic and will therefore not be used (directly) to sample the posterior distribution $P(\theta|\mathbf{O})$ of an HMM. The *stochastic dynamics* (Section 5.2.2) and *hybrid Monte Carlo* (Section 5.2.3) methods meet both of the above conditions and we will therefore be able to use either of these procedures.

5.2.1 Gibbs sampling

Gibbs sampling [82, 7] can be used to sample the distribution of a multi-dimensional parameter. Gibbs sampling is also known as the *heatbath* method in the physics literature.

In Gibbs sampling a Markov chain is simulated, in which the new n -dimensional sample $\theta^{(t+1)}$ is generated from $\theta^{(t)}$ using the following iterative procedure:

Generate $\theta_1^{(t+1)}$ from the conditional distribution of θ_1 given $\theta_2^{(t)}, \theta_3^{(t)}, \dots, \theta_n^{(t)}$.
 Generate $\theta_2^{(t+1)}$ from the conditional distribution of θ_2 given $\theta_1^{(t+1)}, \theta_3^{(t)}, \dots, \theta_n^{(t)}$.
 \vdots
 Generate $\theta_i^{(t+1)}$ from the conditional distribution of θ_i given $\theta_1^{(t+1)}, \dots, \theta_{i-1}^{(t+1)}, \theta_{i+1}^{(t)}, \dots, \theta_n^{(t)}$.
 \vdots
 Generate $\theta_n^{(t+1)}$ from the conditional distribution of θ_n given $\theta_1^{(t+1)}, \dots, \theta_{n-1}^{(t+1)}$.
 Generate $\theta_1^{(t+2)}$ from the conditional distribution of θ_1 given $\theta_2^{(t+1)}, \theta_3^{(t+1)}, \dots, \theta_n^{(t+1)}$.
 \vdots
 and so on.

The transition resulting from the above steps being executed leaves the desired distribution Q invariant and is reversible. Gibbs sampling is, however, not necessarily an ergodic Markov chain and depending on the situation will not converge to its equilibrium distribution if we do not start from the desired distribution.

Djurić and Chun [34] proposed a method of estimating non-stationary (duration modeling) discrete hidden Markov models. The posterior of the model parameters was sampled using a Gibbs sampler. Their implementation was tested on an extremely simple HMM consisting of 3 states and with 5 possible emission variables. Histograms of the posterior samples for certain parameters were presented, showing that the mode of the posterior was reasonably accurate.

Gibbs sampling is dependent on being able to sample the distribution of one parameter conditional on the other parameters. For continuous density HMMs, the conditional distribution of one parameter given the values of the other parameters is non-trivial, where any parameter of a given state is dependent on all the other parameters of all the states. Gibbs sampling is, therefore, not a suitable sampling method that can be used to sample the posterior of an HMM. The posterior of HMMs will be multi-modal and due to the fact that Gibbs sampling is not ergodic, it will also not necessarily converge to the correct distribution under these conditions.

We will not use Gibbs sampling in itself, but rather as part of the stochastic dynamics and hybrid Monte Carlo algorithms described later in this section.

5.2.2 The stochastic dynamics method

The *stochastic dynamics method* [82] for sampling of distributions, otherwise known as the *refreshed molecular dynamics method*, was introduced by Anderson [3] and applied to the field of quantum chromodynamics by Duane and Kogut [36, 38]. It is also sometimes known in the literature as the *hybrid method*, because it contains two standard update steps – uniformly sampling values of variables q and p which have a fixed total energy $H(q, p)$, and sampling states with different values of H . Here, the stochastic dynamics method is used to generate the samples of the posterior distribution $P(\theta|\mathbf{O})$, which will be used to numerically integrate Eq. (5.6).

Let us assume that we wish to sample from a distribution for a variable q , which has n dimensions. In the systems for which the techniques described in this section were developed, q is typically the coordinates of the particles in a physical system. In our work, q will be the HMM parameters, i.e. $q = \theta$. We therefore have a system with continuously valued coordinates q_i , with the probability of the variable q defined as

$$P(q) \propto e^{(-E(q))}, \quad (5.13)$$

where $E(q)$ is the potential energy function. Any non-zero probability function, can be written in this form by defining $E(q) = -\ln[P(q)] - \ln(Z)$ for $Z > 0$.

A *momentum* variable p is introduced which has n components, one for each of the components of q . Here, the kinetic energy is used which is

$$K(p) = \sum_{i=1}^n \frac{p_i^2}{2m_i}, \quad (5.14)$$

with m_i being the “mass” associated with each component. It will be assumed, for the rest of this discussion, that $m_i = 1$ for all parameters. The terms *state* and *configuration* will be used to indicate the combination of the system coordinates ($q = \theta$) and momenta (p), i.e. $state = configuration = (p, q)$. The total energy is $H(p, q) = E(q) + K(p)$ and the probability for q and p is therefore

$$P(p, q) \propto e^{-H(p, q)}. \quad (5.15)$$

Consider Hamilton’s equations of motion for this system,

$$\frac{\partial q_i}{\partial t} = \frac{\partial H}{\partial p_i} = p_i \quad (5.16)$$

$$\frac{\partial p_i}{\partial t} = -\frac{\partial H}{\partial q_i} = -\frac{\partial E}{\partial q_i}. \quad (5.17)$$

Three important properties of the Hamiltonian dynamics necessary for sampling are:

- The motion defined by the above equations leave the total energy H constant.
- The volume of regions of phase space is conserved, i.e. if we follow points in some region of volume V , we find that the region where these points end up after a given time T also has volume V . This is important as the probability of a configuration is really the probability density times the volume element in phase space.
- The motion is reversible. After having simulated Hamilton's equations for a time T , we can change the sign of the momenta, apply Hamilton's equations for the same period of time and end at the original starting point.

An update which consists of generating a random number R , multiplying all the momenta by -1 if $R < \frac{1}{2}$ and then integrating Hamilton's equations for a given time interval satisfies the conditions required by the fundamental theorem [111].

Since each p_i is independent of the q_i and the other p_i , the probability distribution of the momenta can be sampled by using Gibbs sampling and assuming each to be a Gaussian distributed random number. Note that, it is not necessary to randomly reverse the momenta as described above, as the Gibbs sampling is just as likely to generate p_i as $-p_i$. This is known as "refreshing" the momentum variable p .

The length of time T over which Hamilton's equations are integrated is a critical parameter which must be found. Figure 5.1 illustrates the path of the stochastic dynamic method for different integration times. If the time T is small, the coordinates

q will not change much and the result is effectively a random walk of parameter space (Figure 5.1 (a)). Alternately, a large integration time T will result in a path which is periodic in nature and we will waste our time in generating such a long path which could easily end up close to where we started (Figure 5.1 (b)). There is therefore an optimal value of T which lies somewhere between the two extremes (Figure 5.1 (c)). In certain situations this can be done analytically [39], but in realistic systems the simulation time T must be empirically determined. Fortunately, however, there is often a relatively large range of T which give good results. Ultimately, we wish to ensure that the correlation between updates is a minimum.

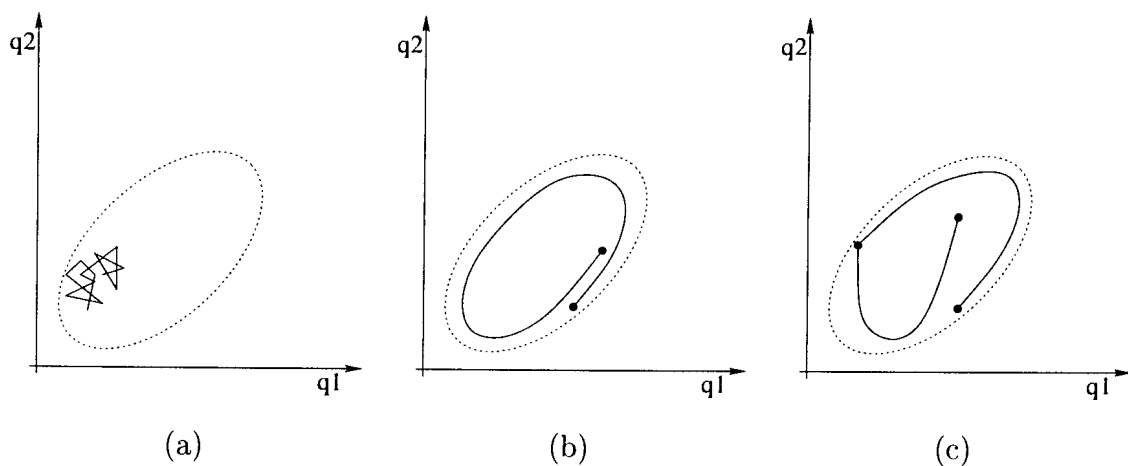


Figure 5.1: Movement through parameter space using the stochastic dynamics sampling method. The underlying distribution being sampled is indicated using a dotted line. (a) Small time period T ; a random walk of parameter space. (b) Large time period T ; a waste of time in “periodic” movement. (c) More optimal simulation time T .

In practice, we cannot integrate the Hamiltonian dynamics exactly. The leapfrog integration scheme described next can, however, be used to approximate the Hamiltonian dynamics.

Leapfrog integration In the leapfrog integration scheme [111] approximations of the position and momentum, $q_i(t + \epsilon)$ and $p_i(t + \epsilon)$ from $q_i(t)$ and $p_i(t)$ are obtained as follows:

$$p_i(t + \frac{\epsilon}{2}) = p_i(t) - \frac{\epsilon}{2} \frac{\partial E(t)}{\partial q_i} \quad (5.18)$$

$$q_i(t + \epsilon) = q_i(t) + \epsilon \frac{p_i(t + \frac{\epsilon}{2})}{m_i} \quad (5.19)$$

$$p_i(t + \epsilon) = p_i(t + \frac{\epsilon}{2}) - \frac{\epsilon}{2} \frac{\partial E(t + \epsilon)}{\partial q_i}, \quad (5.20)$$

where the step size ϵ is a small, finite, positive constant.

In order to follow the Hamiltonian dynamics for a given time T , Eqs. (5.18) to (5.20) are applied in order for $L = \frac{T}{\epsilon}$ steps. When applying the leapfrog step more than once, the last momentum update (Eq. (5.20)) and the first (Eq. (5.18)) of the next step can be combined. All but the very first and very last momentum half-steps can be merged. Figure 5.2 illustrates the leapfrog integration scheme of q and p over a time interval $[0, T]$ using the energy and momentum steps defined in Eqs. (5.18) to (5.20). If preferred, Eqs. (5.18) to (5.20) can be rewritten such that we start with an energy half-step followed by a full momentum step.

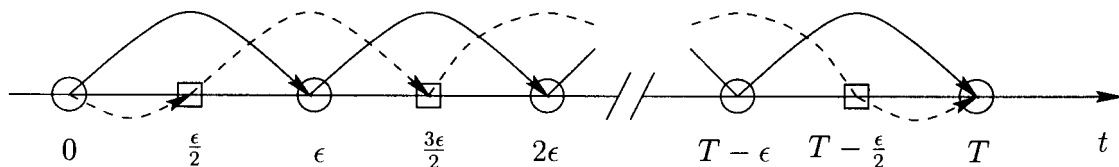


Figure 5.2: The leap frog integration scheme of q and p over a time interval $[0, T]$. The position q is evaluated at the points marked with circles and the momentum p at the points marked with squares. The dashed line indicates the momentum update and the solid line indicates the energy update.

Each of the energy and momentum steps has an error of order ϵ^3 [111]. Integrating for L leapfrog steps will therefore result in an error of order $L\epsilon^3$. A step size ϵ is chosen which is small enough to give an acceptable error. If we choose $L \approx \frac{1}{\epsilon}$, which is often done, then the error is of order ϵ^2 . The hybrid Monte Carlo algorithm discussed next, is based on the stochastic dynamics method and was developed to eliminate these systematic errors introduced by leapfrog integration.

5.2.3 The hybrid Monte Carlo algorithm

The hybrid Monte Carlo algorithm of Duane *et al.* [37] for sampling of distributions eliminates the systematic errors of the stochastic dynamics method resulting from the finite integration step size, where the total energy is not conserved during leapfrog integration. The Hybrid Monte Carlo algorithm is an extension of the stochastic dynamics method and uses the Metropolis algorithm [79] to eliminate the bias introduced by the errors resulting from the leapfrog integration of the Hamiltonian dynamics.

The Metropolis algorithm The Metropolis algorithm was introduced in 1953 by Metropolis *et al.* [79]. It has since been extensively used, and is the basis for the simulated annealing optimization method: Letting A be the current configuration and B a potential configuration (q, p) generated such that $P(A \rightarrow B) = P(B \rightarrow A)$, where $P(A \rightarrow B)$ is the probability of generating the trial configuration B given the current configuration A . If $P(B) > P(A)$ then configuration B is accepted. If $P(B) \leq P(A)$, then the configuration B is accepted with probability $\frac{P(B)}{P(A)}$. If B is rejected, the next configuration is A .

The hybrid Monte Carlo algorithm uses the Metropolis algorithm to determine whether or not to accept a new configuration generated using the refresh and integration step of the stochastic dynamics method. A hybrid Monte Carlo algorithm step can therefore be described as follows:

1. Refresh the momenta using Gibbs sampling.
2. Starting with the current state (p, q) , perform L leapfrog steps to generate a trial next configuration (p', q') .
3. Accept the trial next configuration with probability $\min(1, \exp(H(p, q) - H(p', q')))$, otherwise choose the new state to be the same as the old.

If Hamilton's equations were simulated exactly, the change in H would be zero and the trial configuration would always be accepted. When an approximation is used (leapfrog), H will change and a trial configuration will sometimes be rejected. This exactly eliminates the bias introduced by leapfrog. Note that it is important to maintain a relatively high acceptance rate, so as to minimize correlation between consecutive states.

It is expected that the work required in simulating Hamilton's equations for a fixed time will grow with the volume (V) of the system as $T \propto V^{\frac{5}{4}}$ [26, 50], as compared to growth proportional to V for the stochastic dynamics method. Although this is a relatively slow growth, as pointed out by Toussaint [111], the stochastic dynamics method will be eventually be preferred for systems which are very large. Due to current computational resources, and the large size of HMM systems, the hybrid Monte Carlo algorithm turned out to be infeasible for our work. This will, however, undoubtedly change as computers become faster with time.

We have, up until now, only investigated the theoretical aspects of implementing Bayesian learning. The next section will discuss implementation issues for a Bayesian learning approach for HMM speech recognizers.

5.3 Implementation of Bayesian HMM learning

The training process for the implementation of Bayesian learning described in this chapter generates the N_m samples of the posterior $P(\theta|\mathbf{O})$ required to numerically integrate Eq. (5.6) during recognition of an unknown utterance. This section will discuss several aspects of the implementation for the training and recognition processes.

The prior used here is the same as in the gradient-based MAP algorithm described in Section 4.4. The gradient of the energy function $E(\theta)$ is required here for the leapfrog integration. This gradient (derivative) of $E(\theta)$ with respect to the individual

parameters of both HMM and prior are therefore exactly the same as in Chapter 4 (Eq. (4.47) to Eq. (4.64)) and are as a result, not reproduced here.

5.3.1 HMM constraints

In Chapter 4 we used transformations to ensure that the original HMM constraints were maintained. Unfortunately, we cannot use transformations when we wish to sample a distribution. These are second order sampling techniques and using transforms will deform the resultant distribution being sampled. To understand this, let us look at a simple example.

Let us assume that we wish to sample the distribution of the standard deviation σ of a Gaussian model, and that the gradient of the energy function with respect to σ is a constant (k) for a given region in parameter space (i.e. $\frac{\partial E}{\partial \sigma} = k$). Here, the coordinate variable q consists of only the variable σ , i.e. $q = \{\sigma\}$. Given the momentum at time t , $p(t)$, a single leap frog momentum step (Eqs. (5.18) to (5.20)) will be

$$\begin{aligned} p\left(t + \frac{\epsilon}{2}\right) &= p\left(t - \frac{\epsilon}{2}\right) - \epsilon k, \\ q(t + \epsilon) &= q(t) + \epsilon p\left(t + \frac{\epsilon}{2}\right). \end{aligned} \tag{5.21}$$

We would therefore expect the coordinate variable q to accelerate at a constant rate ($\propto -\epsilon k$) in this region of parameter space. If, however, we use the following transform $\tilde{\sigma} = \ln \sigma$, then the derivative of E with respect to $\tilde{\sigma}$ becomes $k\sigma$ and the leapfrog step becomes

$$\begin{aligned} \tilde{p}\left(t + \frac{\epsilon}{2}\right) &= \tilde{p}\left(t - \frac{\epsilon}{2}\right) - \epsilon k \sigma \\ \tilde{q}(t + \epsilon) &= \tilde{q}(t) + \epsilon \tilde{p}\left(t + \frac{\epsilon}{2}\right). \end{aligned} \tag{5.22}$$

Note that $\tilde{p}(t) = p(t) \frac{\partial \tilde{\sigma}}{\partial \sigma}$. Transforming these update equations back from the transformed domain results in the following update

$$\begin{aligned} p(t + \frac{\epsilon}{2}) &= p(t - \frac{\epsilon}{2}) - \epsilon k \sigma^2, \\ q(t + \epsilon) &= q(t) e^{\epsilon p(t + \frac{\epsilon}{2}) / \sigma}. \end{aligned} \tag{5.23}$$

Note that neither the coordinate update $q(t + \epsilon)$ or momentum update $p(t + \frac{\epsilon}{2})$ are as in Eq. (5.21), as they should be. The distribution sampled will therefore be a distorted version of the true distribution.

Transformations were used in the gradient based MAP algorithm (GMAP) as there we were only searching for the mode of the posterior and the position thereof will not change when using transformations. But now it is not possible to use transformations with the Markov chain Monte Carlo methods described. A constrained version of the sampling algorithms is therefore implemented, in which the HMM constraints are applied. In particular, the variances are not allowed to be below a predefined value which is greater than zero (0.0001 has been used throughout).

It is non-trivial (if not impossible) to implement a constrained algorithm which maintains the constraints for the transition probabilities and mixture weights (i.e. $\sum_j a_{ij} = 1$ and $\sum_i c_i = 1$), without using transformations. These parameters are therefore treated as being fixed (at some estimate such as ML) and only the posterior distribution of the means and variances of the HMM Gaussian mixtures are therefore sampled.

5.3.2 HMM prior and hyperparameters

The prior distribution for the HMM Gaussian mixture mean and variances is again chosen to be a normal-Wishart distribution. Prior distributions are not required for the mixture weights or transition probabilities, as they are assumed to be fixed (as

motivated above) and their posterior is not sampled. The normal-Wishart prior was introduced in Section 4.2 and is reproduced here for convenience,

$$g_{\text{gaussian}}(r_{jk}, \mu_{jk} | n_{jk}, \nu_{jk}, m_{jk}, \tau_{jk}) = (2\pi)^{-\frac{D}{2}} |\nu_{jk} r_{jk}|^{\frac{1}{2}} e^{-\frac{1}{2} \nu_{jk} (\mu_{jk} - m_{jk})^T r_{jk} (\mu_{jk} - m_{jk})} \\ c |\tau_{jk}|^{-\frac{n_{jk}}{2}} |r_{jk}|^{\frac{(n_{jk} - D - 1)}{2}} e^{-\frac{1}{2} \text{tr}(\tau_{jk} \tau_{jk}^{-1})}, \quad (5.24)$$

where $(n_{jk}, \nu_{jk}, m_{jk}, \tau_{jk})$ are the prior distribution parameters associated with mixture component k of state j . See Table 4.1 for a summary of the HMM parameters and their associated priors and prior parameters. The value c is a normalizing constant which ensures that the integral of the prior is equal to one.

It can be easily shown that choosing the parameters m_{jk} and τ_{jk} in Eq. (5.24) to be

$$m_{jk} = \mu' \quad (5.25)$$

$$\tau_{jk} = (n_{jk} - D) \Sigma', \quad (5.26)$$

results in the mode of the prior being at the point $\mu_{jk} = \mu'$ and $\Sigma_{jk} = \Sigma'$. The parameters n_{jk} and ν_{jk} determine the degree to which the prior is peaked about its mode.

Although not necessary, it makes sense (simplifies calculations) to reduce the number of variables in the prior by expressing the parameters n_{jk} and ν_{jk} in terms of a common parameter C_{jk} ,

$$n_{jk} = C_{jk} + D \quad (5.27)$$

$$\nu_{jk} = C_{jk} + 1, \quad (5.28)$$

with $C_{jk} > 0$.

The parameters C_{jk} , m_{jk} and τ_{jk} are now given their own distributions. These distributions and their parameters must be chosen in a meaningful way, such that the hyperparameters contain *a-priori* knowledge (albeit vague). The choice of the distributions and parameters are discussed next.

We would expect the parameters for the mixtures of a given state to be highly correlated. This is often not the case with HMMs, where a single state can represent multiple acoustic units, which can result in a higher confusability.

Let us, for example, look at a three state HMM, with each state having several Gaussian mixtures, used to represent the word “four”. If, assuming, certain of the training utterances are pronounced such that the “r” is not heard, then the last state will probably represent both the consonant r and the vowel ou which is the end of some of the training utterances. The result is a model which is easier to confuse with other utterances (or part thereof). For example, this HMM will give a high output for the phone sequence $f-ou-r-ou-r$.

Illina and Gong [55] investigated this phenomenon, which they called trajectory folding. One solution to this problem is to change the topology of the model. It was found that using a trajectory mixture HMM (TMHMM) [107] and segment-based mixture stochastic trajectory HMMs (MSTM) [49] reduced the effect of this problem.

Trajectory folding, is not only caused by differences in pronunciation, but can also be caused by the training procedure. This is especially the case when the model is complex and little training data is available. From the above discussion, it is evident that states representing multiple acoustic units are not desirable. In an attempt to avoid trajectory folding, the prior parameters of a given state have therefore been given common distributions and hyperparameters. Note that the prior parameter distributions and the associated hyperparameters which have been chosen are not necessarily the only (or best) possible choice. The manner in which the parameters are grouped is a modeling choice, which is made on the basis of prior knowledge. Many possibilities

remain, and it is more than likely that the performances reported in this chapter can be improved upon by making other choices.

Our choice of the priors are as follows:

The prior mean m_{jk} of state j is given a normal distribution with mean ω_j and standard deviation ς_j , i.e.

$$P(m_{jkl}) = \mathcal{N}(\omega_{jl}, \varsigma_{jl}), \quad k = 1, \dots, M, \quad (5.29)$$

where M is the number of mixture components, and ω_{jl} and ς_{jl} are the hyperparameters for the prior mean m_{jk} .

Given that we wish to keep the hyperparameter distributions as intuitive as possible, the distribution of the prior mode Σ' in Eq. (5.26) has been chosen to be a gamma distribution (Section A.4) with parameters ϕ_j and ψ_j . The distribution of the prior parameter τ_{jk} (precision of μ_{jk}) of state j can, using Eq. (5.26), therefore be written as

$$P(\tau_{jkl}) = (n_{jk} - D)\mathcal{G}(\phi_{jl}, \psi_{jl}), \quad k = 1, \dots, M, \quad (5.30)$$

where ϕ_{jl} and ψ_{jl} are the hyperparameters of τ_{jkl} . Note that from Eq. (5.25), the prior mean m_{jk} is also the mode of the prior with respect to the mean parameters. The distribution for the prior parameter C_{jk} is also chosen to be a gamma distribution with parameters ν and ϱ , i.e.

$$P(C_{jk}^{(i)}) = \mathcal{G}(C_m, C_\nu), \quad k = 1, \dots, M, \quad (5.31)$$

for every state $j = 1, \dots, N$ in every HMM $i = 1, \dots, N_h$. A summary of the HMM

parameters, their prior distributions, the prior parameters and their distributions can be found in Table 5.1.

Table 5.1: Summary of the HMM parameters, prior distributions, hyperparameters and their distributions. A Wishart distribution is represented with \mathcal{W} and a Gamma with \mathcal{G} .

Parameter	Prior distribution	Prior parameter	Distribution
$\mu_{jk} R_{jk} = r_{jk}$	$\mathcal{N}(m_{jk}, (C_{jk} + 1)r_{jk})^{-1/2}$	m_{jk}	$\mathcal{N}(\omega_{jl}, \varsigma_{jl})$
$R_{jk} (\Sigma_{jk}^{-1})$	$\mathcal{W}((C_{jk} + D), \tau_{jk})$	C_{jk}	$\mathcal{G}(v, \varrho)$
		τ_{jk}	$(n_{jk} - D)\mathcal{G}(\phi_{jl}, \psi_{jl})$

Determining hyperparameters The ML estimate, which is used as the starting point for the stochastic dynamics sampling, is used to estimate the hyperparameters ω_j, ς_j (prior mean m_{jk} in Eq. (5.29)) and ϕ_j, ψ_j (prior parameter τ_{jk} in Eq. (5.30)).

The sample mean and variance (over the mixtures) of the mixture means μ_{jk} for state j of the ML estimate are reasonable estimates for the parameters ω_{jl} and ς_{jl} , i.e. ω_{jl} and ς_{jl} are chosen to be

$$\omega_j = \frac{1}{M} \sum_{k=1}^M \mu_{jk} \quad (5.32)$$

$$\varsigma_j = \frac{1}{M} \sum_{k=1}^M (\mu_{jk} - \omega_j)^2. \quad (5.33)$$

Likewise, the sample mean and variance of the mixture variances for state j of the ML estimate are reasonable estimates for the mean and variance of distribution $\mathcal{G}(\phi_{jl}, \psi_{jl})$. Letting the sample mean of the mixture variance be $\bar{\sigma}_j = \frac{1}{M} \sum_{k=1}^M \sigma_{jk}$, and using Eqs. (A.14) and (A.15), the parameters ϕ_j and ψ_j are therefore estimated as follows

$$\phi_j = \frac{\bar{\sigma}_j^2}{\frac{1}{M} \sum_{k=1}^M (\sigma_{jk} - \bar{\sigma}_j)^2} \quad (5.34)$$

$$\psi_j = \frac{\bar{\sigma}_j}{\frac{1}{M} \sum_{k=1}^M (\sigma_{jk} - \bar{\sigma}_j)^2}. \quad (5.35)$$

The hyperparameters C_m and C_v (defined in Eq. (5.31)) are determined by the user and express our trust in the ML estimate. Large values of C_m will result in prior distributions which are peaked around the mode of the posterior, while small values of C_m will result in a relatively non-informative prior distribution. The effect of these hyperparameters on the performance of the system will be determined in the Section 5.4.

5.3.3 Refreshing hyperparameters

In the stochastic dynamics method (Section 5.2.2), before using the leapfrog integration scheme, the momenta are “refreshed” using Gibbs sampling. It is at this same step in the stochastic dynamics method, where we will also update or “refresh” the prior parameters using Gibbs sampling.

As a result of using hyperparameters, we have to obtain the prior using Eq. (5.10). This is accomplished by Gibbs sampling of the hyperparameters γ after each transition (i.e. before leapfrog integration) of the stochastic dynamics or hybrid Monte Carlo algorithm, which results in the numerical integration of Eq. (5.10).

The hyperparameters are easily sampled, as they are independent of other parameters and hyperparameters and since their distributions have been chosen to be normal and Gamma, standard techniques for generating normal and Gamma distributed variates can be used.

5.3.4 Implementation of stochastic dynamics method

Writing the posterior in the required form of an energy function as discussed in Section 5.2.2, we get the following “potential energy” function,

$$E(\theta) = -\log[P(\theta)] - \sum_{i=1}^n \log[P(\mathbf{O}_i|\theta)]. \quad (5.36)$$

The implementation of the stochastic dynamics method is then as follows:

1. Refresh momenta and prior parameters using Gibbs sampling.
2. Starting with the current state perform L leapfrog steps to generate a new configuration. Here we require $\frac{\partial E(t)}{\partial \theta}$ (Eqs. (5.18) and (5.20)). The energy function is the same as that used for the gradient-based MAP algorithm in Chapter 4 (Eqs. (4.47) to (4.55)), and the gradient of the energy function with respect to the HMM parameters is therefore exactly the same. The derivation thereof is not repeated here.
3. Keep current configuration, repeat steps 1 through 3 until N samples are generated.

A segmental approach is used, in which we use the best state sequence, as opposed to the sum of all possible sequences, i.e. the probability $\max_q P(\mathbf{O}, \mathbf{q}|\theta)$ is used as an approximation of $P(\mathbf{O}|\theta) = \sum_{\text{all } \mathbf{q}} P(\mathbf{O}, \mathbf{q}|\theta)$. An embedded version has also been implemented, which uses the modified-Viterbi trellis search (Section 2.1.5) to obtain the best HMM and state sequence.

5.3.5 Recognition

In a Bayesian framework, we evaluate Eq. (5.6) for each class and classify using the following decision rule

$$C(\mathbf{O}) = C_i \text{ where } i = \underset{j}{\operatorname{argmax}} P_j(\mathbf{O}_{\text{unknown}}|\mathbf{O}_1, \dots, \mathbf{O}_n). \quad (5.37)$$

Implementation of this decision rule is straightforward when applied to discrete (or label-based) phoneme or word classification: use Eq. (5.11) to numerically determine $P_j(\mathbf{O}_{\text{unknown}}|\mathbf{O}_1, \dots, \mathbf{O}_n)$ for each class (phoneme or word).

The implementation of Eq. (5.37) for continuous speech recognition is not a trivial task. Here, a class is a string of phonemes or words and there are many possible combinations thereof. The direct implementation of Eq. (5.11), although relatively simple to implement, would therefore not be viable in terms of computation complexity.

An alternative approximation would be to use an N -best search to generate the N_s best strings and then evaluate Eq. (5.11) for only these combinations. Although this is a reasonable approximation, it is still somewhat computationally expensive. This approximation, although considerably less computationally expensive, is at best $N_s \cdot N_m$ times slower than the equivalent standard HMM recognizer. Given that N_s will typically have to be relatively large, this technique will unfortunately not be feasible either.

The following is a simple, yet reasonable approximation of the classification process described by Eq. (5.37). Although we wish to obtain $P(\mathbf{O}_{\text{unknown}}|\mathbf{O}_1, \dots, \mathbf{O}_n)$, it is not necessary to do it directly using Eq. (5.11). Formulating the problem in terms of a single HMM, and not the sum of several sampled HMMs, we can write $P(\mathbf{O}_{\text{unknown}}|\mathbf{O}_1, \dots, \mathbf{O}_n)$ as follows:

$$P(\mathbf{O}_{unknown}|\mathbf{O}_1, \dots, \mathbf{O}_n) = \sum_{\text{all } \mathbf{q}} P(\mathbf{O}|\mathbf{q}, \mathbf{O}_1, \dots, \mathbf{O}_n)P(\mathbf{q}|\mathbf{O}_1, \dots, \mathbf{O}_n), \quad (5.38)$$

where \mathbf{q} is a state sequence $(q_1 q_2 \dots q_T)$. The probability of a new observation \mathbf{O} given the state sequence is

$$P(\mathbf{O}|\mathbf{q}, \mathbf{O}_1, \dots, \mathbf{O}_n) = \prod_{t=1}^T P(\mathbf{o}_t|q_t, \mathbf{O}_1, \dots, \mathbf{O}_n), \quad (5.39)$$

where independence of observations has been assumed. The probability of a single observation \mathbf{o}_t given a state sequence can be written as

$$P(\mathbf{o}_t|q_t, \mathbf{O}_1, \dots, \mathbf{O}_n) = \int P(\mathbf{o}_t|q_t, \theta)P(\theta|\mathbf{O}_1, \dots, \mathbf{O}_n)d\theta, \quad (5.40)$$

where we have integrated with respect to the posterior of the Gaussian mixture parameters. We now numerically integrate Eq. (5.40) using N_m samples of the posterior distribution $P(\mathbf{o}_t|\theta^{(m)})$,

$$\begin{aligned} P(\mathbf{o}_t|q_t, \mathbf{O}_1, \dots, \mathbf{O}_n) &\approx \sum_{m=1}^{N_m} P(\mathbf{o}_t|\theta^{(m)}) \\ &= \sum_{m=1}^{N_m} b_{q_t}^{(m)}(\mathbf{o}_t), \end{aligned} \quad (5.41)$$

where $b_{q_t}^{(m)}(\mathbf{o}_t)$ is the observation probability of \mathbf{o}_t for state q_t and the m th sample of the posterior distribution.

It remains to decide on $P(\mathbf{q}|\mathbf{O}_1, \dots, \mathbf{O}_n)$ in Eq. (5.38). Note that we have not allowed the posterior of the transition probabilities to be sampled (as discussed in Section 5.3.1) – a reasonable assumption would be that the posterior of the transition probabilities is reasonably peaked and the ML estimate is a reasonable approximation, therefore

$$P(\mathbf{q}|\mathbf{O}_1, \dots, \mathbf{O}_n) \approx P(\mathbf{q}|\theta) = \pi_{q_0} \prod_{t=1}^T a_{q_t q_{t+1}}. \quad (5.42)$$

The probability of the observation sequence $\mathbf{O}_{unknown}$ given the training data $\mathbf{O}_1, \dots, \mathbf{O}_n$ for a given state sequence \mathbf{q} can therefore be written as

$$P(\mathbf{O}_{unknown}|\mathbf{q}, \mathbf{O}_1, \dots, \mathbf{O}_n) \approx \pi_{q_0} \prod_{t=1}^T \left[\left(\sum_{m=1}^{N_m} b_{q_t}^{(m)}(\mathbf{o}_t) \right) a_{q_t q_{t+1}} \right]. \quad (5.43)$$

Therefore, the procedure described above results in integration with respect to the posterior of the Gaussian mixtures each time an observation probability is required. This is, as opposed to integrating with respect to the full HMM posterior as in Eq. (5.6).

Viterbi implementation Here, where we would have used a single observation probability in the normal Viterbi algorithm, we now calculate the sum of the observation probabilities for the posterior samples (Eq. (5.41)). The implementation of the Viterbi algorithm for this approximation is therefore almost trivial. To find the single best state sequence, $\mathbf{q} = (q_1 q_2 \dots q_T)$, for the given observation sequence $\mathbf{O} = (\mathbf{o}_1 \mathbf{o}_2 \dots \mathbf{o}_T)$, we use the following modified Viterbi algorithm:

1. Initialization

$$\delta_1(i) = \pi_i \left(\sum_{m=1}^{N_m} b_i^{(m)} \right) \quad 1 \leq i \leq N \quad (5.44)$$

$$\psi_1(i) = 0 \quad (5.45)$$

2. Recursion

$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] \sum_{m=1}^{N_m} b_j^{(m)} \quad (5.46)$$

$$\psi_t(j) = \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] \quad (5.47)$$

3. Termination

$$P^* = \max_{1 \leq i \leq N} [\delta_T(i)] \quad (5.48)$$

$$q_R^* = \arg \max_{1 \leq i \leq N} [\delta_T(i)] \quad (5.49)$$

4. Path backtracking

$$q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T - 1, T - 2, \dots, 1. \quad (5.50)$$

Extension to the trellis search, which is a frame-synchronous, modified implementation of the Viterbi search, is trivial. This approximation has the advantage over direct implementation of Eq. (5.11), and the N -best approximation, that it is only approximately N_m times slower than a standard HMM recognizer, where N_m is the size of the sample.

5.4 Experiments

The goal of this section is to experimentally determine the utility of the Bayesian learning approach described in this chapter. The hypothesis that will be tested here, is that “Bayesian learning will improve performance markedly in situations where little data is available for training purposes”. In situations where sufficient training data is available, the posterior will be peaked sharply about the MAP point and there will

therefore be little advantage in using Bayesian learning rather than MAP estimation under such conditions.

The three datasets SUNSpeech, TIMIT and TIDIGITS are once again used here. However, the assumption here is that there is only a small amount of data available with no non-task-specific data available for adaptation purposes. The results obtained using Bayesian learning will be compared with that obtained using the adaptation techniques described in Chapter 4 where non-task-specific data was assumed to be available.

Unless otherwise stated, the stochastic dynamics method was run for 100 iterations, i.e. 100 samples were generated. The last N_m samples were used for recognition tasks. For example, a recognition experiment using $N_m = 30$ will therefore use the last 30 samples, i.e. samples 71 through to 100.

Bayesian learning is not a deterministic process and there will therefore be a certain variance in performance for the resultant systems. Each experiment was therefore repeated 10 times, so as to provide an indication of the variance in accuracy that results. Error bars are given for each result, indicating the minimum, mean and maximum accuracy for a given configuration.

5.4.1 SUNSpeech

This section evaluates the Bayesian learning algorithm under the same conditions as that used to test the adaptation techniques in Section 4.6.1, except that here, there is no English training data available for cross-language adaptation purposes. The effect of the sample size, parameters C_m , C_v , simulation time T and leapfrog step size L will be experimentally determined for the SUNSpeech dataset in this section. Results will be presented for both the small (A_S) and full (A) Afrikaans training sets. The speaker independent Afrikaans testing set is used for evaluation purposes.

Table 5.2 presents the Afrikaans test set accuracy for the two training sets available,

namely the full Afrikaans training set A in the SUNSpeech dataset and the reduced training set A_S .

Table 5.2: Base system accuracy for SUNSpeech Afrikaans test set

Training set	Mixtures	
	5	10
Afrikaans train (A)	48.6	51.5
Afrikaans adapt (A_S)	42.5	41.2

The following results are for a 3 state, 5 mixture HMM system, when only the small Afrikaans training set is available. Selected results for the full dataset, as well as 3 state, 10 mixture HMMs will be presented later in this section.

Effect of sample size Figure 5.3 shows the effect of the sample size N_m on the performance of the system. Firstly, it is noticeable that the performance for a single sample, is already better than that attained by the ML estimate. This improvement is primarily due to the regularization effect of the chosen prior and associated hyperparameters. The choice of prior and hyperparameters described in Section 5.3.2 therefore seems justified.

Considerable improvements result from increasing the sample size, with 46.1% being attained when using a sample size of 90 (a relative improvement in error rate of 6.2%). The sudden decrease in recognition rate for 100 samples, is indicative of the fact that the algorithm had not converged within the first 10 transitions (samples) of the stochastic dynamics method.

The variance on the resulting systems is relatively high, but as expected decreases as the number of samples used is increased. We would expect the variance in the solution to tend towards zero as the number of samples tended towards infinity, i.e. the numerical integration becomes the exact integration in Eq. (5.6).

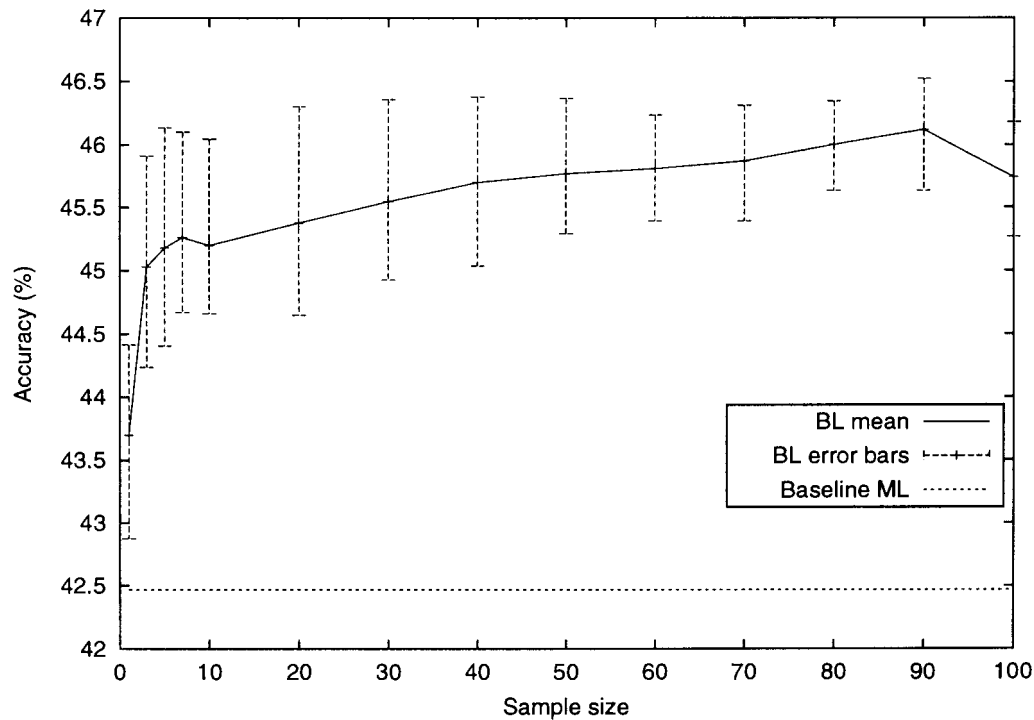


Figure 5.3: Performance of the Bayesian system (BL) versus the size of the sample used to numerically integrate with respect to the posterior.

Effect of the hyperparameter C_m As mentioned in Section 5.3.2 the parameter C_m , along with C_v , indirectly determines the degree to which the prior is peaked about its mode. It contains our *a priori* trust in the ML model used as starting point.

Figure 5.4 presents the effects of C_m on the performance of a 3 state, 5 mixture HMM when the reduced training set is used. As can be seen, an optimal value for C_m exists. Values of C_m which are too low do not force the regularization of the prior to be of any consequence. Alternatively, values which are too large result in a prior which is too restrictive and therefore unnecessarily restricting the posterior.

Figure 5.5 shows the results for a 3 state, 10 mixture system versus the parameter C_m . Note that, here the optimal value of C_m is lower ($C_m = 5$ or $C_m = 10$) than that obtained for the 3 state, 5 mixture system. The ML estimate for the more complex 10 mixture HMM system, is worse when little data is available (Table 5.2). The ML estimate for the 10 mixture system can therefore not be “trusted” as much as the 5

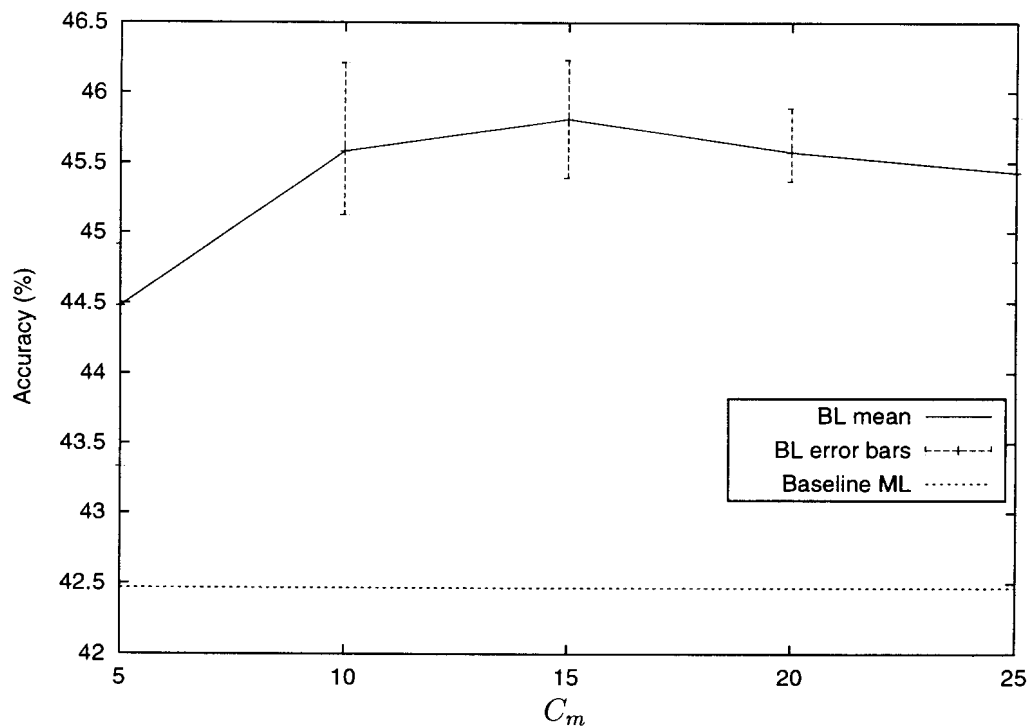


Figure 5.4: Performance of the Bayesian system (BL) versus the hyperparameter C_m , for $C_v = 4$, $L = 100$ and using a sample size of 60.

mixture system. A smaller value of C_m which makes the prior less informative, therefore results in the peak performance for this system.

Effect of the hyperparameter C_v Figure 5.6 presents the effects of the C_v on the performance of a 3 state, 5 mixture HMM when the reduced training set is used. Again, extreme values of this parameter result in degradation in performance. Fortunately, the performance of the system is relatively invariant to reasonable variations in this parameter.

Simulation time T A detailed discussion of the effect of the simulation time T was given in Section 5.2.2, and is critical to the optimal usage of the stochastic dynamics method. We therefore need to experimentally determine the effect of the number of leapfrog steps used on the performance of the system. The step size is kept constant

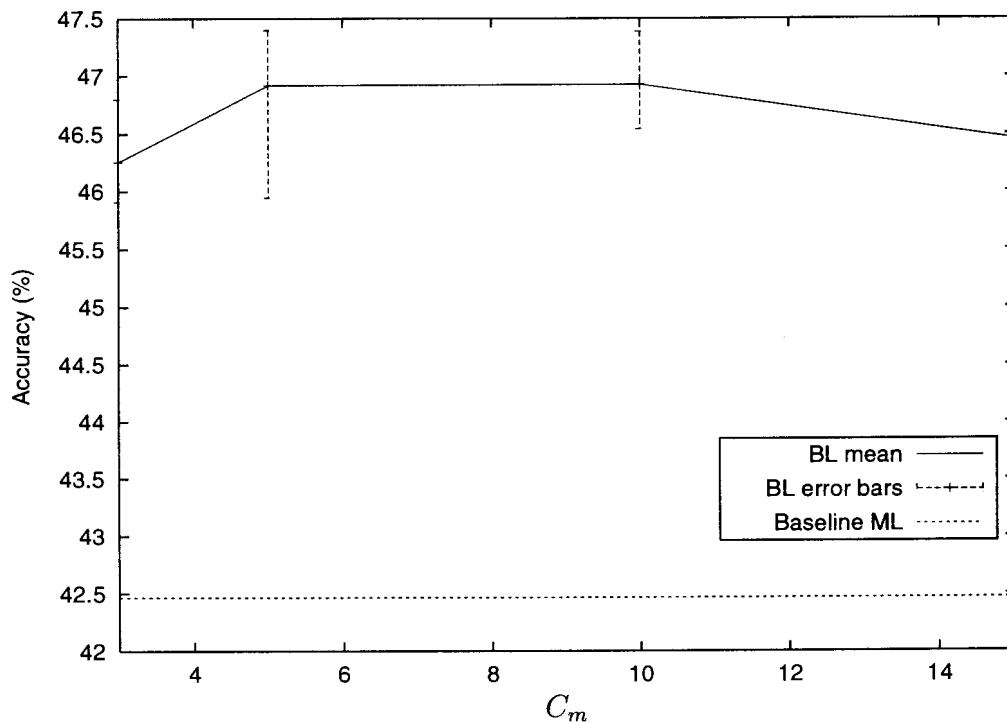


Figure 5.5: Performance of the Bayesian system (BL) versus the hyperparameter C_m , for a 3 state, 10 mixture HMM and $C_v = 4$, $L = 100$ and using a sample size of 60.

for this experiment, at $\epsilon = 0.001$.

Figure 5.7 shows the performance of the Bayesian system versus the number of leapfrog steps L , for the configuration $C_m = 15$, $C_v = 4$ and using a sample size $N_s = 60$. Remember that the simulation time is approximately equal to the step size ϵ multiplied by the number of leapfrog steps, i.e. $T \approx \epsilon L$. We expect that too few leapfrog steps result in samples which are highly correlated, which results in a sample which is not representative of the posterior. This can be seen in Figure 5.7, where using 10 leapfrog steps results in relatively poor system performance, with the minimum accuracy being less than that of the ML baseline system.

System accuracy increases as the number of leapfrog steps increase, and peaks at $L = 100$, where an average accuracy of 45.8% is attained. However, as mentioned in Section 5.2.2, values of L which are too large will not only waste time, but could worsen performance due to higher correlation between samples, resulting from the

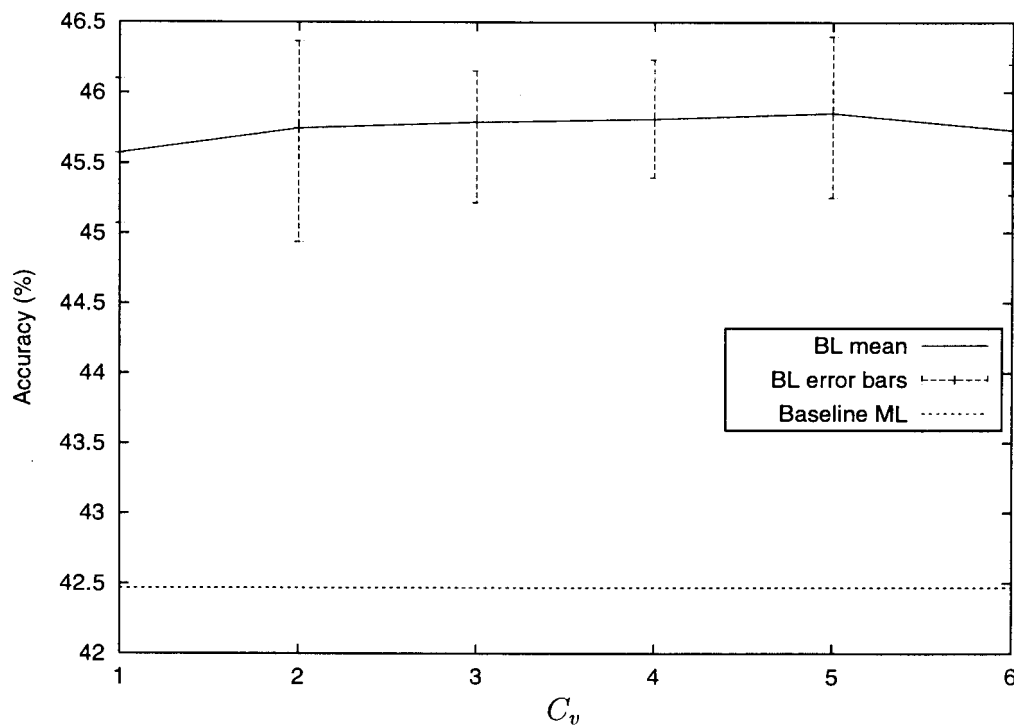


Figure 5.6: Performance of the Bayesian system (BL) versus the hyperparameter C_v , for $C_m = 15$, $L = 100$ and using a sample size of 60.

periodic nature of Hamilton's equations. This effect is also observed in Figure 5.7 where accuracy at $L = 200$ (45.75%) is slightly worse than that at $L = 100$.

Leapfrog step size Table 5.3 presents the results for the Bayesian system using step sizes of $\epsilon = 0.001$ and $\epsilon = 0.0002$. The results presented until now have used a learning rate of $\epsilon = 0.0001$. The number of steps used have been adjusted such that the simulation time T is approximately the same. This ensures that the simulation time T does not affect the results, as seen in the previous experiment. Using a step size of $\epsilon = 0.0002$ results in an average absolute improvement of approximately 0.3%. Using the smaller step size $\epsilon = 0.0002$, however, requires that 500 leapfrog step be used (for equivalent simulation time). The extra 0.3% in absolute performance is, however, not enough to justify the increase in simulation time (5 times).

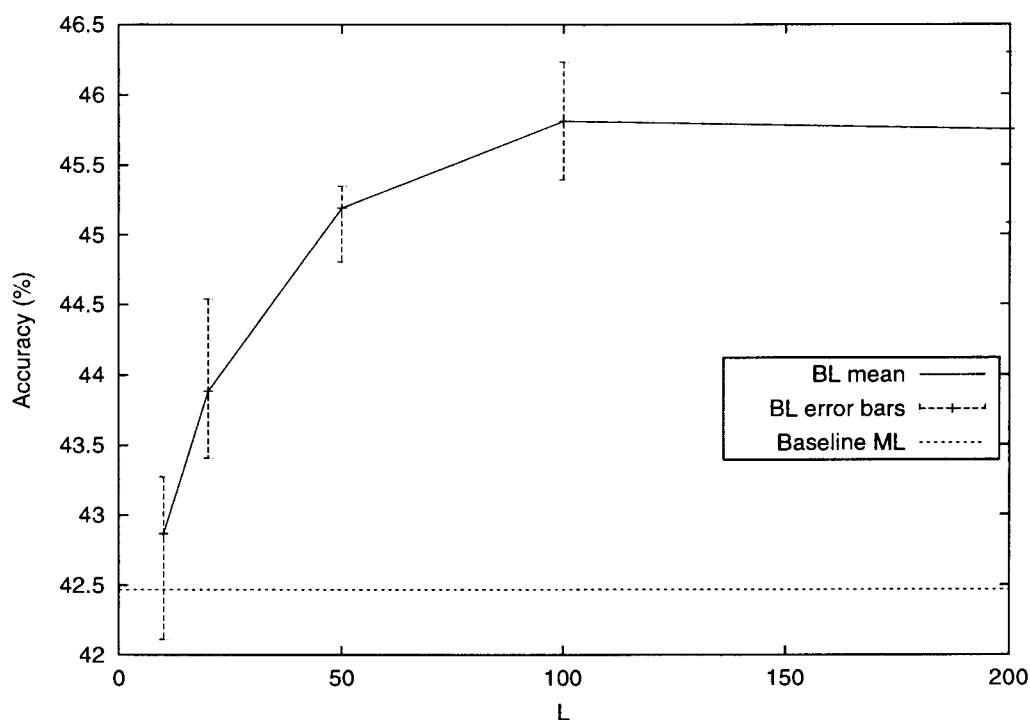


Figure 5.7: Performance of the Bayesian system (BL) versus the number of leapfrog steps L per transition, for $C_m = 15$, $C_v = 4$ and using a sample size of 60.

Summary of results Table 5.4 gives the best results obtained for the SUNSpeech dataset when using the Bayesian learning algorithm described in this chapter. Note that here, the 10 mixture model always performs best. This is as opposed to the ML estimate, where the 10 mixture system performed worse than the 5 mixture system when the small training set is used (Table 5.2).

Peak performance of 47.0% accuracy was attained using 3 state, 10 mixture HMMs and

Table 5.3: Comparison of test set accuracy for leapfrog integration using differing step sizes ($\epsilon = 0.001$ and $\epsilon = 0.0002$) but having equivalent simulation time $T \approx \epsilon L$

Leapfrog configuration	Sample size									
	10	20	30	40	50	60	70	80	90	100
$\epsilon = 0.001, L = 100$	45.2	45.4	45.5	45.7	45.8	45.8	45.9	46.0	46.1	45.7
$\epsilon = 0.0002, L = 500$	45.6	45.8	46.0	45.9	46.1	46.1	46.3	46.3	46.4	46.0

Table 5.4: Summary of the results obtained using Bayesian learning with the SUNSpeech dataset, for $L = 100$, $\epsilon = 0.001$. Relative improvement in error rate expressed as a percentage, relative to the associated baseline performance given in Table 5.2 is given in brackets.

Training set	Mixtures	Configuration	Sample size		
			10	50	90
A_S	5	$C_m = 15, C_v = 4$	45.2 (4.7)	45.8 (5.7)	46.1 (6.2)
A_S	10	$C_m = 5, C_v = 4$	46.5 (9.0)	46.9 (9.7)	47.0 (9.9)
A	5	$C_m = 10, C_v = 4$	50.5 (3.7)	50.3 (3.3)	50.3 (3.3)
A	10	$C_m = 10, C_v = 4$	53.7 (4.5)	54.0 (5.2)	54.1 (5.4)

a sample size of 90 for the small training set. This equates to a relative improvement in error rate of 9.9%. The same system using only 10 samples attained a 46.5% testing set accuracy (9.0% relative improvement in error rate).

When using the full training set, it is once again the 10 mixture system which performs best. Peak performance of 54.1% is attained for the 90 samples system, which is a relative improvement in error rate of 5.4%. The considerably smaller 10 sample system gives a 4.5% relative improvement in error rate (53.7% accuracy).

Whether one can justify using the 90 sample system, which is 9 times slower than the 10 sample system and approximately 90 times slower than that of the equivalent standard HMM system, is currently unlikely.

5.4.2 TIMIT

This section evaluates the Bayesian learning algorithm under the same conditions as that used to test the MCE algorithm in Section 3.4, except that here, only the small TIMIT training set T_S will be used. The procedure followed in selecting this dataset has been described in Section 2.4.1. The full TIMIT testing set has been used for all experiments presented in this section. A 3 state, 5 mixture HMM is used for all

experiments in this section. The base system test set performance for the 3 state, 5 mixture HMM when using the small TIMIT training set is 52.6%.

Here, so as not to unnecessarily duplicate results, only the effect of sample size and parameter C_m will be experimentally evaluated – to show that performance is similar for another dataset and language.

Effect of sample size Figure 5.8 shows the effect of the sample size N_m on the performance of the Bayesian learning system. The performance for a sample size of 1 ($N_m = 1$) is, once again, in general better than that of the Baseline ML system – further evidence that the choice of prior and hyperparameter is justified.

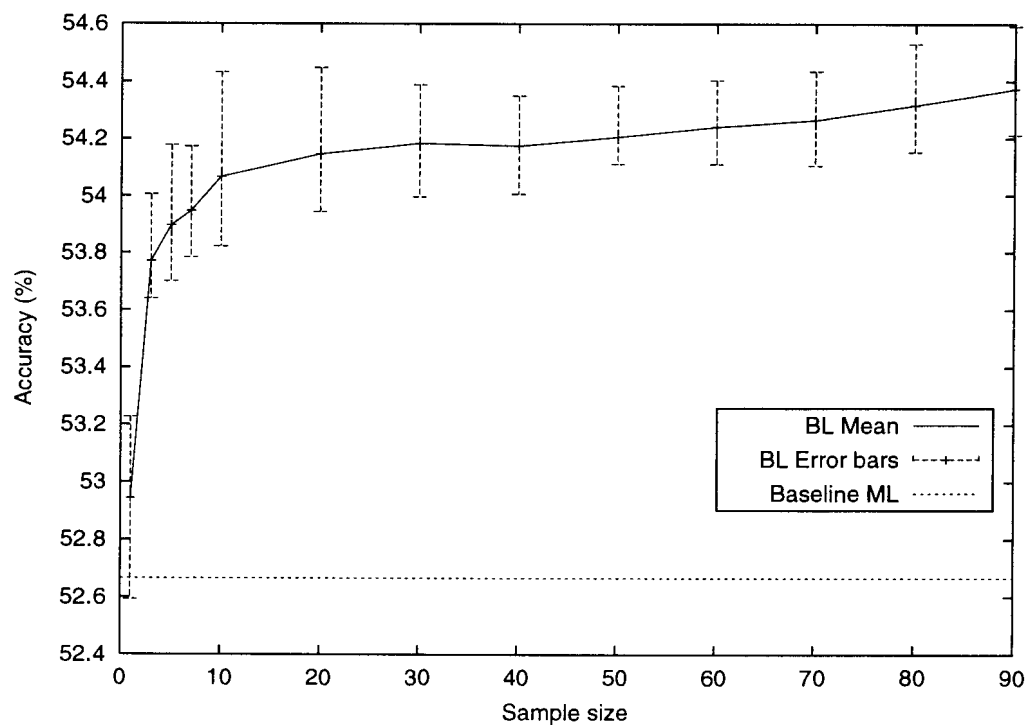


Figure 5.8: Performance of the Bayesian system (BL) versus the size of the sample used to numerically integrate with respect to the posterior.

As with the SUNSpeech experiments, by far the greatest improvements resulted when using 10 or fewer samples. Using 10 samples results in a system accuracy of 54.1%. The performance increases to a reasonable degree for sample sizes greater than 10,

though whether the extra computational effort can be justified is once again somewhat doubtful. Maximum performance of 54.4% is attained when using a sample size of 90, though better performance would probably result when using a larger sample size. As expected, the variance of the system accuracy does, in general, tend to decrease as the sample size grows.

Effect of the parameter C_m Figure 5.9 presents the effect of the parameter C_m on the performance of a 3 state, 5 mixture HMM based Bayesian system. Once again, an optimal value exists ($C_m = 15$), with the performance of the system deteriorating for values smaller or larger than this optimal value.

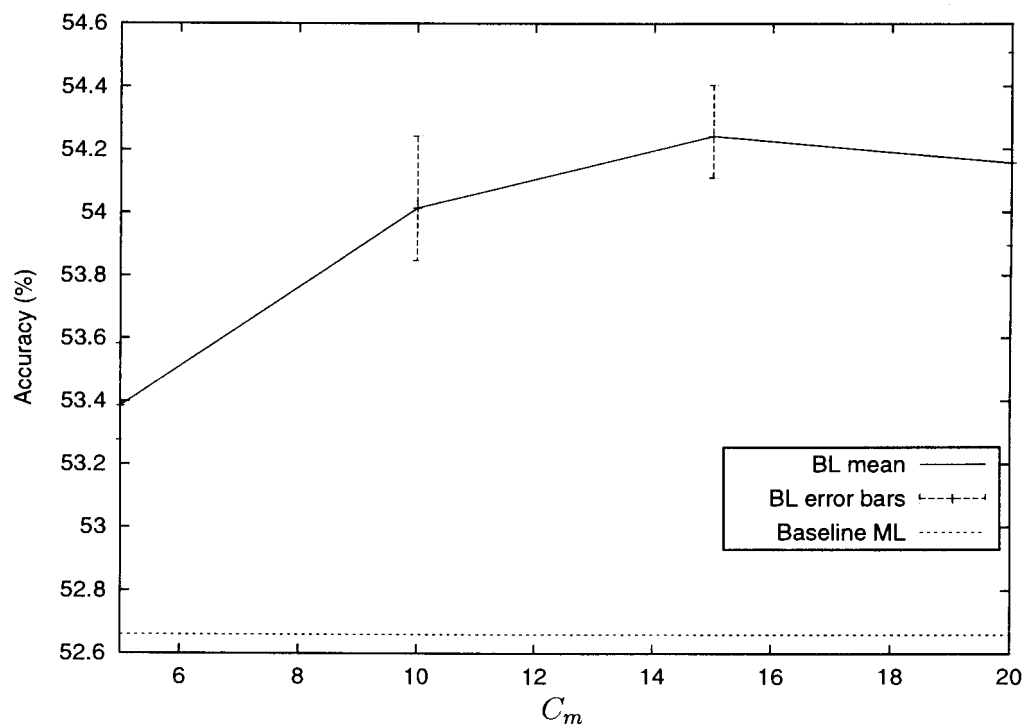


Figure 5.9: Performance of the Bayesian system (BL) versus the hyperparameter C_m , for $C_v = 4$, $L = 100$ and using a sample size of 60.

5.4.3 TIDIGITS

This section evaluates the Bayesian learning algorithm for a connected digit recognition task. The TIDIGITS database, as described in Section 2.4.2, is used for this purpose, where a small gender-independent dataset was created. The reduced speaker set (T_S) was created, using five randomly chosen female speakers and five randomly chosen male speakers, using all 77 digit sequences per speaker. The relevant training and testing sets are summarized in Table 5.5.

The full TIDIGITS testing set has been used for the experiments presented in this section. An 8 state, 5 mixture HMM is used for all experiments. The base system test set performance for the 8 state, 5 mixture HMM when using the small TIDIGITS training set is 95.8%. Once again, so as not to unnecessarily duplicate results, only the effect of sample size and parameter C_m will be evaluated.

Table 5.5: Training and testing sets used with the TIDIGIT database

Description	Label	Number of speakers			Digit sequences	Duration (minutes)
		Male	Female	Total		
Training sets:						
Full (standard)	T	55	57	112	8623	253.4
Man	T_M	55	0	55	4235	121.9
Small	T_S	5	5	10	770	21.5
Testing set		55	57	113	8623	254.4

Effect of sample size Figure 5.10 shows the performance of the Bayesian learning system versus the sample size N_m , with $C_m = 15$, $C_v = 4$ and $L = 100$. As was the case for the previous two datasets, the performance for a sample size of 1 ($N_m = 1$) resulted in better performance (mean of 96.3%) than that of the Baseline ML system (95.8%). A small improvement in accuracy (96.5%) is observed when a sample size of two is used. However, larger sample sizes result in degradation of system accuracy,

with accuracy dropping back down to 96.3% when 30 samples are used.

A connected digit task is relatively simple compared to continuous phoneme recognition. We can therefore expect that the posterior distribution for the HMM we are estimating will not be as complex as that of the two scenarios previously used. For a very simple task, one expects that a single mode posterior would result, and if the posterior is peaked no significant improvement in performance will be attained when using the Bayesian learning approach described in this chapter. With the connected digit task, however, we do observe a small, though significant, improvement in performance when using a sample size which is greater than one. The posterior of this task is therefore slightly more complex than a single mode. The improvement in performance may, for example, be attributable to a bi-modal posterior, with one mode representing male speakers and the other mode representing female speakers.

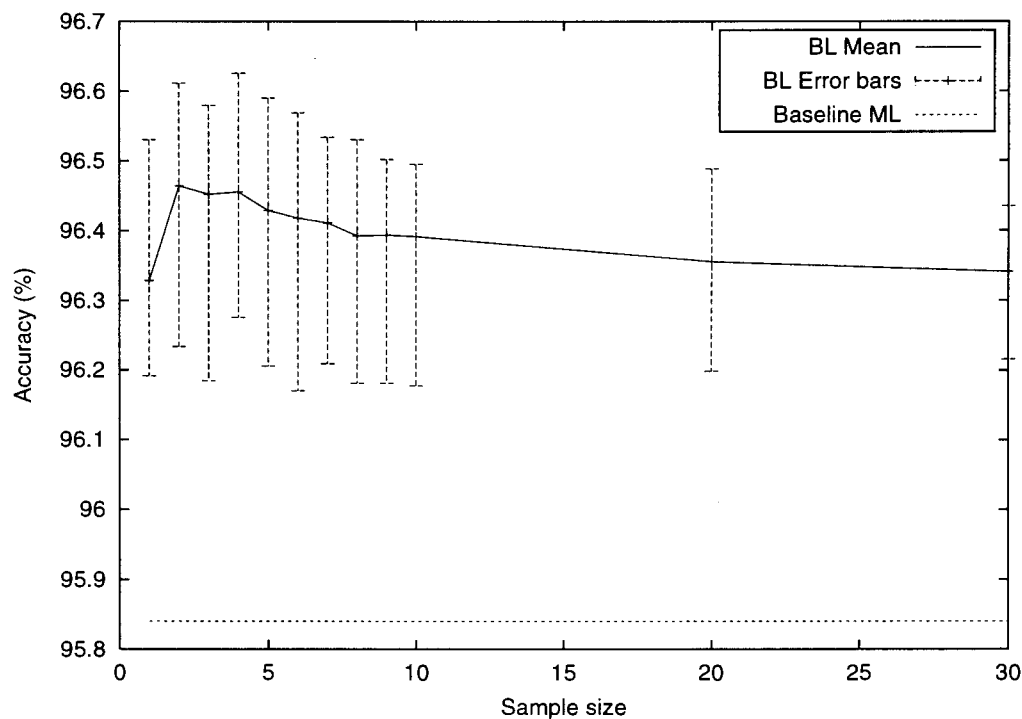


Figure 5.10: Performance of the Bayesian system (BL) versus the size of the sample used to numerically integrate with respect to the posterior.

One would, however, not expect the performance to deteriorate as the sample size

increases. A decrease in performance indicates that the distribution of the samples generated using the Markov chain Monte Carlo method is not that of the desired distribution. In Section 5.4.1 an example of this phenomenon was observed, when samples from a Markov chain were used before it had converged. Here, however, the Markov chain has converged. The decrease in accuracy observed here is probably due to the systematic errors introduced by the leapfrog algorithm. A finite step size will result in a posterior sample which is slightly different from that of the true distribution. As the sample size increases, the probability of including a sample which is not representative of the true or desired distribution therefore increases, and as a result, the performance of the system will tend to decrease as larger samples are used. This phenomenon has only been observed here, and its effect is relatively small compared to the improvements resulting when the posterior is relatively complex.

Effect of the parameter C_m The performance of the system (not shown graphically) is somewhat invariant to the parameter C_m , probably due to the fact that little or no improvement results from using the Bayesian learning procedure. Peak mean performances of 96.47%, 96.46% and 96.50% are attained for C_m values of 10, 15 and 20 respectively.

5.4.4 Summary of results

In Section 5.4.1 the Bayesian learning procedure introduced in this chapter was tested in a situation where a limited amount of training data was available for a new language specific recognizer. The SUNSpeech dataset was used for this purpose. The effects of the sample size, parameters C_m and C_v , simulation time T and leapfrog step size ϵ were determined for this scenario. Significant improvements in accuracy were attained when using the Bayesian learning procedure. Relative improvements in mean error rates of up to 9.9% and 5.4% were attained when using the small and full Afrikaans training sets respectively.

A small gender independent training set was created using the TIMIT dataset, which was used to test the system under much the same conditions (Section 5.4.2). A 3.8% relative improvement in mean error rate was attained for this scenario. The effect of sample size and the parameter C_v were found to be very similar to that observed when using the SUNSpeech database.

In Section 5.4.3 the system was tested using a connected digit recognition task. The TIDIGITS database was used to create a small gender-independent dataset, so as to simulate the scenario where little training data was available for training a connected digit recognizer. A relative improvement in error rate of 16.7% was attained when using a sample size of two. Using only one sample resulted in an 11.9% relative improvement, due to the choice of prior and hyperparameters. A decrease in performance was observed for sample sizes greater than two, which was ascribed to the systematic errors resulting from the leapfrog integration algorithm.

If one was to have the luxury of having an independent cross-validation set, then the resultant performances can be expected to be even better. Under these conditions, one could create several samples of the posterior and then choose the best sample using the cross-validation set. Although this will not necessarily give the truly best sample, it will probably improve results such that they are closer to the maximum results attained in this section.

5.5 Summary

This chapter introduced Bayesian learning and its usage within a continuous speech recognition framework. Markov chain Monte Carlo methods, which are used to sample the posterior, were described. Implementation specifics such as maintaining the HMM constraints, choice of the prior distributions and parameters, and an efficient implementation for recognition were also discussed.

Section 5.4 experimentally evaluated the Bayesian learning procedure introduced in this chapter. The SUNSpeech, TIMIT and TIDIGITS databases were used for this purpose. The effects of sample size, hyperparameters C_m and C_v , simulation time T and leapfrog step size ϵ were determined. Significant improvements in system accuracy were attained for the scenarios created using all three datasets. The improvements realized for the connected digit task (TIDIGITS) were, however, mainly due to the regularization effect of the prior, with a sample size of two resulting in peak performance. This was attributed to the fact that the connected digit task is relatively simple and the posterior thereof is therefore probably also relatively simple.

The results attained show that there is much promise in the usage of the Bayesian learning procedure proposed. Although considerably more computationally expensive than most HMM training techniques, I believe, however, that the improved performance warrants the additional computational expense.

Chapter 6

Conclusion

6.1 Overview

The problem addressed in this thesis is the design of a speech recognizer when only sparse training data are available. The recording and subsequent segmenting and annotation of large speech databases is an expensive and labour-intensive process. The aim of the techniques developed in this thesis is therefore to improve recognition performance of hidden Markov model (HMM) systems when training data is limited.

This thesis started by considering the minimum classification error (MCE) parameter estimation procedure. The MCE criterion is closely related to the goal of ideal Bayes classification. Overspecialization is, however, prevalent when using MCE, especially when training data is limited. Several modifications to the standard MCE procedure were therefore proposed which limited the effect of overtraining. The MCE algorithm was used within a training framework (Section 1.2).

The usage of the maximum *a-posteriori* (MAP) estimation algorithm was investigated for situations where little task-specific training data is available, but a reasonable quantity of non-task-specific training data can be used. The MAP estimation procedure

was used within an adaptation framework (as described in Section 1.1) in this thesis, where non-task-specific data was used in conjunction with limited task-specific data to improve recognition performance for a specific task.

Finally, Bayesian learning was investigated and an implementation for HMM speech recognition was proposed and implemented. Here no assumption is made about the availability of non-task-specific data, and only the task specific data is used in the training process. The Bayesian learning procedure as proposed and implemented in this work is a training algorithm (Section 1.2), where only task-specific data is used.

The algorithms proposed and tested in this thesis have been found to be extremely effective in situations where little training data is available. Significant improvements in recognition performance have been attained for the sparse data scenarios used for evaluation purposes. The algorithms proposed are, in general, considerably more computationally expensive, with execution times for the GMAP and MAPMCE being of the order of 3 hours and 10 hours respectively (depending on the task and computer used) as compared to the 13 minutes required by MAP (or ML). The Bayesian learning procedure is extremely computationally expensive with the execution time being around one to three days, depending on the task and computer used. The training of speech recognition systems is typically done offline and the time taken to estimate the model parameters is therefore often not critical. The additional computational expense can therefore be justified by the significantly improved recognition performance resulting from the use of these algorithms.

The contributions of this thesis to the field of speech recognition are,

1. New modifications to the MCE algorithm are proposed in Chapter 3. These modifications limit the effect of overtraining which is prevalent when using MCE training.
2. A new gradient-based MAP adaptation algorithm (GMAP) that does not make any assumptions concerning the form of the prior distribution was proposed in

Chapter 4. This algorithm was shown to outperform the standard MAP approach of Gauvain and Lee [45] for the conditions tested.

3. A new MCE based MAP adaptation algorithm was proposed and tested in Chapter 4. This algorithm too was shown to work better than the standard MAP approach, as well as being better than standard MCE.
4. Bayesian learning was introduced. An original implementation of Bayesian learning for hidden Markov model speech recognition was introduced and discussed. This is, to the author's knowledge, the first time that Bayesian learning using Markov chain Monte Carlo methods has been used for hidden Markov model speech recognition.

6.2 Summary by Chapter

This thesis centered on the following topics:

In *Chapter 1* the topic of speech recognition was introduced, along with the topic of data sufficiency and the effects thereof. The algorithms currently used to alleviate the problem of insufficient data were briefly reviewed.

In *Chapter 2* the basic theory of hidden Markov models and the speech recognition system used in this thesis were described. Overtraining was discussed in terms of the bias/variance dilemma. The experimental procedure common to all experiments conducted in this thesis was described and details of the three speech databases used for experimental evaluation were given.

In *Chapter 3* the minimum classification error criterion was introduced and its usage within a continuous speech recognition framework was discussed. The MCE criterion realizes optimal decision boundaries in a Bayes sense. The MCE criterion focuses on minimizing the probability of error. Classifiers trained using the MCE criterion

are therefore able to attain the goal of Bayes classification even when the modeling assumptions are incorrect. Embedded or string-level MCE, where the recognition error for entire strings is minimized, was discussed.

The effect of a smoothed zero-one loss function was discussed and experimentally determined for the TIMIT dataset. Furthermore, the need for a zero-one loss function was questioned and the conclusion was reached that there is little evidence to suggest that there is an advantage or disadvantage to using a smoothed zero-one loss function. This result was used later in modifying the MCE criterion, where the modification could not be mathematically justified when a non-linear loss function was used.

Overtraining within the MCE framework was discussed and three modifications were proposed. The first modification stops the mixture variances from becoming very small, which we expect to happen when little data is available. The second added a weighted likelihood term to the MCE criterion, thereby reinforcing correct substrings, as well as improving discrimination for incorrect substrings. Lastly, a word-based string-level MCE algorithm was proposed, in which smaller word-based substrings were used, instead of the the entire string. Significant gains in performance resulted when using these modifications with the TIMIT database. Improvements of up to 10% in relative error rate on the test set were achieved for the TIMIT dataset.

In *Chapter 4* Bayesian adaption and its usage within a continuous speech recognition framework was introduced. The MAP algorithm of Gauvain and Lee [45] was described. A gradient-based MAP algorithm (GMAP) which makes no assumption about the form of the prior was proposed and the implementation thereof was discussed. A Bayesian inspired MCE-based adaptation algorithm (MAPMCE) was also proposed. The MAPMCE algorithm is a simple extension of the MCE algorithm and its implementation is relatively simple.

The three algorithms were experimentally evaluated using the SUNSpeech database for language adaptation, and using the TIMIT and TIDIGIT databases for gender

adaptation experiments. Figure 6.1 shows the best relative improvements in error rate obtained using the MAP, GMAP and MAPMCE algorithms for the SUNSpeech, TIMIT and TIDIGITS datasets. Relative improvements in error rate of up to 10.2% were attained for the SUNSpeech dataset (using the GMAP algorithm). For the gender adaptation task using the TIMIT dataset relative improvements in error rates of up to 25.5% were attained (using the MAPMCE algorithm). Finally, relative improvements in error rate of up to 66.0% (MAPMCE) were reported for the TIDIGITS based gender adaptation task. When using duration modeling with the connected digit task (TIDIGITS), the GMAP algorithm performed best (relative improvements in error rate of 54.9%). The MAP algorithm, in general, performed somewhat worse.

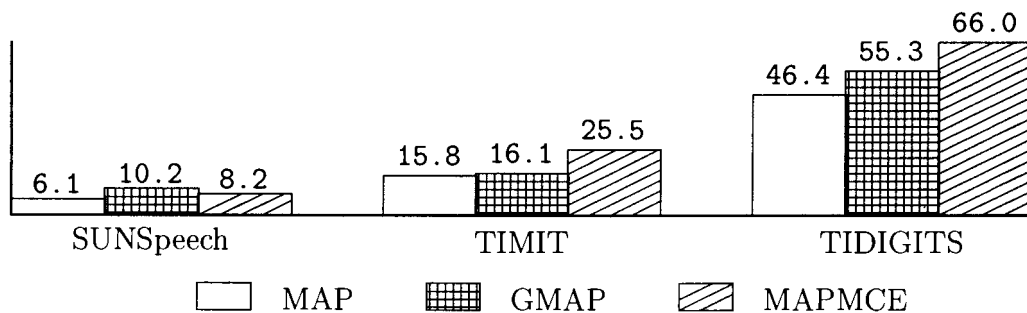


Figure 6.1: Comparison of best relative improvements in error rates obtained using the MAP, GMAP and MAPMCE algorithms for the three datasets used (SUNSpeech, TIMIT and TIDIGITS).

In *Chapter 5* Bayesian learning and its usage within a continuous speech recognition framework was discussed. Markov chain Monte Carlo methods, which are used to sample the posterior, were described. Implementation specifics such as maintaining the HMM constraints, the prior distributions and an efficient implementation for recognition were also discussed.

The SUNSpeech, TIMIT and TIDIGITS databases were used to experimentally evaluate the proposed Bayesian learning procedure. The effects of various algorithm configurations were determined. Significant improvements in system accuracy were attained for the scenarios created using the three datasets. Relative improvements in error rates of up to 9.9%, 3.8% and 16.7% were attained for the SUNSpeech, TIMIT

and TIDIGITS databases respectively. The results attained show that there is much promise in the usage of the Bayesian learning procedure proposed. Although considerably more computationally expensive than standard HMM training, the improved performance warrants the additional computational expense for certain situations.

6.3 Future research

The research performed in this thesis focused on the sparse data problem. The algorithms described were applied to a relatively limited set of tasks and languages. Future research could therefore incorporate a wider variety of typical scenarios where training data is limited, such as cross-database adaptation, speaker adaptation and dialect adaptation. The application of the techniques described here, to context dependent modeling, i.e. bi-phone or tri-phone modeling, will be an important extension of this work. A detailed comparison of Bayesian adaptation techniques with transformation based adaptation algorithms (MLLR, MAPLR) for a wide variety of sparse data scenarios will also be of much interest.

Bayesian methods can also be used for purposes other than the sparse data problem. These include using Bayesian methods to choose the optimal model configuration for a point estimate paradigm, as done by MacKay [72] in the field of neural networks.

The Bayesian learning procedure proposed can be improved upon and several aspects thereof must therefore still be investigated, they include:

- sampling the complete posterior of all HMM parameters including transition probabilities, mixture weights and if applicable, duration modeling parameters,
- using the Hybrid Monte Carlo method to sample the posterior distribution,
- making the modified Viterbi search (Section 5.3.5) more efficient, and

- other prior and hyperparameter configurations may result in even better system performance.