# Fitting of survival functions for grouped data

## on insurance policies

by

## Elizabeth Magrietha Louw

Submitted in fulfilment of part of the requirements

for the degree of Doctor of Philosophiae

in the Faculty of Natural and Agricultural Sciences

University of Pretoria

Supervisor: Professor N.A.S. Crowther

January 2002

# Preface

I am extremely grateful to my supervisor, Professor N.A.S. Crowther, for his careful supervision, positive guidance and constructive comments throughout.

I appreciated the enthusiastic encouragement of my family and friends. The patience and understanding and continuous support of my husband, Johann, and daughters, Hannelie and Annemarie, kept me going.

# Contents

# List of Tables

# List of Figures