# Table of Contents

# Part I:  Data Mining Foundations

# Part II:   Web Mining