# Machine Learning Fundamentals

## A Concise Introduction

Hui Jiang

# Machine Learning Fundamentals

**A Concise Introduction**

Hui Jiang

*York University, Toronto*

# Contents