# Real-Time Large-Scale Data Analytics and Information Retrieval in Practice

Aleksandar Bradic, Igor Bogicevic

December 25, 2009

# Contents

# Chapter 1

# Introduction

# Chapter 2

# The nature of large-scale data

# Chapter 3

# The challenges of real-time information processing

## 3.1 Problem description

# Chapter 4

# The nature of real-time data

# Chapter 5

# Fundamental Algorithms in Data Analytics and IR

## 5.1 Statistical analysis framework

### 5.1.1 Regression analysis

### 5.1.2 Forecasting

### 5.1.3 Parameter estimation

### 5.1.4 Non-parametric methods

# Chapter 6

# Advanced Algorithms

# Chapter 7

# Software toolkits for large-scale data analysis

# Chapter 8

# Large-scale IR Cookbook

**8.1** Building AVMs on vertical data

**8.2** Model selection in the real world

# Chapter 9

# Moving from batch to real-time

## 9.1   Paradigm shift

# Chapter 10

# Concurrency : a new frontier

## 10.1   New challenges

# Chapter 11

# Real-world real-time applications

## 11.4   Quantitive Finance

In this chapter we describe computational aspects related to Quantitative Finance applications :
    - security pricing
    - stochastic process as a central concept in quant finance, as well as the central object in real-time
analytics
    - drawing analogies
    - equivalents of financial concepts in fields such as web analytics
    - continuous vs discrete variables
    - real-time discrete-time
    - real-time continuous-time
    notes :
    - Markov process - variances of the changes in sucessive time periods are additive
    - analogies
    - Twitter topic process
    - Trend detection and keyword bidding
    - Economy of online auctions
    - Towards efficient online marketplaces
    - Prediction markets
    - Financial software deals with real-time, but not large-scale data

## 11.5 Online collaboration

# Chapter 12

# Algorithms and Data Structure in support of large-scale real-time framework

## 12.1 Convolutional procedures

### 12.1.1 Example : Viterbi algorithm

## 12.2 Convolutional representation of fundamental algebraic operations

### 12.2.1 Average,Mean,Median,Variance

### 12.2.2 Matrix operations

## 12.3 Randomized Algorithms

### 12.3.1 Fast vs. Convolutional

## 12.4 Queue-based structures

# Chapter 13

# voidbase : queue-based computing framework

## 13.1  Overview

## 13.2  Paradigms

# Chapter 14

# voidbase cookbook

# Chapter 15

# Future challenges in Real-Time Large-Scale analytical processing

## 15.1 Representation problem

## 15.2 Fundamental limits