#### Welcome to

### DS504/CS586: Big Data Analytics --Introduction & Logistics

Prof. Yanhua Li

Time: 6:00pm –8:50pm THURSDAY Location: AK 232 Fall 2016

## What is DS504/CS586 about?

- A second Level DS/CS course (primarily) for graduates
  - CS/DS Ph.D students in big data analytics and related areas;
  - then other Ph.D students or MS students with
    - Experience in databases and/or in data mining, or equivalent knowledge.
    - Sufficient programming experience is expected so that you are comfortable to undertake a course project.

## **Big Data Analytics**



## techniques and tools for managing, analyzing and extracting knowledge from "big data"



## Introduction

#### What is "Big Data"?

## Big Data – What is it?

- A "big" buzzword ...
- No single standard definition...
- Talk to 1000 people, there will be 1000 "definitions" ...

"*Big Data*" is data whose scale, diversity, complexity, and/or quality require new architectures, techniques, algorithms, analytics, and interfaces to manage it and extract value and hidden knowledge from it...

## Why Now?



#### **Big Data and Big Challenges**



## Big Data

- Volume
- Variety
- Velocity
- Veracity





-Thanks: http://www-01.ibm.com/software/data/bigdata/images/4-



As of 2011, the global size data in healthcare was estimated to be

**150 EXABYTES** [ 161 BILLION GIGABYTES ]



By 2014, it's anticipated there will be

#### 420 MILLION WEARABLE, WIRELESS HEALTH MONITORS



Variety DIFFERENT **FORMS OF DATA** 

30 BILLION PIECES OF CONTENT

are shared on Facebook every month

4 BILLION+ HOURS OF VIDEO

are watched on YouTube each month



400 MILLION TWEETS

are sent per day by about 200 million monthly active users

Thanks: http://www-01.ibm.com/software/data/bigdata/images/4-Vo of big data ing

(u')<sup>2</sup> in m<sup>2</sup>s

POLYTECHNIC BUILD BUILD

The New York Stock Exchange captures

#### 1 TB OF TRADE INFORMATION

during each trading session



#### Modern cars have close to 100 SENSORS

that monitor items such as fuel level and tire pressure



**Velocity** ANALYSIS OF

STREAMING DATA

By 2016, it is projected there will be

#### 18.9 BILLION NETWORK CONNECTIONS

 almost 2.5 connections per person on earth

#### 

Thanks: http://www-01.ibm.com/software/data/bigdata/images/4-V





Thanks: http://www-01.ibm.com/software/data/bigdata/images/ 4-\/s-of-big-data ing

## 4Vs



## The Model Has Changed...

## Old Model of Generating/Consuming Data has Changed

#### **Old Model:**

Few privileged companies are generating and "owning" data, all others are consuming data (in controlled packages)



## The Model Has Changed...

14

 New Model of Generating/Consuming Data has Changed

#### **Producers :**

• Everyone - Man, Woman and Child, and Devices

#### **Consumers:**

- Professionals
- Businesses
- Scientists
- And us
- Everyone wants a piece of this pie ...



## What Sectors Can Benefit?

- Businesses
- Transportation
- Science & Engineering
- Governments
- Energy
- Healthcare
- Education
- Entertainment

Utilize data to improve people's life quality

# Done with the high level introduction

# Begin with application stories

## **Big Challenges in Big Cities**

















Urban Computing: concepts, methodologies, and applications. Zheng, Y., et al. *ACM transactions on Intelligent Systems and Technology*.



Zheng, Y., et al. Urban Computing: concepts, methodologies, and applications. ACM transactions on Intelligent Systems and Technology.

## **Urban Sensing**

#### A sample of data $\rightarrow$ An entire dataset

Biased distribution



• Data sparsity and missing



Air quality monitoring stations

Inferring Gas Consumption and Pollution Emission of Vehicles throughout a City. KDD 2014. Zheng, Y., et al. U-Air: when urban air quality inference meets big data. KDD 2013

## **Urban Sensing**

#### A limited resource (budget, labors, land...)

- Static sensing: Where to deploy sensor to maximize the gain?
- Crowdsensing: How to arrange the incentives dynamically?



Suggesting locations for monitoring stations, KDD 2015

#### Improving Medical Emergency Services using Big Data



- Select locations for Ambulance Stations
- Dynamic ambulance allocation



Yilun Wang, **Yu Zheng**, et al. <u>Travel Time Estimation of a Path using Sparse Trajectories</u>.. KDD 2014 Location Selection for Ambulance Stations: A Data-Driven Approach, ACM SIGSPATIAL 2015



Zheng, Y., et al. Urban Computing: concepts, methodologies, and applications. ACM transactions on Intelligent Systems and Technology.

## **Urban Data Management**

- Managing multi-modality data
  - Categorical and numeric data
  - Different scales, densities, updating frequency, and ST properties

- Dynamic and big volume
  - Group query strategy
  - Computing in parallel







Zheng, Y., et al. Urban Computing: concepts, methodologies, and applications. ACM transactions on Intelligent Systems and Technology.

### **Data Integration vs Knowledge Fusion**



Yu Zheng. Methodologies for Cross-Domain Data Fusion: An Overview. IEEE Transactions on Big Data, 1, 1, 2015.

## **Multi-View-Based Learning**



## Urban Computing for Urban Planning

#### Best Paper Nominee Award at UbiComp 2011 The Most Cited Paper



## **City-Wide Traffic Modeling**

- Partition a city into regions with major roads
- Regions are root causes of the problem



Yu Zheng, et al. Urban Computing with Taxicabs, In Proc. Of UbiComp 2011



#### Shanghai Big Data Hotpot Restaurant



## When Urban Air Meets Big Data

KDD 2013

http://urbanair.msra.cn/



## Air Pollution: A Global Concern !

PM2.5, PM10, NO<sub>2</sub>, SO<sub>2</sub>, CO, O<sub>3</sub>

• Air quality monitor station







## We do not really know the air quality of a location without a monitoring station!



## Inferring Real-Time and Fine-Grained air quality throughout a city using Big Data



Zheng, Y., et al. U-Air: when urban air quality inference meets big data. KDD 2013

## Urban Air System



Zheng, Y., et al. <u>U-Air: When Urban Air Quality Inference Meets Big Data</u>. KDD

## Multi-View Learning Framework

- Features: Non-overlapping features providing different views
- Models: Model extrapolation and trend regression respectively
- Training: Combination of small models vs. a big model







## **Revisit Big Data**

- NOT a single data source which is very big
- NOT mean full data
- NOT mean very dense data
- May need less domain knowledge

•

- Data across different domains
- Sample of (label) data
- Data sparsity always exists
- More understanding of data itself and data science
- Many unsolved problems

Tools are ready Big Data ≠ Mining Single Dataset ≠ Simple Statistics Big Data ≠ Machine learning ≠ Deep Learning Big Data ≠ Cloud Computing ≠ Hadoop

Big Data needs comprehensive capabilities to deliver end-to-end services!



## Take Away Messages

- 3B: *B*ig city, *B*ig challenges, *B*ig data
- 3M: Data Management, Mining and Machine learning
- 3W: Win-Win-Win: people, city, and the environment

#### 3-BMW

Zheng, Y., et al. Urban Computing: concepts, methodologies, and applications. ACM transactions on Intelligent Systems and Technology.

Yu Zheng. Trajectory Data Mining: An Overview. ACM Transactions on Intelligent Systems and Technology. 2015

Yu Zheng. Methodologies for Cross-Domain Data Fusion: An Overview. *IEEE Transactions on Big Data*, 1, 1, 2015.

#### What data available in our course?

## **Road Map in Shenzhen**



20,656 Road Segments

## **Subway Lines**



## **Bus Routes**



8,875 Buses serve 814 Bus Routes

## **Bus Stop Distribution**





## **Transportation Billing Data**







## **Urban Issues**



#### Regional Weather-Traffic Sensitivity





Traffic Estimation & Prediction



Smart shuttle service

## **Urban Issues (cont.)**



Logistic Planning

Low Sample Rate Map Matching

## Questions?