

IST718

Big Data Analytics

Unit 1: A Course Introduction

About me

- Daniel Acuna, Assistant Professor, iSchool
 - Ph.D. Computer Science, University of Minnesota, Twin Cities
 - Postdoctoral Researcher, Northwestern University & RIC
 - Member of Metaknowledge Research Network, University of Chicago
 - Affiliated to the Center for Computational and Data Science @ SU
- Research interests
 - “Science of science”, human-AI collaboration

Recommendation system



$$u = (1 + \alpha) \frac{\sum_{i \in \text{Relevant}} d_i}{|\text{Relevant}|} - \beta \frac{\sum_{j \in \text{Not Relevant}} d_j}{|\text{Not Relevant}|}$$



Activity recognition: in the past

We have very good quality but invasive features:



- GPS
- Audio
- WiFi location
- App activity

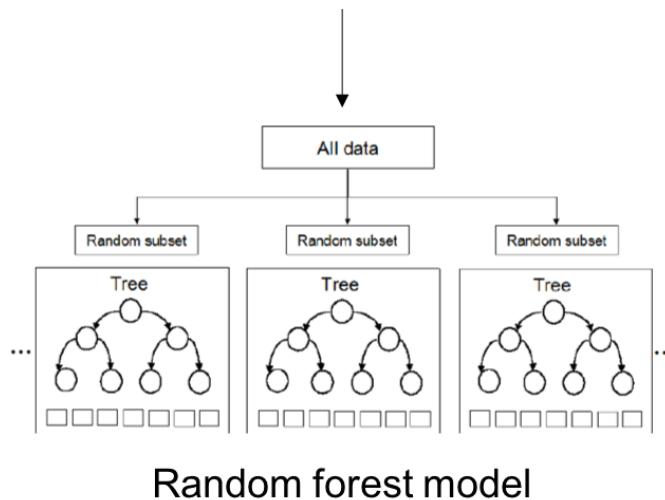


Predict location and activity

Activity recognition: big data

With big data we might not need such high quality features but just a ton of data

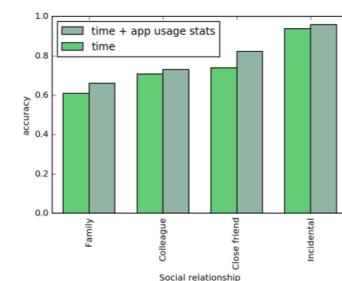
- App open and close activity
- And clock 2 billion seconds of phone activity



Show Me Your App Usage and I Will Tell Who Your Close Friends Are: Predicting User's Context from Simple Cellphone Activity

Alain Shema
School of Information Studies
Syracuse University, Syracuse,
USA
sralain@syr.edu

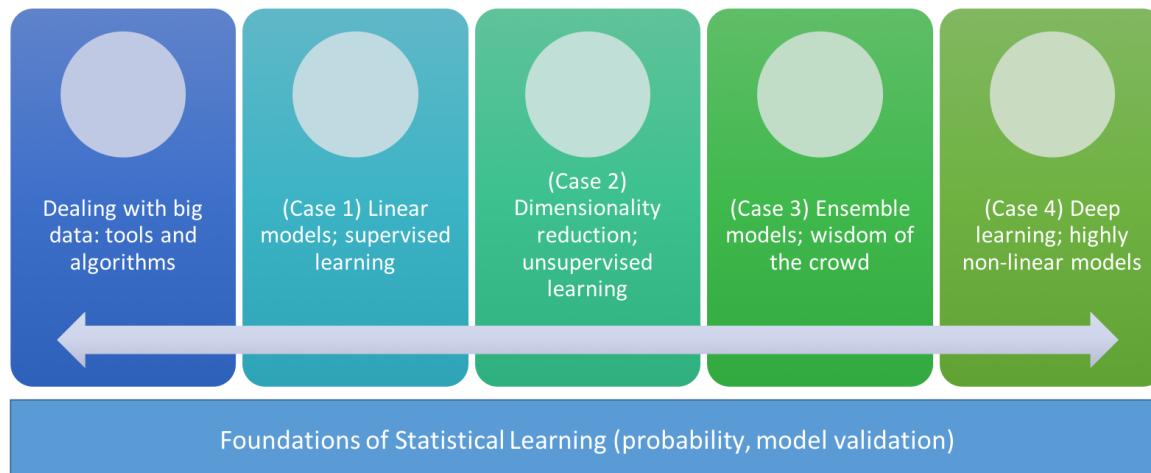
Daniel E. Acuna
School of Information Studies
Syracuse University, Syracuse,
USA
deacuna@syr.edu



About the course (go over Syllabus)

- 1/4 of the course covers prerequisite skills required for big data analytics
 - Python programming, understanding machine learning
- 1/4 of the course covers the Spark and Hadoop, providing skills required to perform big data analytics.
- 1/4 of the course consists of case studies, where you apply your skills and knowledge to real-world applications.
- 1/4 of the course consists of a project, where you work in groups in a real world application

Course roadmap



- Small/medium data
- Low model complexity
- High interpretability
- Low computational power
- Big data
- High model complexity
- Low interpretability
- High computational power

Course preview (1)

- Slides

Results

- 60% training, 20% validation, 10% testing
- Regularized logistic regression: 63% AUC
- Random forest: 84% AUC

#	Rank	Team Name * <small>In the money</small>	Score	Entries	Last Submission UTC	(Best - Last Submission)
1	↑1	Perfect Storm <small>↓*</small>	0.869558	128	Thu, 15 Dec 2011 05:35:00 (-3.2d)	
2	↑4	Gxav *	0.869295	54	Thu, 15 Dec 2011 09:41:23 (-26.9h)	
3	↑14	occupy *	0.869288	9	Thu, 20 Oct 2011 00:40:05	
4	↑16	D'yakonov Alexander (MSU, Moscow, Russia)	0.869197	64	Thu, 15 Dec 2011 22:08:19 (-5.1d)	

Course preview (2)

- Notebooks and labs to be done on your own

The screenshot shows a Databricks notebook interface with the following content:

- Header:** lab-sentiment_analysis (PySpark)
Detached File View: Code Permissions Run All Clear
- Text:** After fitting, the Pipeline becomes a transformer:
- Diagram:** PipelineModel (Transformer) → Tokenizer → HashingTF
PipelineModel .transform() Raw text → Words
- Text:** (Images from <http://spark.apache.org/docs/latest/ml-pipeline.html>)
- Text:** Importantly, transformers can be saved and exchanged with other data scientists, improving reproducibility
- Section:** Loading packages and connecting to Spark cluster
- Code (Cmd 8):**

```
1 import numpy as np
```

Command took 0.64 seconds -- by deacuna@syr.edu at 11/30/2017, 8:47:48 AM on My Cluster
- Code (Cmd 9):**

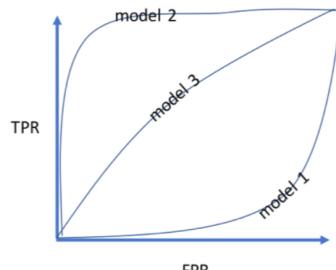
```
1 # dataframe functions
2 from pyspark.sql import functions as fn
3 from __future__ import division
```

Command took 0.67 seconds -- by deacuna@syr.edu at 11/30/2017, 8:47:48 AM on My Cluster
- Section:** Transformers and Estimators

Course preview (3)

- Quiz

8 (2 pts) Multiple choice: In terms of area under the curve (AUC), from best to worse, rank the following models based on their ROC curve



- a) Model 1, model 2, and model 3
- b) Model 2, model 1, and model 3
- c) Model 1, model 3, and model 2
- d) Model 2, model 3, and model 1

Course preview (4)

- Project

Project Proposal, IST 718
Predicting the trends of cryptocurrencies based on historic data

Team members: Neha Humbal, Shikhar Agrawal, Smit Udani, Suchitra Deekshithula

Our project's aim is to look at different cryptocurrencies prevailing in the market. For example, Bitcoin is a digital currency which is decentralized. It is a peer-to-peer network which handles the generation of units of currency and verify the transfer of funds, operating independently of a central bank. Since evolution, cryptocurrencies like Bitcoin have grown more than 1000%, trading at over \$10000. We are trying to predict the price of Bitcoin by taking into account several factors like price of other cryptocurrencies, news, hype, buying and selling on exchanges.

The broader idea is to understand the trends in each of these currencies to strategize investments. The main challenge is that there is no central bank or authority which regulates them. Without understanding factors associated with the fluctuation, it is difficult to make future predictions. Experts are considering Bitcoin for funding. Considering there is neither an optimal nor a permanent solution to the problem of inflation, we can consider the price of Bitcoin to be determined in total about 15 million which have been mined for now which shows inflation levels are decreasing. Cryptocurrency has the potential to revolutionize the world and act as a digital gold card which could help centralized payments and transactions and reduce fraud. Our objective is to identify and analyze the factors associated with each of the several cryptocurrencies and understand the trends.

Goal: Predict the value of a particular cryptocurrency

Tasks

- **Data Acquisition:** Data has been obtained from Kaggle
- **Data Preparation:** Clean, format and identify relevant attributes from the data
- **Feature Extraction:** Check which attributes contribute better to the model and help in better predictions
- **Model Selection and Testing:** Create different models, evaluate models based on validation data test and compare prediction accuracy on test data
- **Visualization:** Present the analysis and trends using appropriate visualization techniques

Experimentation: Designing a model to predict future values and seasonal trend in price fluctuation based on historical prices/ market capitalizations of various currencies

Expected problems: Considering the data we have is of stock market past trends, the data could exhibit unexpected drops. We are taking several currencies into account for handling this inconsistency and this would be something we would have to be careful about

Dataset: This dataset has been obtained from Kaggle and has historical data of several cryptocurrencies. Data can be retrieved from:
<https://www.kaggle.com/sudarshanraj/cryptocurrencyhistory>

Tools: Databricks, PyCharm, Spark, Scikit-learn, Seaborn, Pandas, Numpy, Matplotlib, etc

Models: Linear regression, Polynomial regression, Knn, Svm, Random forest, CART

Criteria: Final model will be selected by comparing the performance of different models using k-fold cross validation and checking accuracy.



Course preview (5)

- Discussion of current trends throughout the semester

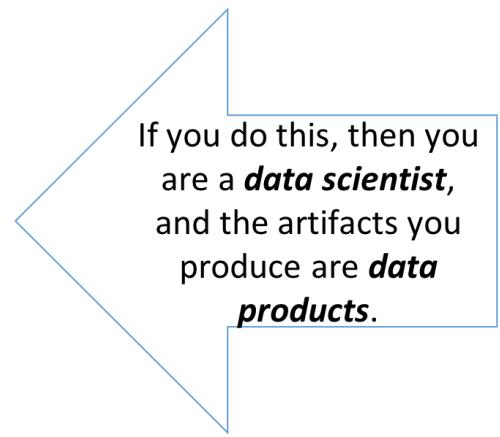
What is one of Google's most important assets?

- 1) Location
- 2) Cash
- 3) Data
- 4) Engineers

Question: What is data science?

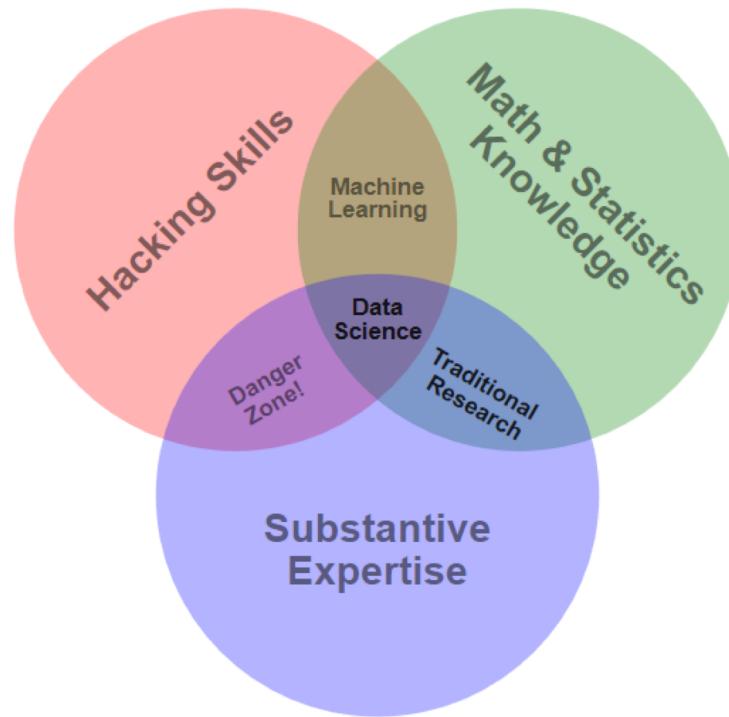
Data Science is

- A combination of disciplines:
 - Information / Computer science
 - Mathematics
 - Statistics
 - Research / Management Science
 - Domain Knowledge
- With the goal of:
 - Using data to make decisions and drive actions



If you do this, then you are a ***data scientist***, and the artifacts you produce are ***data products***.

Data Science Venn Diagram*



*Drew Conway:

https://s3.amazonaws.com/aws.drewconway.com/viz/venn_diagram/data_science.html
[\(https://s3.amazonaws.com/aws.drewconway.com/viz/venn_diagram/data_science.html\)](https://s3.amazonaws.com/aws.drewconway.com/viz/venn_diagram/data_science.html)

"Classic" data science

- Expert is in charge of creating model
- Expert is in charge of providing features that describe or predict new data
- Expert typically produces small models
- Expert typically produces very transparent and easy to understand models

Big Data

What is Big Data?

Data with the following characteristics:

- Data Volume too large to store on a single system.
- Data Velocity too fast for processing by a single computer.
- Data Variety too complex for traditional processing techniques.

These are known as the "three V's" of big data.

"A new kind of data science"

Classic data science



- Expert is in charge of creating model
- Expert is in charge of providing features that describe or predict new data
- Expert typically produces small models
- Expert typically produces very transparent and easy to understand models

Big-data data science

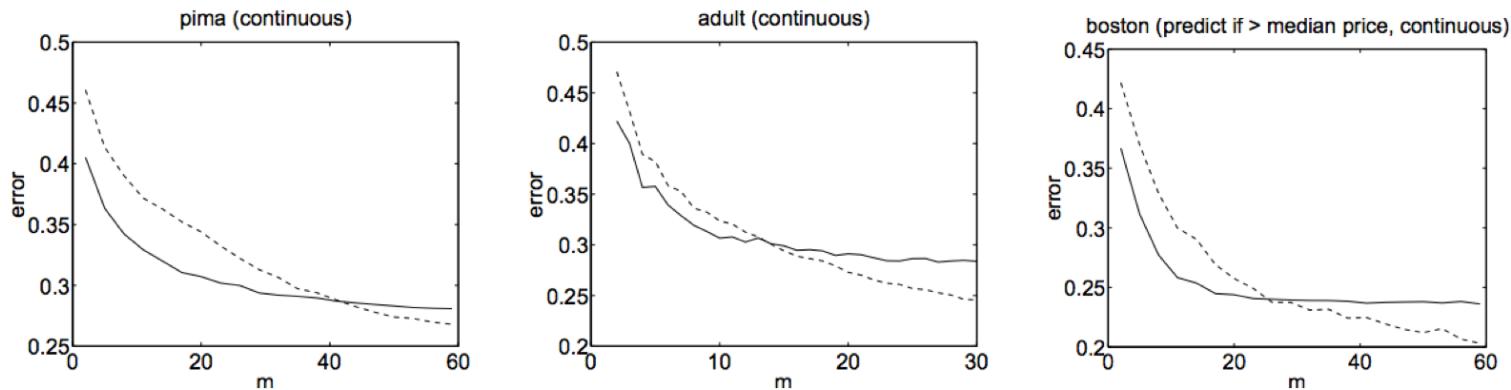
- Generally, there is no expert
- Model used is very general or there is no model at all!
- Features are very low level (e.g., raw transactions vs credit scores)
- Models are very large when fit (e.g., Baidu speech recognition is several terabytes)
- Models are black boxes and almost impossible to understand

On Discriminative vs. Generative classifiers: A comparison of logistic regression and naive Bayes

(NIPS 2001)

Andrew Y. Ng
Computer Science Division
University of California, Berkeley
Berkeley, CA 94720

Michael I. Jordan
C.S. Div. & Dept. of Stat.
University of California, Berkeley
Berkeley, CA 94720



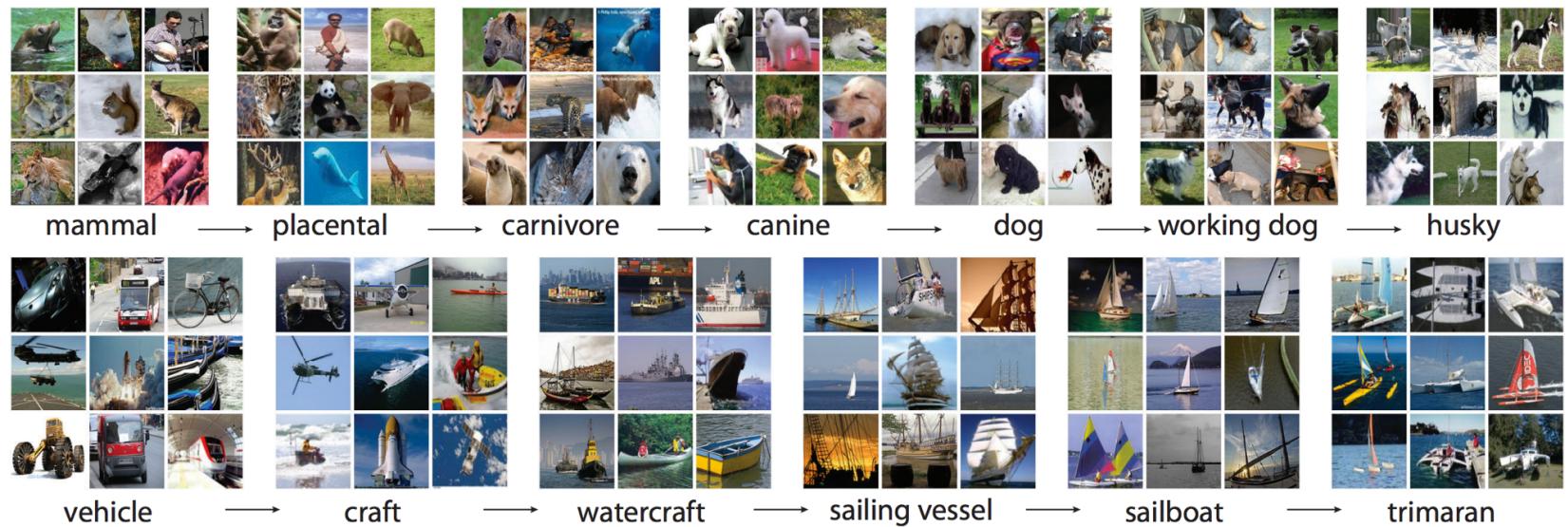
Proposition 1 Let h_{Gen} and h_{Dis} be any generative-discriminative pair of classifiers, and $h_{\text{Gen},\infty}$ and $h_{\text{Dis},\infty}$ be their asymptotic/population versions. Then¹ $\varepsilon(h_{\text{Dis},\infty}) \leq \varepsilon(h_{\text{Gen},\infty})$.

Image recognition



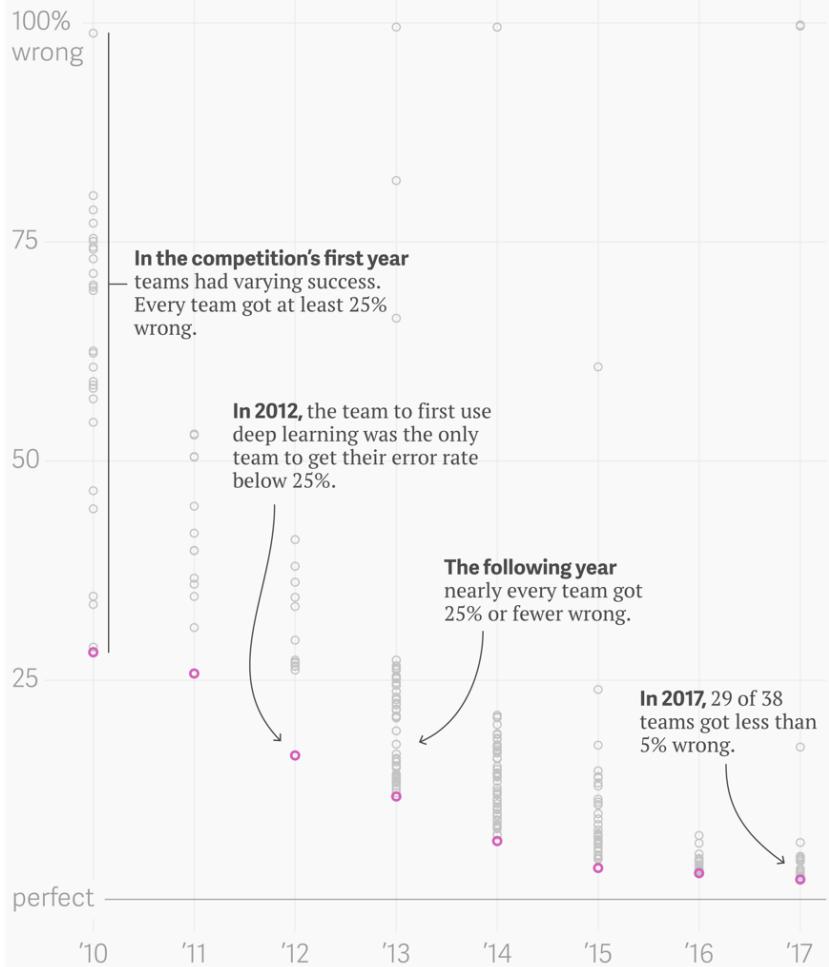
A task that was incredibly challenging and a benchmark for computer algorithms

Image recognition: ImageNet



First models were based on Nearest Neighbor and Naïve Bayes algorithms

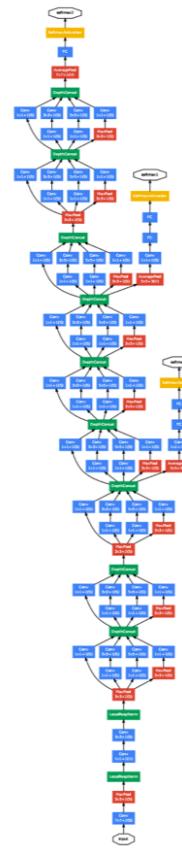
ImageNet Large Scale Visual Recognition Challenge results



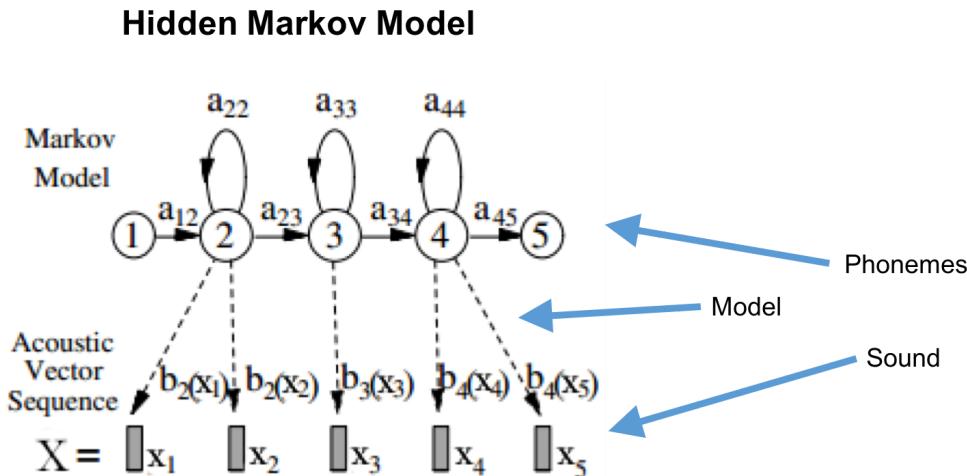
State of the art ImageNet model

GoogleLeNet neural network model

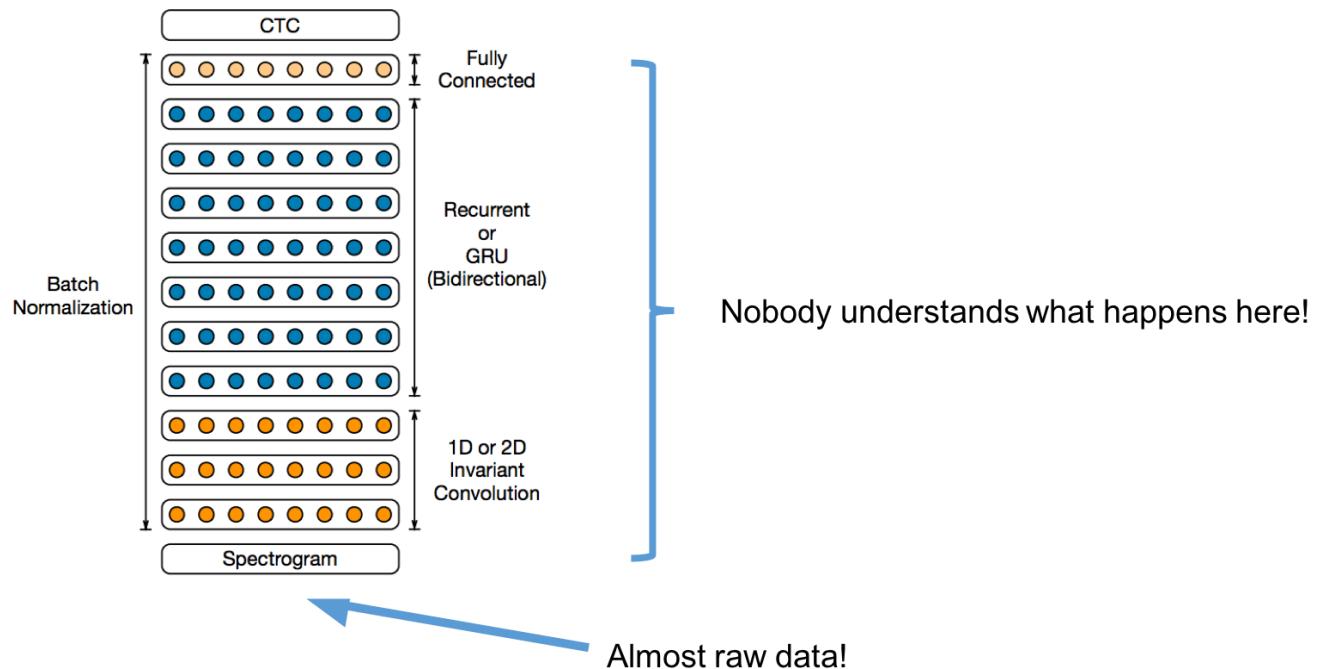
Team	Year	Place	Error (top-5)	Uses external data
SuperVision	2012	1st	16.4%	no
SuperVision	2012	1st	15.3%	Imagenet 22k
Clarifai	2013	1st	11.7%	no
Clarifai	2013	1st	11.2%	Imagenet 22k
MSRA	2014	3rd	7.35%	no
VGG	2014	2nd	7.32%	no
GoogLeNet	2014	1st	6.67%	no



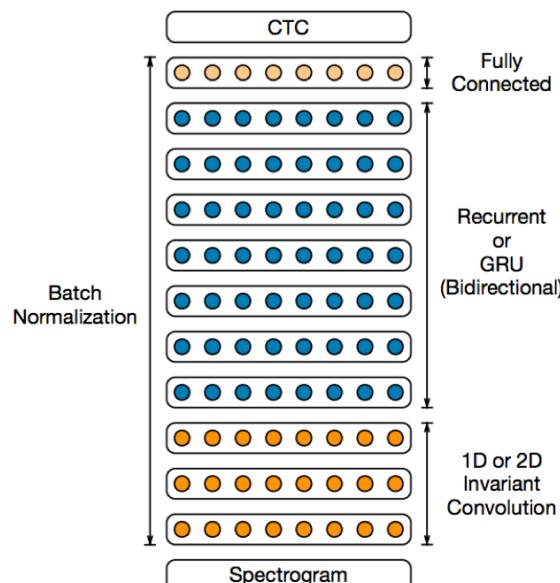
Speech recognition: Generative model



Speech recognition: Deep Speech 2 (Baidu) (1)

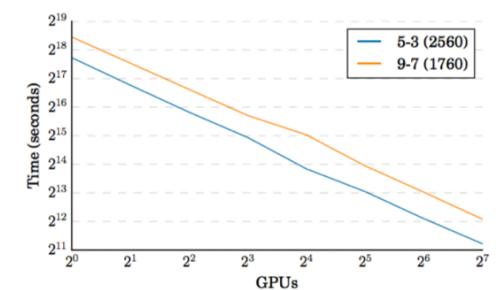


Speech recognition: Deep Speech 2 (Baidu) (2)



Test set	Read Speech		
	DS1	DS2	Human
WSJ eval'92	4.94	3.60	5.03
WSJ eval'93	6.94	4.98	8.08
LibriSpeech test-clean	7.89	5.33	5.83
LibriSpeech test-other	21.74	13.25	12.69

Model size	Model type	Regular Dev	Noisy Dev
18×10^6	GRU	10.59	21.38
38×10^6	GRU	9.06	17.07
70×10^6	GRU	8.54	15.98
70×10^6	RNN	8.44	15.09
100×10^6	GRU	7.78	14.17
100×10^6	RNN	7.73	13.06



Music: Implicit-feedback ALS (Spotify)

Users

$$\begin{pmatrix} 10001001 \\ 00100100 \\ 10100011 \\ 01000100 \\ 00100100 \\ 10001001 \end{pmatrix} \approx \underbrace{\begin{pmatrix} X \\ Y \end{pmatrix}}_f f$$

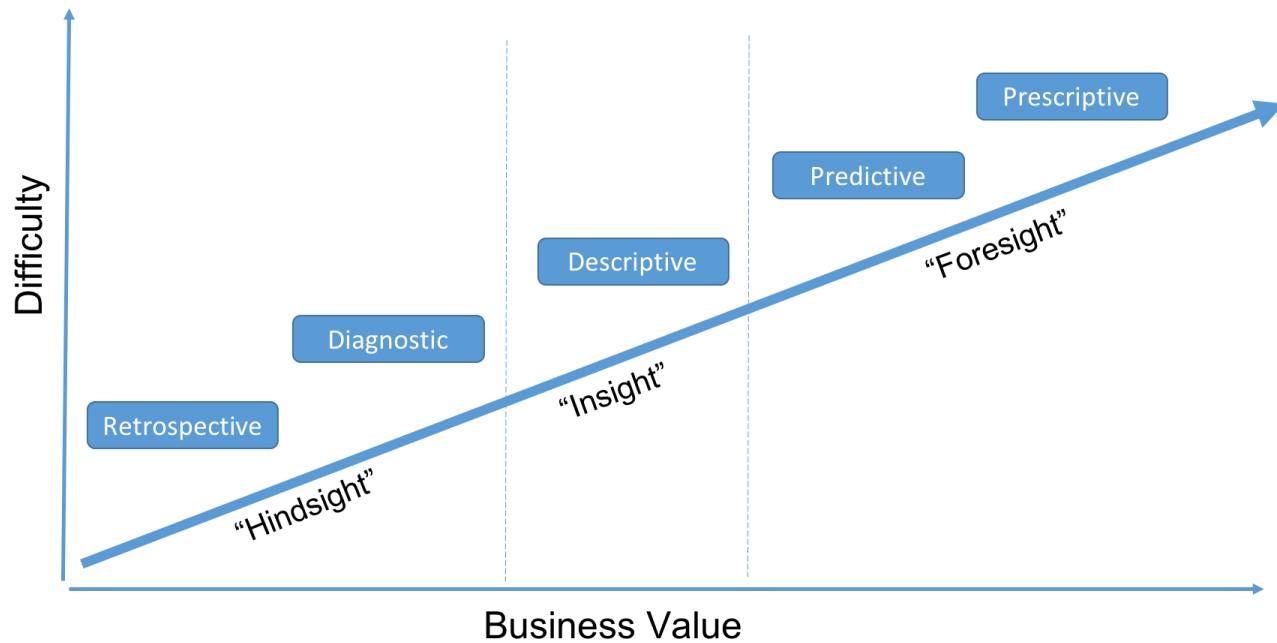
Songs

Fix tracks

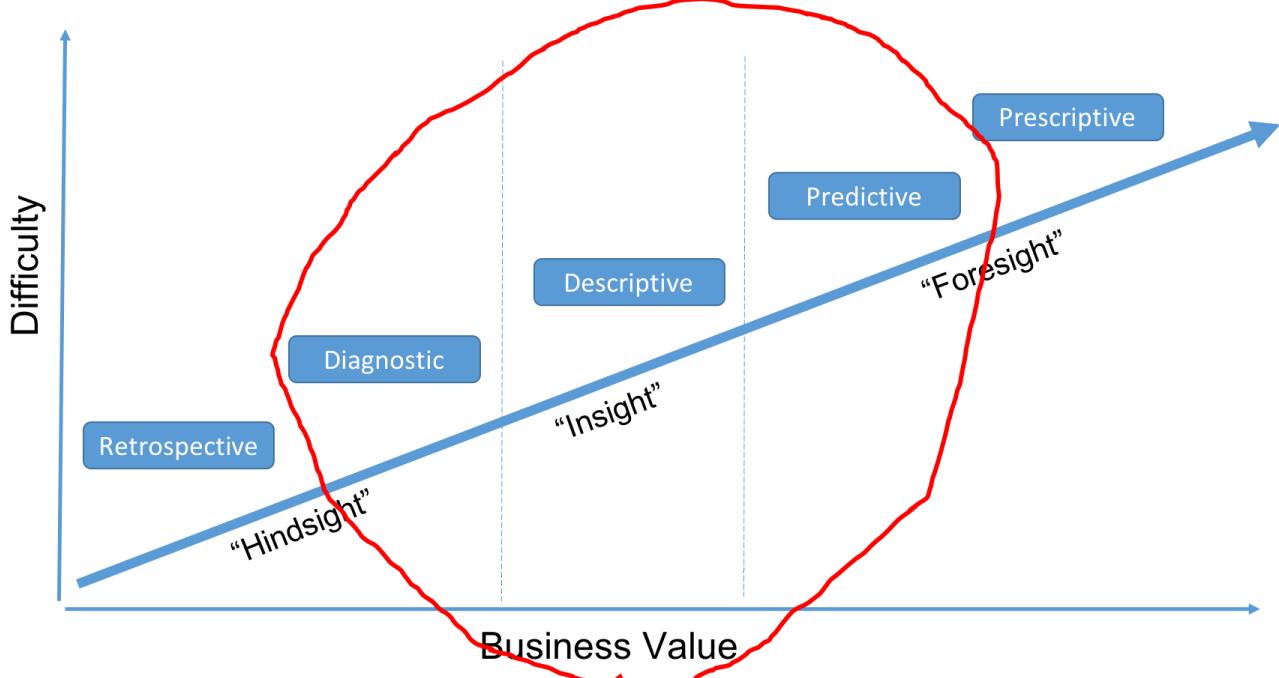
$$\min_{x,y} \sum_{u,i} c_{ui} (p_{ui} - x_u^T y_i - \beta_u - \beta_i)^2 + \lambda (\sum_u \|x_u\|^2 + \sum_i \|y_i\|^2)$$

- p_{ui} : 1 if user u streamed track i else 0
- $c_{ui} = 1 + \alpha r_{ui}$
- x_u = user u 's latent factor vector
- y_i = item i 's latent factor vector
- β_u = bias for user u
- β_i = bias for item i
- λ = regularization parameter

What will this course cover? (1)



What will this course cover? (2)



First steps

- The sooner you start playing with Python and Spark, the better
- Create a Github account <http://github.com> (<http://github.com>)
- Give me your Github username by completing the "Github username" test assignment on Blackboard
- Go to <http://notebook.acuna.io> (<http://notebook.acuna.io>) and use Github to login
- Optional
 - Install Anaconda Python 3.6 <https://www.anaconda.com/download/> (<https://www.anaconda.com/download/>)
 - Install PySpark through pip <https://pypi.org/project/pyspark/> (<https://pypi.org/project/pyspark/>)
- Go to <http://spark.apache.org> (<http://spark.apache.org>) and read the documentation
 - Do “Quick Start”, and maybe “SQL, DataFrames, and Datasets”
 - See how the API is structured
<http://spark.apache.org/docs/latest/api/python/index.html>
(<http://spark.apache.org/docs/latest/api/python/index.html>)

Questions?

- Contact me at deacuna@syr.edu
- Leave a message in Blackboard class forum
- Visit office hours often