

LECTURE #1.

January 28 2014.

①

Administrivia.

Social Network X
CSCI-GA.3033-004.

Rm 1002, 715 Broadway, NY. 10003.

Phone: 212.998.3464

Email: mishra@nyu.edu.

Office Hrs: By appt. +


Class Hrs. + Class Room:

Tuesdays, 5¹⁰ pm - 7⁰⁰ pm. EST.

CIWW (251 Mercer Street). Room 317.

Text Book: John Hopcroft & Ravi Kannan.
Foundations of Data Science.
Unpublished Notes, 2013.
Available on Course Home page.

+ Mark Newman
Networks: An Introduction
Oxford University Press, 2010.



◇ Introduction

◇ Syllabus:

- a) Statistical Inference / Data Science
- b) Graph Theory / Network Analysis.
- c) Game Theory / Signaling
- d) Entrepreneurship.

◇ Motivation: Startup 2.0.

• The Economist, January 18 2014.

Tech Startups: A Cambrian Moment.
(A 14-page Special Report on Tech Startups)

Coming to an office near you...

What today's technology will do to
tomorrow's jobs.

- Creating jobs / Venture Capitalists / Accelerators / Building Companies / Business Communities / The Dark Side / Hardware Startups / Platforms

- A Cambrian Moment.



Snowball Earth



Oxygen ↑
Ozone ↑

→ Multicellularity (Germline Evolution)



Haplo-diploidy



Sex → Mate Selection

↓ Deception

Social Grooming

↓ Loss of Hair

Language (State of Nature)

↓ Deception

Migration



Trust
Verification

(Society - Social Contract)

Incredible
Threat

(Religion)

Doctrine of Two-Swords
Monotheism

Prophets

↓ Conflict/Deception

Information

Stability

[Social Media
Complexity
→ Low Error]

[Low Complexity
Stable]

↓
High Cost/Energy

↓
Non-adaptive

Global Warming
Carbon Dioxide ↑

Ozone ↓

↓
Extinction
(Water World.)

CAMBRIAN RADIATION:

(3)

◇ The first geological period of the Paleozoic Era
542 mya ~ 486 mya.

◇ Evolution of Multicellularity:

- Animals, Plants, Fungi, Brown, Red- and Green-Algae.
- Signaling: Tyrosine Kinases.
MicroRNA ...
- Symbiosis: Mitochondria, Chloroplasts, ...
- Cytophagy:
- Sex: Haplo-diploidy.
- Morphology: Organelle, Organs, Organisms ...
- Eco-Systems: Eco-System Engineering.

◇ Causes:

- Environmental Changes

~ Snowball Earth (Metastable Climate)

~ Increase in Oxygen Level (Result of photosynthesis)
Area of oxygen absorbing surface

~ Ozone Formation

~ Volcanic Activity of Mid-ocean Ridges

- Calcium + Phosphorus.

- Exoskeleton, skeleton, Teeth, Morphology

~ Ecology

- Food Chain, Arms race among Prey-Predators

⇒ Signaling

⇒ Ecological Networks

⇒ Competition & Cooperation..



ANALOGY.

④

- 1) Signaling → Data + Platforms
 - 2) Ecological Networks → Social Networks, P2P, B2B networks, Communities
 - 3) Competition → Game Theory, Deception, Security, Privacy.
- Desiderata

"Why Software is Eating the World?"

Marc Andreessen. WSJ, Aug 20, 2011.

- a) Understanding Data Science (Hopcroft + Kannan)
- b) Understanding Technology Communities (Silicon Valley vs. Silicon Alley)*
- c) Understanding Deception (Cyber Security, Bitcoins, Market Microstructure)

Data

Datum vs Desideratum.

Dare → Desiderare

Give → Desire

Example: Omne Datum Optimum

Pope Innocent II 1139

Order of the Poor Knights of Christ
and of the Temple of Solomon
(Knights Templar)

"As for the things that you will receive from the spoils, you can put confidently^{put} to your own use, and we prohibit that you be coerced against your will to give anyone a portion of these."
Every Perfect Gift.

(5)

Datum \equiv Given, Gift.

→ Input, what is given, known, ...

Gift and Social Contract:

Marcel Mauss (1872-1950) "The Gift"

"Essai sur le Don"

→ Open Source

Gifts → Reciprocal Exchange

Inalienable Bond between Giver & Gift.

Failure to reciprocate →

Results in loss of Social Status.

Social Grooming

Data Economy → Based on Social Contracts

→ Social Networks.

◊ Social Status / Social Rank / Reputation

◊ Signaling / Reciprocal Exchange / Trust

◊ Properties (Alienability) / Valuation, Recommendation

~ Freemium, Open Source, Potlatch Economy,

Crowd-Sourcing

~ Deception.

Mapping: Data → Desiderata

↓

Information

(Asymmetric Information Game
Between an Informed Sender
and an Uninformed Receiver)

↓

Deception.

{ Machine Learning
from
Data, MetaData,
Provenance Data.

(6)

Two Examples for Data Science:

- a) House Price : Data: a) Lot Size, x_i
(amount of land being conveyed from a buyer to seller)
b) House Price, y_i
(Transaction deed data).

$$D = \{ \langle x_i, y_i \rangle \mid i = 1, \dots, n \}$$

← Training Set.


Desiderata: Given an unpriced/new house with lot size x_0 , estimate y_0 .
 y_0 = Your bid price.

- b) Cancer Survival: Data: a) Tumor Size, x_i
(E.g. Tis (DCIS), Breast cancer, Ductal Carcinoma in situ TisX_i, code records only greater dimension.)
→ At detection.
b) Survival Time, y_i
(Time to relapse)
Right-censoring Problem.

$$D = \{ \langle x_i, y_i \rangle \mid i = 1, \dots, n \}$$

← Training set

Desiderata: Given a newly diagnosed breast cancer (DCIS) with size x_0 , estimate y_0 .
 y_0 → Determines whether to select mastectomy?



Simple Models. (Linear Regression)

(7)

→ Find the equation of a straightline that provides the best fit for the data point.

$$y = \alpha + \beta x + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$y | \alpha, \beta, x, \sigma^2 \sim \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(y - \alpha - \beta x)^2}{2\sigma^2}}$$

Maximum Loglikelihood Estimator:

$$\arg \min_{\alpha, \beta} \sum_{i=1}^n \hat{\epsilon}_i^2 = \arg \min_{\alpha, \beta} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

Heuristic Soln:

$$y_i = \alpha + \beta x_i$$

$$\sum y_i = \alpha \sum 1 + \beta \sum x_i \rightarrow \bar{y} = \alpha + \beta \bar{x}$$

$$\sum x_i y_i = \alpha \sum x_i + \beta \sum x_i^2 \rightarrow \bar{x}\bar{y} = \alpha \bar{x} + \beta \bar{x}^2$$

$$\therefore \begin{aligned} \bar{x}\bar{y} &= \alpha \bar{x} + \beta \bar{x}^2 \\ \bar{x}\bar{y} &= \alpha \bar{x} + \beta \bar{x}^2 \end{aligned}$$

$$\beta [\bar{x}^2 - (\bar{x})^2] = \bar{x}\bar{y} - \bar{x}\bar{y}$$

$$\hat{\beta} = \frac{\bar{x}\bar{y} - \bar{x}\bar{y}}{\bar{x}^2 - (\bar{x})^2} = \frac{\text{Cov}[x, y]}{\text{Var}[x]}$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} = \bar{y} - \frac{\bar{x}\bar{y} - \bar{x}\bar{y}}{\bar{x}^2 - (\bar{x})^2} \cdot \bar{x}$$

$$= \frac{\bar{y} \bar{x}^2 - \bar{y} (\bar{x})^2 - \bar{x} \bar{x}\bar{y} + (\bar{x})^2 \bar{y}}{\bar{x}^2 - (\bar{x})^2}$$

$$= \frac{\bar{y} \bar{x}^2 - \bar{x} \bar{x}\bar{y}}{\bar{x}^2 - (\bar{x})^2}$$

Are these models good?

(8)

◊ Linear relationship between a scalar dependent variable y and a single independent/explanatory variable x may not suffice.

- More explanatory variable?

$$y = \theta^T x + \epsilon \quad \left\{ \begin{array}{l} \vec{y} = X \vec{\theta} \rightarrow x^T \vec{y} = (x^T X) \vec{\theta} \\ \hat{\theta} = (X^T X)^{-1} X^T \vec{y} \end{array} \right.$$

- How many more?
Problem with underfitting/overfitting.

$y = \theta^T x + \epsilon \rightarrow$ Sparse Solution

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \|\vec{y} - X \vec{\theta}\|_2^2 + \lambda \|\vec{\theta}\|_0$$

Note $\|x\|_p = \sqrt[p]{\sum x_i^p}$

$$\|x\|_0 = \sqrt[0]{\sum x_i^0} = \sum \mathbb{1}_{x_i \neq 0} = \#(i | x_i \neq 0)$$

$$\|x\|_1 = \sum |x_i|$$

$$\|x\|_2 = (\sum x_i^2)^{1/2}$$

⋮

Another Formulation:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \|\vec{y} - X \vec{\theta}\|_2^2 + \lambda \|\vec{\theta}\|_1$$

Lasso Regression

→ Shrinkage + Continuous Model Selection.

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \|\vec{y} - X \vec{\theta}\|_2^2 + \lambda \|\vec{\theta}\|_2$$

Ridge Regression

→ Shrinkage to origin

No sparsity is forced.

Computational Complexity: NPC, LP, Linear Algebra...

Are these models sound?

9

◇ Gauss-Markov Theorem:

In a linear regression model in which errors have
{ expectation zero,
no correlation, and
equal variances,
the best linear unbiased estimator (BLUE) of the
coefficients is given by ordinary least square (OLS)?

- GLS (Generalized Least Squares)

$$\vec{y} = X\vec{\theta} + \epsilon \quad \epsilon \sim \mathcal{N}(\vec{0}, \Sigma)$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} (\vec{y} - X\vec{\theta})^T \Sigma^{-1} (\vec{y} - X\vec{\theta})$$

$$\hat{\theta} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} \vec{y}$$

- Knowledge about Relationship among
explanatory variables (e.g. gene network)
or dependent variables (e.g. neighbour's prices)

~ Graph Laplacians

~ Kriging / Cokriging

⋮

Model Selection?

◇ Model Complexity (Rate) vs Sample Error (Distortion)

Connection to Rate Distortion Theory (RDT)
Information Theory.

- Regression, Logistic Regression, Ridge -, Lasso, ...
- Support Vector Machines
- Ensemble Methods
- PCA, ICA, ...

There will still be some issues with our housing and cancer models ...

Housing Model:

◇ Least square error → ?

Need to model mortgage loans

→ Default, Prepayment

→ Credit rating, Interest Rates

→ Subprime, Regular, Jumbo

→ Derivative Market

Risks, Liquidity, hedonic

◇ Do houses have an intrinsic ~~hedonic~~ price?

◇ Can it be manipulated?

Gentrification.

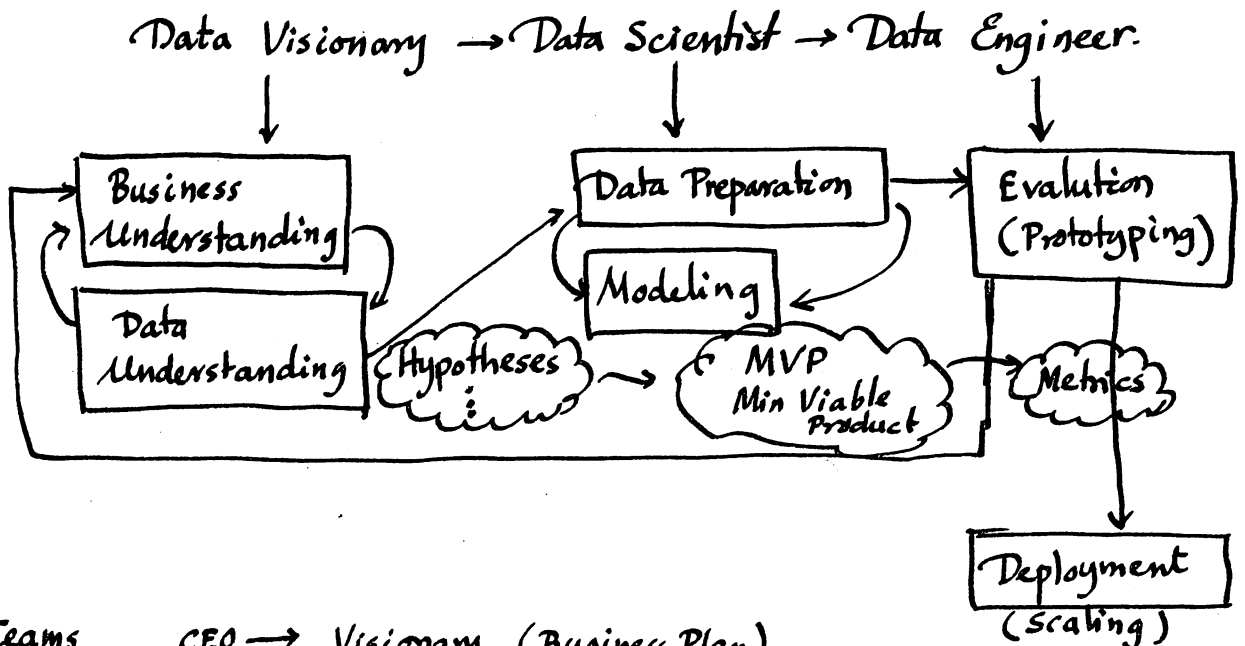
Securitization ...

Cancer Model ?

- ◊ Gene Network , gene-gene interaction.
- ◊ Genomic Instability , copy number, syntenic correlations.
- ◊ Multiple Hypothesis Testing
Driver mutations.
Oncogene vs Tumor Suppressor Genes.
Dominant vs Recessive.
- ◊ Progression Models. Causality.
- ◊ Population Stratification - / Age / Covariates
P73 vs P53.

Examples BRCA1
Age effect
LOH, 2hit.

Projects : Structure:



- Teams
- CEO → Visionary (Business Plan)
 - CSO → Scientist (IP/Design)
 - CTO → Engineer
 - CIO → Deployment
 - CFO → Deployment