Introduction aux sciences des données

Kavé Salamatian





Data Science – A Definition

Data Science is the science which uses computer science, statistics and machine learning, visualization and human-computer interactions to collect, clean, integrate, analyze, visualize, interact with data to create data products.

Turn data into data products.





Data Science – A Visual Definition









Why computers ?

- Capture and storage of massive amount of data
 - Diverse, heterogeneous, imperfect
- Data processing
 - Sanitisation
- Statistical modeling
 - Co-occurrences/correlations
 - Uncertainty structure
- Specialisation
 - Application of the model to specific individuals
- Action
 - Intelligent behaviour

DATASphere AND AnthropocenE

The Economist

HAY 287H-JUNE 34D 2011 Economist.com

me to the Anthro

Getting Spain's protesters off the plazas Obama, Bibi and peace The costly war on cancer How the brain drain reduces poverty A soft landing for China

Every 60 seconds

98,000+ tweets

695,000 status updates

11million instant messages

698,445 Google searches

168 million+ emails sent

1,820TB of data created

217 new mobile web users

Intermediation

5 Vs of Big Data

- Raw Data: Volume
- Change over time: Velocity
- Data types: Variety
- Data Quality: Veracity
- Information for Decision Making: Value

Evolution of statistics and statistics

Peter Luhn

1997: "Machine Learning"

World Economic Forum 2011

2007:"The Fourth Paradigm"

NIVERSITÉ

BLANC

Epistemiology: Types of scientific inferences

<u>Deduction</u>	Induction	Abduction
All the beans in this bag are white These beans are from this bag	These beans are white These beans are from this bag	(Peirce: "Hypothesis") All the beans in this bag are white These beans are white
These beans are white	All the beans in this bag are white	These beans are from this bag
Properties:	Properties:	Properties:
 The truth of the premises warrants the truth of the conclusion. 	- The truth of the premises does not warrant the truth of the conclusion.	- The truth of the premises does not warrant the truth of the conclusion.
 In a deductively valid argument, it is impossible that the premises are true and the conclusion is 	 It is possible that the premises are true and the conclusion is false. 	 It is possible that the premises are true and the conclusion is false.
false.	- Not necessary	- Not necessary
- Logically necessary	- Ampliative ("synthetic")	- Ampliative ("synthetic")
- Not ampliative ("analytic")		 The conclusion "explains" the premises.
- A deductive argument is called		- Inference to the best explanation

UNIVERSITÉ

NT BLANC

be true.

"sound" if its premises happen to

Inference to the best explanation -

Epistemic steps

- Going from reality to model and back
 - how to not get lost in translation
- Building model steps
 - Separate what is « in » from « out »
 - Characterize the model
 - Simplify it
 - Evaluate the model
- Issues
 - Separation
 - Characterization
 - Evaluation
- We will focus on evaluation

As proposed by the project sponsor.

As specified in the project request.

As designed by the senior architect.

As produced by the engineers.

As installed at the user's site.

What the customer really wanted.

On statistics

- There is no probability but probabilistic models
 - Some are useful, all are wrong
- At beginning there is a concrete reality
 - We move from the reality to realm of the probability model
 - We need to come back to the real world
 - Statistics defines how to not get lost in translation
- No free lunch !
 - How much I pay ?

Should we regress when we do regression ?

- Rationality conflict
 - Is the model logic dominating or the problem (real world) logic
- Oprah event
 - Discrimination
- Evaluation and transparency
 - COMPAS case

ISTIC

The quest for objectivity

- Finally what is objective in statistics
 - What can be transferred back in real world from the realm of statistics
- We throw a coin
 - The probability of the model is p \checkmark
 - Pr{10 th throw is face}=p X
- We construct a model of customer interests
 - The average interests of customers is characterized by given parameters \checkmark
 - The customer X will be interested to Y with probability pX
- Application of statistical models to individuals can only be subjective

Correlations vs. Causation example

- "The correlation between workers' education levels and wages is strongly positive"
- Does this mean education "causes" higher wages?
 We don't know for sure !
- Recall: Correlation tells us two variables are related BUT does not tell us why

Correlation vs. Causation

- Possibility 1
 - Education improves skills and skilled workers get better paying jobs
 - Education causes wages to Λ
- Possibility 2
 - Individuals are born with social environment A which is relevant for success in education and on the job
 - Social environment (NOT education) causes wages to Λ

Consider the following research undertaken by the University of Texas Health Science Center at San Antonio appearing to show a link between consumption of *diet* soda and weight gain.

- The study of more than 600 normal-weight people found, eight years later, that they were 65 percent more likely to be overweight if they drank one diet soda a day than if they drank none.
- And if they drank two or more diet sodas a day, they were even more likely to become overweight or obese.

- A strong relationship between two variables does not always mean that changes in one variable causes changes in the other.
- The relationship between two variables is often influenced by other variables which are lurking in the background.
- There are two relationships which can be mistaken for causation:
 - 1. Common response
 - 2. Confounding

• Common response refers to the possibility that a change in a lurking variable is causing changes in both our explanatory variable and our response variable

• **Confounding** refers to the possibility that either the change in our explanatory variable is causing changes in the response variable OR that a change in a lurking variable is causing changes in the response variable.

Example Correlation v.s. Causation

One study during the polio epidemic in the 1920s showed a strong correlation between ice cream consumption and cases of polio. As a result, the public was warned to avoid eating ice cream as it increased the risk of contracting the disease.

Thoughts?

- Again, there was a strongly confirmed correlation. However, it turned out that there was NO causation. With a properly controlled experiment, it could have been easily shown that increased ice cream consumption did NOT increase the risk of polio.
- Again, there was a lurking variable hiding in the background. It turns out that the virus that causes polio (a virus of the *picornoviridae* family for anyone who cares) thrives in warmer weather. So the lurking variable here was temperature!

Example

STUDY: In a study, babies of women who bottle feed and women who breast feed are compared, and it is found that the incidence of gastroenteritis, as recorded in medical records, is lower in the babies who are breast-fed.

RANDOM ERROR

By chance, there are more episodes of gastroenteritis in the bottle-fed group in the study sample. (When in truth breast feeding is not protective against gastroenteritis).

Or, also by chance, no difference in risk was found. (When in truth breast feeding is protective against gastroenteritis).

MISCLASSIFICATION

Lack of good information on feeding history results in some breast-feeding mothers being randomly classified as bottle-feeding, and vice-versa.

If this happens, the study *underestimates either of the two groups.*

BIAS

The medical records of bottle-fed babies *only* are *less complete* (perhaps bottle fed babies go to the doctor less) than those of breast fed babies, and thus record fewer episodes of gastro-enteritis in them only.

This is called bias because the observation itself is in error.

In this case the error was not conscious.

CONFOUNDING

The mothers of breast-fed babies are of higher social class, and the babies thus have *better hygiene, less crowding* and perhaps other factors that protect against gastroenteritis.

Less crowding and better hygiene are truly protective against gastroenteritis, but we mistakenly attribute their effects to breast feeding.

This is called confounding, because the observation is correct (breast-fed babies have less gastroenteritis), but its explanation is wrong.

Overfitting or something else?

- Cars climbing up trees (at CMU)...
- Road sides look like parallel lines.
- But, unfortunately, so do trees!
- Related "incident":
- Task: Given pictures of wooded areas,
- find pictures where tanks are hidden
- (late 1960s, DARPA challenge)

- Then, general noticed, pictures containing tanks were taken in the late afternoon. Without tanks, in the morning. ⁽³⁾
- What is the issue here? Good or bad learning? Overfitting? Or something else?
- ML is "good" but: The training data itself "flawed"!!!!

 Also, problem with "algorithmic bias." ML methods learn / reinforce "unwanted" biases eg in hiring or loan decisions. But you may not realize it!

LISTIC

WIKIPEDIA

Who is a Data Scientist?

In addition to advanced analytic skills, this individual is also proficient at integrating and preparing large, varied datasets, architecting specialized database and computing environments, and communicating results. A data scientist may or may not have specialized industry knowledge to aid in modeling business problems and with understanding and preparing data.

Creating value from data requires a range of talents: from data integration and preparation, to architecting specialized computing/database environments, to data mining and intelligent algorithms

An individual responsible for modeling complex business problems, discovering business insights and identifying opportunities through the use of statistical, algorithmic, mining and visualization techniques.

UNIVERSITÉ

BLANC

Data scientists can be invaluable in generating insights, especially from "big data;" but their unique combination of technical and business skills, together with their heightened demand, makes them difficult to find or cultivate.

LISTIC

Figure 3. Core Data Scientist Skills

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

MarketingDistillery.com is a group of practitioners in the area of e commerce marketing. Our fields of expertise includemarketing strategy and optimizations customer tracking and on-site analytics: predictive analytics and econometrics: data warehousing and thig data systems: marketing channel insights in Paid Search.SUS, Oscial, CMR and brand.

PROGRAMMING & DATABASE