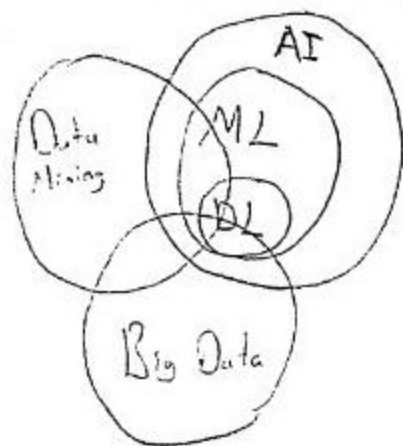


Introduction to Machine Learning - Beck's Part

Lecture (1)

1.1 The Place of ML and AI & DS

Different Parts in Data Science



The Field of Data Science



Big Data: Collecting and processing any data which is huge in volume, arrival/processing rate or inconsistent in structure.

Data Mining: Process of finding out hidden patterns in the structured data and find hidden information in the data.

Data Analytics: It is a process which is one step above data mining. Data analytics identifies the type of the analysis to be performed within which data mining techniques will be performed.

Data Analyzing: It is a more general approach of finding insights out of the raw data by forming a hypothesis and proving them using statistical tests.

Data Science: It defines the process of understanding the business problem to deliver the solution.

Machine Learning: It is a tool used in data analytics to predict/find out a hidden layer of information in data.

1.2 Data Science Steps

Step 1: Collecting Data

- ① Generally, data is not ready to consume!
- ② Data comes from many sources such as client database, web logs, etc.
- ③ 3 V Problem (mostly in Big Data)
 - Volume: Data in TB.
 - Velocity: Streaming data with high throughput.
 - Variety: Data with varying structures.

Step 2: Cleaning Data

- ① Removing discrepancy from data.
- ② Outlier Analysis, Missing Value Reduction, transforming, modelling, visualization.
- ③ Data Analysts fit here.

Step 3: Analyzing Data

- ① Create a plan to do analysis on data.
- ② May include Descriptive Analytics, Predictive Analytics, Prescriptive Analytics.
- ③ Data Mining & Data Analytics fit here.

Step 4: Sharing Insights & Business Intelligence Reports

- ① We make future predictions and validate our previously defined hypothesis.
- ② Machine Learning fit here.

Step 5: Taking Action!

③ The Terms: "Learning" and "Machine Learning"

(Machine) Learning: is programming computers to optimize a performance criterion using examples or past experience.

① Learning general models from a data of particular examples.

② Data is cheap & abundant; knowledge is expensive & scarce.

Aim: Build a model that is a good and useful approximation to the data.

Machine Learning's sub-roles

Role of Statistician

- Inference from a sample.

Role of Computer Science

- Solve the optimization problem.
- Representing and Evaluating the model for inference.

1.4 Data Types

Data: Collection of factual information based on numbers, words, observations, measurements which can be utilized for calculation, discussion and reasoning.

① Data is plural for Datum in Latin Language.

② Data is kind of postulate/axiom in math, they are accepted as bare truths, which we don't use some or other.

Categorical/Qualitative Data: Based on descriptive information (e.g., "He is a clever boy")

Binomial/Boolean: Variable data with only 2 options (e.g., good-bad, true-false)

Nominal/Unordered: Variable data which is in unordered form (e.g., red-green-blue, single-married-widow)

Ordinal: Variable data with proper order (e.g., short-medium-tall, bad-good-excellent)

Numerical/Quantitative Data: Based on numerical information (e.g., we have 4 lectures)

Discrete: Countable data (e.g., ID data, # of people)

Continuous/Masurable Data: (e.g., height, width, length)

- **Interval**: No true zero (e.g., absence of temperature, decibel of sounds)

- **Ratio**: Absolute zero (e.g., height, duration)

Special Data: Data which is not classifiable like numerical or categorical (E.g., date-time, network, sound, image)

1.5 Data Representation in Computer world

Main Data Formats: Frequently-used data formats in file-system.

CSV: A text format. Comma-separated-values, column separator can vary, e.g., ',', '!', '-' character can be used as delimiter.

ARFF: A text format. Mainly used for WEKA. Expanded form is Attribute-Relation File Format. Data + Meta-data. office formatting odr, xlsx. Not in text format, not easy to directly use.

Other Data Formats: Image, sound, video, a database table
JSON, XML

Structured vs. Unstructured Data

Structured Data: Data whose elements are addressable for effective analysis. Simplest way to manage information (e.g., a relational database (SQL))

Semi-Structured Data: Information that does not reside in a relational database but that have some organizational properties that make it easier to analyze (e.g., XML data, JSON)

Unstructured Data: A data which is not organized in a pre-defined manner or does not have a predefined data model (e.g., pdf, media logs, office format)

Meta-Data: Data about data. It's computer's! What is 101001? "101001" or (101001)₁₀ or (101001)₂

1.6 Data Table Example & Feature Terms

Obere Dataset

ID	Height (in cm)	Weight (in kg)	BFP Percentage	Age	Gender (F=1, M=0)	Parent's Obesity	isObere
1	165	97	0.26	19	0	NO-O	Yes
2	182	126	0.35	41	1	O-O	Yes
3	178	98	0.29	18	1	NO-1p	No
4	182	110	0.30	25	0	NO-NO	Yes
5	176	65	0.21	32	2	NO-O	No

Feature Engineering/Inspection

ID: Index/Primary-key/Meta-data

Height: Input Variable, Numerical & Discrete

Weight: " " & " "

BFP: Input Variable, Numerical & Continuous & Ratio

Age: Input Variable, Numerical & Discrete

Parent's Obesity: Input Variable, Categorical & Ordinal

Gender: Input Variable, Categorical & Binomial (if there is other option)

isObere: Output Variable, Categorical & Binomial (also class)

Terms

* Features = Attributes = Variables = Column

* Input Variable = Independent Variable (if variables are independent from each other, depending among input variables is discouraged)

* Instance = Observation = Row

* Output Variable = Dependent Variable = Class

Postulate/Axiom of the Data

* Postulates (like in math) are present in data science

Our Main Postulate: IV Input Variable, OV Output Variable

There is a function W for every variables, such that

$$W(IV_1, IV_2, \dots, IV_N) = OV \text{ for every instance}$$

* Our main purpose to find best W , i.e., the function gives the best output value with the given input values

Without Class (no Supervised Learning)

* Without class, we have only features and can only compare with each other!

* Generally, such cases ambiguous and needs more data.

Note That: If no repetition and no ordinality in feature, this feature is useless!

Brainstorming: Why do we use (machine) learning in order to calculate instead of finding a formula like Body-Mass Index?

- Data can change.
- Norm for class (e.g. obesity) can change.
- Data teaches us what the truth is.

1. Extra

Where to find datasets?

- Kaggle
- GitHub
- Google Datasets
- UCI Machine Learning Repository

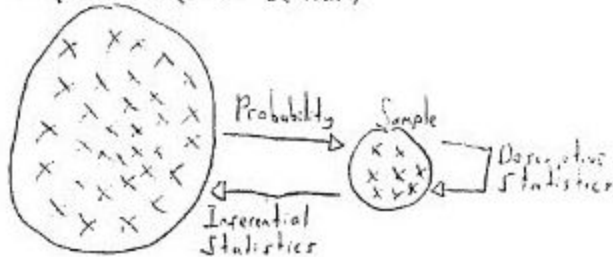
Data Science Tools

- WEKA
- Scikit Learn
- TensorFlow
- PyTorch

1.7 Descriptive Statistics

Statistic in a nutshell:

Population (o.k. Census)



Population Data: Collection of all items of interest

* Denoted by "N".

* The numbers we obtain = PARAMETERS

Subset Data: Subset of population

* Denoted by "n"

* The numbers we obtain = STATISTICS

Frequency Distribution Table: A graphical representation of variables.

Grade	Frequency (Abs)	Frequency (Relative)
Prep School	14	14/94
1	36	36/94
2	32	32/94
3	5	5/94
4	5	5/94
4+	1	1/94
Master/PhD	1	1/94

$$\Sigma = 94$$

$$\Sigma = 94/94 = 1$$

Measures of Central Tendency

* A single value that explains a set of data by identifying the central position within that set of data.

* Also called "Measures of Central Location".

* 3 Measures: Mean, mode, median

Mean: Most popular, used with both discrete & continuous data.

* Sample Mean: $\bar{x} = \frac{\Sigma x}{n}$

* Population Mean: $\mu = \frac{\Sigma x}{N}$

Median

Why we use Median: Suppose we have a dataset.

(Yearly Income of Employee at a Company)

18.5 22.8 23.7 24.6 26.8 25.2 132.5 (K\$)

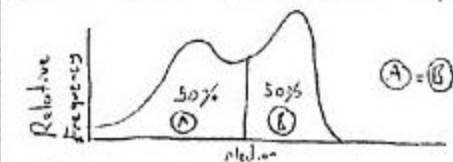
→ $\bar{x} = 47.4$, which is misleading!

→ $\bar{x} = 26.6$, better "center".

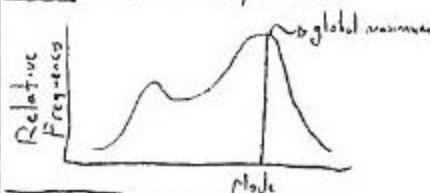
* Resilient to outliers!

* Sample Median & Population Median: Md

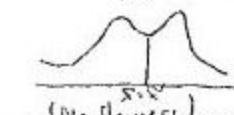
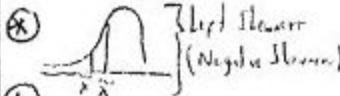
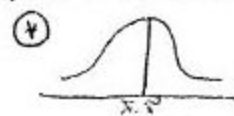
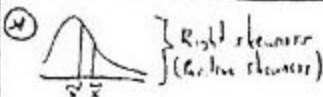
* Median divides area of Relative Frequency into 2.



Mode: Mo, most frequent item



Measures of Skewness: o.k. Measures of Asymmetry



Measures of Variability : The range, the variance, the standard deviation

Range : The simplest method

$$R = X_{\max} - X_{\min}$$

Variance :

For sample : $s^2 = \frac{\sum (x - \bar{x})^2}{n-1}$

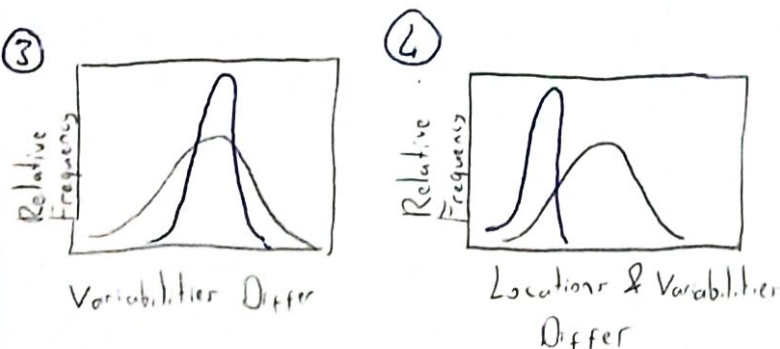
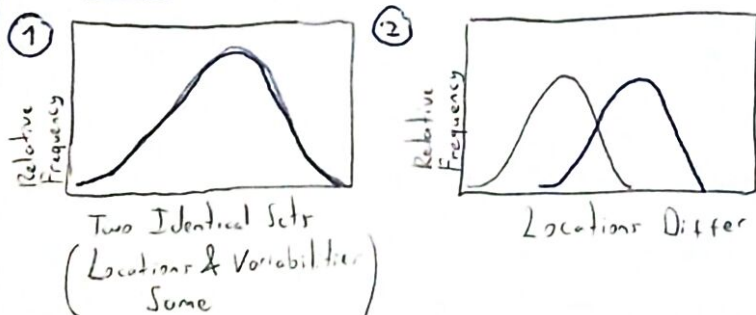
For population : $\sigma^2 = \frac{\sum (x - \mu)^2}{N}$

Standard Deviation :

For sample : $s = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$

For population : $\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$

Difference Between Two Datasets :



Covariance : A statistical measure which is defined as a systematic relationship between a pair of random variables wherein change in one variable responded by an equivalent change in another variable.

Range of Cov. $S_{xy} \sim (-\infty, +\infty)$

$Cov_{xy} = 0$: Independent variables

$Cov_{xy} > 0$: Two variables move together in same direction

$Cov_{xy} < 0$: Two variables move together in opposite direction

* Doesn't show correlation directly. For example; if one variable doubles then covariance doubles!

Sample Covariance Formula : $S_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$

Population Covariance Formula : $\sigma_{xy} = \frac{\sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)}{N}$

Correlation (Coefficient) : Same as covariance, but in different range

Range : $r_{xy} \sim [-1, 1]$

$r = 0$: Independent variables (no correlation)

$r = 1$: Perfect positive correlation

$r = -1$: Perfect negative correlation

Sample Correlation Formula : $r_{xy} = \frac{S_{xy}}{S_x S_y}$

Population Correlation Formula : $\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$

* r and $\rho \rightarrow$ greek letter "rho"

Covariance & Correlation Example :

