

NOTES

Commonly used ML algorithms

- Linear Regression → estimate real value based on continuous variables.
- Logistic Regression → estimate discrete values (0, 1) based on given set of independent variables.
- Decision Tree → classification problem.
- SVM - classification method, plot graph
- Naive Bayes → assumes presence of particular feature is unrelated to presence of any other feature.
- kNN - classification, regression
- k-means - clustering problem.
- Random forest - ensemble of decision trees.
- Dimensionality reduction.
- Gradient boosting algorithm.
 - ↳ GBM → used while dealing with plenty of data, boosting algo
 - ↳ XGBoost → high predictive power -- boosting algo
 - ↳ Light GBM → faster training speed & higher efficiency, treebased learning algo.
 - ↳ Cat Boost → open sourced ML

Content - 30 hrs

5. Python 2:55 hrs

- ✓ Introduction
- ✓ Variable & Data types
- ✓ Basic Ps Python syntax } 1:01 hrs
- ✓ Other Python operator
- ✓ Conditional statements }
- ✓ Python Functions }
- ✓ Sequences }
- ✓ Iterations }
- ✓ Advance Python tools. } 0:46 hrs

6. Advanced Statistical method in Python 4:54 hrs

- * Linear regression with Statsmodel
- * Multiple Linear regression with Statsmodel
- * Linear regression with Sklearn } 1:39 hrs
- Practical example
- * Logistic Regression
- * Cluster Analysis } 1:17 hrs
- * K-mean clustering
- * Other type of clustering

7. Mathematics 0:51 hrs.

8. Deep learning 13:06 hrs.

- 46m - Intro to Neural network
- 21m - How to build a neural netw from scratch-Numpy
- 28m - Tensor Flow Intro
- 26m - Digging deeper into NNs - Deep Neural Netw
- 20m - Overfitting
- 8m - Initialization
- 21m - Digging into Gradient Descent & learning rate selection
- 16m - Preprocessing
- 33m - Classify on the MNIST dataset

1. Introduction

- ✓ Field of Data Science - The various DS disciplines
- ✓ The field of Data Science - Connecting the DS discipline

2. Introduction 2:09 hrs.

- ✓ The field of Data Science
- ✓ Various DS disciplines
- ✓ Connecting the DS disciplines
- ✓ Benefits of each disciplines
- ✓ Popular DS techniques
- ✓ Popular DS tools
- ✓ Career in DS

3. Probability 3:37 hrs.

- ✓ Combinatorics
- ✓ Bayesian Inferences
- ✓ Distributions
- ✓ Probability in other field.

4. Statistics 3:49 hrs

- * Descriptive Statistics } 1:04 hrs
- ✓ Practical example
- * Inferential statistics fundamentals } 0:22 hrs
- * Inferential Statistics: Confidence Intervals } 0:56 hrs
- ✓ Practical Example
- * Hypothesis testing } 0:55 hrs
- ✓ Practical example

7:39 - Bushers Case Example

7:41 Conclusion

2h 8. Appendix

- ✓ 29m - Deep learning Tensorflow - Intro
- ✓ 30m - Classify an MNIST dataset
- ✓ 31m - Business Case

30m 10. Software Integration

31m 11. Case Study

- ✓ 10m - What's next

- ✓ 12m - Preprocessing the Absenteeism data
- ✓ 13m - Applying ML to create the absenteeism model
- ✓ 14m - Loading the absenteeism module
- ✓ 15m - Analyzing predicted O/P in Tableau

1.1h 12. Appendix

- 40 - Additional Python tools
- 59 - Pandas Fundamentals

13 Bonus Lecture

3. Process the Data

(Examine data at a high-level \Rightarrow Understand every column; identify errors, missing values (or corrupt records))

(Clean the data \Rightarrow Throw away, replace, or filter corrupt/error prone/malformed values)

Scripting language Python/R

Data wrangling (Cleaning) Python Pandas library
Distributed + Iterative Hadoop MapReduce
Spark.



Data Science Process

1. Frame the problem (translate ambiguous req. into a concrete, well-defined problem. Identify business priorities to strategy decisions that will influence your work.)

Domain knowledge (needs), Product Indication (metrics), Business Strategy (priorities), Teamwork (people), Colleagues (resources)

2. Identify Colleagues (collect raw data)

(Identify all available datasets, extract data into usable format)

Database management systems: MySQL, PostgreSQL, Oracle, MongoDB, Query Structured Database (SQL)

Retrieving unstructured info: (Information retrieval / text)

Distributed Storage (Hadoop HDFS, Spark, Flink)

4. Explore the data

(Play around with the data → Split, segment, plot the data in different ways)

(Identify patterns & extract features → use statistics to identify & test significant variables)

Scientific computing (Python, numpy, matplotlib, scipy, pandas) Inferential Statistics (hypothesis testing, correlation vs causation) Experimental design (A/B tests, controlled trials)

5. Perform In-Depth Analysis

(Create a Predictive model & use feature vectors from step #4)

(Evaluate & Refine model & Perhaps return to # 2, 3 or 4)

Machine Learning (Supervised / Unsupervised algo) contextual processing ML tools library (Python: scikit-learn) Advanced Math (linear algebra & multivariate calculus)

6. Communicate results

(Identify business insights → Return back to the business problem) (Visualize your findings & keep it simple (or priority driven)) (Tell a clear Actionable story & Effectively communicate to non-technical audiences)

Business Acumen (Non-technical terminology)

Data Visualization tools (Tableau, D3.js, Google visualization, matplotlib, ggplot, seaborn)

Data storytelling (presently & speakily, repeatly, visually)

Introduction

Scatter Plot

k-mean clustering

Alienated

low loyalty & low satisfaction

Fan

high loyalty & high satisfaction

Support high loyalty & low dissatisfaction

Improve shopping experience

Roamer Low loyalty & high satisfaction

Loyalty cards, stoppath

Discount, raffle

→ Cluster analysis - states the problem

Dendrogram

Heat map

Free tableau training - 50% completion message

"Hi! I completed 50% of the course! Can I have my gift, please?" → Q/A section.

Field of Data Science

Data Data team Big data team

Business Intelligence Data Science Business analytics

Data Analytics.

Statistics → Datamining → Predictive analytics
Data Science.

* Analysis

Past

Explain

How?

why?

Qualitative

Quantitative

Data

Show sale decreased

* Analytics

Future

explore potential future events.

Qualitative

Intuition + analysis

Quantitative
formulas + algorithm

8/2/22:

Exploratory Data Analysis.

- ↳ Find out the weak areas where you can work to make more profit.

- * Data science - discipline relevant on data availability while business analytics does not completely rely on data
- * Data science can be used to improve the accuracy of predictions based on data extracted from various activities typical for driving efficiency
- * Business intelligence aims to explain past events using business data
- * Machine learning
 - ↳ creating
 - ↳ implementing algo
 - ↳ let ml/c receive data to use data to → make prediction
 - ↳ analyze patterns
 - ↳ give recommendations
- * Artificial Intelligence.
Simulate human knowledge & decision making with computer
- * Symbolic Reasoning → AI
 - ↳ based on high level human-readable representations of problem & logic

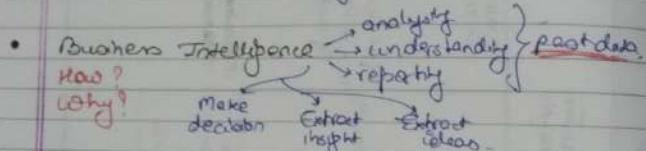
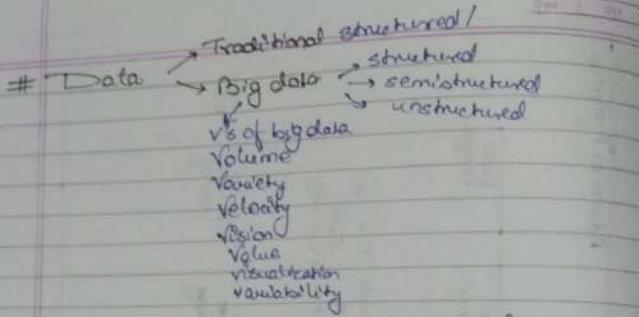
Past

Present

Future.

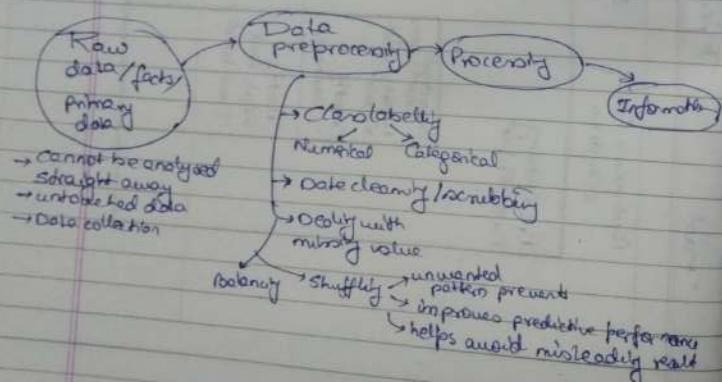
Advanced Analytics

Business Case Studies	Quantitative Analytics	Business Analytics	Data Science	Data Analytics
Preliminary Data Report			Sale Forecasts	Optimization of Drilling Operation
<p>Business Intelligence Reporting with Visuals</p> <p>Creating Dashboard</p> <p>Creating Real-time Dashboard</p>		<p>Client Retention</p> <p>Fraud Prevention</p>	Digital Signal Processing	Machine Learning
				AI



- Traditional methods
 derived from statistics
 Regression Cluster Factor analysis

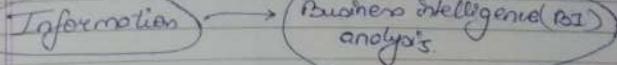
- Data → Raw facts
 → Processed data
 → Information



Big Data → financial trading data record

Data mining

Data masking
confidentiality preserving data mining.



Data skills + business knowledge intuition

explain past performance
 what when how which
 many regla

Quantification - the process of representing observations as numbers

Observations

Quantifications

Measures

Measure → related to a simple descriptive statistics of past performance

Metric = Measure + Business meaning
comparison.

Key performance (KPIs) = Metric + Business objectives

within timeframes generated only from user who have clicked on a link - ad campaign

Metric KPI

traffic of a page
weighted by a type of user

traffic generated only from users who have clicked on a link provided

5/2/22

Exploratory Data analysis.

↳ understand various aspects of data.

Objectives.

↳ identify faulty part.

↳ understand relationship b/w the variables.

Steps

Understand
data

Clean
Data.

Analyse
relationship b/w
variables.

import seaborn ← visualization.

* Understand data

data.head(), data.tail() (range)

data.shape, data.describe()
(no. of rows, column)

count mean std min. max.

data.columns,

data.nunique() → unique value.

data['__'].nunique()

* Cleaning.

→ Check null data.isnull().sum()

→ Delete redundant data std = data.drop(['year', 'parent1'], axis=1)
variable

→ check outlier - data point that differ significantly from others.

finding std [(std['cgpa'] > 8.80) | (std['cgpa'] < 5.11)]

trimming new std = std [(std['cgpa'] < 8.80) &
(std['cgpa'] > 5.11)]

* Relationship analysis

• Correlation matrix

correlation = student, corr=True

• use heatmap - visualize volume

sns. heatmap (correlation, xticklabels= correlation, columns=y ticklabels= correlation, columns)

annot=True

• pairplot - plot pairwise relationship per dataset

sns. pairplot (student)

• Scatterplot

sns. relplot (x='math score', y='readiness score', hue='gender', data=student)

• Histogram - bar plot

represent data in form of groups

sns. distplot (std ('writing score'), bins=5)

no. of bins

→ sns. catplot (x='writing score', kind='box', data=student)

! matplot lib - interactive with backend & frontend

↳ plotly cannot below o/p cell of prev. plot.

input go

provides funcⁿ for creating to remove a directory, folder's content, change directory, current directory.

16/2/22

Linear regression $y = b_0 + b_1 x_1 + \dots + b_n x_n$

→ side quantity causal relationship among variables

Logistic regression → decision making



Cluster analysis



$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_n x_n$$

$x \rightarrow$ explanatory variable

→ regressor

→ independent variable

→ predictor variable

* Factor analysis

↳ grouping explanatory variable together

* Time series → plotting values against time

Machine learning

improve computer
computational models

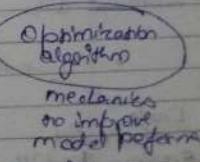
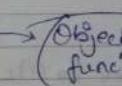
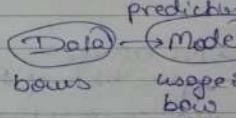
create an algo

↓
find model

fit data

accurate
prediction

eg



Training model → no rules
→ a final goal measure inaccuracy
improve

→ Types of machine learning (ML)

Supervised learning
• training algo.
labelled data
know target prior
sum
NN DL Bayesian
RW

Unsupervised learning
unlabelled data
large data
K-mean → DL

Reinforcement learning
+ reward system

Hadoop → a software - collection of programs.

→ Roles

• Data → DB Administrator.

Handles this control of data

mainly works with transactional data

Data Architect

Data engineer

processes the obtained data

so that it is ready for analysis

Design the way
data will be retrieved,
processed & consumed

• Business intelligence (BI)

→ BI developer

performs analyses
specifically designed
for company

BI Analyst

→ BI consultant

perform analyses
& reporting of past
historical data

external BI

analyst

• Data Science & ML

→ ML engineer

applies state of the
art computational
methods.

Data Analyst

prepares more
advanced types
of analyses

Data Scientist

employs traditional
statistical methods
or unconventional AI
techniques for
making prediction.

Data

Data Science

When it is applied	Traditional	Big	Business Intelligence			Traditional Methods	Machine Learning
	At the beginning of your analysis.			Past	Now	Future	
Why you need it	data driven decisions require well organised & relevant raw data stored in digital format.		use data to create reports to dashboards to gain business insights.	Predictive Analytics access potential future scenario by using advanced statistical method		utilise artificial intelligence to predict behaviour in unprecedented ways	
What techniques involved	DATA COLLECTION Preprocessing <ul style="list-style-type: none"> class labeling (categorical vs numerical) data cleaning dealing with missing values Case specific eg - Balancing & Shuffling datasets ER diagram RDBMS	DATA COLLECTION Preprocessing <ul style="list-style-type: none"> class labeling (no, text, digital images, digital audio data) data cleaning dealing with missing values Case specific Tent data maturity confidentiality - Preserving data mining techniques	ANALYSE THE DATA Extract info, and present it in the form of: <ul style="list-style-type: none"> Metrics KPIs Reports Dashboard 	Regression logistic regression Clustering Factor Analysis Time Series		Supervised learning <ul style="list-style-type: none"> SVMs NNs Deep learning Random forest Bayesian networks Unsupervised learning <ul style="list-style-type: none"> K-means Deep learning Reinforcement learning <ul style="list-style-type: none"> Similar to supervised learning, but instead of minimising losses, one maximises reward 	
Where it is used	Basic customer data Historical stock price data	Social media Financial trading data	Price optimization Inventory management	User experience (UX) Sales Forecasting Time series		Fraud Detection Client retention	
How using what tools	Programming lang. 	Programming lang. 	Prog. lang. 	Prog. lang. 		Prog. lang. 	
Software	 	 	 	 		 	
Who	Data Architect Data Engineer Database Administrator	Big Data architect Big Data Engineer	BI Analyst BI Consultant BI Developer	Data Scientist Data Analyst		Data Scientist Machine learning engineer.	
Big Data - not just volume it defines Variety, Variability, Velocity, Veracity & others	Qualitative analysis SWOT Strength Weakness Opportunity Threats not used		Excel, SPSS Stata - useful		OL is still a debate on why algo used to outperform all conventional		

Probability

→ the likelihood of an event occurring

higher prob. value
0.8, 0.9, 1

→ higher likelihood.

likely
unlikely
relatively more likely

$$P(A) = \frac{\text{Preferred favourable}}{\text{All Sample space}}$$

Independent event $P(A \text{ and } B) = P(A)P(B)$

Expected value → used to make predictions about future based on past data
aug. outcomes we expect if we run an exp.
many times → good approximation
experimental probabilities → easy to compute

$E(A) \rightarrow$ outcome we expect to occur when we run an exp. $E(A) = P(A) \times n$

$$E(x) = P(A) \cdot A + P(B) \cdot B + P(C) \cdot C$$

Probability frequency distributions → collection of prob.
eg sum of 2 dice $\begin{array}{|c|c|c|c|c|c|} \hline & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 1 & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 2 & 2 & 4 & 6 & 8 & 10 & 12 \\ \hline 3 & 3 & 6 & 9 & 12 & 15 & 18 \\ \hline 4 & 4 & 8 & 12 & 16 & 20 & 24 \\ \hline 5 & 5 & 10 & 15 & 20 & 25 & 30 \\ \hline 6 & 6 & 12 & 18 & 24 & 30 & 36 \\ \hline \end{array}$ for each possible outcomes

$P > 1 \leftarrow$ double counting
 $P < 1 \leftarrow$ not counting fully

$$P(A) + P(B) + P(C) = 1$$

$$A' = B + C$$

$$P(A') = 1 - P(A)$$

Combinatorics

- ↳ combination of objects from a diff specific finite sets
- ↳ Repetition
- ↳ Order
- ↳ other

- Permutation → no. of diff. possible ways we can arrange elements
ord & completely up to es

$$P_n = n!$$

→ v.e. no. don't have factorial

$$(n+k)! = n! \times (n+1) \times (n+2) \times \dots \times (n+k)$$

$$(n-k)! = \frac{n!}{(n-k+1)(n-k+2) \times \dots \times 1}$$

$$\frac{n!}{k!} = (k+1)(k+2) \dots n$$

Variation - permutation with repetition

first pick $\bar{V}^n = n^p$ p-many element
then arrange some elements
ABC  $\bar{V} = 3^2 = 9$

Variations w/o repetition

when arranging p elements out of total of n.

$$\bar{V} = \frac{n!}{p!}$$

Combination

↳ don't take into account - double counting
↳ set.

↳ all different permutations of a single combination are different variations

$$C_p^n = \frac{n!}{(n-p)!p!} = \frac{V_p^n}{p!}$$

Permutation? No. of fav. Outcomes

Variations

Combinations

Picking more elements leads to having fewer combinations

$${}^6C_4 = 15 \quad {}^6C_5 = 6$$

↳ Symmetry of combination

$$p > \frac{n}{2} > n-p$$

* Solving combinations with separate sample spaces
eg ${}^3C_1 \times {}^2C_1 \times {}^5C_1$

↳ Determine the appropriate amount of time
it would take for such a task to be completed

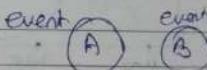
Set

$x \in A$
element set

$A \ni x$
set element.

$x \in \text{part of } A$

$A \subset \text{contain } x$



$\forall i:$
↑ Such that
for all

Mutually exclusive - non overlapping

$$A \cup B = A + B$$

* Independent event.

If $P(A) = P(A|B)$ then the two are independent and
if not, A and B are dependent.
eg getting head after another head.

Conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{if } P(B) > 0$$

$$P(A) = P(A|B)$$

↑↑
independent if any two events are independent
 $P(A \cap B) = P(A) \times P(B)$

$P(A|c) \neq P(A) \leftarrow A \& C \text{ are also dependent}$

If $P(B) = 0 \rightarrow$ event B would never occur
 $A|B \rightarrow$ not h.t.p.eable.

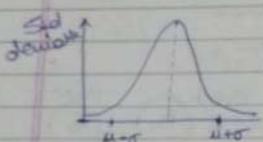
(law of total probability)

$$P(A) = P(A|B_1) * P(B_1) + P(A|B_2) * P(B_2) + \dots$$

Bayes Rule

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Possible values a variable can take & how freq. they occur



More congested the middle of distribution, the more data falls within that interval

Probability distributions

Finite no. of outcomes \rightarrow Discrete \rightarrow True
Infinite many outcomes \rightarrow Continuous \rightarrow False

Variable type characteristics

Discrete distributions \rightarrow finite distinct outcome

• Bernoulli Distribution $[P^2 = p(1-p)]$

events with only two possible outcomes $1 \text{ trial} \rightarrow P$

• Binomial Distribution

Carry out similar exp. several times in a row.

* 2 outcomes per situation

* many iterations

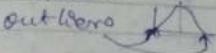
• Poisson Distribution

Test out how unusual an event freq. is for a given interval

$$P(Y \leq y) = P[Y < (y+1)]$$

Continuous distribution.

• Normal distribution



Student's T distribution

\hookrightarrow accommodates extreme values significantly better

c) Any extreme value represents a much bigger part of population



• Chi squared

\hookrightarrow not often mirror real life events

\hookrightarrow used in hypothesis testing

\hookrightarrow Asymmetric

\hookrightarrow non negative values

• Exponential distribution

\hookrightarrow events that are rapidly changing early on

• Logistic distribution

\hookrightarrow useful in forecast analysis

\hookrightarrow useful for determining a cut off point for a successful outcome

Uniform distributions

\hookrightarrow all outcomes have equal probability i.e. each outcome is equally likely

Drawback - No predictive power

* Both mean & variance are uninterpretable

Binomial Distribution

1 T/F question \rightarrow Bernoulli event
Test of 10 T/F ques \rightarrow Binomial distn

$X \sim B(10, 0.6)$ no. of times we expect
 $B(n, p)$ to get specific outcome

Expected values
 $E(X) = p_1 x_1 + p_2 x_2 + \dots + p_m x_m$
 $P(X) = {}^n C_r p^r (1-p)^{n-r}$

$$\text{Var } \sigma^2 = E(y^2) - E(y)^2 = n p(1-p)$$

Poisson Distribution, $y \sim P_0(\lambda)$
 \hookrightarrow freq. with which an event occur
 \hookrightarrow no. of occurrence.

$$P(y) = \frac{\lambda^y e^{-\lambda}}{y!} \quad e \rightarrow \text{euler's no. / Napier's const}$$

$$E(y) = \lambda \quad \sigma^2 = \lambda$$

Continuous distribution

\hookrightarrow Infinitely many consecutive outcomes
 \hookrightarrow Sample space is infinite

$f(y) > 0$
 \hookrightarrow Associated prob. of every possible value y
 \hookrightarrow Likelihood of each individual one would be extremely small equal to 0.

Prob. for any individual value

$$P(y) = \frac{1}{\infty}$$

\hookrightarrow Prob. of every individual value y is 0.

Integral check \rightarrow wolframalpha.com

Cumulative distribution Function (CDF)

$$F(-\infty) = 0 \quad F(\infty) = 1$$

Prob. of intervals

$$P(b < y < a) =$$



$$\int_a^b p(y) dy$$

PDF $\xrightarrow{\text{Interpret}}$ CDF $\int_{-\infty}^y p(y) dy = F(y)$

Derivative
PDF $\xleftarrow{\text{Derivative}}$ CDF

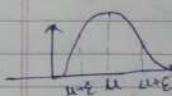
$$P(y) = F(y) \frac{d}{dy}$$

$$E(y) = \int_{-\infty}^{\infty} y p(y) dy$$

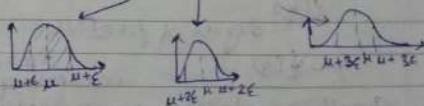
$$\text{Var}(y) = E(y^2) - E(y)^2$$

Normal Distribution \rightarrow Bell shaped curve

\hookrightarrow prob. distribution that is symmetric about the mean

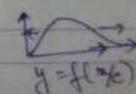
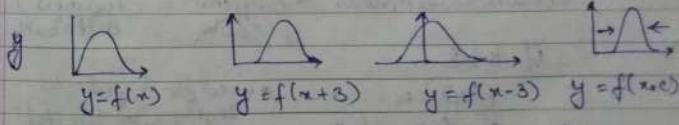


68, 95, 99.7 law:



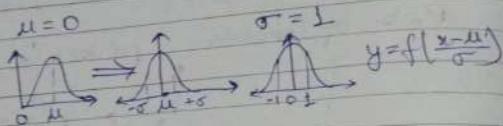
Standard normal distribution

Transformation \rightarrow Way in which we can alter every element of a distribution to get a new distribution



Standardize

$$E(x) = 0 \quad \text{Var}(x) = 1$$



$$z \sim N(0, 1) \quad z = \frac{y - \mu}{\sigma}$$

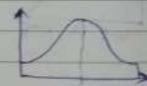
no. of standard deviations away from the mean

Drawbacks

- useful when we have normal distribution
- requires lot of data
- Risk of outliers drastically affecting our analysis

Student's T distributions

$$t(k) \quad \text{degree of freedom}$$



$y \sim t(k)$ To accommodate the occurrence of values far away from the mean
Small sample size approximation of a normal distribution

Certain characteristics + Sufficient data = Normal distribution

Certain characteristics + Sufficient data = Student's T distribution

If $k > 2$

$$E(y) = \mu$$

$$\text{Var}(y) = \sigma^2 k$$

$$k-2$$

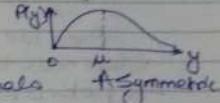
* used for hypothesis testing with limited data

* CDF table (T-table)

Chi squared distribution

$$X^2(k) \quad Y \sim X^2(3)$$

degree of freedom



→ Statistical analyses

↳ Hypothesis testing

↳ Computing confidence intervals

↳ Table of known values: N_3, T

$$E(x) = k$$

$$\text{Var}(x) = 2k$$

Exponential Distribution

$$\text{Exp}(\lambda) \quad x \sim \text{Exp}(\lambda)$$

$\lambda \rightarrow$ Rate parameter.

$$[E(y) = \lambda]$$

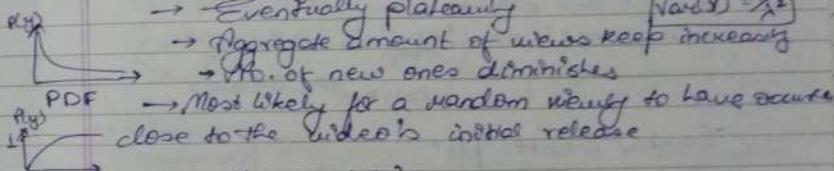
$$[\text{Var}(y) = \lambda^2]$$

→ Start off high

→ Initially ↓

→ Eventually plateau

→ Aggregate amount of views keep increasing
→ No. of new ones diminishes



$$Y \sim \text{Exp}(\lambda) \quad x \sim N(\mu, \sigma^2)$$

$x = \ln(y)$ One of the most common transformations

Logistic Distribution

$$\text{Logistic}(L, \mu, s)$$

location Scale parameter

$$y \sim \text{Logistic}(6, 3)$$

Application - forecasting victory or defeat

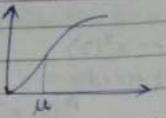


1 mps

mean mps

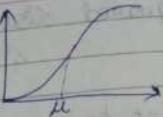
Scale parameter

2 mps



Once we reach values near the mean, the prob. drastically goes up.

$$E(Y) = \mu$$



The later the graph starts to pick up, but the quicker it reaches values close to 1.

$$\text{Var}(Y) = \frac{\sigma^2 n^2}{3}$$

Excel

Histogram Tower/bar of histogram \rightarrow bin
 1. Select column
 2. Insert \rightarrow Statistical chart \rightarrow Histogram
 Bin width \rightarrow X-axis

Charts/Scatter plot

Select \rightarrow Insert \rightarrow chart \rightarrow scatter plot
 column

expanding Trendline \leftarrow +

Sum \rightarrow new column
 $= B2 + \dots \rightarrow = B2 + C1$

Head Project manager responsibilities
 1. Development of the game

- a) Balanced
 - b) Realistic
2. Supervise the release
- a) Social media marketing
 - b) Competitor analysis

3. Customer analysis.

Probability in Finance

\hookrightarrow Option pricing

Option - An agreement b/w two parties for the price of a stock or item at a future point \rightarrow in time

\hookrightarrow Allows one side to decide

\hookrightarrow How much we are willing to pay to receive that part? (Highest premium)

Expected payoff

$E(P) < 0$ Disadvantageous (Avoid buying option)

$E(P) = 0$ Fair Deal

You expect to make as much as you paid

$E(P) > 0$ Favourable deal

Go through the deal

Pricing on option $\rightarrow E(P) < 0 \rightarrow$ unfavourable deal
 $\rightarrow E(P) > 0 \rightarrow$ fair deal

Probability with statistics

Statistics

based on sample of data

Uncertainty

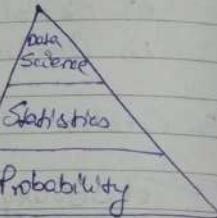
Prediction intervals

analyses numeric & categorical data

Mean, variance & expected value

Characteristics

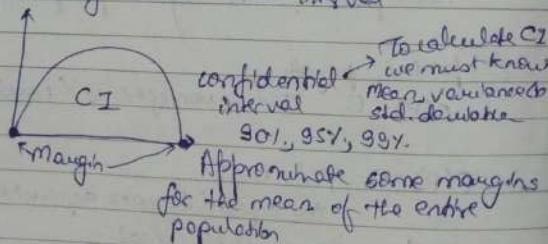
based on population



Experimental data

↳ Many useful concepts based on probability theory.

↳ Expresses the likelihood of the population mean being within that interval



Stats used for hypothesis testing

Hypothesis → An idea that can be tested.

3 crucial requirements

↳ Mean, Variance & type of distribution
↳ validate similar stats to a specific degree of certainty

Distribution → Shape
Characteristics } Type
B, S²

- * Any distribution predicts a value for all points within our datasets.
- * Distribution anticipates the actual data points

Mathematical modeling (Supervised ML)

→ extension of stats that data scientists deal with

Probability in Data Science

Data Analysis → Analyse past data

→ find insight

→ Make reasonable predictions

Mathematical

modeling → We run artificial simulations to see how well our predictions match up to various possible future outcomes.

Monte Carlo simulation

→ we generate artificial data to test the predictive power of our mathematical models

↳ Not completely random data
follows restrictions

→ we can break restrictions

→ Many models look well when we train them

→ Perform poorly with new data

Expected values

ML → An extremely fast forced trial & error process.
The more predictors it makes, the more precise they become.

Forecasting

- ML & DL have very high predictive power → not 100% certain
- Even data scientists assign probability to their predictions

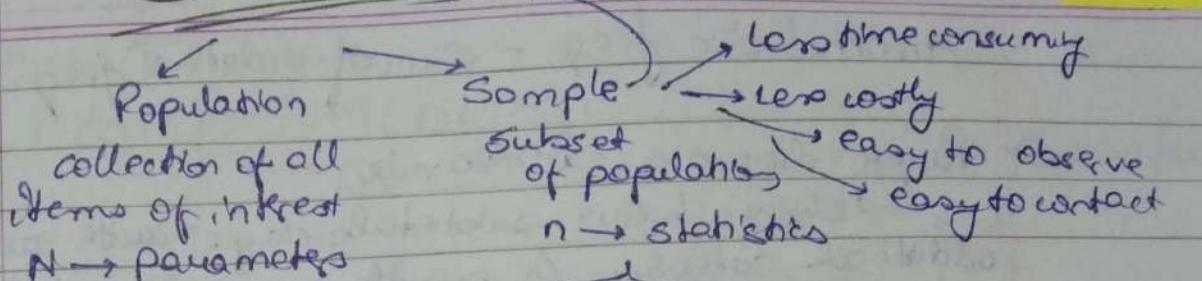
Data Science

An expansion of probability, statistics & programming that implements computer technology to solve more advanced questions

Statistics

Page No.

Date: / / 201



Descriptive Statistics

chosen strictly by chance

Representative
accurately reflect members of entire populations

Types of data

Categorical or group

- car brand
- yes/no

Numerical

Discrete

• finite

eg grades
no. of objects

continuous

• infinite

weight
height

area, distance, time

Measurement levels

Qualitative

Nominal categories

Ordinal
order

Quantitative

Interval
does not have zero

Ratio

has a true zero

Categorical data

Frequency distribution tables

Bar charts

Pie charts → Market shares

Pareto diagram

Frequency distribution table

Relative frequency - $\frac{1}{\text{total freq.}}$ for each category

Pareto diagram

→ special type of bar chart, category shown in descending order of freq.

Page No. _____ Date _____ / _____ / _____

cumulative freq - sum of relative freq.

* Pareto diagram 80/20 rule
↳ shows how subtotals change with each additional category & provide us with a better understanding of our data.

Excel (Take a tour)

lookups & functions

* Add (i) $=\text{SUM}(\text{D4:D7})$

(ii) "Alt = " then enter to select all values from

(iii) $=\text{SUMIF}(\text{D11:D15}, >50)$ enter
→ add only if number > 50

w) $=\text{SUM}(\text{D48, G48:G51, 100})$

v) $=\text{Sumif}(\text{D73: D77}, >50)$

Excel → Data → from other sources

Home → auto sum

* Fill → select cursor (drag down / up / left / right)
Same to copy same item in rows

* Split

Ctrl E → Flash fill.

Home → F4 → Flash fill.

* Based on delimiter

Step Select cells → Data → Text to col

Next ← Tab / Comma ← Next ← Delimiter

* Split a columns with formulas

$=\text{LEFT}(\text{C5}, \text{FIND}(":", \text{C5}) - 1)$

* Transpose Rows ↔ Columns

Ctrl C → Ctrl + Alt + V

Ctrl C → select cell → Home → Paste special
transpose

* With formula

$=\text{Transpose}(\text{C33:H54})$

* Sort & filter with ease Right click
Select column schema Home > Sort & Filter → sort A to Z
→ largest to smallest
→ custom filter

* Filter → select all → Home → Sort & filter
OK ← Select filter ← column ← filter
clear filter → number filter

* Table

Select All → Insert → Table
value

Total rows in tables

Select any cell → Table design → Total row
or Ctrl + Shift + T

* Drop downs

Select cell → ~~ctrl~~ Data → Data Validation
OK ← List ←

Select cell → Ctrl O → quickly make a chart

Bias - Variance Tradeoff

Bias → error b/w our model prediction & ground truth
 $\text{bias} = E[f'(x)] - f(x)$
 tells us the capacity of underlying model to predict the values.

Variance → tells how much "the func" can adjust for the change in the data set.

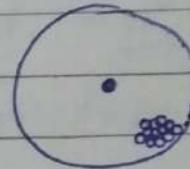
$$\text{Variance} = E[(f'(x) - E[f'(x)])^2]$$

Low Bias



Low Variance

High Bias



over simplified model
Under fitting
represents both test & train data.

High Variance



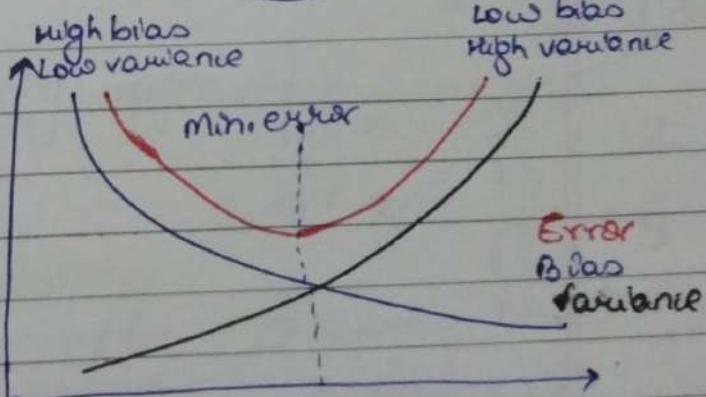
Overly complex model

↪ overfitting

↪ lower on train data

↪ higher on test

↪ starts modelling the noise in the input



$$\text{Relative freq} = \frac{\text{Freq}}{\text{Total freq}}$$

Page No. _____
Date. / / 2013

15

15

Histogram - unequal interval

Outliers → go against the logic of whole dataset

Variables → Categorical → cross-table
 Numerical → Scatterplot

Catter plot → somewhat vertical → two variables not related

Mean mode median.

Mean - average

$$\text{Median } \left(\frac{n+1}{2} \right)$$

mode - frequently appearing

500 7
 515 10
 515 130
 515 130
 515 130

10 12

$\sum f_i x_i$ = mode
 53
 54
 55
 62

5

Measure of central tendency

Measures of asymmetry - skewness

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3$$

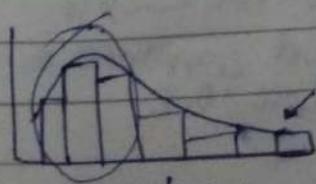
$$\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}^3$$

Skewness

→ checks whether

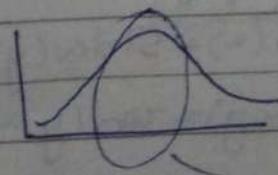
The data is conc. or one side

mean > median



mean = median = mode

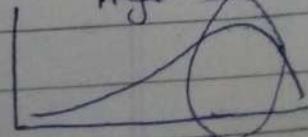
zero skew



mean < median

negative skew

where the data is situated



Variability

Sample data $\sum_{i=1}^n x_i^2$
 Data → Sample data
 Population data $n = \text{size of sample}$

Variance → dispersion around mean

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

$$\sum_{i=1}^n n_i$$

$N = \text{size of pop}$

$$\begin{aligned} &= \text{Var. S} \\ &= \text{Var. P} \end{aligned}$$

Standard Var.
 deviation → variability of single data

Coeff. of var. / relative std deviation → compare 2 or more
 Standard deviation
 mean

Covariance - How two var are correlated
 → indicates rel. of 2 vars whenever one var changes

$$\text{Sxy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

n-1 ← sample

> 0 two var move together

< 0 " " move opposite

= 0 " " are independent

Correlation - describe degree to which two var. move in coordination with one another

$$\frac{\text{cov}(x, y)}{\text{std}(x) \cdot \text{std}(y)}$$

Disregard correlations below 0.2

$$\text{correl}(x, y) = \text{correl}(y, x)$$

Titanic

What sorts of people were more likely to survive?
 Trainees → Testees
 48 rows entries

Solution → colouring

$$\begin{matrix} \checkmark & \checkmark \\ \text{Percentage} & \text{Survived} \end{matrix}$$

$\frac{2+2}{2} = 4$

Causality - direction of causal relationship

Correlation does not imply causation

Age - continuous variable

Level of measurement → Nominal

→ Interval

Data → Categorical

→ Numerical

Inferential Statistics Fundamentals

→ Probability theory

→ Distribution theory

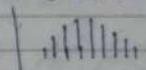
Distribution = providing prob. of occurrence of diff. possible outcomes in an exp.

→ How values are distributed in relation to each other

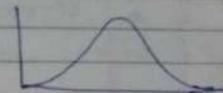
Uniform distn



Binomial distribution



Normal distribution



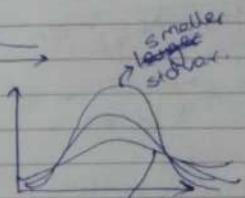
Student's t-distr



smaller
variance

smaller
std. dev.

larger std. var. (σ)



smaller
mean (μ)

larger
mean (μ)

Gaussian distribution
or Bell curve.
used to represent
a noisy data

$N \sim (\mu, \sigma^2)$

mean
var.

$$Z = \frac{x - \mu}{\sigma} \quad N \sim (0, 1)$$

origin

Standardization

- compare differently normally distributed datasets
- detect normality
- detect outliers
- create confidence intervals
- test hypotheses
- perform regression analysis

Types of Data

Categorical

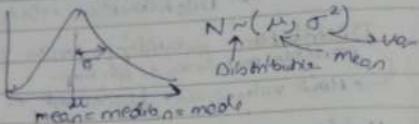
Discrete

Numerical

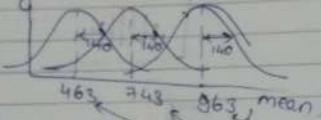
Continuous

Normal Distribution

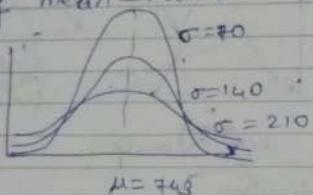
Gaussian distribution - bell curve.



- * Consistently std. deviation \rightarrow constant



Consistently mean \rightarrow want



Standardization



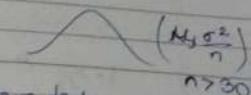
$$\sim N(\mu, \sigma^2) \rightarrow \sim N(0, 1)$$

$$Z = \frac{x - \mu}{\sigma}$$

Central Limit theorem
Original distribution



Sampling distribution



- * CLT allows to
 - perform tests
 - solve problems
 - using Normal distribution
 - even when the population is not normally distributed

Standard errors - The std. deviation of distribution formed by the sample means

$$= \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$$

→ shows variability of sample means
, it shows how well you approximated the true mean

Std error ↑ Sample size ↓

Estimates → points
are like judges

Unbiased estimator

expected value = population parameter

smallest
var

confidence level $1 - \alpha$

[Point estimate - Reliability \Rightarrow Std error \Rightarrow Point estimate + rel. factor \times std. error]

$$[\bar{x} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$$

Confidence interval

$$90\% \rightarrow z = 1.07 \text{ or } 0.10$$

$$\frac{\alpha}{2} = 0.05$$

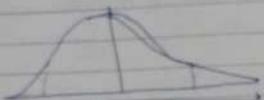
$$1 - \frac{\alpha}{2}$$

$$0.95$$

Vt stat

When $1-\alpha$ is lower CI is smaller
" 1- α is higher CI is larger

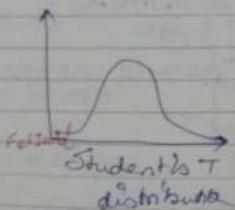
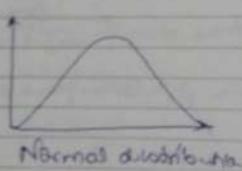
95% → accepted norm



100% width
as it comprises a wider range

Student's T distributions

↳ inference through small sample
↳ unknown population variance



$$t_{n-1, \alpha/2} = \bar{x} - \mu$$

$$s/\sqrt{n}$$

Degrees of freedom = $n-1$
 $n \rightarrow \text{sample size}$

For large enough samples → Student's T distribution converges with normal dist.

$$\bar{x} \pm t_{n-1, \alpha/2} \frac{s}{\sqrt{n}}$$

Population var. unknown

Margin of error ME

Known

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}}$$

Unknown

$$\bar{x} \pm t_{n-1, \alpha/2} \frac{s}{\sqrt{n}}$$

confidence interval : $\bar{x} \pm ME$

biggest ME

smallest ME
narrower CI

↑ statistics ↑ ME

↑ std dev ↑ ME

↑ Sample size ↑ ME ↓

Samples

$$CI = (\bar{x} - \bar{y}) \pm z_{\alpha/2} \sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

Independent

• Before & after situation

• cause & effect

Pooled variance formula.

$$s_p^2 = \frac{(n_x - 1)s_x^2 + (n_y - 1)s_y^2}{n_x + n_y - 2}$$

$$CI = (\bar{x} - \bar{y}) \pm t_{n_x + n_y - 2, \alpha/2} \sqrt{\frac{s_p^2}{n_x} + \frac{s_p^2}{n_y}}$$

no. of variable

$$v = \left(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y} \right)^2$$

$$\left(\frac{s_x^2}{n_x} \right)^2 / (n_x - 1) +$$

$$\left(\frac{s_y^2}{n_y} \right)^2 / (n_y - 1)$$

Hypothesis Testing

→ Ideas that can be tested

- Steps in data driven decision making
- Formulate a hypothesis
 - Find the right test
 - Execute the test
 - Make a decision

→ ^{hypothesis} we try to reject

Null hypothesis H_0 - one to be tested

Alternative hypothesis - other than null hypothesis

Significance level

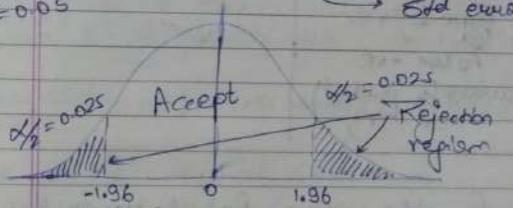
α - prob. of rejecting the null hypothesis if it is true $0.01, 0.05, 0.1$

$$Z\text{-test}$$

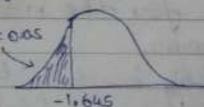
$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

Sample mean
hypothesized mean
 σ/\sqrt{n} → Std error

$\alpha = 0.05$



One-sided test $\alpha = 0.05$



Types of error.

Type I error:

- ↳ Reject a true null hypothesis
- ↳ False positive

Type II error

- ↳ accept a false null hypothesis
- ↳ False -ve

Rejecting a false null hypothesis $P = 1 - \beta$.

e.g. H_0 - She doesn't like you.

The truth

	H_0 : True She doesn't like you	H_0 : False She likes you
Accept (Do nothing)	✓	Type 2 error (False -ve) missed your chance
Reject (Indict her)	Type 1 error False +ve (incorrectly indicted her)	✓

Test for the mean, Population variance is known.

Testing is done by standardizing the variable at hand & compare it to thez

$$Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

$$Z \sim N(\bar{x} - \mu_0, 1)$$

$$\sim N(0, 1)$$

P-value universal concept works with every distribution smallest level of significance at which we can still reject the null hypothesis, given the observed sample statistic

Rule You should reject the null hypothesis if $p\text{-value} < \alpha$

Test at 90% : $0.0001 < 0.1$

Test at 95% : $0.0001 < 0.05$

Test at 99% : $0.0001 < 0.01$

} Reject null hypothesis

Rule You should reject the null hypothesis if $p\text{-value} < \alpha$

Test at 90% : } Reject.

Test at 95% : }

Test at 99% : } Cannot reject

How to find p-value manually p-value online calculator, One sided p-value $\{ 1 - \text{no. from the table} \} \Rightarrow 1 - 0.983 = 0.017$

$$\text{Two-sided p-value } (1 - \text{no. from the table}) \times 2 \Rightarrow (1 - 0.983) \times 2 = 0.034.$$

Where & How p-values used.

→ Most statistical software calculates p-values for each test.

→ Researcher decides significance post-hoc

→ p-values are usually found with 3 digits after the dot $x.xxx$

→ The closer to 0.000 the p-value, the better significance result

* p-value lower than level of sign - you reject null hypothesis

Test for the mean, population variance unknown

Like confidence intervals with var., unknown & a small sample, the correct stats to use is the t-statistics

$$T = \frac{\bar{x} - \mu}{S/\sqrt{n}}$$

Decision rule

Accept if: The absolute value of the T-score < critical value
Reject if: " " " " " T-score > critical value

Accept if - p-value > α
Reject if p-value < α

Test for the mean, dependent samples.

$$H_0: \mu_B - \mu_A \geq 0$$

Do.

$$H_0: D_0 \geq 0 \quad T = \frac{\bar{D} - D_0}{\text{st. error}} \xrightarrow{\text{meandif}}$$

$$H_1: D_0 < 0 \quad \text{st. error} \xrightarrow{\text{or pooled std.}}$$

Test for the mean, independent sample known variance

Z-test {Big sample
known variance
can also be used with
unknown var.

$$Z = \frac{\bar{x} - \mu_0}{\sqrt{\frac{s_e^2}{n_e} + \frac{s_m^2}{n_m}}}$$

{Small sample}
Unknown var. } $t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$

Test for the mean, independent samples unknown variances pooled in variance

$$\text{pooled var. } s_p^2 = \frac{(n_x - 1) s_x^2 + (n_y - 1) s_y^2}{n_x + n_y - 2}$$

$$\text{pooled std dev } \sqrt{\frac{s_p^2}{n_x} + \frac{s_p^2}{n_y}}$$

Rule of thumb: Reject the null hypothesis when the T-score is bigger than 2

Generally for z > T, a value higher than 4 is extremely significant

Steps

- Find Sample Var
- State null hypothesis
- Calculate pooled var
- Find t score
- Find p-value
- Interpret result.

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_p^2}{n_x} + \frac{s_p^2}{n_y}}}$$

$$T = (\bar{x} - \bar{y}) - (\mu_m - \mu_l)$$

$$\sqrt{\frac{s_p^2}{n_x} + \frac{s_p^2}{n_y}}$$

Once we have surpassed 50 degrees of freedom the student's T ~ Normal distribution

Point of Ad. Stats.

Overfitting & Underfitting

- * Underfitting - The model has not captured the underlying logic of the data
low train accuracy, low test accuracy
- * Overfitting
Our training has focused on particular training set
so much it has "injected the points"
high train accuracy, low test accuracy
- * Good model
captures the underlying logic of the dataset
high train accuracy
high test accuracy

1. We will use our model to make predictions based on the test data.
2. We will compare those with the actual outcome
3. Calculate the accuracy
4. Create a confusion matrix.

Misclassification = $\frac{\# \text{misclassified}}{\# \text{all elements}}$

We use map funcⁿ Yes → 1 No → 0.

Confusion matrix

```
def confusion_matrix(data, actual_values, model):
    pred_values = model.predict(data)
    bins = np.array([0, 0.5, 1])
    cm = np.histogram2d(actual_values, pred_values,
                        bins=bins)[0]
```

$$\text{accuracy} = \frac{cm[0,0] + cm[1,1]}{cm[\cdot,\cdot]}$$

Cluster Analysis

The goal of clustering is to maximize the intra-class variance.

Programming

- Problem solving → abstract thinking
- Mechanistic thinking
- Good style of coding

Programming involves

- formulating problems
- Breaking them down into meaningful steps
- Communicate those to steps to computer

Python

- technical advantage
- case sensitive
- practical applicability
- open-source - free
- General-purpose - interoperability
- high level lang. with C, Java, R

Jupyter → facilitates communication tremendously

Ctrl + Enter - code execute

Shift + Enter Execute cell

Cut → Select cell → Press X

Paste → " → Press V

Copy → " → . → Press C

Cell above → " → Press A
a cell

Cell below → " → Press B

Delete → " → Press D twice

Markdown cell documentation

→ Press M

Markdown → code cell
Press Y

Variable

$(a, b) = (1, 2)$

`type(-6)` → int

$16/3 \rightarrow 5$

floats

$\text{float}(16)/3 \rightarrow 5.33$

$16.0/3 \rightarrow 5.33$

$16 \% 3 \rightarrow 1$

$y = 12.5$

$y = 12.5$ True

$y = 12.6$ False

Comment # Comment

Line continuation This is

↓ Rahul

'Friday'[0] → a

logical operator → and
or
not

Precedence
Not > And > Or

Identity operator is, is-not

String
'George'

$y = 10$
`print(str(y) + ' Dollar')`

'I\'m fine'
↑ escape

function

def function-name (parameters):

function body

Return can be used only once in a func

Print → does not affect
calculation of o/p

Return → does not
visualize the o/p

it specifies what a
certain func is suppose
to give back

Built-in functions

`int(5.0) → 5`

`max(10, 20, 30) → 30` `min(10, 20, 30) → 10`

* `abs()` → allows you to obtain the absolute value of its argument

$z = -20$ `abs(z) → 20`

* `sum()` → sum of all elements

* `round(x, y)` → returns float of its argument (x)
rounded to a specified no. of digits (y)
after decimal.

$2**10 \rightarrow 1024$

`pow(2, 10) → 1024`

`len() ← length`

Methods

→ name, op method()

→ .append()
→ .extend()
→ .index()

• List slicing

`list_name[5:0]`

→ from to position in list
→ reverse = True
→ opposite.

BiggerList = [List_1, List_2]

• Tuple - immutable (, ,)
→ .split()

Dictionary
dict = {^{'k₁}: 'cat', ^{'k₂}: 'dog'}

eg Team = {}
Team[^{'T₁}] = 'Ran'
Team[^{'T₂}] = 'Shyam'
Team.get(^{T₂}) → Shyam.

Iteration.
for n in even:

even = [0, 2, 4, 6, 8, 10, 20]

Range (start, stop, step)
↑
(start + 1)

stop > start > stop

list(range(10)) → [1, 2, 3, 4, ..., 10]

Advance Statistical methods

* OOPs

Object = Data + Manipulation operations
Instance → float, integer, string, lists

Class defines the rules for creating an object
↓
Object → attributes / properties
→ methods

Function
can have many parameters
exists on its own
function()

Method
the object is one of its parameters
belongs to a certain class
object.method()

Package - directory of a collection of modules

from math import *

from math import sqrt

For more info type → help(math)

Regression analysis cause → effect ↗
Factor analysis

Linear regression with stats model
linear approximation of a causal relationship b/w two or more variable.

$$y = \beta_0 + \beta_1 x_1 + \epsilon \rightarrow \text{error}$$

↓ ↓ ↓
independent variable
Dependent variable quantity & effect
constant value irrespective of n

$$\hat{y} = b_0 + b_1 x_1$$

estimated/predicted value

Correlation
Relationship

Movement together
 $P(x, y) = P(y, x)$

Single point

Regression
one variable affects the other
cause & effect
one way

* Regression line \rightarrow best fitting line through the data points

* Package \rightarrow import numpy - work with multidimensional array
import pandas - allows us to organize data in tables form & attach descriptive label to results
import scipy - scientific calculator - matlib, M_L
import statsmodels.api - very good summaries
import matplotlib - visualization of numpy comp.
import seaborn - based on matplotlib & very attractive
* import sklearn - statistical graphics & ML libraries. \rightarrow better visualization

data = pd.read_csv('filename.csv')
pandas

data.describe - features of data

plt.scatter(x, y)

plt.xlabel('SAT', fontsize=20)

plt.ylabel('GPA', fontsize=20)

plt.show.

OLS - ordinary least squares regression

fit() - specific estimation to obtain the fit of the model

Regression {
x = sm.add_constant(x)
result = sm.OLS(y, x).fit()
result.summary()

OLS Regression \rightarrow Model summary
 \rightarrow coeff. table \rightarrow values of b_0, b_1 ,
 \rightarrow some additional test.

Decomposition of variability
ANOVA framework.

Sum of square
of total SS T
 $\sum (y_i - \bar{y})^2$ $\frac{SS_T}{T}$

dispersion of the observed variable around mean
- measure total variability of dataset

Sum of square
regressor SSR
 $\sum (\hat{y}_i - \bar{y})^2$ $\frac{SS_R}{T}$

- measures the exp. variability by our line
 \rightarrow measure how well our line fits the data

Sum of square error
 $\sum e_i^2$ $\frac{SSE}{n}$
RSS

measures the unexplained variability by the regression

$SS_T = SSR + SSE$
 $\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum e_i^2$
Total Variability Explained Var.

Ordinary least squares → min SSE
 Lower error \Rightarrow better explanatory power
 depict the closest to all points

$$S(b) = \sum (y_i - x_i^T b)^2 = (y - Xb)^T (y - Xb)$$

↪ OLS estimator of β for a simple linear reg.

Other method

↪ Generalized least sq.

↪ Max. likelihood estim.

↪ Bayesian regression

↪ kernel regression

↪ Gaussian process regression

R squared
 measure goodness of fit

$$R^2 = \frac{SSR}{SST}$$

your reg.
 explains none
 of the variability

your ref.
 explains the entire
 variability

It measures how much of the total variability is explained by our model

Steps linear regression

→ Import relevant libraries
 Load data

Create regressor

↪ Declare independent & dependent variables

Multiple linear regression

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon$$

↪ inferred value of intercept

Adjusted R² available

↪ penalizes excessive use of variables

↪ Basis for company mode

↪ same dependent var

↪ same dataset

↪ usually smaller than R²

$x = \text{data}[[\cdot, \cdot \text{SAT}], [\cdot, \cdot \text{year}]]$

when a new parameter is added if it decrease R-squared
 but ↓ adjusted R-squared → variable can be omitted
 since it holds no predictive power

F-statistics
 testing the overall significance of the model

H₀: $\beta_1 = \beta_2 = \beta_3 = \dots = \beta_k = 0$

H₁: at least one $\beta_i \neq 0$

if all betas are 0, then none of Xs matter & our model has no merit

* Lower the f-statistics, the closer to a non significant model

↪ larger f value \rightarrow significant, small p-value

OLS assumptions

- (i) Linearity $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon$
- (ii) No endogeneity $\text{E}\epsilon = 0 : \forall i \in E$
- (iii) Normality & homoscedasticity $\epsilon \sim N(0, \sigma^2)$
- (iv) No autocorrelation $\text{E}\epsilon_i \epsilon_j = 0 : \forall i, j$
- (v) No multicollinearity $\text{P}(\text{corr} \neq 1) \text{ if } i \neq j ; i \neq k$

• Linearity $y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \epsilon$

 If straight line
linear regression is suitable.

 Linear reg. not appropriate.
Fix:
- run new linear reg.
- exp. transformation
- log transformation

• No endogeneity

$$\text{E}\epsilon_i = 0 : \forall i \in E$$

Error is correlated with independent variable

Omitted variable bias → easiest way to detect is Durbin-Wu-Hausman test

↳ you forgot to include a relevant variable

↳ always diff.

↳ always sneaky

↳ only exp. to advanced knowledge can help

• Normality & homoscedasticity

$$\epsilon \sim N(0, \sigma^2)$$

Normality t-test & F-test work because we have zero mean

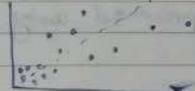
assumed normality

Central limit theorem

$$\epsilon \sim N(0, \sigma^2)$$

$$\text{Homoscedasticity } \sigma_{\epsilon_1}^2 = \sigma_{\epsilon_2}^2 = \dots = \sigma_{\epsilon_n}^2 = \sigma^2$$

↳ to have equal variance



Prevention

↳ look for omitted variable bias

↳ look for outliers

↳ transform - log transformation

↳ take log of the variable

↳ perform regression

Semi-log model $\hat{y} = b_0 + b_1 (\log x_1)$

or $\log \hat{y} = b_0 + b_1 x_1$

Log-log model

$$\log \hat{y} = b_0 + b_1 (\log x_1)$$

elasticity

• No autocorrelation - no serial correlation

$$\text{E}\epsilon_i \epsilon_j = 0 : \forall i, j$$

auto-correlation - not observed in cross-sectional data.

Day of the week effect → high returns on Friday, low returns on Monday

Detection

* If there are no patterns to be seen → no autocorrelation

* Durbin-Watson falls b/w 0 & 4

2 → no autocorrelation

<1 & >3 → cause an alarm.

No remedy → when in the presence of autocorrelation avoid the linear regression model.

Alternatives

↳ Auto-regressive model

↳ moving avg. model

↳ Auto-reg. moving avg. model

↳ Auto-reg. integrated moving avg. model

- No multicollinearity - easy to spot & easy to fix
 $r_{x,y} = 1 : x \perp\!\!\!\perp y$

$$\text{eg } a = 2 + sb \quad b = \frac{0.2}{s}$$

$r_{ab} = 1 \rightarrow \text{perfect multicollinearity}$
 If a can be represented using b , there's no point in using both.

Fixes

- Drop one of the two variable
- Transform them into one
- Keep them both

Prevention

Find the correlation b/w each two pairs of independent variables

$$r_{x,y} \text{ for } x \perp\!\!\!\perp y$$

- Dummy variable

Imbalance of categories with no.

$$\text{SAT} \rightarrow \text{GPA}$$

Attendance

dummy var.

$$\text{data}['Attendance'] = \text{data}['Attendance'].map(\{'Yes': 1, 'No': 0\})$$

- Making prediction - linear rep.

Add new data

```

new_data = pd.DataFrame({'const': 1, 'SAT': [1200, 1320],
                         'Attendance': [0, 1]})

new_data = new_data[['const', 'SAT', 'Attendance']]
    
```

Clearing index { new_data.rename(index={0: 'Bob', 1: 'Alice'}) }

predict_predictions = results.predict(new_data)

Predictions

Linear regression with sklearn

Numpy, scipy, matplotlib → fast & efficient
 prefer working with array

pandas data frame → statsmodel
 great for learning

Numpy array → scikit learn
 for preferred

Advantages → Incredibly documented

Variety → Representa

Classification
 Clustering
 Support vector m/c
 Dimensionality red.

Import from sklearn.linear_model import LinearRegression

GPA ← SAT
 dependent var. independent var.
 O/P target I/P feature

Regresson itself reg = LinearRegression()
 reg.fit(x, y) reg.coef_, reg.intercept_

Residual matrix x_matrix = x.values.reshape(-1, 1)
 matrix in 2D

Calculate reg.score(x_matrix, y)
 the R-squared

L_1 norm

(least absolute deviations) \rightarrow least squares
or least absolute error (LAE):

$$S = \sum_{i=1}^n |y_i - f(x_i)|$$

L_2 norm

(least squares)

$$S = \sum_{i=1}^n (y_i - f(x_i))^2$$

find { neg. intercept
intercept }

find the coefficient { neg. coef - }

making predictors { neg. predict() }

Required slightly larger than adj. R^2
 $R^2 > R_{adj}^2 \rightarrow$ not penalized a lot
for the inclusion of 2 independent variables.

\Rightarrow Multiple Linear regression - Sklearn

$$R_{adj}^2 = 1 - \frac{(1 - R^2) * (n - 1)}{n - p - 1}$$

$n \rightarrow$ no. of observations

$p \rightarrow$ no. of predictors

def adj_r2(x, y):

$r^2 = \text{reg. score}(x, y)$

$n = x.shape[0]$

$p = x.shape[1]$

$$\text{adjusted_r2} = 1 - (1 - r^2) * (n - 1) / (n - p - 1)$$

return adjusted_r2

\Rightarrow Feature selection

Simplifies models, improves speed $\&$ prevents a series of unwanted issues arising from having too many features

If a variable has a p-value > 0.05 we can disregard it.

feature_selection, f_regression:

→ creates simple linear regression of each feature

& the dependent variable

+ Note that for a simple linear regression the p-value of F-stat = the p-value of the only independent variable

Feature selection:

from sklearn.feature_selection import f_regression

f_regression(x, y)

array(_____, ____), array(_____, ____)

P-values = f_regression(x, y)[1]

P-values

P-values = round(3)

+ Create summary table

reg_summary = pd.DataFrame(data=x.columns)

columns = ['Features']

reg_summary['Coefficients'] = reg.coef_

reg_summary['P-values'] = reg.pvalues_.round(3)

reg_summary

• Feature scaling / Standardization

transforming data into a standard scale.

$$\text{Standardized Variable} = \frac{x - \bar{x}}{\sigma}$$

10.00 → Intuit
Weights - ML word for coeff.
→ Bigger the weight → Bigger the impact

Standardization:
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler
scaler.fit(x)

x_scaled = scaler.transform(x)
x_scaled

Create summary table

reg_summary = pd.DataFrame([{'f1': 100}, {'SAT': 100},
'Rand 123'}, columns=['Features'])
reg_summary['Weights'] = reg.intercept_ + reg.coef_[0]
reg.coef_[1] * x

Lec 2213

making
pred
with
stand.
wrt.

new_data = pd.DataFrame(data=[[1700, 2], [1800, 1]]),
columns=['SAT', 'Rand 123'])

new_data

new_data_scaled = scaler.transform(new_data)
new_data_scaled

reg.predict(new_data_scaled)

Overfitting & Underfitting

Our training has
focused on the
particular training
set so much that
it has missed a
point

high train accuracy

The model has not captured
the underlying logic of the data

Low accuracy

Low test accuracy

→ We can create the regression on the train data
but Test on the test data

Train - Test Split

import numpy as np
from sklearn.model_selection import train_test_split

generally { a = np.arange(1, 10)
data } → np.arange([start], stop, [step])
returns evenly spaced values
within a given interval.
b = np.arange(50, 60)

split to do { train_test_split(a) } → splits array or matrix
data into train and test
randomly into subsets
a_train, a_test = train_test_split(a, test_size=0.2,
shuffle=False, random_state=42)

explore { a_train_slope, a_test_slope.
its result } { a_train
a_test }

⇒ Linear Regression { Clean data
Relax assumption
log transformation
Create a model
Create dummies }

Data preprocessing

`data.describe(include='all')` categorically defines data.

dropping columns {
`data1 = data.drop(['Model'])`, axis=1}
 → remove less important variable

give no. of {
`data1.isnull().sum()`
 null values
 in variable}

- Rule of thumb:- if you are removing less than 10% of the observations, you are free to just remove all that have missing value

remove the rows that contain
 Null value {
`data_no_mv = data.dropna(axis=0)`

Exploring - the prob. distn. func (PDF)

for optimal results give world b/w
 looking for normal distribution

$$q = \text{data_no_mv}['Price'] \cdot \text{quantile}(0.99)$$

* One way to deal with outliers seemingly is to remove top 1% of observations

,
`data1 = data_no_mv[data_no_mv['Price'] < q]`
`data1.describe(include='all')`

recently to hide after cleaning {
`data_cleaned = data1.reset_index(drop=True)`

* log-transformer especially useful when facing exponential relationship

np.logins → returns the natural log of a new array of no.

$$\text{log_price} = \text{np.log}(\text{data_cleaned}['Price'])$$

$$\text{data_cleaned}[\text{log_price}] = \text{log_price}$$

Multicollinearity

→ one of the best ways to check for multicollinearity is through VIF (variance inflation factor)

→ `data_cleaned.columns.values`

→ from statsmodels.stats.outliers_influence import
 variance_inflation_factor

variables = `data_cleaned[['Mileage', 'Year', 'EngineV']]`

vif = `pd.DataFrame()`

vif['VIF'] = `[variance_inflation_factor(variables, i) for i in range(variables.shape[1])]`

vif['features'] = `variables.columns`.

vif = 1 = no multicollinearity

1 < vif < 5 = perfectly okay

10 < VIF = unacceptable

drop features with no multicollinearity

`data_no_mc = data_cleaned.drop(['Year'], axis=1)`

Create dummy variables.

`pd.get_dummies(df, drop_first=True)` - spots all categorical variables & creates dummies automatically.

If we have N categories for a feature we have to create N-1 dummies

Steps

- * Import the relevant libraries
- * Loading the raw data
- * Preprocessing
 - Exploring the descriptive statistics of variables
 - Determining the variables of interest
 - Dealing with missing values
 - Exploring the pdfs
 - Dealing with outliers
- * Checking the OLS assumptions
 - Relaxing the assumptions
 - Multicollinearity
- * Create dummy variables
 - Rearrange a bit
- * Linear regression model
 - Declare the i/p & the targets
 - Scale the data
 - Train Test split
 - Create the regressor
 - Finding weights & bias

* Scaling has no effect on the predictive power of dummies, once scaled, though, they lose all their dummy meaning

Train test split

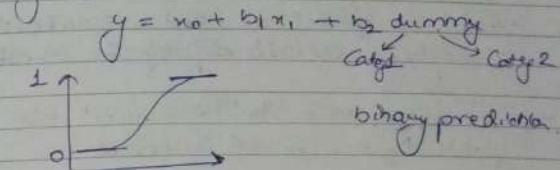
{ from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test =
train_test_split (inputs_scaled, targets, test_size=0.2,
random_state=365)

* Residual = Diff b/w the targets & the predicted values of errors

Logistic Regression

categorical

non-linear data



Assumption

- No endogeneity
- Normality & Homoscedasticity
- No multicollinearity
- No autocorrelation

→ Predict prob. of an event occurring.

$$\text{logistic model } p(x) = \frac{e^{(B_0 + B_1 x_1 + \dots + B_k x_k)}}{1 + e^{(B_0 + B_1 x_1 + \dots + B_k x_k)}}$$

referred.

$$\text{odds } \rightarrow \left(\frac{p(x)}{1-p(x)} \right) = e^{(B_0 + B_1 x_1 + \dots + B_k x_k)}$$

$$\text{logit model } \log \frac{p(n)}{1-p(n)} = B_0 + B_1 x_1 + \dots + B_k x_k$$

Regression

$$n = \text{sm.add_constant}(n)$$

$$\text{reg_log} = \text{sm.Logit}(y, n)$$

$$\text{result_log} = \text{reg_log.fit() } \rightarrow \text{result.logisfit}$$

* Current func. value gives value of objective func
at 10th iteration

* Maximum Likelihood estimation (MLE)
 likelihood function - estimates how likely it is that the model at hand describes the real underlying relationship of the variables.

→ The bigger the likelihood func., the higher the prob. of our model is correct

→ MLE tries to maximize the likelihood func.

* Log-Likelihood almost but not always -ve.
 bigger log-likelihood - higher the probability

* LL Null - (log likelihood - null)
 The log-likelihood of a model which has no independent variables.

* LLR (log likelihood ratio)
 measures if our model is statistically different from LL null aka a useless model

* Pseudo R-squared McFadden's R-squared.
 A good pseudo R-squared is somewhere b/w 0.2 to 0.4.

→ This measure is mostly useful for continuous variables of the same model.

→ Different models will have completely diff. & incomparable Pseudo-R-squares.

Logit model

SAT - Admitance regression

$$\log\left(\frac{x}{1-x}\right) = -69.91 + 0.042x \text{ SAT}$$

$x \rightarrow$ prob. of an event occurring.

$1-x \rightarrow$ prob. of an event not occurring

$$\log\left(\frac{\text{odd}_2}{\text{odd}_1}\right) = 0.042 (\text{SAT}_2 - \text{SAT}_1) \rightarrow \text{1 unit}$$

$$\frac{\text{odd}_2}{\text{odd}_1} = 6.667 \rightarrow 10 \text{ units}$$

Difference of 1 unit of SAT.

$$\log\left(\frac{\text{odd}_2}{\text{odd}_1}\right) = 0.042 \cdot \frac{\text{odd}_2}{\text{odd}_1} = 0.042 \cdot 10 \rightarrow \text{odds}_2 = 10.24$$

$$\frac{\text{odd}_2}{\text{odd}_1} = 1.042$$

When the SAT score increases by 1, the odds of admittance ↑ by 4.2%.

$$\Delta \text{ odds} = e^b x$$

Dummies → linear regression

Binary predictors → logistic regression

→ Accuracy of the model

sm. Logit results.predict() → returns the values predicted by our model

np.set_printoptions(formatter={float: lambda x: f'{(0.0000000000000002*x).format(n)}})

results_log.predict()

sm. Logit Results.pred_table() → returns a table which compares predicted to actual values.

confusion matrix

$\text{cm_df} = \text{pd.DataFrame}(\text{results}, \text{log} + \text{pred}[[\text{y}]]$
 $\text{cm_df.columns} = [\text{'Predicted 0'}, \text{'Predicted 1'}]$
 $\text{cm_df} = \text{cm_df.rename}(\text{index}=\{\text{0: 'Actual 0'},$
 $\text{i: 'Actual 1'}\})$
 cm_df

 $\text{cm} = \text{np.array}(\text{cm_df})$

$$\text{accuracy_train} = (\text{cm}[0,0] + \text{cm}[1,1]) / \text{cm.sum()}$$

 accuracy_train

Underfitting & Overfitting → After inferential
statistics
fundamentals

Cluster Analysis

unsupervised learning.
The goal of clustering is to maximize the
similarity of observations within a cluster &
maximize the dissimilarity b/w clusters.

Market segmentation
Image segmentation

Explore the data → identify patterns

② unsupervised learning → linear & logistic regression
→ Classification probability on a/p category
→ labelled data

Maths prerequisites
→ clusters b/w 2 data points
→ centroids
→ euclidean distance
 $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$
mean position of a group of points

k-mean clustering

↳ no. of clusters

- 1) Choose the no. of clusters
- 2) Specify the cluster needs
- 3) Assign each point to a centroid
- 4) Adjust the centroids.

import {KMeans} from sklearn.cluster import KMeans
libaries

Steps

Import Libraries

Load the data

No. of clusters Plot the data

dataframe.loc (row indices, col indices) etc
Select the features from the data frame, given below to be kept

kmeans = KMeans(n_clusters=k)
kmeans.fit()

Clustering Standardize

Clustering results.

WCSS method.

Sklearn.cluster.KMeans().fit_predict(x)
returns the cluster predictions in an array

How to choose no. of clusters

Elbow method.

within-cluster sum of squares
or WCSS

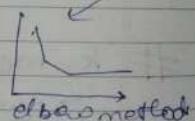
observation = N

clusters = N

WCSS = 0 (min)

N N
1 Small
max low

WCSS → kmeans.inertia_



WCSS = []

for i in range(1, 7):

kmeans = KMeans(i)

kmeans.fit(x)

wcss_iter = kmeans.inertia_

wcss.append(wcss_iter)

wcss

Elbow method

number_clusters = range(1, 7)

plt.plot (number_clusters, wcss)

plt.title ('The elbow method')

plt.xlabel ('No. of clusters')

plt.ylabel ('Within Cluster Sum of Squares')

Pros of k-mean Simple to understand

Fast to cluster

widely available

Easy to implement

Always yield a result (also consider to decide)

Cons

We need to pick k

Sensitive to initialization

Sensitive to outliers

Produce optimal sol.

Standardization

Remedy

The elbow method

K-means++

Remove outliers

Standardization

The ultimate aim of standardization is to reduce the set of higher no.

Standardization is trying to put all variables on equal footing

Market Segmentation

Brand Loyalty → Churn rate

→ Retention rate

Customer lifetime value (CLV)

Standardise the variable

from sklearn import preprocessing
n.scaled = preprocessing.scale(n)

* sklearn.preprocessing.scale(n) -

Standardize with mean 0 & std dev of 1
each variable (column) separately

More Types of Analysis

clustering mostly

Exploratory

- Get acquainted with the data
- Search for patterns
- Plan clustering

Data visualizable

Confirmatory Exploratory

- Explain a phenomenon
- Confirm a hypothesis
- Validate prev. research

Descriptive

→ Plot → k-means

Types of clustering

→ Agglomerative (bottom-up)

Hierarchical

→ Divisive (top-down)

Dendrogram

- The closest b/w the links

Shows similarity b/w observations

Pros

- Shows all possible linkage b/w clusters
- understand data much better
- No need to preset the no. of clusters
- Many methods to perform hierarchical clustering

Cons

scalability

complex

Heatmap

sns.clustermap(~~dist~~)

Mathematics

Linear Algebra

Matrix, $R \times C$

Scalors have 0 dimension. $[1]$, $[2]$ $[1 \times 1]$
like point

Vector
like line. $\begin{bmatrix} 5 \\ -2 \\ 4 \end{bmatrix}$ $[m \times 1]$
 \uparrow

Scalar

$s = s$

vector

$v = np.array([s, -2, 4])$

Matrices $m = np.array([[5, 3, 6], [-3, 0, 14]])$

`m.shape` → returns dimension of matrix

Tensor Flow

→ scalar → Tensor of rank 0 → 1×1
→ vector → P " " " " $1 \rightarrow m \times 1$
→ matrix → " " " " $2 \rightarrow m \times n$

addition $m_1 + m_2$
matrices.

`v1 = np.array[1]`

Transpose $m_1.T$ $\xrightarrow{\text{matrix}}$ transpose

Dot product

$$\begin{bmatrix} 2 \\ 8 \\ -4 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ -7 \\ 3 \end{bmatrix} = [-66] \quad np.dot(m1, m2) \xrightarrow{\text{matrix 1}} \xrightarrow{\text{matrix 2}}$$

Matrix multiplication $m \times n$ with $n \times k$

$$A_{100 \times 300} \cdot B_{300 \times k} = C_{100 \times k}$$

- Application of linear algebra in DS

↳ Vectorized code

↳ Image recognition

↳ Dimensional reductions

Deep learning

Training an algo involves:-

Data → feed

Model → compose

Objective Function

Optimization algs. vary

Pages
Date
Page

1 / 100

100%

cross entropy
classification &
regression loss
Types of ML → Supervised Training model
→ unsupervised clustering
→ Reinforcement.

The linear model.

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b$$

input

w

Intercept

Linear model with multiple I/p & multiple O/p

$$y_1 = x_1 w_{11} + x_2 w_{21} + b_1$$

$$y_2 = x_1 w_{12} + x_2 w_{22} + b_2$$

loss → ↓ loss function ↑ level of
supervised accuracy

Objective funcⁿ → Reward (reinforcement)
used to evaluate how well the model is able
to match the desired correct values.

Loss

Classification

Regression

Cross Entropy

$$\text{L}^2 - \text{norm} = \sum (y_i - t_i)^2$$

OLS

lower the error the
lower the loss

$$L(y, t) = -\sum t_i \ln y_i$$

Optimization algo.

Gradient Descent → can find the minima

$f'(x) \rightarrow \nabla F(x_1, x_2, \dots, x_n)$

$x_{i+1} = x_i - \eta f'(x_i)$

learning rate.

When the min is reached

$x_{i+1} = x_i = 0$ (no longer updates)

Oscillate - repetitive variable around target value

Learning rate to be
high enough so we can reach the closest min. in a natted amount of time
low enough so we don't oscillate around the min.

n - Parameter Gradient descent

$$y_i w + b = t_i \rightarrow t_i$$

$$L(y, t) \rightarrow \text{loss}$$

$$C(y, t) \rightarrow \text{cost}$$

$$E(y, t) \rightarrow \text{err}$$

$$x_{i+1} = x_i - \eta f'(x_i)$$

$$\left\{ \begin{array}{l} w_{i+1} = w_i - \eta \nabla_w L(y, t) \\ b_{i+1} = b_i - \eta \nabla_b L(y, t) \end{array} \right.$$

$$\nabla_w L = \sum_i \nabla_w \delta_2(y - t_i)^2 \rightarrow x, w + b$$

$$= \sum_i \delta_2(y - t_i)$$

$$= \sum_i \delta_i \delta_i$$

weights

Supervised inputs
learning weights
outputs targets

Simple Linear Regression

Import the relevant libraries

Generate random i/p data to train on

Create the targets we will aim at

Plot the training data

Create weights

Create biases

Set a learning rate

from mpl_toolkits.mplot3d import Axes3D

2) Generate random i/p data to train on.

`np.random.uniform(low, high, size)` - draw random values from interval (low, high) where each no. has an equal chance to be selected

Size = $n \times k \rightarrow$ no. of variables
no. of observations

`np.column_stack` (appropriate types) - takes a seq. of 1D arrays & stack them into a 2D array.

3) Create the targets we will aim at.

noise
target = $2 * x_1 + 3 * x_2 + 5 * \text{noise}$.
eg

Train - the model

for i in range(100):
 $y_i = \mathbf{w} \cdot \mathbf{x}_i + b_i$ outputs = np.dot(inputs, weights) + biases
 $\delta_i = \text{outputs} - \text{targets}$

to norm formula

$\text{lens} = \frac{1}{2} \sum (y_i - b_i)^2$
 $\text{lens} = \text{np.sum}(\delta_i^2) / 2 / \text{observations}$
print(lens)

$\delta_i = \delta_i / \text{observations}$

$w_{i+1} = w_i - \eta \sum x_i \delta_i$ (weights = weights - learning rate * np.dot(inputs, deltas - scaled))
 $b_{i+1} = b_i - \eta \sum \delta_i$ (biases = biases - learning rate * np.sum(deltas - scaled))

- * In the memory of computer, the variables, weights, biases & outputs contain their optimized values or those from the last iteration of the loop

Data
Model

Objective func

Optimization algos

Tensorflow → CPUs, GPU, TPU
→ for neural netw

Scikit-learn better for k-mean clustering &
random forest
High-level packages → PyTorch
Tensorflow library → keras

Tensors n-dimensional arrays

np.savetxt (filename, array) since n-dimensional arrays in .npz format, using a certain keyword (label) for each array

tf.keras.Sequential() → func that specifies how the model will be laid down ('stack layers')

Linear combination + O/P = layers

tf.keras.layers.Dense(output_size) takes the I/P provided to the model to calculate the dot product of the I/P & wt & adds the bias

model.compile(optimizer, loss) → configures model for training

SGD = Stochastic Gradient Descent

- i) L2-norm lens = least sum of squared (least sum of squared error)
- ii) Scaling by # observations = avg (mean)

model.fit(inputs, targets) fit (trains) the model

Epoch = Iterations over the full dataset

Verbosity = 0 → silent → no info about training displayed
Verbosity = 1 → progress bar
Verbosity = 2 → one line per epoch

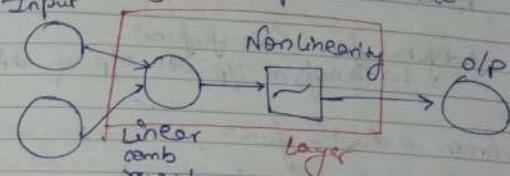
model.layers[0].get_weights() → extract weights

model.predict_on_batch(data) calculates the o/p given i/p/s

tf.keras.optimizers.SGD (learning rate)
Stochastic gradient descent optimizers, include support for learning rate, momentum, decay, etc.

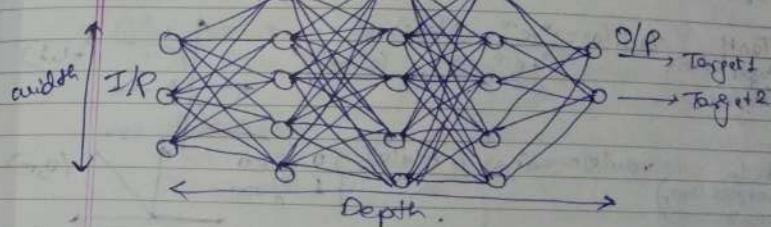
Loss func? → alternative loss function for regression is huber loss

• Layer - building block of neural net



$$\text{Non-linear sigmoid} = \sigma(w) = \frac{1}{1 + e^{-w}}$$

Deep NN



Hyperparameters
Widths Preset by us
Depth
Learning rate (η)

Parameters
wt (w) found after training
biases (b)

Non linearities
↓
Stacking layers
↓
Depth
↓
Deep learning

Batch size
Momentum coeff (α)
Decay coeff (γ)

Non-linearities also called as activation func?

Activation func? / transferfunc?

→ transform i/p into o/p of diff. kinds

Common activation func?

Name	Formula	Derivative	Graph	Range
Sigmoid (logistic) func	$\sigma(a) = \frac{1}{1+e^{-a}}$	$\frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1-\sigma(a))$		(0, 1)

Tanh (hyperbolic tangent)	$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$	$\frac{\partial \tanh(a)}{\partial a} = 4 \cdot \frac{(e^a - e^{-a})^2}{(e^a + e^{-a})^2}$		(-1, 1)
ReLU (rectified linear unit)	$\text{relu}(a) = \max(0, a)$	$\frac{\partial \text{relu}(a)}{\partial a} = \begin{cases} 0 & \text{if } a < 0 \\ 1 & \text{if } a > 0 \end{cases}$		(0, infinity)

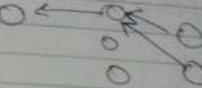
Softmax	$\sigma_i(a) = \frac{e^{a_i}}{\sum_j e^{a_j}}$	$\frac{\partial \sigma_i(a)}{\partial (a_j)} = \sigma_i(a) (\delta_{ij} - \sigma_j(a))$	diff. (a _j) everytho
---------	--	---	-------------------------------------

- All common activation func? are monotonic, continuous & differentiable

- The softmax transformation transforms a bunch of arbitrarily large or small no. into a valid prob. distribution

- Softmax is often used as activation of o/p layer in classification problems

- Back propagation



- Forward propagation

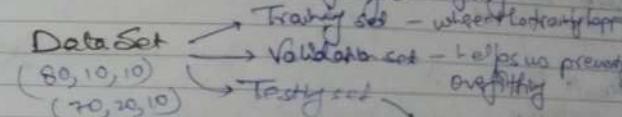
process of pushing i/p through the net

Overfitting - missed the point

Underfitting - not captured (spike)

- Bias-variance trade-off

The balance b/w underfitting & overfitting



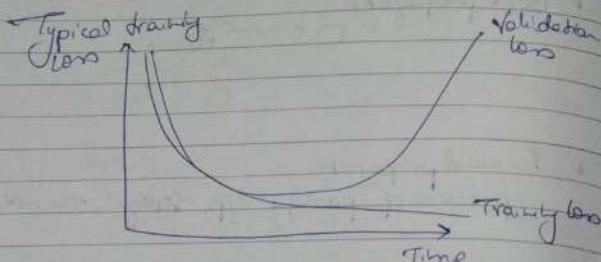
measure first prediction power of the model

N-fold cross validation → Training set + Validation set → Test set

- We train the model until the loss func? is minimized.

- early stopping - prevent overfitting

- preset no. of epochs
- Stop when updates become too small, $X_{t+1} - X_t = 0.00$
- Validation set strategy



	Preset epochs	Updates too small	Validation set
Solves the prob	✓	✓	✓
Contains features	X	✓	✓
↳ initialized			
Don't iterate endlessly	X	✓	X
↳ prevent overfitting	X	X	✓
			→ Best practice

Initialization - Set initial values of weights

Xavier - Glorot initialization

Uniform Xavier initialization - draw each wt w_i from random uniform distribution in $[-\bar{x}, \bar{x}]$

$$\text{for } x = \sqrt{\frac{6}{\sigma^2 + \sigma^2}}$$

Normal Xavier initialization - draw each wt w_i from a normal dist. with a mean of 0 & std dev

$$\sigma = \sqrt{\frac{2}{\sigma^2 + \sigma^2}}$$

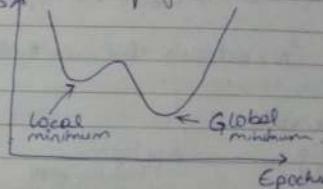
default initializer

Stochastic Gradient Descent

↳ Motivation:

We update the weights after each batch
→ we lose bias accuracy

Gradient descent pitfall



$$w \leftarrow w(t) - \eta \frac{\partial L(t)}{\partial w} - \alpha n \frac{\partial L(t-1)}{\partial w}$$

current update

update a moment ago

$$\eta = 0.9$$

$\eta \rightarrow$ learning rate → small enough - gently descend
instead of oscillate & diverge
big enough - some weight it's a rational amount of time

Learning schedules

1. We start from a high initial learning rate
2. At some points we lower the rate to avoid oscillation
3. Around the end we pick a very small rate to get a precise ans

First 5 epochs

$$\eta = 0.1$$

Next 5 epochs

$$\eta = 0.01$$

Until the end

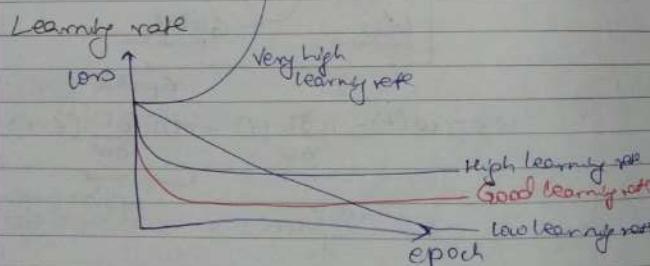
$$\eta = 0.001$$

Exponential schedules

$$\eta_0 = 0.1 \rightarrow \text{current epoch}$$

$$\eta = \eta_0 e^{-\frac{t}{T}} \quad \text{const} \rightarrow \text{hyperparameter}$$

No set rule but same order of mag.



Learning rate schedules

$$G_i(t) = G_i(t-1) + \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

→ Adaptive gradient algo.

$$w_i(t+1) = w_i(t) - \eta \frac{\partial L}{\partial w_i}(t)$$

$$\Delta w_i = -\eta \frac{\partial L}{\partial w_i}(t)$$

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)}} \frac{\partial L}{\partial w_i}(t)$$

Adaptation magic.

• ADAGRAD

→ Smart

adaptive learning rate schedule
based on the training itself
per weight

• RMS prop → root mean square propo.

$$\Delta w_i(t) = -\eta \frac{\partial L}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = \beta G_i(t-1) + (1-\beta) \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

at beginning point $G_i(0) = 0$

β - yet another parameter usually around 0.9

• ADAM (Adaptive Moment estimator)

→ learning rate + momentum

most advanced optimizers (very fast & efficient)

$$\Delta w_i(t) = -\eta \frac{M_i(t)}{\sqrt{G_i(t)} + \epsilon}$$

$$M_i(t) = \alpha M_i(t-1) + (1-\alpha) \frac{\partial L}{\partial w_i}(t)$$

$$M_i(0) = 0$$

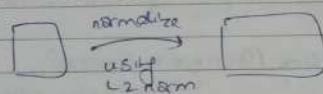
Preprocessing

- compatibility
- order of input
- Data transformation
- Generalization
- `np.savez()`
- Types
 - relative
 - logarithmic

Standardization / Feature scaling

$$\text{Standardized variable} = \frac{x - \mu}{\sigma} \rightarrow \text{mean}$$

Normalization



PCA - principal component analysis

Dimension reduction technique used to combine several variables into bigger (latent) variable.

Whitening

performed after PCA - remove unrelated correlations

Categorical Data

- One hot encoding
- Binary encoding
- Few category
- Many categories

Classifying the MNIST dataset

MNIST → very unusual problem

- Extremely common
- Easy to build up to CNN
- Very big to preprocessed

Each photo 28x28 pixel

- 0 → Black
- 255 → White

MNIST action plan

1. Prepare our data to preprocess it. Create training, validation & test datasets
2. Outline the model to choose the activation function
3. Set the appropriate advanced optimizer (the loss function)
4. Make it learn
5. Test the accuracy of the model.

Package import tensorflow-datasets as tfds

`tfds.load(name)` → loads a dataset from tensorflow datasets
`as_supervised=True` → loads the data in a 2-tuple structure [input, target]

`with_info=True` → provides a tuple containing info about version, features, # samples of the dataset

Preprocess the data

```
mnist_train, mnist_test = mnist_dataset['train'],  
mnist_dataset['test']  
num_validation_samples = 0.1 * mnist_info.splits  
[train].num_examples  
num_validation_samples = tf.cast(num_validation_samples,  
tf.int64)  
    ↳ convert a var into  
    given datatype
```

```
# def scale(image, label):  
#     image = tf.cast(image, tf.float32)  
#     image /= 255.  
#     return image, label
```

dataset.map(function) → applies a custom transform to a given dataset. It takes as i/p a func' which determine the transformation

batch size = 1 - SGD

batch size = # samples - GD

↳ batch size < # sample - mlp: batch GD

dataset.batch(batch_size) → combines consecutive elements of a dataset into batches

when batched - find avg - loss.

Outline the model

input_size = 784
output_size = 10
hidden_layer_size = 50.

- tf.keras.Sequential() → func' that is laying down the model (used to stack layers)
- tf.keras.layers.Flatten() (original shape)
 ↳ transforms (flattens) a tensor into a vector
- tf.keras.layers.Dense(output_size)
 takes i/p, provided to the model & calculates dot product of the i/p's to wts & add the bias
 This is also where we can apply an activation func'

Select loss & the optimizer

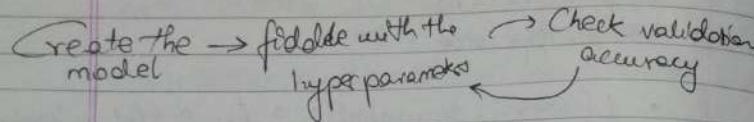
model.compile(optimizer, loss) configures the model.
Mainly

What happens inside an epoch

- At the beginning of each epoch, the training loss will be set to 0.
- The algorithm will iterate over a pre-set no. of batches, all from train data.
- The weights & biases will be updated as many times as there are batches.
- We will get a value for the loss func', indicating how the training is going.
- We will also see a training accuracy.
- At the end of the epoch, the algo will forward propagate the whole validation set.

width - hidden layer size

When we reach the max no. of epochs the training will be over



Testing the model

model.evaluate()

return the loss value & metrics value for the model in test mode

Deep Neural Net for MNIST Classification

MNIST dataset → handwritten digit recognition.

"Hello world!"

Dataset → 70000 images (28×28 pixel)
→ 1 digit per image

Goal To write an algo that detects which digit is written (0, 1, 2, ..., 9)

↳ to build a neural net with 2 hidden layers

Import relevant libraries

Load & preprocess data

→ train dataset
→ validation dataset
→ test dataset

extract train & validation data

batch train & validation data

batch test data

scale
rep
Batch, Buffer size
shuffle

↓
Model

outline the mode.

define { input
output }
Hidden layers } size

then define how model will look.

choose the optimizer & loss function

↓
Trainig

↓
Test model

* Width (hidden layer size) $50 \rightarrow 200 \rightarrow 500$
validation accuracy significantly higher
it takes the algo much longer to train

* Depth ↑ → add another hidden layer to algorithm
→ accuracy of the model does not necessarily improve
↳ Fiddling with the hyperparameters may not be enough
↳ Sometimes deeper nets need to also be wider
in order to have higher accuracy.

* Fiddle with activation functions: relu → sigmoid.
relu → cleans the noise in the data
Sigmoid → does not filter the signals as well as
relu

* Fiddle with activation func → use different act funcs
for different layers
↳ The result should not be significantly diff

* Adjust the batch size

$$100 \rightarrow 10,000$$

↳ Bigger the batch size results in slower training.

Notice that the validation accuracy starts from a low no. & with 5 epochs actually finishes at a lower no. That's because there are fewer updates in a single epoch.

* Adjust the batch size

$$100 \rightarrow 1$$

Batch size 1 → results in SGD. It takes the algorithm very little time to process a single batch (as it's one data point) but there are thousands of batches (say 10,000 to be precise).

Thus the algo is slow. Remember that it depends on the no. of cores that you train.

If you are using a CPU with 4 or 8 cores, you can only train 4 or 8 batches at once.

The middle ground (mini-batch 100 samples) is optimal.

Notice that the validation accuracy starts from a higher no. That's because there are lots of updates in a single epoch. Once the training is over, the accuracy is lower than all other batch sizes (SGD was an approximation).

* Adjust the learning rate 0.0001

↳ Since the learning rate is lower than normal, we may need to adjust the max epochs (say 50).

↳ The result is basically the same, but we reach it much slower.

↳ While Adam adapts to the problem, if the order of mag. are too different, it may not have enough time to adjust accordingly.

* Adjust learning rate 0.02

↳ While Adam adapts to the problem, if the order of mag. are too diff., it may not have time to adjust accordingly. We start overfitting before we can reach a neat solution.

Business Case → Audiobooks

Action Plan

1. Preprocess the data
 - Balance the dataset
 - Divide dataset in training, validation & test
 - Save the data in tensor friendly form.
2. Create a clear dataloader
3. Create the machine learning algo
 - same structure different model.

`np.delete(array, obj to delete, axis)`

`np.ndarray.astype()` - create a copy of the array cast to a specific type

Early stopping → median loss
tensorflow → callbacks

`tf.keras.callbacks.EarlyStopping` (difference)
decides how many consecutive to we can tolerate

Test model.

`model.evaluate()` returns the loss value & metrics values for the model in test mode

- * Test accuracy should be lower & equals to the validation accuracy

Import relevant libraries

↓
Data & preprocessing

↓
Model

outline, optimizing loss, early stopping/decay

↓
Test the model

Convolutional neural netw \rightarrow into flattening
(CNN)

kernel \rightarrow weights - 5×5

pooling 2×2 $\rightarrow 12 \times 12$

No of kernels used \rightarrow hyperparameters

\hookrightarrow mainly used in image recognitions

\hookrightarrow spatial proximities are preserved

\hookrightarrow a detail is looked for everywhere in a photo

Deep mind.

Recurrent neural netw (RNN)

\hookrightarrow Sequential data

\hookrightarrow reading music \rightarrow NLP

* The hidden layer for the second sample is a function of the new I/p & updated weights, and the hidden layer of prev. iteration.

Non neural Netw

* Discriminative models

I/p \rightarrow Model \rightarrow Prob whether o/p is correct

Random forest - not good at classification as it overfits a lot

many bad classifier \equiv Good classifier

* Generative models

I/P \rightarrow Model \rightarrow Prob whether an o/p is correct \rightarrow Target joint probability

\hookrightarrow Hidden Markov model
 \hookrightarrow Bayesian Netw

$$P(H|E) = \frac{P(E|H) * P(H)}{P(E)}$$

Appendix

Page No.
Date: 20/10/2018

- Numpy - 3rd party package used for computation
→ multidimensional array.
- Matplotlib 2D plotting package, especially designed for visualizing python's numpy computation.
- Tensorflow - Machine learning - deep learning.
Nodes → in graph represent mathematical operators
graph edges represent multidimensional data arrays (tensors) communicated b/w them

Data → Preprocess → save .npz file

np.savez(file name, arrays) → used for saving n-dimensional arrays in npz format

tf.placeholder → where we feed the data

tf.variable → preserve their values across steps while placeholders don't

- tf.matmul (A, B) = np.dot(A, B)

* Loss func : $\frac{1}{2} \cdot \text{norm}_{\text{L2}}^2(\text{observed} - \text{predicted})$

tf.losses.mean_squared_error → method equals to any L2-norm loss

Optimizable

tf.train.GradientDescentOptimizer(lr)
minimize
→ minimize op

Prepare for execution

tf.InteractiveSession() → used when we want to execute directly

Initialize variables

tf.global_variables_initializer() → method that initializes all tensor flow objects marked as variable.
run → initialize.

Running → prints

curr_vars → actually return something for optimization (variables)

→ underscore → special symbol to disregard a return value of a function (method)

Plotting the data.

Extremely reusable code of tensorflow

* The huber loss is more appropriate than the L2 norm when we have outliers, as it is less sensitive to them.

plotted

* Change no. of observations 1000 → 100,000
↳ algo will take more time to solve the problem
↳ matplot lib cannot plot the data, because too many points.

* Change learning rate → 0.0001

↳ takes a lot of time to finish working.

↳ loss is not minimized.
↳ wts & biases are far from what we want them to be.

→ More iterations are needed for this learning rate to solve the problem

→ The problem is not solved.

* Change learning rate → 0.1 or 1

↳ takes a lot of time to finish working.

↳ loss diverges to infinity.
↳ wts & biases are completely random / entirely bip.

→ More iterations would not solve the issue, algorithm is not converged

→ The problem is not solved.

* Change lossfunction L2 norm → huber loss

mean_loss = tf.losses.mean_squared_error(labels=targets, predictions=outputs) / 2

↓
mean_loss = tf.losses.huber_loss(labels=targets, predictions=outputs)

↳ Any func. that has to prop. to be least for better results (worse for worse result can be last func). This includes all

↳ Almost everything seems identical.

↳ The values of the loss are generally lower.
↳ Both huber & L2 norm works equally well.
↳ Huber loss is used when we have lots of outliers.

Classifying On the MNIST Dataset

Relevant package

- ↳ `from tensorflow.examples.tutorials.mnist import input_data` → automatically downloads the MNIST dataset to the directory of Jupyter notebook
- ↳ `mnist = input_data.read_data_sets('MNIST_data', one_hot=True)`

Outline the models

- ↳ `tf.reset_default_graph()` → clears the memory of all variables left from previous runs
(reset the computational graph.)
- ↳ `tf.get_variable("name", shape)` → func used to declare variables
- ↳ `tf.nn` → contains neural network support. Among other things, it contains the most commonly used activation functions like `tf.nn.sigmoid`,
`tf.nn.tanh`,
`tf.nn.relu`,
`tf.nn.softmax`
- ↳ `tf.nn.softmax_cross_entropy_with_logits()` → function applies a softmax activation to calculate a cross entropy loss.
Logits → unscaled probabilities.
- ↳ `tf.reduce_mean()` → finds the mean of elements of a tensor across a dimension.

Tensorflow → code → reusable.

Prediction Accuracy

In what % of the cases the o/p of algorithm matched the target

(tf.equal) checks if two values are equal

* To run early stopping will come if the validation loss starts increasing

6 mnist.train.next_batch (size of batch)
function that comes with the MNIST data provider,
which loads the batches one after the other

Batchnorm case

Extract the data from csv

Balence the dataset

Standardize the inputs

Shuffle the data

Split the dataset into train, validation, test
Save the 3 datasets in npz

Preprocess

the data

self → method → instance method.

* A ReLU activation of the first hidden layer will cause many hidden units to be zeroed

Improve model accuracy

- Improve preprocessing
- Fine-tune the model
- Play around with the activation func
- Tiddle with the batch size
- Experiment with learning rate/optim
- Check out tf.contrib
- Take your own data & use it
- kaggle

Software Integration

Date

form of symbols 1s and 0s, other digits, letters, special characters, etc.
 → it can be collected, measured, analysed, processed further

Servers

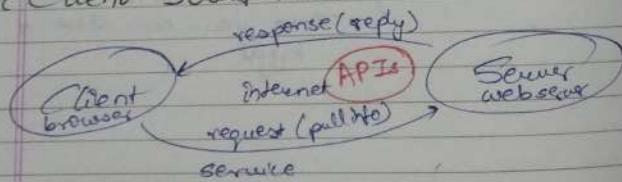
a combinations of hardware & software responsible for storing, managing & processing large amounts of data.

Web servers - web pages

Database servers - database queries

FTP servers - files

+ Client-Server model.



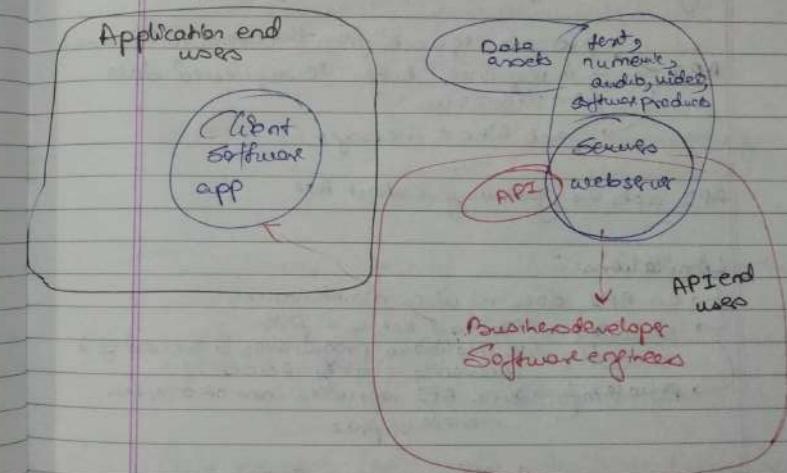
Data connectivity → APIs.

regards the ability to connect clients (browsers) securely the swift & voluminous transfer of information b/w them.

- API (Application Programming Interface)
 - a contract allowing software to share data with each other
 - lets device & software applications communicate in real time

API applications

a collection of endpoints to which developers can attach, and then extract specific information that can be used by those who are working with the apps.



Data assets

- data - that is expected to have some value in the future = forms of data or software related to the processing of information
- Intellectual property → patent
 - Databases
 - Websites
 - codes

• Applications

- a program designed to perform a specific set of operations for the end user, be it a person or another application
- web browsers
- video editing programs
- Database management software

• Programming

- refers to the functions the app performs
- regards the process of converting input to output

API → processes requests from the clients to the servers
 API → in response, brings the requested data from the server
 API → acts like a messenger

API acts as gateway extract front

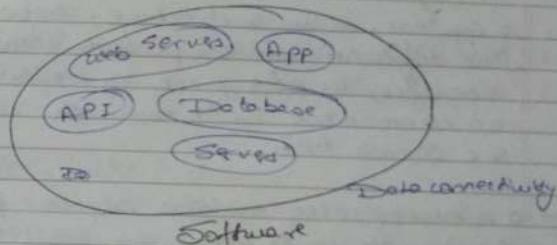
Limitations

- an API does require maintenance
- Security → APIs must act as a filter (vs. malicious programmes & overloading to guarantee quality access)
- private info. → some API services come at a certain mandatory price

Your service → build your own API to improve communication with apps & people
 → Having your own interface helps you improve the functionality of the services you are offering

Benefits of API

- data is changing extremely quickly
- all you need is a specific bit from entire dataset
- don't forget how huge the database of an external source could be
- Saves enormous amounts of your time
- provides access to numerous databases
- is very efficient when delivering precisely the info. you need.



* Different pieces of software (software products) communicate through text files

* JSON → Javascript Object Notation
 → just an ordinary text file

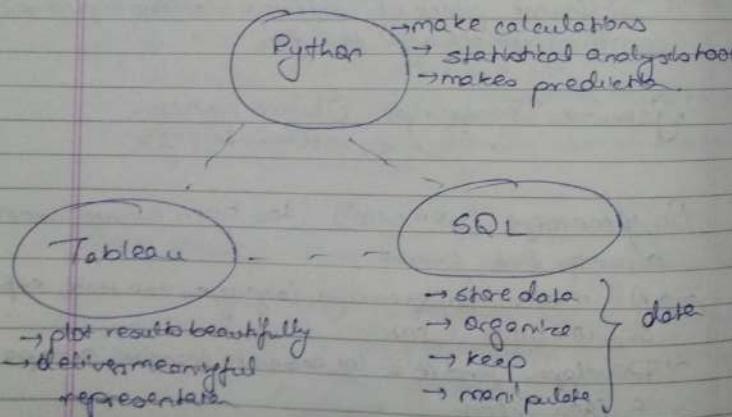
Only message (you need) has been turned into the common text format

- All modern programming languages can read, export or work with JSON
- Developers find it a lot easier to translate, or to parse a JSON string

- APIs act like a messenger b/w the end user of the app & the servers
- APIs are like a gateway to the server (or a bridge b/w the two pieces of software)
- Interfaces help you establish connections

Integration

- a system, or an architecture, composed of a few different software products, various programming languages, or other pieces of software which can communicate with each other via API or common API
- a situation where multiple software products can be set up to work as one tool



- Tableaus - — SQL
- Driver
 - o a software driver can drive a certain software product within another one
 - o It is a piece of software. Its function is to mediate b/w two products
 - o Tableau has built-in drivers
 - * Tableau to Python share the same code to functions through a common API
 - * Tabby server allows you to execute Python code remotely (to on-the-fly thanks to a built-in driver)

- Tableau built-in drivers
- ↳ may leave the user completely unaware of their existence
 - ↳ allows Tableau to centralize the advantages of using SQL relations & Python computation

Case Study

Page No. _____
Date 12/21

Good Business analyst must

SQL → be able to manage information

Python → have a substantial amount of mathematical

De statistical tools

Tableau → present results in the most intuitive way
business logic → be able to tackle the problem from a
business perspective.
understanding

The exercise will address ABSENTEEISM at
a company during work time
↳ absence from work.

→ Problem

higher competitiveness → increased pressure
unachievable business goals → raised stress levels
elevated risk of becoming unemployed → raised stress levels
unemployed

DATA
Primary
data you have created

- Survey questionnaire
- Sales operations
- Supply inventory
- Account ledger
- Business intelligence

Secondary
an already existing datasets
that somebody else has created
NOT YOURSELF

- free data from a website
- paid data from organizations

Data processing
a group of operations that will convert your raw data
(the data you have been given) into a format that is easier
to understand and hence useful for further processing &
analysis

- fix the problems that can inevitably occur with data gathering
- organize your info in a suitable & practical way

Preprocessing

1. Importing the data

import pandas → allows us to work with panel
pandas → data
Panel → handle data in tabular form

• Head_csv() → assign info from initial csv file to this variable

2. Checking the content of the dataset

* You must always make a copy of your initial dataset

df = data.copy()

df → dataframe

→ Display all items

pd.options.display.max_columns = None
pd.options.display.max_rows = None.

display(df)

→ df.info() → gives summary of dataset.

* Variable in economics → characteristics / feature, attribute → quantity that may change.

Variable in programming - storage location.

maths → vector → matrix

pandas → array

Variables → dependent → target

→ independent → predictors

Logistic regression → 0 False
1 True

→ Dropping column in dataframe

df.drop(['ID'], axis=1) → temporary off

df = df.drop(['ID'], axis=1)

→ Analyzing the reasons for absences

df [specify column].min() → find min value
e.g., 'Reason' .max() → find max value of column

(, pd.unique(df['Reason for Absence']))
→ distinct values in the columns

→ or df['Reason'].unique()

len() → no. of elements

to find missing element use
isnull() → Searched list

Quantitative analysis → add numerical meaning to our categorical nominal values → dummy variables

explaining binary variable
1 → certain categorical
absent → 0

get_dummies() → converts categorical variables into
dummy variables
Effect that
converts categorical variables into
dummy variables

• Checkpoints

In programming creating checkpoints refers to saving the current version of your code, not really the content of a variable.

→ Create a temporary save of your work so that you reduce the risk of losing important data at a late stage.

• timestamp - used for values representing dates (the classical date type found in many prog. languages). pd.to_datetime() → converts values into timestamp

When doing this conversion, you must specify the proper format of date values you will be working on

format - 'string' allows you to take control over how python will read the current date, so that it can accurately understand which no. refer to days, month, year

%d → day

%m → month

%Y → Year

%H → Hour

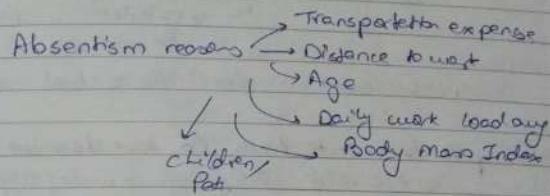
%M → minute

%S → second

(format = '%d %m %Y')

• weekday() → returns an integer corresponding to the day of the week

travel expenses = fuel + parking + meals + transportation + other



• value_counts() → Count no. in different categories

→ The more manual way of doing preprocessing gives you a higher level of control over your data

Applying ML to create absenteesism module

logistic regression

→ import libraries

random forest

→ load preprocessed data

Neural net

Logistic regression → Create targets → above which

extra

below which

np.where (and, value if True, value if False)

median() gives mean of code

value if False)

→ Using median as a cut off line is numerically stable & rigid.

→ Select 2/3 for the representation

→ A balance of 45-55% is almost always sufficient

df.iat(), since data by position of column

→ Standardize the data

absenteeism scalar = StandardScaler()

absenteeism scalar.fit(df)

calculate to store the mean & std dev.

.transform() → does the actual scaling

→ Split data into train & test dataset

↳ to avoid overfitting

80 → 20 split

↳ shuffle

random_state → make shuffle pseudo random

→ Logistic regression with sklearn

* statsmodels are not always numerically stable, for more complicated models.

↳ so use sklearn

Import

Train -> the model

{ * reg. fit(X, y) → fits the model acc. to given training data

ML { LogisticRegression()

* reg. score(X_p, targets)

→ returns the mean accuracy on the given test data & labels.

manually calculate

reg. predict(t) → predicts the class labels

for given % samples

→ Creating a summary tables with coeff & intercept.

reg.intercept_

reg.coef_

→ A feature is not particularly important:

↳ if its coefficient is around 0

↳ if its odds ratio is around 1

→ A lot (coeff) of 0 multiplies that no matter the feature value, we will multiply it by 0

→ For a unit change in the standardized feature, the odds increase by a multiple equals to the odds ratio. (1 = no change)

$$\text{odds} \times \text{odds ratio} = \text{new odds}$$

$$5:1 \times 2 = 10:1$$

$$5:1 \times 1 = 5:1$$

* When we standardized we also standardized the dummy → that's a problem

→ The further away from 0 a coefficient is, the bigger its importance.

→ Standardized models almost always yield higher accuracy

* ML engineers - prefer models with higher accuracy so they normally go for standardization

* Econometricians & Statisticians

prefer less accurate but more interpretable models because they care about underlying reasons behind different phenomena.

+ Data Scientists

may be in either position. Sometimes, they need higher accuracy, other times - they must find the main drivers of a problem.

→ The intercept or the BIAS calibrates the model

→ Backward elimination.

- * The idea is that we can simplify our model by removing all features which have close to no contribution to the model.
- * When we have the p-values, we get rid of all coefficients with p-values > 0.05.
- * If the weight is small enough, it won't make a difference anyway.

→ List comprehension is a syntactic construct which allows us to create a list from existing lists based on loops, conditions, etc.

Testing the model

→ Often the test accuracy is 10-20% lower than the train accuracy (due to overfitting)

`reg.predict_proba(x)` → returns the probability estimates for all possible o/p (classes)

Saving the model

Saving the model means saving the reg object
`pickle[module]` → is a python module used to convert a python object into a character stream

import pickle
from sklearn import linear_model

with open('model', 'wb') as file:
 pickle.dump(reg, file)
 save

workspaces

model

reg

Object to be dumped

pickle → old python tool
→ serialization
→ deserialization

pickling → converting python object → string of characters
unpickling → converting string of characters → python object

• Pickle & python version may differ

• Pickle is slow

• Pickle is not secure

→ Never unpickle data received from an untrusted or unauthorized source
→ reason - anything can be pickled so you can easily unpickle malicious code.

Deployment of the model through a module

→ store code in a module will allow us to reuse it w/o trouble.

Model

the analytical tool applied to solve the business problem

Module

A software component containing the code that will help us execute the model
→ a file containing python definitions & statements, with the suffix .py

Scaler → It contains the statistical parameters needed to adjust the magnitude of all numbers we have in this data set.

load and clean data) will preprocess the entire data set we provide

predicted_outputs - its role is to feed the cleaned data into the model, & deliver the O/pure disbursed

Tables

slope = scatter plot

Additional Python tools

template.format() → very useful tool for pretty printing
{ } {} → text data
in-place → applicable to string values only
other → will have formatted the values as text first, regardless of their initial datatype

We can use both as index values (or keys) at the same time

Iterating over range objects

help(range) range (start, stop, step)

(list comprehension)

i+j for i in range(2) for j in range(5)

lambda func → creating anonymous functions

raise_to_power2 = lambda n: n**2

raise_to_power2(3)

def lambda (n: n/2) (1) Ans 5.5

shift + tab → document

one or many parameters
contains code expressively

Lambda function

can be applied to the larger exp. they have been written

Flask

Functions

Independent entity
square = lambda

Methods

can have accessible objects data

pd.Series()

can manipulate the objects state

.idxmax() → returns index label corresponding to highest value in a series

.max()

.head(n)

.describe()

.len()

.unique() → array → delivers the values in order they have appeared in the dataset

.nunique() → no. of unique elements

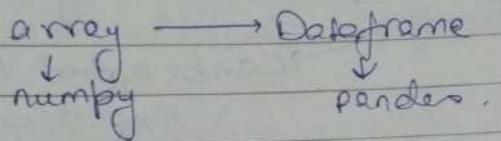
.set_value(s) → append = False

Series
single column data



Dataframe
multi-column data
2-D collection of objects

Every column from the Dataframe is a Series object itself



• `iloc[]` - purely integer location based indexing for selection by position.

• `iloc[3, 0]`

• `iloc[3, :]` → entire row.

• `loc[]`