

Finance & Risk Analytics

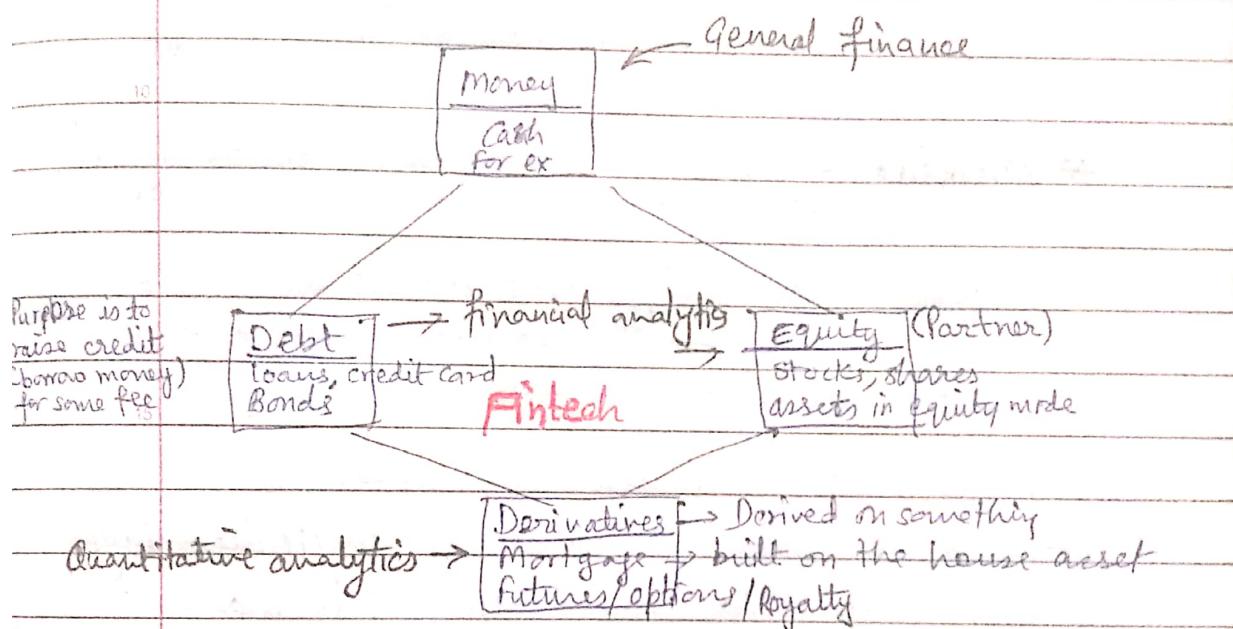
Combining

Abhinanda Parkar

Base is money - shared, traded for other things
money moving around (forex)

Liquidity risk - will you be able to buy/sell anything
- ex need to sell apartment, but not getting buyer

- You wanted cash, but you can't have it



Fintech - technologies that are used to facilitate all the above boxes.

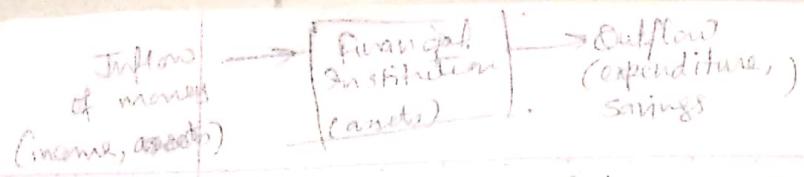
for ex, Paypal, paytm, security of a transaction visa/master. Anything which enables financial transactions, say, internet banking.

Policy nature of finance (rules & regulations)

RBI, CFPB etc

Only govt. has rights to create money, that's why on a note, RBI Governor declaration is done given.

However, its debatable now whether other than govt. anybody else should have rights to create money or not like Bitcoin, blockchain (can money be created this way? Is it valid?)



CamScanner

Risk - Uncertainty, not just one, can be opportunity also (something may happen increased unpredictability)

- Downside *

Hedging - Have another source of income apart from salary to pay the loan.
Salary is uncertain, though another source of income is also uncertain, but both together would help in reducing the risk of not being able to pay the loan.

* Downside - future plans issues due to uncertainty - getting a picture of future, what could happen.

Both uncertainty and downside helps to take the decision → This exercise is done in banks or financial institution as well.

Risks of Financial Institutions

- Risk associated with what an institution buys or sells. For banks they take money from customer (in terms of savings/ED) and sells back to customer again (in terms of loans, card).

Bank's profit = selling - buying.

Banks biggest risk is that the trading business would fail - known as credit risk.

Credit means somebody owing someone something.

2 types kind of credit risk

- 1) As a buyer - bank owes you money
- 2) As a seller - Somebody owes bank money
Bank may transfer its risk to 3rd party but as a consumer, bank owes you money and you can't share your risk with anyone

Default / credit risk is when any of the party is getting default

Market Risk — changes associated with the market, eg. Policy changes, interest rates, state of share market (deals with what a market is). Market — Real/Virtual env where buying / selling / trading is happening

For IT professional — state of IT industry is market risk. Not in control of an individual

Things ~~are~~ that can be controllable —

can be put in credit risk

Things that can ~~be~~ not be controlled —

can be put in market risk

Operational risk — Not associated with any financial institutions. It is just the risk of doing business. ~~Eg.~~ Attrition of employees (could be in control / could not be in control)

Liquidity risk can or can't be treated (often treated) as market risk as most of the time it is generated by market.

However its implications are on credit risk

Main analysis in credit Risk

① Probability of default — If the person is going to default or not

② Exposure at default — When something goes wrong how much is the impact (30% , 40% ...) how much it is going to hurt me.

- LR - whether the person would default or not (0 or 1)
DA - What are the special characteristics of a 0 person or a 1 person so when a new person comes in, does he/she have characteristics of 0 or 1 (statistical method for differentiation)

(3) Loss given default - Default has happened

It has impacted me. So how much is the final value now (could be a number or %age)

Specific kind of market risk -

- Stock market - lot of uncertainty, so you need to know, what goes in there

& Data on multiple companies - 3 yrs data of weekly closing stock price

- Work with returns (not the actual price)

$$R_t = \frac{P_t - P_{t-1}}{P_{t-1}}$$

P_t = Price
 R_t = Return

So convert current price value (weekly values) into returns and then work with distribution of returns further

- Distribution of returns

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

↑
normal distribution

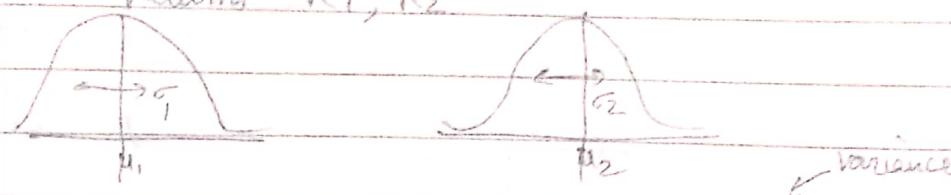
For portfolio optimization, if we understand the distribution of returns, we can take better decisions

like I know average return (μ), 10% of it has 68% possibility, 20% (has higher variance) but possibility of loss is also high and possibility of gain is also high

Portfolio Management

2 stocks in portfolio

Returns R_1, R_2



$$\text{Expected return } E(R_1) = \mu_1 \quad \text{Var}(R_1) = \sigma_1^2$$

$$E(R_2) = \mu_2 \quad \text{Var}(R_2) = \sigma_2^2$$

~~18~~ Correlation b/w $R_1, R_2 = \text{Corr}(R_1, R_2) = \rho$

Now you want to find in which proportion, we need to invest in 1st stock and 2nd stock

Suppose this proportion invested for stock 1 = wt. 'w' on R_1 , weight w on R_1 and $(1-w)$ on R_2

$$0 < w < 1$$

$$\text{Portfolio Return } (R) = wR_1 + (1-w)R_2$$

$$\text{Expected Return on } R \quad E(R) = w E(R_1) + (1-w) E(R_2)$$

$$\text{Variance on } R \quad \text{Var}(R) = \text{Var}(wR_1 + (1-w)R_2)$$

$$= w^2 \text{Var}(R_1) + (1-w)^2 \text{Var}(R_2)$$

$$+ 2w(1-w) \text{Cov}(R_1, R_2)$$

~~20~~ $\rho = \text{Corr}(R_1, R_2) = \frac{\text{Cov}(R_1, R_2)}{\text{sd}(R_1) \text{sd}(R_2)}$

~~Variance of portfolio~~ $\text{Var}(R) = w^2 \sigma_1^2 + (1-w)^2 \sigma_2^2 + 2w(1-w) \rho \sigma_1 \sigma_2$

~~When 2 returns are uncorrelated,~~

i.e. $\rho = 0$, $\text{Var}(R)$ depends on 1st 2 terms

If $\rho = +ve$, var on portfolio increases

If $\rho = -ve$, var on portfolio decreases

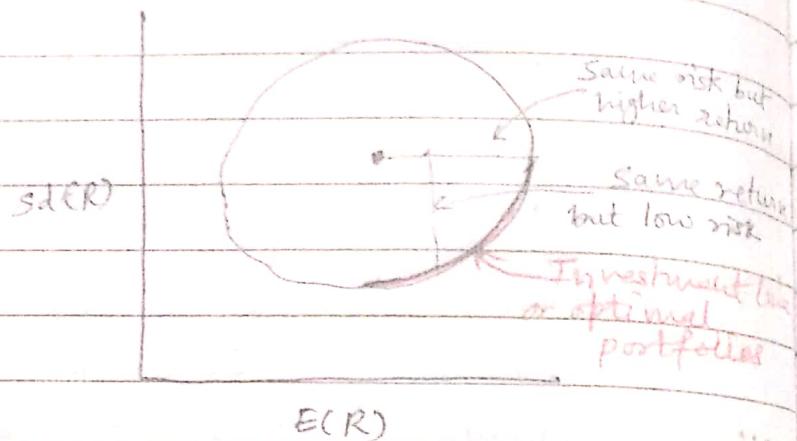
So we should pick -vely correlated stocks

i.e. if one stock goes up, other goes down

this reduces the variability and in turn

Risk

$\rho = +ve$ implies more risk but higher return as well



$$E(R) = w\mu_1 + (1-w)\mu_2$$

$$sd(R) = \sqrt{w^2\sigma_1^2 + (1-w)^2\sigma_2^2 + 2w(1-w)\rho\sigma_1\sigma_2}$$

Being on border is better than being in the circle

If I pick 2 stocks with

$$\mu_1, \sigma_1 = 5, 2$$

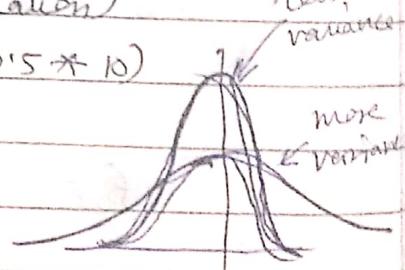
$$\mu_2, \sigma_2 = 10, 5$$

Now if $w = 0.5$ (50% in R₁, 50% in R₂)

$$\rho = 0 \text{ (no correlation)}$$

$$E(R) = 7.5 (0.5 \times 5 + 0.5 \times 10)$$

$$sd(R) = 2.6925$$



Now if $\rho = -0.5$

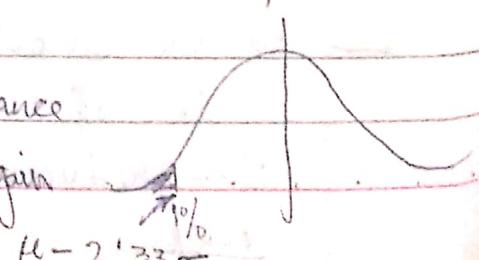
$$sd(R) = 2.179 \text{ (with same return as } \rho=0)$$

So how to do the trade off b/w risk and returns

So I don't want to play very safe but in terms of returns but at the same time I don't want to sink down (safe from worst case)

Say grey area says 1% chance

of losing, 99% chance of gain



$$\text{So } q(\text{norm}) \quad q_{\text{norm}}(0.0) = -2.3263$$

i.e. -2.326 σ point corresponds to lowest 1% risk point. Also known as **Value at Risk (VaR)** — estimate of worst case scenario

So we need to minimize VaR i.e.

$$\mu = 2.33 \sigma$$

Two ways to minimize VaR — either increase μ or decrease σ

VaR — Best of the worst

↳ a measure of "Typically" Market Risk

So in example μ is given, σ is given
you need to find weights of portfolio which would minimize VaR.

To convert price matrix into return matrix

$$R_t = \frac{P_t - P_{t-1}}{P_{t-1}}$$

$$\text{Acc.ts} = \text{as.ts}(\text{Acc}) \quad \leftarrow \text{Convert data into time series data}$$

$$\text{LAcc.ts} = \text{lag}(\text{Acc.ts}, 1)$$

$$\text{Acc.ts} - \text{LAcc.ts} \quad \leftarrow \text{change in Price}$$

$$\text{trans} = \text{Acc.ts} - \text{LAcc.ts}$$

$$\text{mean(trans)} = 0.00606$$

$$\text{sd(trans)} = 0.04113$$

$$\text{mean(BHEL)} = -0.00020$$

$$\text{sd(CBHEL)} = 0.0434$$

$$\text{cor(trans, BHEL)} = 0.3412508 \quad \leftarrow \text{the correlation}$$

It shows high risk

$$\text{Sim} = 0.5 \text{trans} + 0.5 \text{BHEL}$$

$$\text{mean(Sim)} = 0.002927 \quad \text{mean is less than individual stocks but sd also decreases}$$

$$\text{sd(Sim)} = 0.034$$

To find Var $\mu = 2.330$

$$\text{mean(sim)} - 2.33 * (\text{sd(sim)}) = -0.07782$$

* So 7% is worst case scenario on week on week data of 3 yrs horizon, that is 0.5 AEC and 0.5 BHEL. So 1% chance chance that I can loose more than 7%.

Now we need to try with different weights and choose the one which is giving best Var

Get Covariance matrix of whole portfolio (Σ)

- ~~Var~~ var(cbind(trans, BHEL))

Mean Return of portfolio \rightarrow

$$[w_1, w_2, \dots, w_{10}] \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_{10} \end{bmatrix} = \mu_e \text{ (mean return of portfolio)}$$

$$w' \Sigma w = \text{Var}(R) \quad (\text{variance of portfolio return})$$

↑ ↑
transpose of weight vector Covariance matrix

Can be applied for number of shares (as weights) instead of investment proportion.

Compoundo / Dividendo problem

average of ratios is NOT same as ratios of average

Easy way to automate it is by "optimization"

- Convert prices into returns in excel

- Create column of weights (w_1, \dots, w_{10})

- Get average returns and multiply by weights
mean of each column

- for variance, $\text{Var}(w_1 R_1 + w_2 R_2 + \dots + w_n R_n)$
- $\text{sd}(R) = \sqrt{w_1^2 \sigma_1^2 + w_2^2 \sigma_2^2 + \dots + w_n^2 \sigma_n^2 + 2 w_1 w_2 \rho \sigma_1 \sigma_2 + \dots}$

Where σ_i is sd of 1st column

- Maximize return $w_1 \mu_1 + w_2 \mu_2 + \dots + w_n \mu_n$
- Now setup criterias VaR

$$\mu = 2.33\sigma$$

- Objective is VaR (maximize)
- wts b/w 0 and 1
- $\sum(\text{wts}) = 1 (100\%)$

VaR Computation - 2 techniques

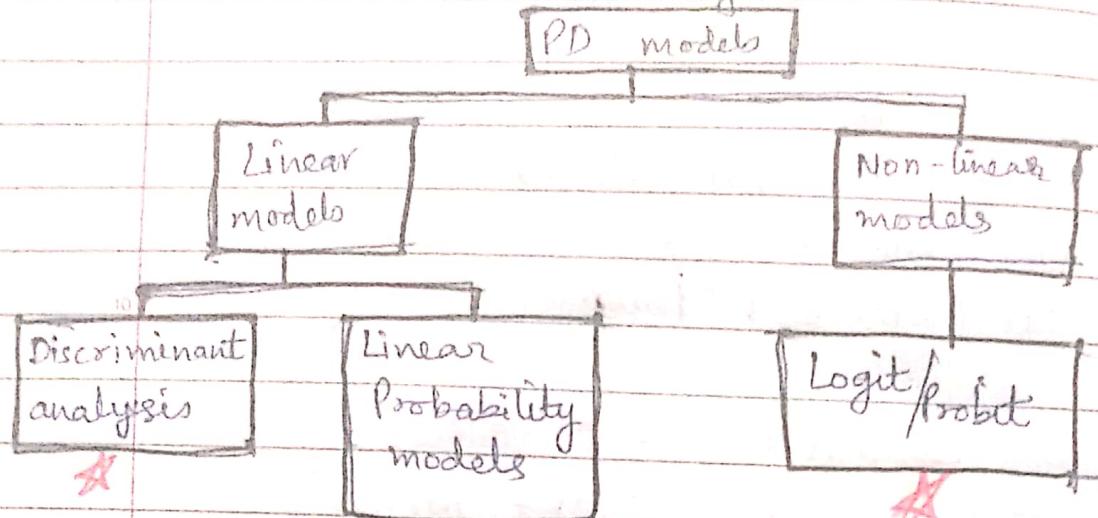
- Parametric - Delta normal
- Historical Simulation

Credit Risk

- How decisions are taken in financial industry
- How ML is adopted in taking the decisions
 - Is the person capable of paying back
 - If yes, how much should I give
 - If not, how much should I give or should not give
 - For what period, what special conditions
- At No wrong activity has happened so far, You need to assess the future based on current scenario and then take a decision

- Look for key attributes
- Check various models

PD models - statistical models in credit risk modelling



15 A person may voluntarily go default which is termed as **Fraud**

For ex. As a ML algorithm, you try to catch the pattern of defaulter/fraudulators. However, fraudulators also at the same time learn to discover new patterns (involc old ones)

for ex. if a rule says that if a person is using ATM card more than 30% beyond 30 km circle, mark the person suspicious.

Now that person would purposefully make the transactions within 30 kms (say 29 kms), if he gets to know about this internal rule.

LR - finds the probability of 0 and 1, In this case probability of default

DA - Gives probability of I given various ratios (posterior probability). So we use prior probability to get posterior probability.

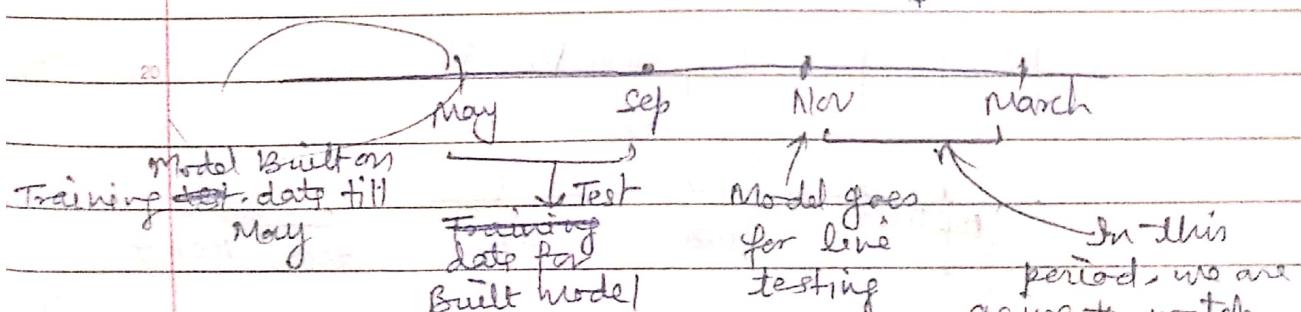
Prior & Posterior Probability

Posterior probability is the probability of assigning to groups "given" the data (condition).

Prior probability is the probability that an observation will fall into a group before you collect the data. For ex, when you know 60% males get default and 40% females get default before running any model (historical), it is called prior probability.

However, after working on different attributes you find out that some different probability, after running models and various algorithms, it is called posterior probability.

Performance window/Pilot testing
↓



here you are going to predict on test data but at the same time, you know actuals (reality) also

In this period, we are going to watch the performance of model for live data
Also, we are going to build model for next window

Generally, analytics team keep building the model every quarter. (If old one is failure, to build again, If success, to build better)

PD - linear models

① Altman's z score model (1968)

$$z(\text{Public}) = 1.2X_1 + 1.4X_2 + 3.3X_3 + 0.6X_4 + 1.0X_5$$

where $z > 2.99$ is healthy

$z < 1.81$ is unhealthy

$1.81 \leq z \leq 2.99$ is intermediate

$X_1 = \text{Working capital} / \text{total assets}$

$X_2 = \text{Retained earnings} / \text{Total assets}$

$X_3 = \text{Earnings Before Interest and Taxes} / \text{Total assets}$

$X_4 = \text{Market value of equity} / \text{Total assets}$

$X_5 = \text{Net Sales} / \text{Total assets}$

Based on above, he did DA to separate out a good company and a bad company.

→ Deep dive on variables →

clear separation in mean values of variables

b/w defaulters and non-defaulters

	3.66		1.99						
20				.77	.67				
				.07	.34				
		PBITINT	TDTA	QR	NCATA				
	3.66	1.99	.77	.67	.34				
	20		.07						
	-1.93								

• - default

• - non-default

② Discriminant function (2000)

India z score model

$$z(\text{Public}) = 1.06 + 0.01 \text{ PBITINT} - 3.12 \text{ TDTA}$$

$$+ 0.48 \cdot QR + 4.61 \cdot NCATA$$

PBITINT - Profit before interest & taxes / Total assets

TDTA - Total borrowings / Total assets

AR - Current assets - Inventories / Current liabilities
+ provisions

NCATA - Profit after tax & depreciation / Total assets

However these kind of clear cut rules (1 and 2)
might not work as ppl can easily break them
(by understanding and working to satisfy it)

10

Challenges of DA

- Zone of indifference (Intermediate)
- Sensitivity to industry
- Dealing with new companies
- Loss of predictive power across time - use of penalty functions
- Multi collinearity of variables
- Some variables provide the same information and are highly correlated
- Assumption of multivariate normality of variables.

③ Logistic model

$$\text{Logit(score)} = -6.9965 + 4.8879 * \text{Total Borrowings} /$$

$$\text{Total assets} - 0.455 * \text{Net working capital} /$$

$$\text{Current Liabilities} - 6.8605 * \text{Net cash accruals} /$$

$$/ \text{Total assets} + 1.6317 * \text{Current Liabilities} /$$

$$\text{Current assets} + 0.1978 * \text{Total borrowings} /$$

$$\text{Total Liabilities}$$

When to use logit and when to use RF.

- logit gives probability of default but at the same time, also gives the coefficients of variables. So answers questions like Why this variable is important or 1 unit increase may impact how much.
- RF gives imp variables (based on max splits) but does not answer to what extent or why this variable is imp. However, gives good prediction
- SVM, xgboost

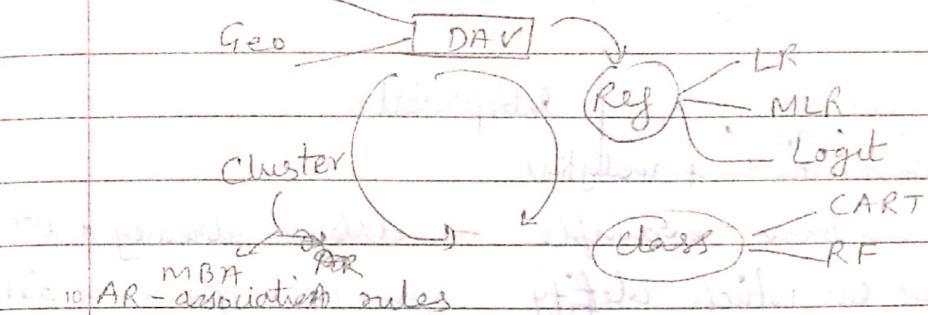
Classification - What is the probability of dependent var being 0 or 1

Discrimination - Given that the dependent variable is 0 or 1, what are the significant characteristics in discriminating 2 groups.

Web and Social media Analytics

- Text mining (Text as data) - Images, text, video, audio
- Google Analytics
- Campaign management

Time



Google to determine traffic congestion - analyses number of android phones on that road, ppl tweeting, sending messages on social media and much more. Analysis of all of this gives true picture of traffic.

20-75-5 rule for Corporate data

Numerical Text data Tacit (Intuition Intuition)

Problems with Text data

- Loose structure
- Poor spellings
- Non traditional grammar
- Multi lingual (E = action, एक्शन)
- Ambiguity (I put my bag in car. It is big & blue)
- Context - Homonyms, metaphors what? car or bag?
Sarcasm

LINK by Michael Atadewell - art of thinking without thinking
GLM allows binomial, trinomial and multinomial also

Carnegie Page
Date

P1 Extract the data
clean it
Get Text

Corpus

P2 → Corpus
- Stopwords
- Stem document
- sparse matrix

DTM/TDM

P3 DTM/TDM
logistic, CART (use various methods of analytics)

Interpretation

Phase 3 Spam base example - data is already a DTM
based on which identify Spam / No spam mail
As we have DTM, as data set

1) Perform logistic regression (In reality Text mining
and logit does not work)

spamModelLR = glm(Spam ~ ., data = spam, family = "binomial")

spamPred = predict(spamModelLR, data = spam,
type = "response")
table(Spam ~ Spam, spamPred > 0.5)

↳ will get confusion matrix

2) Built classification tree (for small data, it works)

spamCART = rpart(Spam ~ ., data = spam, method = "class")
prp(spamCART, extra = 2)

However in general industry use case there would be > 3000 columns, so rpart is not a very good idea for classification. Best is RF.

3) Random forest

spamRF = randomForest(Spam ~ ., data = spam)
varImpPlot(spamRF)

Phase 1 - Text mining

Shambhu C

Word cloud - visualization

SpiralKit - Automated Sentiment Discovery
(Byuzhat)

Correlation

So idea is to get Important variables from RF and use those in CART to fine tune the model further

Suppose there are 3000 columns in DTM, RF will

pick $\sqrt{3000}$ columns each time (randomly)

on a Random sample (with replacement) and

build a tree. Say it will build 200 such

trees and find the top 20/30 variables which

gives best Gini index (Imp variables). It means

these are the imp variables which are

Causing various splits. Once we are sure of this causal effect, we can use these 20/30 variables to fine tune CART. (No point in using these selected variables and build RF again as it would further reduce imp vars to say 5)

However, we can build RF of 3000 vars, get 30 or imp vars and again build RF to give 20/30 vars (In iteration)

Phase 2

Text mining applications

- Web and Social media, Discovery of

- Sentiments

- Opinions

- Emotions

- Topics

- Fraud and Irregularity detection

- Spam identification

- Business Intelligence

- Content enrichment

Phase 2 Take clean set of data and create DTM for further text analysis

example - Apple tweets (columns Tweet and avg rating)

Till 2015, apple used to extract say 1000 tweets and 5-6 linguists used to manually analyze the tweets and give rating from -5 to +5 on a scale of and then average of all 5% ratings were taken for a particular tweet. Now it is generated by machines

To get pure character set of tweets

If apple wants to analyze only -ve tweets could have been classified as +ve, -ve and neutral. Now avg column is not of any interest for us. Only -ve columns along with tweets

stringAsFactors = FALSE while importing or

tweets \$ Negative = as.factor(tweets \$ Avg < -1)

table(tweets \$ Negative)

Create corpus (a bag of words after extracting Text only from tweet to apply further techniques)

corpus = Corpus(VectorSource(tweets \$ Tweets))

ex - corpus of automobile companies chairman speech or say Harry potter books review

(Polarity)

wordcloud(corpus, colors = rainbow(7), max.words = 50)

However better visualization is to use Net diagram instead of wordcloud as it tells connections b/w words instead of plotting unigrams.

Wordcloud just gives an idea of corpus.

Converting corpus into Lower using "toLower" is optional. It is useful for unedited text where user can write in any form.

However, whenever date is fully edited, for ex chairman speech where every capitalization/lower case has a meaning or say TOI news,

- `corpus = tm_map(corpus, tolower)`
similarly remove punctuation, numbers, special characters etc.
- `stopwords ("english") [1:10]` (169 words in english)
- `corpus = tm_map(corpus, removeWords, c("apple", "iphone"))`

10
However, again in some cases stopwords are required like Healthcare industry

Another approach is, for my iphone use replace with my - iphone.

- 15
- Stem the document by finding root of words

`corpus = tm_map(corpus, stemDocument)`

- Now create DTM

`frequencies = DocumentTermMatrix(corpus)`

- Analyze frequencies

`inspect(frequencies [1000:1005, 505:515])`

- check for sparsity

`findFreqTerms(frequencies, lowFreq = 20)`

20
if DV is in column, go for DTM

if DV is in row, go for TDM and take a transpose of it

Sparsity - say there are 1000 tweets, and if a word is not appearing in atleast 5 tweets, we will remove them, will not consider those words which are appearing in < 5 tweets

- `sparse = removeSparseTerms(frequencies, 0.995)`
that means consider terms which are in atleast 5 terms out of 1000 (not in 995 terms)

- Now in sparse we have got 309 columns
so - Ground rule is after removing sparsity, columns should not be more than 400.

- `tweetsSparse = as.data.frame(as.matrix(sparse))`

Imp - If there are any Col names which starts with number or any special character which R does not approve, it would throw error as we process `tweetsSparse`

So better is to convert them into R friendly names (function would append `X - < col name>`)

- `colnames(tweetsSparse) = make.names(colnames(tweetsSparse))`

- Add dependent variable again

`tweetsSparse $ Negative = tweets $ Negative`

→ Now use analytics models as step 1

Example of topic identification based on various feeds (Need to find out crude feed)

- Its not as simple as finding crude word in the feed.

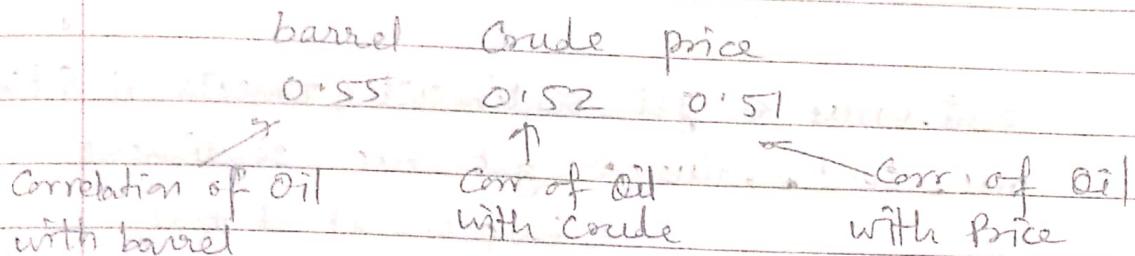
- There are other words / text which model learns to classify it as crude

An important function to find associated words after creating DTM/TDM is "findAssoc".

`findAssoc(frequenciesRT, "Oil", 0.5)`

DTM → word whose association with others need to be found where correlation is more than 0.5

Output → \$ oil



Phase 1 To get data from fire twitter and process it, packages `twitterR`, `RCurl`, `httr`, `syuzhet`, `tm`, `wordcloud`

Get API key, API secret key, Access token and Access secret token

→ Extract data

Setup twitter oauth (authorization)

Once it is established using above 4 Keys, data can be extracted

20. `tweets = searchTwitter("CB1", n = 200, lang = "en")`

`tweets_df = twListToDF(tweets)`

→ Clean data (Using gsub function)

`gsub` = global substitute

`tweets_df$text = gsub("&", "", tweets_df$text)`

similarly other strings like

"(RT) via ((?:\\b\\w*@[\\w]))" } replace with ""
 "http://wt"

"[[:punct:]]"

→ `tweets_df$text <- iconv(tweets_df$text, "UTF-8", "ASCII", sub = "")`

Convert everything which is ~~not~~ UTF or ASCII so that R can handle.

Very Important - To get sentiment

automatically (which earlier linguists were doing and then average reading was being taken)

sent.value ← get_sentiment(tweets_df \$ text)

nrc_sent_value ← get_nrc_sentiment

(tweets_df \$ text)

- Sent.value gives sentiment score of each tweet (earlier we did manually on a scale of -5 to +5 5/6 linguists and then averaging out). It gives a weighted average value say 0.3

- nrc_sent_value gives score of various human emotions (anger, frustration, trust, ...) (weighted avg), it gives scores of these individual emotions

- 2 more values possible with Syuzhet package. (not covered here)

- diff b/w sentiments and emotions

Sentiments = conscious value. If you are writing a tweet, you consciously think what to write, how to write (may be direct, may be indirect)

Emotions = sub-conscious. Your emotion comes from within in your comment. You purposefully don't write it

Syuzhet - How it works - stanford.py

As there was no stability in human analysis of tweets (may depend on mood, fatigue etc) automation of sentiment score analysis was required. Stanford analyzed 4M documents

CR Package in R to read images

CamScanner

linguistically
from 1940 to 2003 → manually analysed
→ extracted sentiment score → built
package which can now identify sentiments
by analysing words in various context

- Now as we have got sentiment score, create corpus and go to Phase 2

corpus = Corpus (VectorSource(tweets_df\$text))

Read Complete Documentation of Syuzhet Package

Bangalore's Police CCTNS - crimes and
criminals tracking Networking System
developed by wipro

Dont use words like CB1 (otherwise any sentence
where CB1 appears would come)

- Use # (hashtags for specific topics)
Use @UserName for specific users

R (Packages) and Python (libraries) are not
very suitable for text mining.

There are various other softwares / Tools
available

- Wordij (Univ of California Product)

Wordij

- Suite of products, a combination of multiple softwares

Wordlink

* Takes data and prepares DTM

* Also finds out bigrams

* Gives their frequencies in the text

* Generates 8 files everytime

VISIJ

* takes care of visualization

* Creates N/w diagram of word cloud

* Gives the frequency that which word is followed by which word and in what frequency

* Arrows will vary in thickness which will directly tell which words are always followed by other like "Wells", "Fargo".

* Link would be visible which would tell Wells is followed by Fargo but Fargo is never followed by Wells (Directional nature)

OptiComm

* Seed word and target words are used to find connection b/w 2 words which do not appear just next to each other (otherwise VISIJ would give less info)

word1 - - - - - word2 (Any connection)
some words b/w 2 words

BAPNet

Analyzes data in different time frames
for ex, tweets of day 1, tweets of day 1+7,
tweets on day 1+15 - - -

and we want to find out the variation in tweets (in corpus) over a period of time, we can use QAPNet.

- Z Utilities

* Gives Statistical inferences

* Coefficients of correlation b/w 2 terms, or 2 N-grams etc

* for ex, we used FindAssoc in R

* Gives z statistics of any term

- Utilities

- Proper Nouns

- identify Proper Nouns so that you can convert it into strings (say Harry and Potter is converted into Harry Potter)

- Now they can be used in wordlink or any N/W diagram.

- TimeSlice

- Suppose out of complete text mining data, we want to analyze data b/w specific time (say Jan 1998 to Jan 1999), we can extract / filter it using TimeSlice. However, we need timeSeries data for it. (based on text collection timestamp)

Real Example using Tweet, tact

WordLink - pick source file as tweet.txt

pick "stoplist" from the documentation folder where you have kept "wordij" folder. It gives stop words for english language.

default setting says "Drop words appearing less often than" 3 and "Drop pairs appearing less often than" 3

= Analyze Now

for huge data also, it takes few seconds

8 files are generated

Gives details about frequency of various words/pairs, their proportion in the dataset

10 VISij - generate N/w diagram by adding tweets.txt.net (in folder) file

Here you can play with # of nodes and changing node settings to picking nodes to analyze

15 * Analyze connections to word "not" or "dont". Here you can easily get idea that sentence why not apple iphone is being build up.

* Also it helps in doing customized stemming by finding out similar words

- Can build list of N/ws of different bags
- N/W diagram for jumbled up paragraph would be identical,

Opticom - select wfg file and ptg file give seed and target words

(ex. apple and freak) gives connection of these 2 words no matter whether they are n words apart

[0.217] apple → new → iphone → break. (28,3333)

Above all gives a picture whether to use any word as stop word (for deletion or not), whether we can combine words, whether we can utilize N/w and opt/cam info for stemming specific words etc.

Utilities - finds Proper Nouns

- Select tweets.txt input file
- Mention an output file say trial.txt
- Output string Replace file - gives hyphenated file
- Wherever in b/w the sentence, it finds capital it considers it as "Proper Noun"
- It does not consider first word of the sentence which by default is capital.

Word Tree Tool (IBM's)

www.jasondavies.com/wordtree

Take any text (here copied complete text from tweets.txt) and paste in Paste Text box
→ Generate word tree

Search "break" in search box →
gives connection of all words with
break in one go

→ It was part of IBM Watson (now orphaned
product)

stopped now **IBM Watson** - main feature of it is
Team b "Ask a question". Not mainly a chatbot but
if u suppose give a HR data, it would
and if you ask "why attrition is so high"
it would do the regression in the background
do some factor analysis and give you
the reasons of attrition being high.

However chatbot is just an
application of AI which gives predefined
answers.

Big ML . com

A tool developed by IBM Watson team, very
user friendly, can do all sort of supervised/
unsupervised techniques (results match very closely
to the manual versions of analysis)

CamScanner	Date
------------	------

- 1) Data is uploaded as "Sources"
- 2) Cloud icon gives option to ready "1-click dataset" or more various other options. 1-click dataset analyses Outliers, Null values, do data processing provides information (! sign) for non-recommended columns etc
- 3) Here you get "DataSource" which has ready data with distribution of various columns.
- 4) Choose "Model" for CART
- 5) Choose "Ensemble" for Random Forest
CART and Random Forest are proprietary and can be used by R only and not any other tool or language
- 6) Choose Model → Objective field (dependent variable)
→ create ML

Complete model can be exported in tableau

- 7) Go back to dataset, choose "Logistic Regression" → Build model → Evaluate (It would give confusion Matrix)
- 8) For text mining, this tool creates only Word cloud

other tools are also available
not to build a new feature
new feature has received from 100+ users

Google Analytics

Current Page
Date

Timeline feature of Google

→ timeline.google.com

google.co.in/maps/timeline

gives all history of your travel

Google analytics certified - level 1

- Who are the people coming on your website
- What are they doing on your website
- Are they looking for specific product/page
- What kind of mobile / Service provider they are using
- Browser, its language they are using
- All information about a user coming on your website

analytics.google.com

googlemerchandisestore.com

get tracking ID after after registering on google analytics.

Now place this tracking ID on every page of your website so that google can track every action of website (put that in CSS).

Now google saves a small cookie on the machine and captures date from your machine.

This is how data for a specific website is accumulated over a period of time (using TID) and browser and other user

specific information is collected (using cookies)

Webserver Master - Means you are admin of your website

Real time

- who are on my website right now
- from which part of the world they are coming
- what are they doing
- what browser they are using etc.
- helpful is creating chat box where you can directly utilize this information

Audience

- who are the people who have ever come to your website from date of inception till today
- where they have come from, their gender & age (if declared by them), which country, city, day-time they visited
- Active users
- Lifetime value (if they have bought anything from you)
- Cohort analysis (what type of ppl coming in)
- Their Demographics
- Their interest levels (Google asks for specific interest)
- Geography details
- Behaviour (how much time spending on page, images or text they are interested in etc)
- Technology, Service Provider, M/c they use

Acquisition

- If users end up buying something, how did they end up,
- How did they come from youtube, google, any organic search, referral, adword which they have put somewhere etc
- As websites pay google \$5 per click, it helps in determining where the traffic came from.

Behaviour

- How users interact with my system
- What is the flow, from which page to which page
- Any page showing 404 error or 505 error
- Is there any page which is becoming exit page (after users not liking it)

Conversion

- About money, how sales is happening
- Conversion rates

Bounce Rate - You spent money to bring ppl to your website (to acquire them), they came to website and without doing anything they disappeared. (Your site's bounce rate is very low)

Analysis of google merchandise site (1st Oct - 31st Dec 2017)

- Shows user peak on Dec 12
- So we check user vs bounce rate, it shows high peak of users but high peak of bounce rate too
- That means customers are not attracted (engaged) they did not traverse through a single page.

⑩ type of customers/users

- Suspect - I don't know them b/c, they don't know them (sending mail to testers by B school for BA course)
- Prospect - I do something, and they come to me
- Customer -
- Lead
- opportunity

⑪ - Conversion

- No of users coming from India = 6.67%
- Out of which MH is highest buyer as state (17%) and Bangalore is highest buyer city (11.9%)
- However conversion (Revenue) is coming more from mumbai (though it is less in terms of total number of users)
- Let's see based on revenue, In tableau create scatter plot and correlation b/w revenue and session. Perform cluster analysis

In analysis tab → uncheck "Aggregate"

"Measures" to create scatter plot of individual values.

- Which age group has higher % of new users
- Age pt of Indians showing highest number of users? Its revenue? What can be said about this?
- Compare Chinese users (not based on country but those who have set up their browser to Simplified Chinese) on the basis of country, what revenue patterns can be seen?
in audience → geo → go to language
select zh-cn

zh-cn = 357.99

china = 298.07

Text mining

Challenges

- Gender Identification - Big challenge in the field of text mining. Finding Proper Name from text (Capitalize letters etc.) but how to find out gender from names?

limit for free

1000 identifd An useful API - genderize.io

R Package - Genderizer (It ab is an R

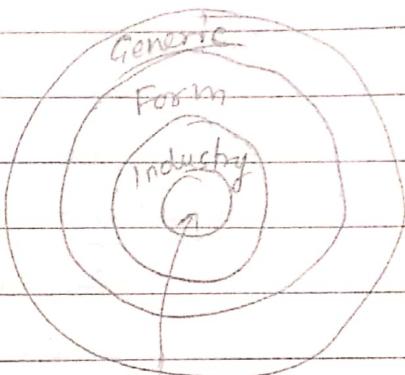
Package which goes to above API only to get info)

Campaign Management

- Nature of competition - 4 types

Brand - Direct substitutes

- ex of buying a car, people suggesting various cars based on their brand value (maruti, hundai, honda etc)



Industry - same product or class

- of product.

ex. first hand maruti car vs

Mercedez benz of 1964 (2nd hand) as a possibility so competition is not just b/w various brands of that product but within other class of products in industry

Form - Same end benefit

eg. want to buy a car but looking for other options like bullet or harley davidson which is very different from basic need but solving the purpose.

Generic - Compete for same sales volume/money

McFarlan Grid (1984) - tells about relationship b/w IT strategy and business & operations

Factory	Strategic
Support	Turnaround

Support - applications that improve management & performance but are not critical to business.
e.g. Many companies open their account on google analytics, google supports but then they don't help business much (no focus on proper usage)

factory - applications that are critical to sustaining existing business. For ex google taking care of your automated data managing instead of manual train

Strategic - applications that are critical for future success. e.g. placing strategic ad-words and attracting people

Turnaround - applications that may be future strategic decisions.

Strategic Alignment Model

Have the right ppl, bring the right job, getting the right reservation and using the data as currency.

- Generating reports is ok but what to do, how to interpret those reports

Humans of internet

- Historians - they search only intended things
- Mavericks - they find something different (adventurous)
- Confused - they do not what they are searching but still keeps searching
- Gamblers - No specific goal, just find what interests you.
- Time wasters

- Facebook data is also possible to download but is paid.
- So better to download your own facebook data in a JSON file and analyse using tableau (facebook settings download my data)

Google adwords - Specific Keywords which someone searches and your website is shown as "Ad" at the top of google. If someone clicks on that ad (URL of my website), I will pay google some \$

If there are 2 or more websites with same adwords/search words, google auctions them. Whoever pays more price will get Top Ad slot.

Women clothing example

- Based on customer reviews, perform sentiment analysis
- Check if there is any relation b/w comments and rating i.e. if somebody has really talked bad about the brand, is rating also low.
- -ve ratings are seen by more number of people and +ve ratings are seen by less number of people?
- Which products have got -ve ratings only?
Which products have got mixed ratings (+ve, -ve, neutral)
-

Python for Data Science and Machine Learning

Camlin	Page
	Date

Popular DS Libraries

- NumPy
- SciPy
- Pandas
- Seaborn
- SciKit - learn
- Matplotlib
- Plotly
- PySpark

To convert notebook to .py → use nbconvert
or save as .py file

Anaconda Virtual environments [allows to set up virtual installation]

- allows to set up virtual installations of python and libraries on your computer
- can have multiple versions of python or libraries
- These env can be easily activated/deactivated
- There is a "virtualenv" library for normal Python distribution
- Anaconda has built in virtual environment manager for easy process

Numpy

- Numpy is Linear Algebra library for Python
- All libraries in PyData ecosystem rely on Numpy as one of the building blocks
- Incredibly fast as it has bindings to C libraries
- If you have anaconda, go to command prompt
canda install numpy or
pip install numpy

- Numpy arrays are the main way we will use Numpy throughout the course.

- Numpy arrays → Vectors { Two flavors of Numpy arrays
- Matrices. } Both can be called arrays

Numpy Arrays

- are the main way we're going to be using the Numpy library
- Creating Numpy array from Python object like list

```
my_list = [1, 2, 3]
```

```
import numpy as np:
```

```
np.array(my_list)
```

array([1, 2, 3]) → array as container for list object

or → Casting of normal python list into array

```
arr = np.array(my_list)
```

```
arr
```

```
array([1, 2, 3])
```

2d list or matrix

```
my_mat = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
```

```
np.array(my_mat)
```

array([[1, 2, 3],
 [4, 5, 6],
 [7, 8, 9]])

→ list of list

Python

} casting of list of lists to
numpy array

* However, we will be using numpy own built in generation methods to actually create arrays a lot faster and simpler

```
np.arange(0, 10)
```

```
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

```
np.arange(start, stop, step)
```

from numpy.random import
randint(2, 10)
↳ 3

functions

to generate specific type of arrays

Camlin Page

Date _____

① np.zeros(3)

array([0., 0., 0.]) → rows → columns.

② np.zeros((5,5)) → tuple np.ones((5,5))

array([[0., 0., 0., 0., 0.],

[0., 0., 0., 0., 0.],

[0., 0., 0., 0., 0.],

[0., 0., 0., 0., 0.],

[0., 0., 0., 0., 0.]]

③ np.linspace(1, 5, 10)

array([1, 1.5, 2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5]) → 10 evenly spaced points

④ Identity matrix - useful matrix when dealing

with linear algebra problems.

np.eye(2) array([[1., 0.]]) → It is two dimensional square

[0., 1.]] matrix

→ diagonal of 1 and all 0's

→ Single digit as argu

⑤ Arrays of random numbers.

np.random.rand(5)

array([0.9225, 0.0715, 0.6829, 0.9956, 0.23256])

→ creates array of random no. b/w 0 to 1

* If we want to return a sample / samples from a standard normal distribution / Gaussian instead of using rand, we can use randn

np.random.randn(2) → 1d

array([0.1748905, -0.917685])

L, numbers from a stdn. normal dist.

np.random.randn(2,2) → 2d

import random np.random.randint(1, 100, 10) → gives 10 nos

start ↓ end (not inclusive)

alternatively

array([80, 10, 12, 50, 60, 66, 70, 88, 13, 93]).

attributes and methods of array

arr = np.arange(4)

[arr

→ array([0, 1, 2, 3])

[arr.shape

→ (4,)

→ 4, indicates it is 1d array

arr.reshape(2,2) → reshape in 2d array

array([[0, 1],
 [2, 3]])

② ranarr = np.random.randint(0, 10, 8)

[ranarr

array([1, 6, 9, 5, 4, 3, 7, 2])

ranarr.max() → returns max value
g or np.max(arr)

ranarr.min() → returns min value

ranarr.argmax() → gives index locⁿ of max value

2 idmax()

20 Data type in array

arr.dtype

dtype('int32')

Numpy Indexing and Selection

import numpy as np:

arr = np.arange(0, 11)

arr

array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10])

arr[8] gives 8

We can use normal python methods to get specific element in numpy

`arr[:6]` gives array `[0, 1, 2, 3, 4, 5]`

`arr[5:]` gives array `[5, 6, 7, 8, 9, 10]`

Numpy array differ from normal Python list because of
`arr[0:5]=100` their ability to broadcast

`arr`

`array([100, 100, 100, 100, 100, 5, 6, 7, 8, 9, 10])`

Slice of arr = arr[0:6]

Slice of arr

↳ `array([0, 1, 2, 3, 4, 5])`

Slice of arr[:]=99 → assign 99 to all values.

Slice of arr

↳ `array([99, 99, 99, 99, 99, 99])`

2 formats for grabbing elements from a 2d array

Matrix

`arr_2d = np.array([5, 10, 15], [20, 25, 30], [35, 40, 45])`

`arr_2d`

`array([5, 10, 15],
[20, 25, 30],
[35, 40, 45])`

`arr_2d[0][0]` → gives 5 * This is double bracket notation

Ex. 2 ways to access 30

① Double bracket notation

`arr_2d[1][2]`

② Comma Single bracket notation → recommended.

`arr_2d[1, 2]`

To get sub matrix from matrix / 2d array slicing

Ex1 `arr_2d[1:2, 1:2]`

Ex2 `arr_2d[1:, 1:]`

↳ `array([10, 15],
[25, 30])`

`array([5, 10, 15],`

`[20, 25, 30])`

`arr = np.arange(1,11)`

`arr`

`array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])`

`bool_arr = arr > 5` → comparison of array ē single digit using boolean operator
`bool_arr[Step1]`
`array([False, False, False, False, False, True, True, True,`
`True, True])`

boolean masking
or masking

`dtype = bool`)

Above can be used to actually do conditional selection

(Step 2) `arr[bool_arr]`

`array([6, 7, 8, 9, 10])` → get results where this boolean array is true

OR

Combination of Step 1 & 2 `arr[arr > 5]` also gives

`array([6, 7, 8, 9, 10])`

Numpy Operations.

Array with Array

Array with Scalars

Universal Array Functions

① Array with Array

`arr = np.arange(0,3)`

`arr`

`array([0, 1, 2])`

`arr + arr`

— Addition of arr

`array([0, 2, 4])`

`arr - arr`

`arr * arr`

Array with Scalars

→ scalar mean just a single number and Numpy does in broadcast that number to every element in array

$\text{arr} + 100$

array ([100, 102, 104])

works same with '-' '*' etc

% in python gives error. However, in numpy it does not give error but warning.

arr/arr → notice % here gives null and warning in
 $\text{array}([1/\text{nan}, 1, 1, 1, 1, 1])$ RuntimeWarning: invalid
 (%) value encountered in true-
 divide.

Similarly, $1/\text{arr}$ gives

array ([inf, 1., 0.5, 0.33, ...])

IMP :-

Universal Array Func'

① np.sqrt(arr).

② np.exp(arr)

np.max(arr) or arr.max() arr.max()

np.sin(arr)

np.log(arr)

→ for all universal functions.

docs.scipy.org/doc/numpy/reference/ufuncs.html

Ques Mat

array ([[1, 2, 3, 4, 5],
 [6, 7, 8, 9, 10],
 [11, 12, 13, 14, 15],
 [16, 17, 18, 19, 20],
 [21, 22, 23, 24, 25]])

array ([f[2],
 mat[1:3, 1:2] gives f[7],
 [12]])

Although removing 2 will give
 Should also give same result
 but it gives array ([2, 7, 12])

PANDAS

- ① Open source library
- ② Build on top of Numpy.
- ③ Allows for fast analysis, data cleaning and preparation.
- ④ It excels in performance and productivity.
- ⑤ It also has built-in visualization features.
- ⑥ It can work with data from a wide variety of sources.
- ⑦ Panda installation
 - conda install pandas
 - pip install pandas.

Series

- ① Data type
- ② Similar to numpy array.
- ③ Build on top of numpy array object.
- ④ Difference b/w numpy array & panda series is that can access labels meaning it can be indexed by label.

```

import numpy as np
import pandas as pd
labels = ['a', 'b', 'c']
my_data = [10, 20, 30]
arr = np.array(my_data)
d = {'a': 10, 'b': 20, 'c': 30}

```

pd.Series(data = my_data) → only passing 1 for data

0	10
1	20
2	30

} looks like numpy array except here we have an index 0, 1, 2 and actual data 10, 20, 30

* In panda series, we can actually specify what we want index to be

E.g. pd.Series(data=my-data, index=labels)

label or index a 10
 b 20 → Label index series
 c 30 → data points

pd.Series(my-data, labels)

Another way:

↳ Since data and index are actually in order as parameters

pd.Series(a, labels)

a 10
 b 20
 c 30

pd.Series(d) → a 10
 b 20
 c 30

* Can hold variety of objects as its data point

pd.Series(data=labels)

0 a
 1 b
 2 c

pd.Series(data=[sum, len, print])

0 <built-in function sum>

1 <, " ", len >

2 <, " ", print >

dtype = object

ser1 = pd.Series([1, 2, 3, 4], [USA, USSR, INDIA, SL])

USA	1	USA
USSR	2	U
INDIA	3	
SL	4	

dtype = INT64

ser1['USA'] = 1

Basic operations in series done based on index

Ser 2 = pd.Series([1, 2, 4, 5], ['US', 'ITALY', 'USSR', 'NORWAY'])

Ser 1 + Ser 2 gives

US 2.0

ITALY NaN

INDIA NaN

USSR 6.0

NORWAY NaN

SL NaN

dtype: float 64

Imp while performing operations with Panda series, integers are converted into floats in order to retain all possible inform

DataFrame

bld of dataframe using series object.

import numpy as np

import pandas as pd

from numpy.random import randn

np.random.seed(101)

→ to get the same random numbers

df = pd.DataFrame(randn(5, 4), ['A', 'B', 'C', 'D', 'E'], ['W', 'X', 'Y', 'Z'])

axis=0	W X Y Z				W, X, Y, Z → is panda's series
	A	B	C	D	
A					df → is just bunch of
B					series that share an
C					same index.
D					
E					

df['W'] or df.W

gives A 2.70

B 0.65

C -2.01

D 0.18

E 0.19

→ returns series

Name: W, dtype: float 64

`df.asmatrix(df)` → would convert dataframes back to
numpy 2D array.

Camlin	Page
Date	/ /

↳ `type(df['w'])`

↳ pandas.core.series.Series

↳ `type(df)`

↳ pandas.core.frame.DataFrame

`df[['w','x']]` → returns data frame with
multiple columns.

`df['new'] = df['w'] + df['y']` → creates 'new' column

`df.drop('new')` → gives Key Error as it does not
'new' column in axis=0

`df.drop('new', axis=1)` → * deletes 'new' column

* This does not happen in
place.

`df.drop('new', axis=1, inplace=True)`

* actually remove
'new' column.

`df.drop('E', axis=0)` → drops row 'E'

Why rows are referred to as zero axis and
column as 1 because it is directly taken from
the shape.

Df are just fancy index markers on top of
a numpy array.

`df.shape`

↳ (5, 4) 0 index has rows

↳ 0 index 1 index 1 " " columns.

Selecting rows - 2 ways to select rows in df

① `df.loc['A']`

N 2.7

X 0.62

Y 0.90

Z 0.50

→ returns series.

Imp. Con - Not only are column
series but rows are
also series.

Name: A, dtype: float64

② `df.iloc[2]` → numerical based index.

Subset of row & column.

`df.loc['B', 'Y']`

-0.8542

`df.loc[['A', 'B'], ['w', 'Y']]` → gives subset of rows and columns.

CONDITIONAL SELECTION USING DATAFRAME

`df > 0`

	w	x	y	z
A	T	F	F	T
B	T	F	T	T
C	F	T	F	F
D	T	F	F	T
E	T	T	T	T

Step 1 \rightarrow booldf = `df > 0`

booldf

Step 2 \rightarrow `df[booldf]`

	w	x	y	z
A	2.7	0.6	0.9	0.5
B	0.6	NAN	NAN	0.6
C	NAN	0.74	0.52	NAN
D	0.188	NAN	NAN	0.9
E	0.19	1.9	2.6	0.68

Step 1 and Step 2 can be combined to get

same result as above `df[df > 0]`

`df['w'] > 0`

A True

B True

C False

D True

E True

Name: w, dtype: bool

`df[df['w'] > 0]`

	w	x	y	z
A	2.7	0.6	0.9	0.5
B	0.6	-0.31	-0.8	0.6
C	0.18	-0.75	-0.9	0.9
D	0.19	1.9	2.6	0.6

On data frame

Stacking command such as bracket notation on top of that such as between x column where that's true

`df[df['w'] > 0]['x']`

A 0.6

B -0.3

D -0.7

E 1.9

Name: x, dtype: float64

`boolser = df['w'] > 0`

`result = df[boolser]`

`mycols = ['Y', 'X']`

`result[mycols]`

Pt. plot ('scatter')

for 2 conditions, df [(df['w'] > 0) & (df['y'] > 1)]

df.reset_index() → reset index to a column
and you'll get the actual
index then to be numerical.

	index	w	x	y	z
0	A	2.7			
1	B				
2	C				
3	D				
4	E				

newind = 'CA NY NY OR CO'. split()
newind

['CA', 'NY', 'NY', 'OR', 'CO']

df['States'] = newind.

	w	x	y	z	States
A					CA
B					NY
C					NY
D					OR
E					CO

If we want States column to be new index,

df.set_index('States'). we need to add

in-place = 'True' to
retain these changes.

new index	States	w	x	y	z
	CA				
	NY				
	NY				
	OR				
	CO				

Multi index and index hierarchy

```
import numpy as np
```

```
import pandas as pd
```

```
outside = ['G1', 'G1', 'G1', 'G2', 'G2', 'G2']
```

```
inside = [1, 2, 3, 1, 2, 3]
```

```
hier_index = list(zip(outside, inside))
```

```
hier_index = pd.MultiIndex.from_tuples(hier_index)
```

```
df = pd.DataFrame(np.random(6, 2), hier_index, ['A', 'B'])
```

	Num	A	B
Groups	1		
G1	2		
1	3		
2	1		
G2	2		
3	3		

df.loc['G1']

	A	B
1		
2		
3		

df.loc['G1'].loc[1]

A
B

Name.

To give names to index

```
df.index.names = ['Groups', 'Num']
```

Cross section

df.xs

To access G1, there are two ways -

df.loc['G1'] or df.xs('G1')

df.xs(1, level='Num')

A B

Groups

	A	B
G1	0.3	1.69
G2	0.16	1.18

{ gives row 1 of both G1 & G2

MISSING DATA

in your dataset

→ Panda will fill missing value with NaN value

import numpy as np

import pandas as pd

d = { 'A': [1, 2, np.nan], 'B': [5, np.nan, np.nan],
'C': [1, 2, 3] }

df = pd.DataFrame(d)

	A	B	C
0	1.0	5.0	1
1	2.0	NAN	2
2	NAN	NAN	3

df.dropna() → * drops complete rows with missing

A B C values. In case you want to

0 1.0 5.0 1 delete missing values from dataset.

* default axis=0, operation occurs

along the row.

df.dropna(axis=1) → drops column with null values

	A	B	C
0	1	25	44
1	2	37	3000
2	3	43	1134

It would not delete those
df.dropna(thresh=2) → drops rows

A B C rows which does not have

0 1.0 5.0 1 at least 2 NonNa values
1 2.0 NaN 2

`df.fillna(value = 'FILL VALUE')`

	A	B	C
0	1	5	1
1	2	FILL VALUE	2
2	FILL VALUE	FILL VALUE	3

`df['A'].fillna(value=df['A'].mean())`

	1.0	o/p	1
1	2.0		2
2	NaN		1.5

Name : A, dtype: float64

GROUP BY

- to tie a group rows of data together and call aggregate functions
- Groupby allows you to group together rows based off of a column and perform an aggregate function on them

`byComp = df.groupby('Company')`

< pandas.core.groupby.DataFrameGroupBy object at 0x0000

`byComp.mean()` → gives mean sales

	FB	GOOG	MSFT
Sales	75	56	152
Count	3	2	2
Name: FB, dtype: int64	1	2	1

`sum().loc['FB']` or `df.groupby('Company').sum().loc['FB']`

Sales

Count

Name: FB, dtype: int64

Basic Data

df.groupby('Company').count()

	Person	Sales
FB	2	2
GOOG	2	2
MSFT	2	2

Company Person Sales

0	GOOG	Sam	200
1	GOOG	Charlie	120
2	MSFT	Amy	340
3	MSFT	Vanessa	124
4	FB	Carl	243
5	FB	Sarah	350

df.groupby('Company').max()

	Person	Sales
FB	Sarah	350
GOOG	Sam	200
MSFT	Vanessa	340

NOTE :- It also returns the person because Python is able to store things in alphabetical order

df.groupby('Company').min()

13	Carl	243
14	Charlie	120
15	Amy	124

df.groupby('Company').describe()

	Sales
Company	
FB	count
	mean
	std
25	min
	25%
	50%
	75%
	max

df.groupby('Company').describe().transpose()

Merging, Joining and Concatenating data frames

$df_1 = \{ 'A': ['A_0', 'A_1', 'A_2'],$ $df_2 = \{ 'A': ['A_3',$
 $'B': ['B_0', 'B_1', 'B_2'],$ $'A_4'],$
 $'C': ['C_0', 'C_1', 'C_2'] \},$ $'B': ['B_3', 'B_4',$
 $index = [0, 1, 2, 3] \},$ $'C': ['C_3', 'C_4', 'C_5'],$
 $index = [3, 4, 5] \}$

$df_3 = \{ 'A': ['A_6', 'A_7', 'A_8'],$
 $'B': ['B_6', 'B_7', 'B_8'],$
 $'C': ['C_6', 'C_7', 'C_8'] \},$
 $index = [6, 7, 8] \}$

Concatenation - basically glues together dataframes.

Dimension should match along the axis

you are concatenating on-

`pd.concat([df1, df2, df3])`

By default concatenation is done on axis=0, M
however same can be done with axis=1.

`pd.concat([df1, df2, df3], axis=1)`

A B key.
0 A0 B0 K0

left
1 A1 B1 K1
2 A2 B2 K2
3 A3 B3 K3

Merging

Page
Date

Merge function
allows you to
merge Dataframes
together.

A B key.

Right
0 C0 D0 K0
1 C1 D1 K1
2 C2 D2 K2
3 C3 D3 K3

pd.merge(left, right, how='inner', on='key')

	A	B	key	C	D
0	A0	B0	K0	C0	D0
1	A1	B1	K1	C1	D1
2	A2	B2	K2	C2	D2
3	A3	B3	K3	C3	D3

on combination of keys

Merging, In case of multiple key column, innerjoin is based

pd.merge(left, right, on=['key1', 'key2'])

	A	B	key1	key2	C	D
0	A0	B0	K0	K0	C0	D0
1	A2	B2	K1	K0	C1	D1
2	A2	B2	K1	K0	C2	D2

left →	key1	key2	A	B	Right	key1	key2	C	D
0	K0	K0	A0	B0	=	K0	K0	C0	D0
1	K0	K1	A1	B1		K1	K0	C1	D1
2	K1	K0	A2	B2		K1	K0	C2	D2
3	K2	K1	A3	B3		K2	K0	C3	D3

Outer Join pd.merge(left, right, how='outer', on=['key1', 'key2'])

	A	B	key1	key2	C	D
0	A0	B0	K0	K0	C0	D0
1	A1	B1	K0	K1	NaN	NaN
2	A2	B2	K1	K0	C1	D1
3	A2	B2	K1	K0	C2	D2
4	A3	B3	K2	K1	NaN	NaN
5	NaN	NaN	K2	K0	C3	D3

Right join

`pd.merge(left, right, how='right', on=['key1', 'key2'])`

A	B	Key1	Key2	C	D
A0	B0	K0	K0	C0	D0
A2	B2	K1	K0	C1	D1
A2	A2	K1	K0	C2	D2
NAN	NAN	K2	K0	C3	D3

Left join

`pd.merge(left, right, how='left', on=['key1', 'key2'])`

A	B	Key1	Key2	C	D
A0	B0	K0	K0	C0	D0
A1	B1	K0	K1	NAN	NAN
A2	B2	K1	K0	C1	D1
A2	B2	K1	K0	C2	D2
A3	B3	K2	K1	NAN	NAN

JOINING - method for combining the columns of two potentially differently-indexed dataframes into single result DataFrame.

left	A	B	right	C	D
K0	A0	B0		K0	D0
K1	A1	B1		K1	D2
K2	A2	B2		K2	D3

`left.join(right)` → by default it is inner join

	A	B	C	D
K0	A0	B0	C0	D0
K1	A1	B1	NAN	NAN
K2	A2	B2	C2	D2

`left.join(right, how='outer')`

	A	B	C	D
K0	A0	B0	C0	D0
K1	A1	B1	NAN	NAN
K2	A2	B2	C2	D2
K3	NAN	NAN	C3	D3

Operations.

col1 col2 col3

0	1	444	abc
1	2	555	def
2	3	666	ghi
3	4	444	xyz

① df['col2'].unique()

→ array([444, 555, 666], dtype=int64)

② df['col2'].nunique()

→ 3

③ df['col2'].value_counts()

→ 444 2

555 1

666 1

Name: col2, dtype: int64.

④ selecting data

⑤ Applied method

def times2(x):

return x*2

df['col1'].apply(times2)

0 2

1 4

2 6

3 8

Name: col1, dtype: int64

⑥ Built-in functions can also be applied.

def['col3'].apply(len)

0 3

1 3

2 3

3 3

Apply is powerful when you combine it with lambda expression as we do not need to define entire function.

df['col3'].apply(lambda x: x*2)

0 888

1 1110

2 1332

3 888

Name: col2, dtype: int64

To drop column

`df.drop(['col1'], axis=1)`

	Col1	Col2	Col3
0	1	444	abc
1	2	555	def
2	3	666	ghi
3	4	777	xyz

	Col1	Col2
0	1	444
1	2	555
2	3	666
3	4	777

`df.columns` → returns index object which is list of Index(['col1', 'col2', 'col3'], dtype='object') columns.

★ This is going to be useful when you are trying to index a column and maybe the spacing is weird or you can't figure out what the string is.

`df.index` → gives list of the index.
since it is range index it actually just reports back the start, stop and step.

`RangeIndex(start=0, stop=4, step=1)`

Setting & Ordering data frame

`df.sort_values('col2')` or `df.sort_values(by='col2')`

	Col1	Col2	Col3
0	1	444	abc
1	2	555	def
2	3	666	ghi
3	4	777	xyz

`df.isnull()` → gives no values in your data frame

	Col1	Col2	Col3
0	False	False	False
1	False	F	"
2	F	F	"
3	F	F	"

Pivot Table

CamScanner

Date

df	A	B	C	D
0	foo	one	x	1
1	foo	one	y	3
2	foo	two	x	2
3	bar	two	y	5
4	bar	one	x	4
5	bar	one	y	1

df.pivot_table(values='D', index=['A', 'B'], columns=['C'])

	C	x	y
A	B		
bar	one	4.0	1.0
	two	NaN	5.
foo	one	1.0	3.0
	two	2.0	NaN

Data Input and Output

Data Sources

html & sql

- ① CSV → conda install BeautifulSoup
- ② Excel → conda install lxml
- ③ HTML → conda install html5lib
- ④ SQL → conda install sqlalchemy

* Below command assumes files are at same location
where jupyter notebook is installed

Read / Write csv files

pd.read_csv('example.csv')

a	b	c	d
0	0	1	2 3

df = pd.read_csv('example.csv')

1	4	5	6 7
2	8	9	10 11

df.to_csv('My_output')

3	12	13	14 15
---	----	----	-------

pd.read_csv('My_output') → Unnamed: 0 a b c d

not read_csv('My_output', index=False) 0 0 0 1 2 3

pd.to_csv('My_output', index=False) 1 1 4 5 6 7

pd.read_csv('My_output') 2 2 8 9 10 11

pd.read_csv('My_output') 3 3 12 13 14 15

→ 0 0 1 2 3
→ 1 4 5 6 7
→ 2 8 9 10 11

R/W Excel files

- Pandas can only import the data.
- It can't import formulas, images or things like macros. It will crash pandas
- You might have to install conda install xlrd
- pd.read_excel('Excel-Sample.xlsx', sheet_name=0)
- df.to_excel('Excel-Sample.xlsx', sheet_name=0, index=False)

R/W HTML files

```
data = pd.read_html('https://www.fdic.gov/..../xxx.html')
```

type(data)

list

↓ Panda find every table

element that is in this html file

data[0]

* Pandas read_html function will read tables off of a webpage and return a list of Dataframe object

data[0]

Bank Name City ST CERT Update Date

data[0].head() → gives top 5 rows

R/W SQL files

→ pandas.io.sql provides collection of query wrappers to both facilitate data retrieval and to reduce dependency on DB-specific API.

→ Database abstraction is provided by SQLAlchemy

→ psycopg2 for PostgreSQL } driver library for MySQL for MySQL } database

→ SQLite - included in Python's standard library by default.

Sal. Job Title [sal['EmployeeName'] == 'JOSEPH DRISCOLL']

CamScanner

Date

- We can find an overview of supported drivers for each SQL dialect in the SQLAlchemy docs.
- If SQLAlchemy is not installed, a fallback is only provided for sqlite (and for mpp/mysql for backwards compatibility, but this is deprecated and will be removed in a future version). This mode requires a Python database adapter which respects the Python DB-API.

```
from sqlalchemy import create_engine  
engine = create_engine('sqlite:///memory:')  
df.to_sql('my-table', engine)  
sqldf = pd.read_sql('my-table', con=engine)
```

sqldf	index	a	b	c	d
	0	0	0	1	2
	1	1	4	5	6
	2	2	8	9	10
	3	3	12	13	14
					15

df.info() → gives number of entries against individual columns.

* Employee with job title 'JOSEPH DRISCOLL'

sal[sal['EmployeeName'] == 'JOSEPH DRISCOLL']⁽³⁾

① Condition:

This gives that specific dataframe for employee name Joseph Driscoll

③ Adding this gives the job title.

* Person with highest salary.

sal[sal['TotalPayBenefits'] == sal['TotalPayBenefits'].max()][^{0%} 'EmployeeName']

sal.loc[sal['TotalPayBenefits'].idxmax()]

Matplotlib

matplotlib.org

Camlin Page
Date

Gallery

- Most popular plotting library
- gives control over all aspects of figure
- designed to have similar feel to Matlab's graphical plotting
- Seaborn is built over Matplotlib

How to install matplotlib

'import matplotlib.pyplot as plt'

% matplotlib inline → allows to see visualization in jupyter notebook

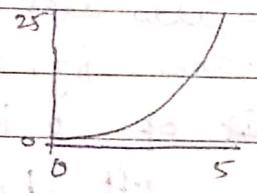
2 ways of plot creation → functional & object oriented

x = np.linspace(0, 5, 11)

y = x * x 2

① plt.plot(x, y)

② plt.show()



mainly useful if you are not plotting in jupyter

① would give plot as output in cell

② would print the plot as "kind of markdown"

plt.plot(x, y, 'r-')

adds colour (Red)

plt.xlabel('X Label')

plt.ylabel('Y label')

plt.title('Title')

(x=0, y=0, 0, 5) sets the plot range

Multiplot on same canvas

plt.subplot(1, 2, 1) ← 1st plot of the row

of rows # of columns ↑ ↑ ← plot now you are referring to

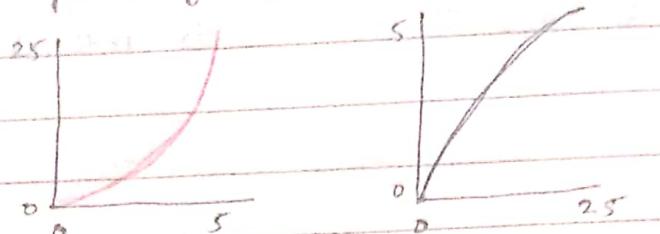
plt.imshow (mat, cmap = 'coolwarm')

Learning Page
Date / /

plt.plot (x, y, 'r')

plt.subplot (1, 2, 2) ← 2nd plot of the row

plt.plot (y, x, 'b') blue



Object-

oriented

API method

means we are going to instantiate figure objects and then call methods or attributes from that

object.

fig = plt.figure () ← creates figure object
blank at this moment,
better way → adds axes to fig.
can add attributes now.

axes = fig.add_axes ([0.1, 0.1, 0.8, 0.8])

list of 4 inputs

left height

bottom width

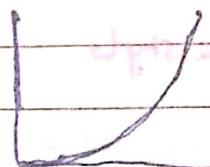
left, bottom, width, height (ranges from 0 to 1)

Basically the % of canvas you want to take

axes.plot (x, y)

x = np.linspace (0, 5, 1)

y = np.sqrt (x)



axes.set_xlabel ('X Label')

axes.set_ylabel ('Y Label')

axes.set_title ('Set Title')

fig = plt.figure ()

axes1 = fig.add_axes ([0.1, 0.1, 0.8, 0.8])

axes2 = fig.add_axes ([0.2, 0.5, 0.4, 0.3])

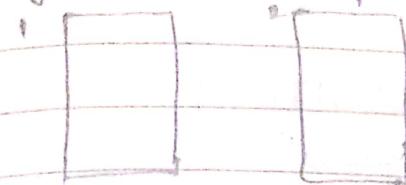
Subplots

→ tuple packing.

CamScanner

Date:

fig, axes = plt.subplots(nrows=1, ncols=2)



fig, axes are defined in previous page.

* plt.tight_layout() → fixes overlap in layout

* axes → is a list of axes object meaning you can iterate through it.

for current_ax in axes:
 current_ax.plot(x, y)

 axes[0]
 axes[1]

plotting by
iterating.

Since we can iterate three axes, we can also index it.

axes[0].plot(x, y)

axes[1].plot(y, x)

plotting by
indexing.

axes[0].set_title('First Plot')

axes[1].set_title('Second Plot')

Figure Size, Aspect Ratio and DPI, dots per inch

fig = plt.figure(figsize=(3, 2), dpi=100)

ax = fig.add_axes([0, 0, 1, 1]) → inches

ax.plot(x, y)

For subplots,

fig, axes = plt.subplots(nrows=2, ncols=1, figsize=(3, 2))

axes[0].plot(x, y)

axes[1].plot(y, x)

* Matplotlib can be used to generate high quality outputs in a number of formats png, jpeg, jpg, spg,

Tif, Pdf

fig.savefig ('my-picture.png')
.jpg

```
# fig = plt.figure()
ax = fig.add_axes([0, 0, 1, 1])
ax.plot(x, y) label = 'x squared'
ax.plot(y, x) label = 'x cubed'
```

ax.legend()

or
ax.legend(loc=0)

ax.legend((loc=0.1, 0.1))

x squared
- x cubed

Plot Appearance

ax.plot(x, y, color='green')

= '#FF8C00')

ax.plot(x, y, color='purple', linewidth=2)

ax.plot(x, y, color='green', linewidth=2, alpha=0.5)

front end or transparency of line

linestyle = '--'

or
ls = '-.'

or
ls = ':'

= 'steps'

markeredgecolor=green, marker='o', markersize=1, markerfacecolor='yellow'

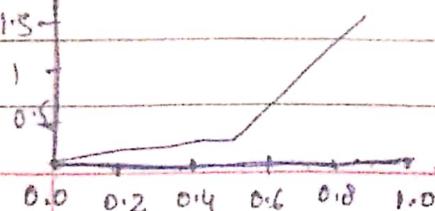
markeredgecolor='green', marker='+' or '+', markersize=10

= 'x', markersize=10

= 't', markersize=10

ax.set_xlim([0, 1]) → sets lower and upper

ax.set_ylim([0, 2]) bound of x axis



Special Plot Types

- Boxplots, histograms, scatter plots can also be created.
- Most of these plots will actually be created using **Seaborn**, a statistical plotting library for Python.

IMP - write tips plotting in jupyter notebook.

- fig should All commands for plotting fig should go in same cell.

SEABORN → very nice visualization library

- Statistical plotting library and built on top of `matplotlib` lib
- has beautiful default styles
- designed to work very well with pandas `Dataframe` objects.

Distribution Plots → Diff. plot types with Seaborn that allow us to visualize the distribution of a data set

import seaborn as sns

% matplotlib inline

tips = sns.load_dataset('tips')

	total_bill	tip	sex	smoker	day	time	size
0	16.99	1.01	F	No	SUN	Dinner	2
1	10.34	1.61	M	No	SUN	"	3
2	21.01	3.50	M	No	SUN	"	3
3	23.68	3.31	M	No	SUN	"	2
4	24.59	3.61	F	No	SUN	"	4

sns.distplot(tips['total_bill'], kde=False)

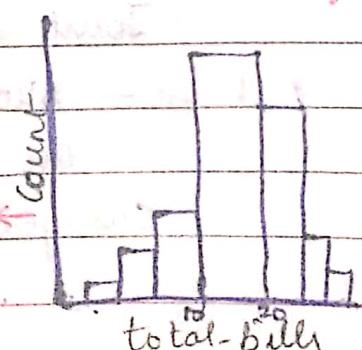
→ Distplot allows us to show the distribution of a univariate set of observations

→ kde → Kernel density estimation

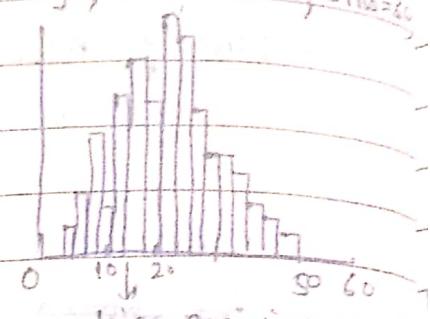
1. This plot is without kde.

2. kde = False

→ Most of total-bills are b/w 10\$ & 20\$



To get little more information
`sns.distplot(tips['total_bill'], kde=False, bins=100)`



bins are increased

If we keep `bins=100`, i.e. too high, we will get weird scenario where we are beginning to plot every instance of total bill for every single price point.

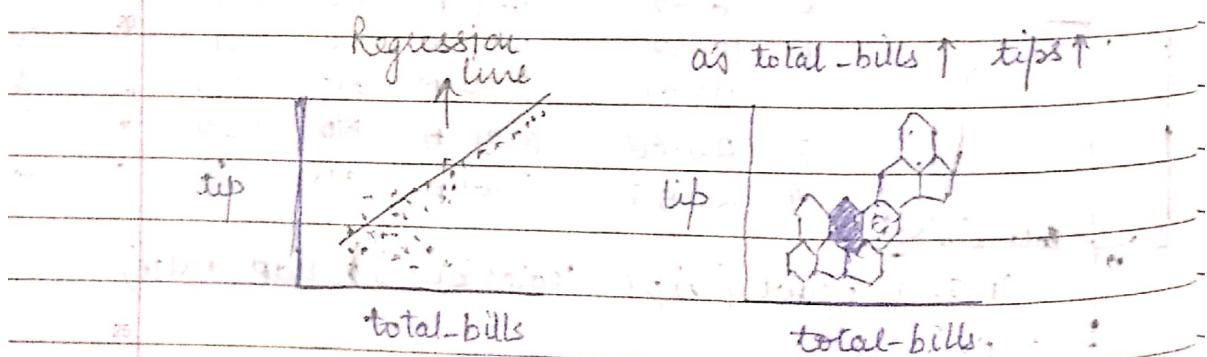
Joint plot - allows to match up 2 dist plots for bivariate data.

`sns.jointplot` meaning we can combine two different distribution plots

kind parameter - allow us to choose how we actually want to compare these two distributions

To compare distribution of total_bill Vs tip size

`sns.jointplot(x='total_bill', y='tip', data=tips)`



Joint plot gives additional argument para 'kind'

kind - Allows you to affect what's actually going on inside of this joint plot. By default it's scatter

`sns.jointplot(x='total_bill', y='tip', data=tips, kind='hex')`

helps you make hexagon distribution representation

- * If hexagon has a certain no. of points it gets darker.
" " " less " " " " lighter
- If kind = 'reg', which stands for regression
In this case Seaborn, is actually going to draw a regression line on it.

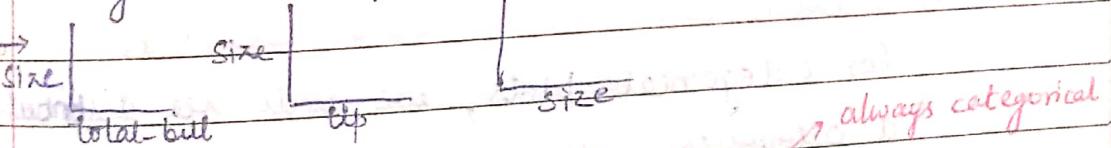

If kind = 'kde' — allows to have 2d 'kde', which shows you the density where points match up most.
→ This joint plot is usually used with scatter, as it is easiest to read and gives quite information.

Pair plot — Plots pairwise relationship across an entire data frame at least for the numerical columns.

- pair plot creates joint plot for every single combination of the numerical columns in this data frame

`sns.pairplot (tips)`

→ Larger the data frame, longer the pairplot takes



- Nice way to quickly visualize data.

`sns.pairplot (tips, hue = 'sex')`

'hue' arguments color data points based on column you put in as hue. Green points are female and blue points are male

`sns.pairplot (tips, hue = 'sex', palette = 'coolwarm')`

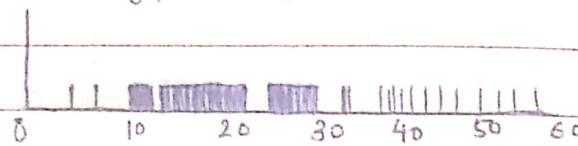
Allows to actually color this with some specific color palette.

Rug Plot

`sns.rugplot (tips['total_bill'])`

- It just draws a dash mark for every points on this uniform or unique variant distribution.
- Rug plot are used to explain kde plot using coding.

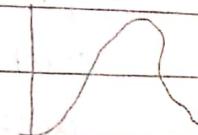
`sns.rugplot (tips['total_bill'])`



- Diff. b/w histogram and rug plot is that the histogram essentially has bins and it counts how many dashes were in that bin and then shows it as a number up here

To KDE Plot.

`sns.kdeplot (tips['total_bill'])`



Categorical Plot using Seaborn

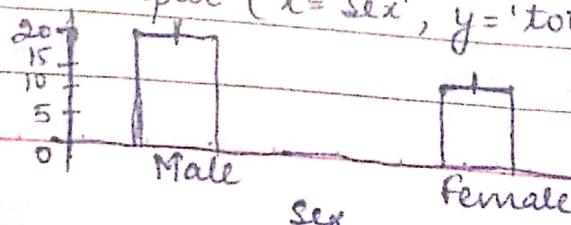
- For categorical plots, we will see distribution of a categorical column with numerical columns or another categorical column.

Bar Plot

- General plot that allows to aggregate the categorical data based of some function (default mean)

`sns.barplot (x='sex', y='total_bill', data=tips)`

mean
(total
- bill)



- Default value of estimator object is mean, however we can put in our own aggregate functions.

import ~~np~~ numpy as np

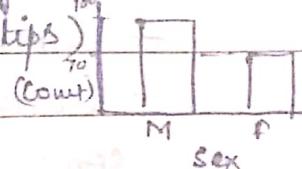
```
sns.barplot(x='sex', y='total_bill', data=tips,
             estimator=np.std)
```

Estimator - is a function that is aggregate func

Count Plot

same as bar plot except the estimator is explicitly counting the number of occurrences.

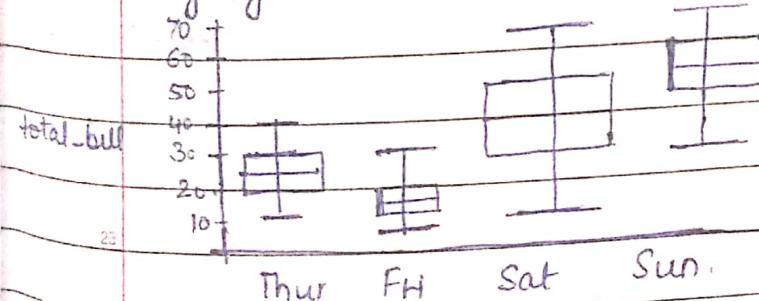
sns.countplot(x='sex', data=tips)
Wherein in bar plot estimator is mean.



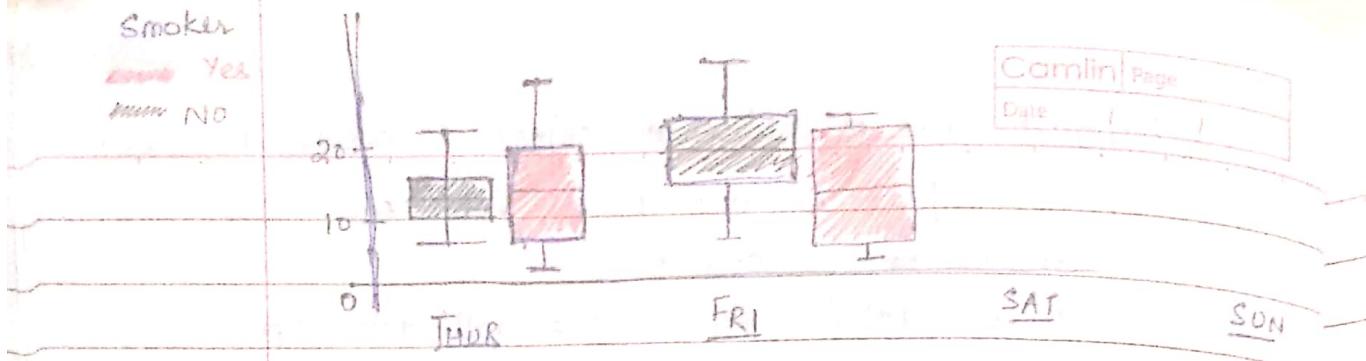
Box Plot and Violin Plot

- Used to show distribution of categorical data
- Box Plot is sometimes known as Box and Whisker plot
- Shows distribution of quantitative data in a way that hopefully facilitates comparison b/w variables.

sns.boxplot(x='day', y='total_bill', data=tips)
trying to see distribution of total-bills per day.



sns.boxplot(x='day', y='total_bill', data=tips, hue='smoker')
By adding another categorical column, we can split up these box plots even further.

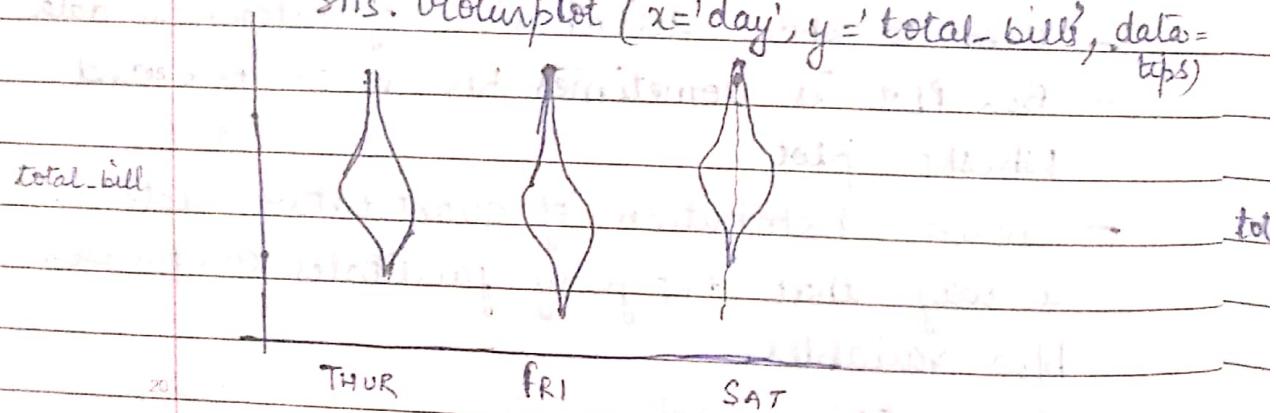


- * **IMP** → Above shows day by day behavior for total bill but also whether or not this person is a smoker.

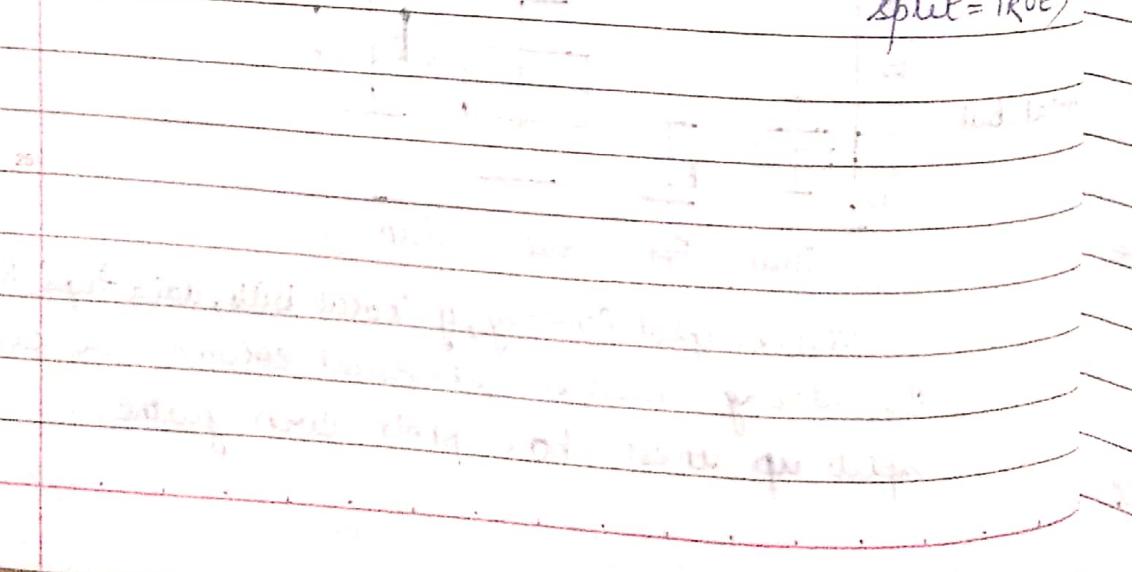
VIOLIN PLOT

- Similar to Box plot
- Used to show the distribution of data across some category

`sns.violinplot(x='day', y='total_bill', data=tips)`



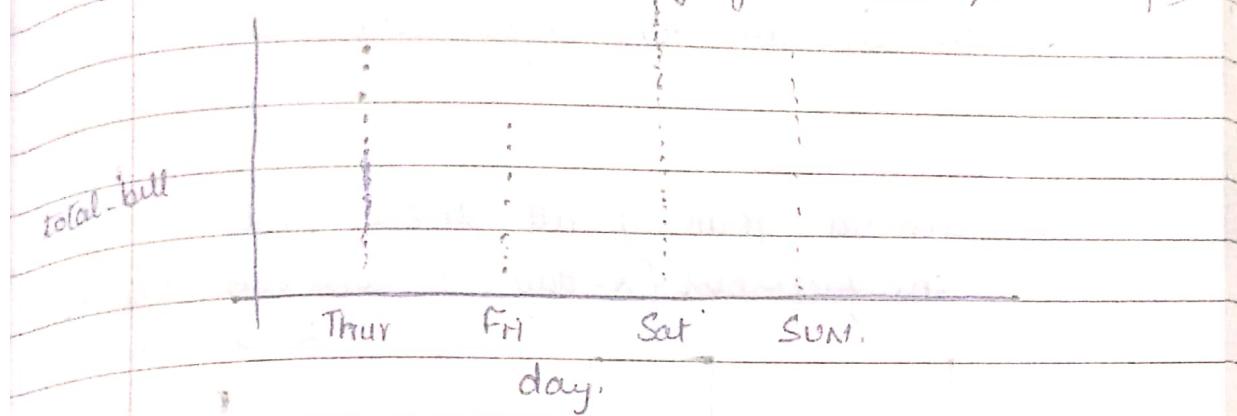
`sns.violinplot(x='day', y='total_bill', data=tips, hue='sex', split=True)`



STRIP PLOT

- Draw scatterplot where 1 variable is categorical

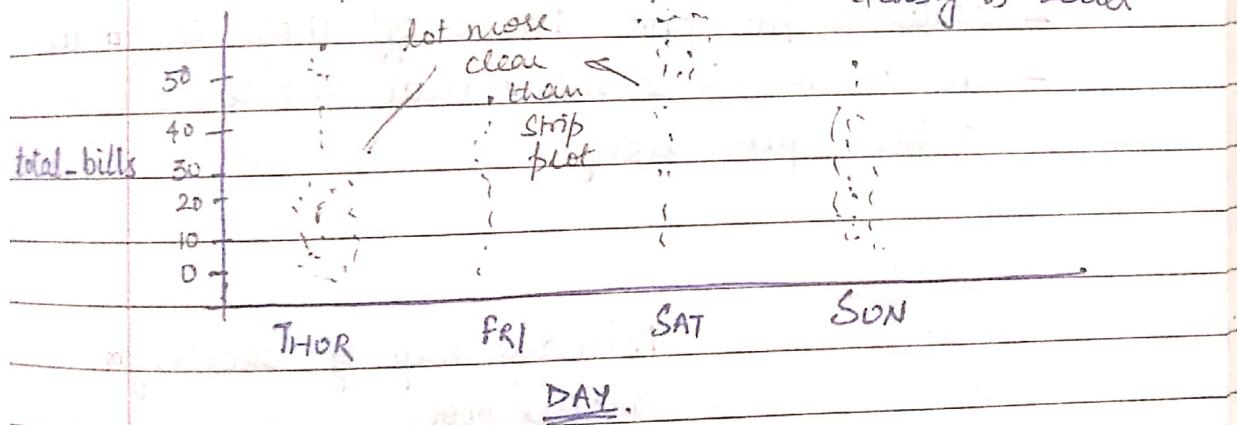
`sns.stripplot(x='day', y='total_bill', data=tips)`



- Issue with Strip plot - can't tell how many points are stacked on top of each other

`sns.stripplot(x='day', y='total_bill', data=tips, jitter=True)`

JITTER - Add a little bit of random noise to separate some of these stacked points as density is better



`sns.stripplot(x='day', y='total_bill', data=tips, jitter=True, hue='sex', split=True)`

SWARM PLOT - COMB^N OF STRIP AND VIOLIN PLOT

- Similar to strip plot except points are adjusted so that they don't overlap. This gives better representation of values

`sns.swarmplot(x='day', y='total_bill', data=tips)`

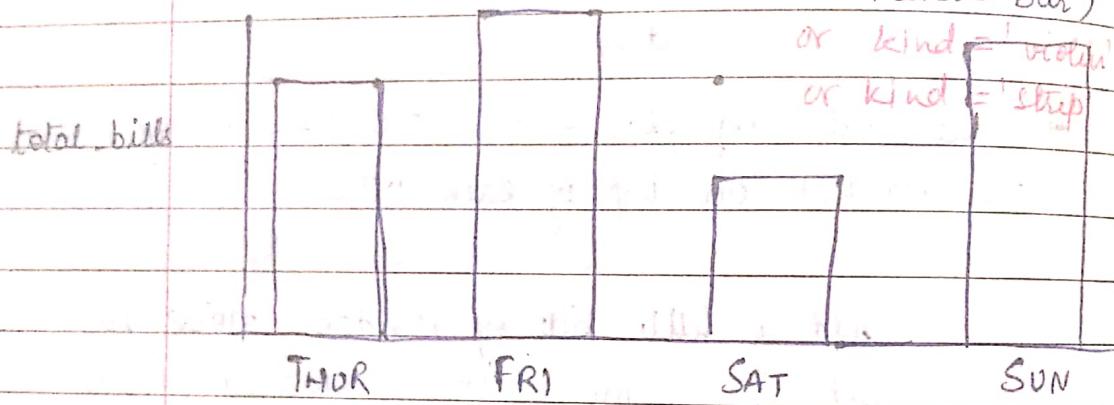
Disadv

- Sometimes they don't actually scale to very large numbers both in terms of the ability to show all the points.
- Not advisable for large data sets.

Factor Plot

General form of all these plots:

`sns.factorplot(x='day', y='total_bill', data=tips,
kind='bar')`



- Creates a bar plot based off those categories
- ~~This~~ However, it is individual choice to call actual plot itself.

Matrix Plots

Heat Map - Primary way of showing matrix plot

Condition - Data should be in matrix form

Matrix Form - Index name and column name match up so that the cell value actually indicates something that is relevant to both of those names. i.e., I need to have both variables on the columns and the same rows. This can be done either thru a pivot table or correlation data as row is not an actual variable so, we need to get it in matrix form.

Creating matrix $\text{tc} = \text{tips_corr}()$

Camlin Page
Date / /

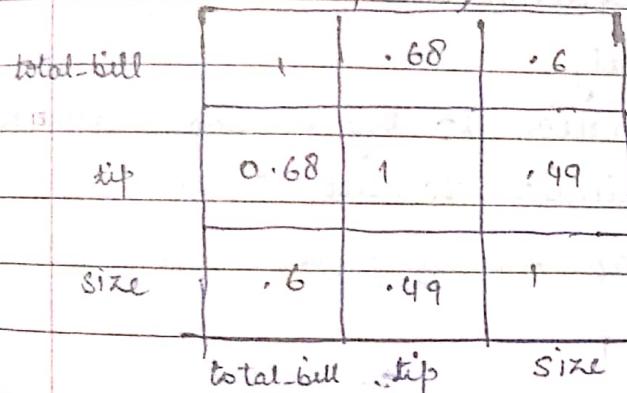
	total_bill	tip	size	
total_bill	1.000	.67	.59	
tip	0.67	1.0	.489	matrix form of tips dataset
size	0.59	.48	1.00	

flight.head()

year month passengers.

0	1949	Jan	112
1	1949	Feb	118
2	1949	Mar	132
3	1949	Apr	129

sns. heatmap (tc , annot=True, cmap='coolwarm')



matrix form of flights data

$\text{fp} = \text{flights.pivot_table}(\text{index}=\text{'month'}, \text{columns}=\text{'year'}, \text{values}=\text{'passenger'})$

Year	1949	1950	1951	1952	1953	1954	1955	1956	1957	1958	1959	1960
Month												
Jan	112	115	145									
Feb												
Mar												
Apr												
May	121	125	172									

Above, if we see no's. are correlated now. Year vs Month

Heat Map can be used now to quickly visualize the data

sns. heatmap(fp)

- * looking at heat map we can easily tell which month of which year has most flights.

sns. heatmap(fp, cmap='magma', linecolor='white', linewidth=1)

CLUSTER MAP (SECOND MATRIX TYPE PLOT)

- Cluster map uses hierarchical clustering to produce a clustered version of heat map.

sns. clustermap(fp)

- This ~~has~~ cluster column and rows together based off their similarity.
- Clusters of information to try to show columns and rows that are similar to each other.

sns. clustermap(fp, cmap='coolwarm', standard=
scale=1)

GRIDS

- How to use sideboards grid's capability to automate subplots based on features of our data

import seaborn as sns

%matplotlib inline

iris = sns.load_dataset('iris')

iris.head()

sepal-length sepal-width petal-length petal-width specie
setosa

0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	"
2	4.7	3.2	1.3	0.2	"
3	4.6	3.1	1.5	0.2	"
4	5.0	3.6	1.4	0.2	"

Simplified version of pair grid

→ Does lot of stuff automatically

PAIR PLOT, sns.pairplot(iris)

PAIR GRID → We have to do little more but gives lot more control.
sns.PairGrid(Iris) → It takes all numerical columns and grids them up.

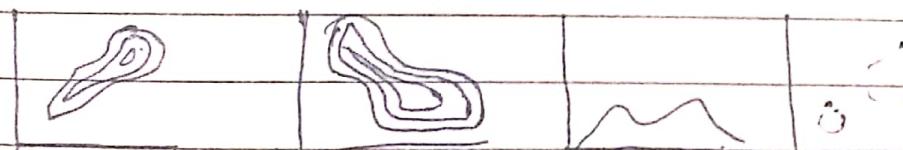
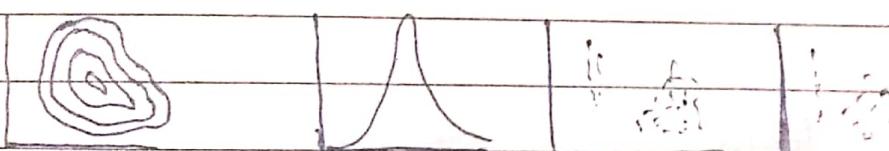
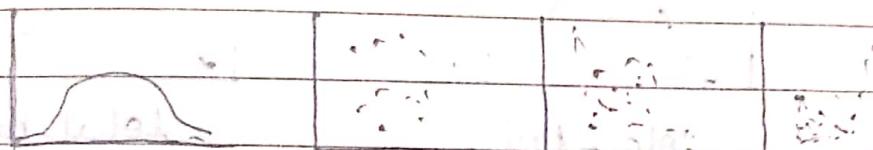
g = sns.PairGrid(Iris) → creates empty grid.

g.map(plt.scatter) → map funcⁿ maps scatter funcⁿ of plt library to complete dataset grid obj.
Since, pair grid gives more control, you can plot as below,

g.map_diag(sns.distplot)

g.map_upper(plt.scatter)

g.map_lower(sns.kdeplot)

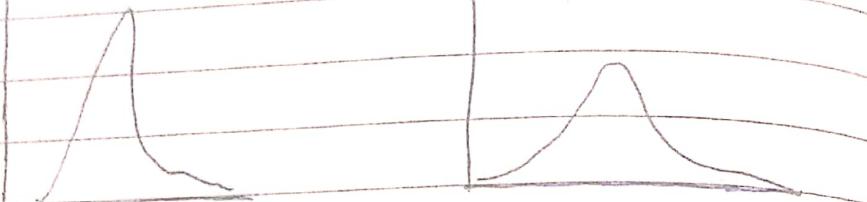


FASET GRID

`g = sns.FacetGrid(data=tips, col='time',
row='smoker')`

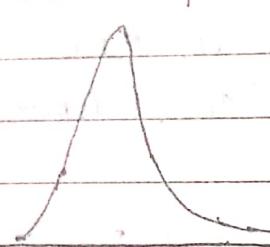
`g.map(sns.distplot, 'total_bill')`

time:
smoker=Yes | Lunch Smoker=Yes | time=Dinner



Smoker=Yes | time=Lunch

Smoker=Yes | time=Dinner



total_bill

total_bill

REPO

15

16

17

18

19

20

21

22

23

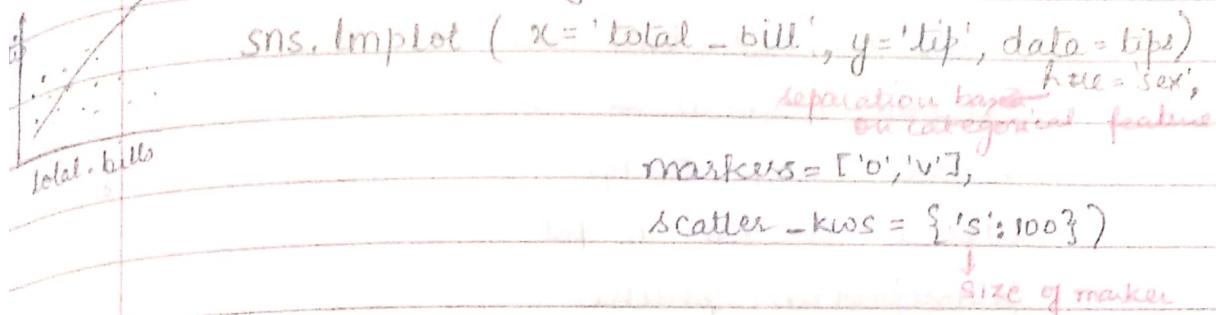
24

25

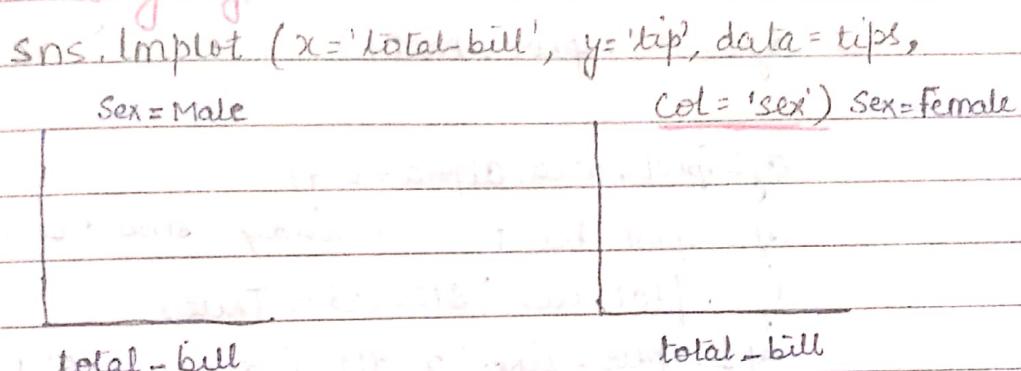
REGRESSION PLOTS WITH SEABORN

#1 LM PLOT

Allows to display linear models with Seaborn.



Same can be done using grid instead of separating by hue:



`sns.lmplot(x='total_bill', y='tip', data=tips, col='day', hue='sex', aspect=0.6, size=8)`

Style and Color (`sns.set_style('whitegrid')`)

`sns.countplot(x='sex', data=tips)`

`sns.despine(left=True, bottom=False)`

Pandas In-built Visualization

To get a histogram of all the values in A column of DF1

✓ df1['A'].hist(bins=30)

import seaborn as sns.

import numpy as np

import pandas as pd

%matplotlib inline

df1 = pd.read_csv('df1', index_col=0)

✓ df1['A'].plot(kind='hist', bins=30)

✓ df1['A'].plot.hist()

df2.plot.area(alpha=0.4)

df2.plot.bar(1) → index should be categorical

df2.plot.bar(stacked=True)

df1.plot.line(x=df1.index, y='B', figsize=(12, 8), loc=1)

A B C D.

2000-01-01

2000-01-02

2000-01-03

2000-01-04

df1.plot.scatter(x='A', y='B')

We can use style sheet using following commands.

✓ plt.style.use()

→ plot-bmh

→ plot-fivethirtyeight

→ plot-ggplot and many more

Import matplotlib.pyplot as plt

plt.style.use('ggplot')

df1['A'].plot.hist(c=-).as_usual

To select specific number of rows
`df3.ix[0:30].plot.area(alpha=0.4)`

PLOTLY AND CUFFLINKS

PLOTLY

- for more details visit github
- for financial analysis

- Interactive visualization library.
- Cufflinks connects plotly with pandas.
 - ↳ It is productivity tool that allows you to quickly call plot's just as you would with pandas for instance a data frame dot plot method.
 - Instead you would actually create an interactive visualization using plotly.

import pandas as pd

import numpy as np

from plotly import version

print(version)

↳ 1.12.9

import cufflinks as cf

from plotly.offline import download_plotlyjs, iplot,

init_notebook_mode, plot,

init_notebook_mode(connected = True) → connects

javascript to notebook

as plotly essentially connects pandas and python to its interactive javascript library. And this is going to allow your notebooks to access these visualizations

cf. go_offline()

`df = pd.DataFrame(np.random.randn(100,4), columns = 'ABCD'.split())`

`df.head()`

Using dict, for creating another dataframe

`df2 = pd.DataFrame({ 'category': ['A', 'B', 'C'],
'values': [32, 43, 50] })`

df2	Category	Values	A	B	C	D
			0	1	2	3
	A	32	0.83	2.4	-2.4	-4.3
	B	43	-2.8	-0.7	-2.6	-1.27
	C	50	-1.2	1.86	0.99	0.50
			-0.4	-0.34	0.99	-1.7
			0.001	1.41	-0.20	-1.5

`df.plot()` → pandas will by default take matplotlib

`df.iplot()` → matplotlib is converted to plotly

interactive, Zoom in, save & download

the plot as png, filter data.

How to use iplot to create other plots:-

~~Scatter~~ `df.iplot(kind='scatter', x='A', y='B')`

→ By default it will join points with lines

`df.iplot(kind='scatter', mode='markers', size=20)`

Bar `df2.iplot(kind='bar', x='category', y='value')`

`df.count().iplot(kind='bar')`

`df.sum().iplot(kind='bar')`

Box `df.iplot(kind='box')`

3D surface plot

`df3 = pd.DataFrame({ 'x': [1,2,3,4,5],`

`'y': [10, 20, 30, 20, 10],`

`'z': [500, 400, 300, 200, 100]`

`df3.iplot(kind='surface')`

`, colorscale='rdylbu')`

`df['A'].iplot(kind='hist', bins=25)`

`df.iplot(kind='hist', bins=25)`

Spread → for stock data

`df[['A', 'B']].iplot(kind='spread')`

→ used in world GDP comparison to population, happiness factor; United Nations Report
Bubble Plot - similar to scatterplot except that it will change size of the points based off another variable

```
df.iplot ( kind = 'bubble', x = 'A', y = 'B', size = 'C' )
```

Scatter matrix plot → similar to seaborn pair plot

→ it just creates a scatter matrix of all columns it can
df.scatter_matrix () → IMP - All columns should be numerical to make this work

- * There are some technical analysis capabilities with cufflinks that are built in. → (under beta)

Geographical Plotting

- challenging due to various formats data can come in
- We will use plotly for plotting
- Matplotlib also has a basemap extension.

Choropleth Maps

- This allows to plot out information either on a global scale or any nation wide scale.
import plotly.graph_objs as go
- import plotly.plotly as py
from plotly.offline import download_plotlyjs, iplot
init_notebook_mode, plot,
- data = dict (type = 'choropleth',
variable ↓ locations = ['AZ', 'CA', 'NY'],
location mode = 'USA-States',
colorscale = 'Portland',
text = ['Text1', 'Text2', 'Text3'],
z = [1.0, 2.0, 3.0],
colorbars = { title: 'Colorbar Title (0 to
Here)' })

layout variable

↑
layout = dict (geo = { 'scope': 'usa' })

choromap = go.Figure (data = [data], layout = layout)
iplot (choromap)

import Pandas as pd

df = pd.read_csv('2011-US-AIRPORT-EXPORTS')

<https://plotly/python/reference/#choroplot>

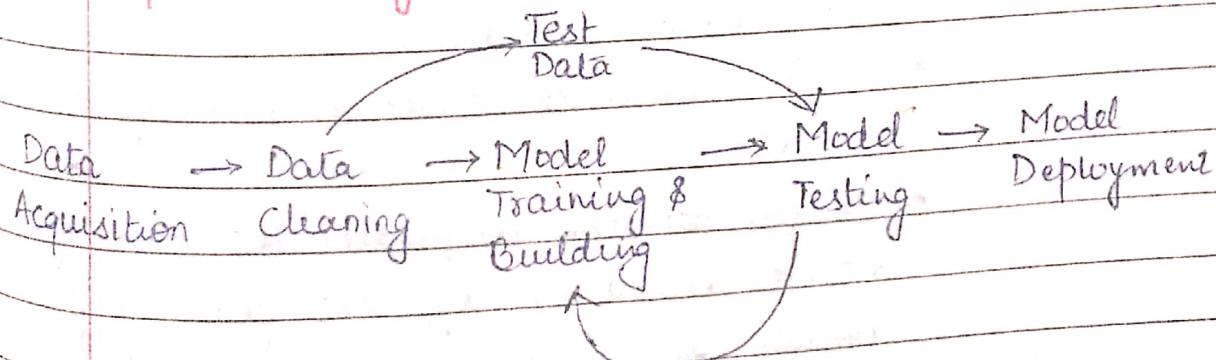
Machine learning:

- Is a method of data analysis that automates analytical model building.
- Using algorithms that iteratively learn from data, machine learning allows computers to find hidden insights without being explicitly programmed where to look.

Usage of ML

1. Fraud detection.
2. Web search results.
3. Real time ads on web pages.
4. Credit scoring and next-best offers.
5. Prediction of equipment failures.
6. New pricing models.
7. Network intrusion detection.
8. Recommendation engines.
9. Customer Segmentation.
10. Text sentiment analysis.
11. Predicting customer churn.
12. Pattern and image recognition.
13. Email spam filtering.
14. Financial modeling.

ML learning Process



3 main types of ML algorithms

1. Supervised Learning.
2. Unsupervised "
3. Reinforcement "

1. we have labeled data and are trying to predict a label based off known features.
2. You have unlabeled data and are trying to group together similar data points based off of features.
3. Algorithm learns to perform an action from experience.

1.a. Algorithm learns by comparing its actual output with correct outputs to find errors.

1.b. It then modifies the model accordingly.

1.c. Through methods like classification, regression, prediction and gradient boosting, supervised learning uses patterns to predict the values of the label on addition unlabeled data.

1.d. Commonly used in applications where historical data predicts likely future events.

1.e. Ex → it can anticipate when c.c. transactions are likely to be fraudulent.

* Which insurance customer is likely to file a claim:

* It can attempt to predict the price of a house based on different features for houses for which we have historical price data.

Unsupervised learning -

- used against data that has no historical labels.
- System is not told the 'right answers'. The algorithm must figure out what is being shown.
- The goal is to explore the data and find some structure within or it can find the main attributes that separate customer segments from each other.
- Techniques includes self-organizing maps, nearest-neighbor mapping, k-means clustering and singular value decomposition.
- These algorithms are also used to segment text topics, recommend items and identify data outliers.

Reinforcement Learning

- used for robotics, gaming and navigation.
- With reinforcement learning, the algorithm discovers through trial and error which actions yield the greatest rewards.
- It has 3 primary components:
 - ① the agent (the learner or decision maker)
 - ② the environment (everything the agent interacts with)
 - ③ actions (what the agent can do)
- The objective is for the agent to choose actions that maximize the expected reward over a given amount of time.
- The agent will reach the goal much faster by following a good policy.
- So, the goal in reinforcement learning is to learn the best policy.

Scikit Learn package - has lot of algorithms built.
Every algorithm is exposed in scikit-learn
via an "Estimator"

from sklearn.linear_model import Model

Example

Estimator obj

from sklearn.linear_model import LinearRegression

Next, step is to instantiate that model or estimator

Estimator parameters - All the parameters of
an estimator can be set when it is
instantiated, and have suitable default values

Example:-

model = LinearRegression(normalize=True)

print(model)

LinearRegression(copy_X=True, fit_intercept=True,
normalize=True)

Splitting Train and Test data

X-features

y-labels

import numpy as np

can be applied
on features and
labels

from sklearn.cross_validation import train_test_split

X, y = np.arange(10).reshape((5, 2)), range(5)

X

list(y)

array([[0, 1],
 [2, 3],
 [4, 5],
 [6, 7],
 [8, 9]])

[0, 1, 2, 3, 4]

X-train, X-test, y-train, y-test = train_test_split(X, y)

test_size=0.3)

> X-train

array([[4, 5],
 [0, 1],
 [6, 7]])

y-train
[2, 3]

x-test

array([[2, 3],
 [8, 9]])

y-test

[1, 4]

Train / Fit our model

Carolina Page

After splitting data, train / fit model on training data

`model.fit(x-train, y-train)`



Model is ready to predict labels or values on the test set.

Predicted values

`predictions = model.predict(x-test)`

We can now evaluate our model by comparing our predictions to the correct values.

The evaluation methods depends on sort of ML algorithm (Ex. Regression, classification, clustering etc.)

Available in all Estimators

- `model.fit()`: fit training data.

- For SL apps, `model.fit(x, y)`

- For USL "", `model.fit(X)`

In supervised estimation

- `model.predict(x-new)`

- `model.predict_proba()` → for classification

problems, some estimators also provide this method, which returns the probability that a new observation has each categorical label. In this case, the label with the highest probability is returned by `model.predict()`

- `model.score()` → for classification or regression problems,

most estimators implement a score method.

Scores are b/w 0 and 1, with a larger score indicating a better fit.

In Unsupervised estimators

- + model.predict() → predict labels in clustering algorithms.
- + model.transform() → given an unsupervised model, transform new data into the new basis.
- + Accepts one argument X_new, and returns the new representation of the data based on the unsupervised model.
- + model.fit_transform → some estimator implement this method, which more efficiently performs a fit and a transform on the same input data.

Choosing an Algorithm (download cheat sheet - scikit-learn algorithm)

```
from sklearn.preprocessing import MinMaxScaler  
scaler_model = MinMaxScaler()  
scaler_model.fit(dataframe)  
scaler_model.transform(dataframe)
```

Linear Regression

CamScanner Page

Date

import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

%matplotlib inline

USAhousing = pd.read_csv('USA_Housing.csv')

USAhousing.head()

USAhousing.info() → gives column details (no. & data type)
→ no. of entries

USAhousing.describe() → count

mean

std

min

25%

50%

75%

max.

USAhousing.columns

FDA sns.pairplot(USAhousing) → Does plotting b/w

sns.distplot(USAhousing['Price']) numerical variables.

↳ creates dist. blot of Price column

sns.heatmap(USAhousing.corr())

↳ draws correlation plot

Training a Linear Regression Model

X = USAhousing[['...features column...']]

y = USAhousing[['Price']] →

Train Test Split

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = % of test data

train_test_split(X, y, test_size=0.4,

random_state=101)

optional - This essentially ensures a specific set of random splits on your data because this train split is going to occur randomly in order to make sure that you train your model is same as mine

Creating and Training the model

from sklearn.linear_model import LinearRegression

lm = LinearRegression()

lm.fit() → to train or fit my model on

lm.fit(x-train, y-train) training data

LinearRegression(copy_X=True, fit_intercept=True,
n_jobs=1, normalize=False)

print(lm.intercept_)

→ -2640159.79.

lm.coef_

array([2.15, 1.64, 1.22, 2.233, 1.5157])

x_train.columns

Index(['column names'], dtype='object')

cdf = pd.DataFrame(lm.coef_, x.columns,
columns=['coeff'])

cdf

Avg. Area Income

coeff
21.52

Avg. Area House Age

164803.28

Avg. area no. of Rooms

122368.67

Avg. area no. of Bedrooms

2233.80

Area Population.

15.15

Ex: Importing real data from sklearn

from sklearn.datasets import load_boston

boston = load_boston()

boston.keys()

dict_keys(['data', 'feature_names', 'target', 'DESCR'])

print(boston['DESCR'])

↓
Column
name
data of
housing
prices

↓
target prices

Predictions from Test set

`predictions = lm.predict(x_test)`

`predictions`

`array([1260960.7, ..., 7])` → predicted prices of the house

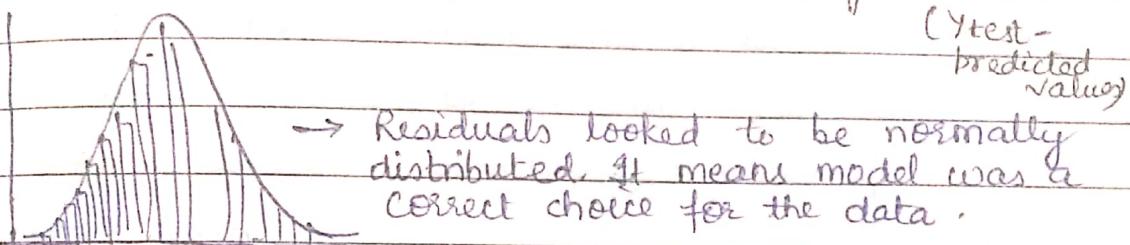
`y-test` → contains correct prices of house

Now, we want to know how far off are the predictions from the test prices the actual prices. This can be done by scatter plot.

`plt.scatter(y_test, predictions)`

`sns.distplot((y-test - predictions))`

→ Histogram of the distribution of our residuals



→ Residuals looked to be normally distributed. It means model was a correct choice for the data.

→ If not, normally distributed, in that case we might have to see if a linear regression model was the correct choice for the dataset.

Regression Evaluation Metrics

There are 3 evaluations metrics for regression problems:

Mean Absolute Error (MAE)

- is the mean of the absolute value of the errors:

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Mean Squared Error (MSE)

is the mean of the squared errors:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE)

is the square root of the mean of the squared errors:

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

MAE

MSE

RMSE

is the easiest to understand, because it's the average error.

is more popular than MAE, because MSE "punishes" large errors, which tends to be useful in the real world.

is even more popular than MSE, because RMSE is interpretable in the "y" units.

* All are loss functions, because we want to minimize them.

```
from sklearn import metrics
metrics.mean_absolute_error(y-test, prediction)
metrics.mean_squared_error(y-test, prediction)
np.sqrt(metrics.mean_squared_error(y-test,
                                     predictions))
```

r^2 can be calculated by
 ↳ a measurement of how much variance model explains.

① metrics.explained_variance_score(y-test, prediction)

or ② r^2 in this case is. 0.98 i.e. 98% of variance is explained by our model on test data.

③ lm.score(y-test, predictions)

LOGISTIC REGRESSION (See 4-403)

- It is method for 'Classification'

- Ex:- ① Spam versus "Ham" emails

② Loan default (yes/no)

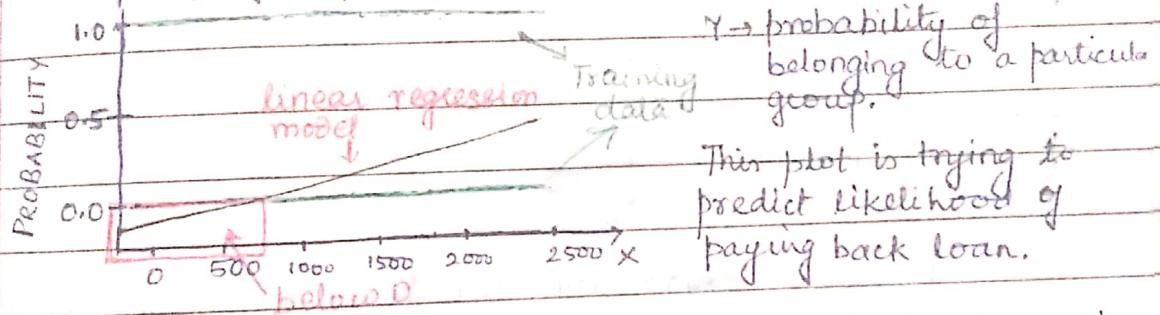
③ Disease diagnosis

Binary classification

- ML and statistics classification is the problem of identifying to which of a set of categories a new observation belongs to based off of training data.

- It allows to solve classification problems where we are trying to predict discrete categories.

IMP - We can't use a normal linear regression model on binary groups. It won't lead to a good fit.

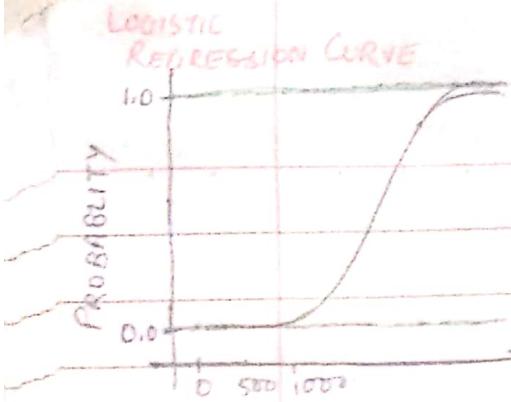


0 → 0% probability as defaulting on their loan meaning they have a 0% probability of being able to pay back their loan.

1 → 100% probability as fully paying back their loan

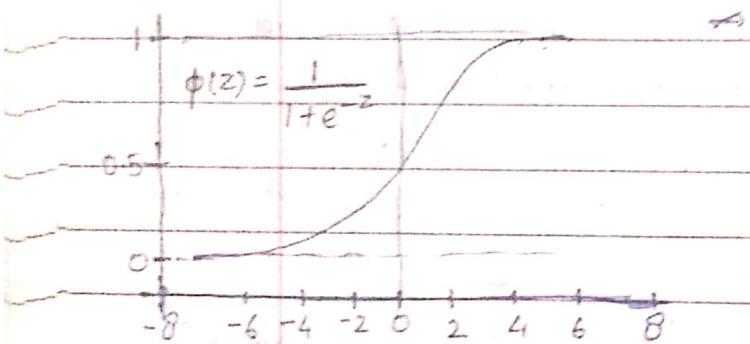
IMP - As paycheck goes lower you have a closer to zero % probability that you're going to be able to pay back your loan and as your paycheck value gets higher you didn't have closer to 100% probability of paying back loan.

On training data (marked green) if we try to use a linear regression model we will get bad fit and end up predicting probabilities below zero % which does not make any sense.



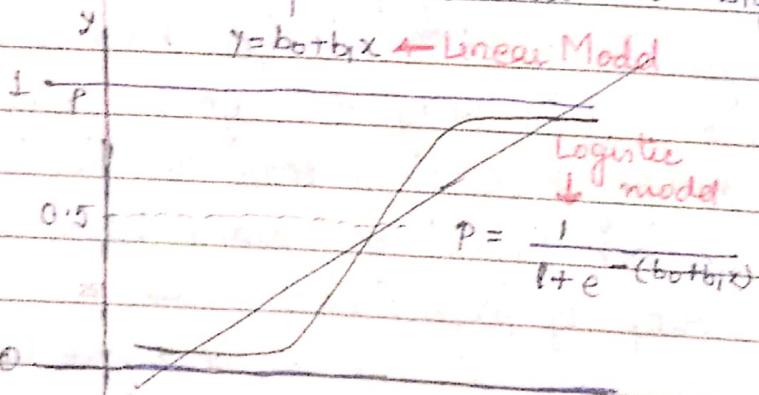
Therefore we can transform our linear regression to a logistic regression curve and logistic regression curve can only go b/w 0 and 1

SIGMOID FUNC^N OR LOGISTIC → takes in any value and outputs it to be between 0 and 1.



Q&A - Logistic regression function is going to be key to understanding using logistic regression to perform a classification

IMP This means we can take our Linear Regression Solution and place it into the Sigmoid funcⁿ



If we take linear model and put it into the sigmoid funcⁿ, we are able to transform linear regression to a logistic model.

* We can set a cutoff point at 0.5.

$$P < 0.5 \Rightarrow 0$$

$$P > 0.5 \Rightarrow 1$$

However, in some cases 0.5 can't be right cut off.
(In imbalance data)

Model evaluation

CONFUSION MATRIX — is used to evaluate classification models. Model performance is evaluated on some test data for which the true values are already known.

$n=165$	PREDICTED: NO		PREDICTED: YES	
	ACTUAL: NO	$TN = 50$	$FP = 10$	$TYPE I ERROR$
ACTUAL: YES	$FN = 5$	$TP = 100$	105	110

Ex:- Test for presence of disease

NO = negative test = False = 0

YES = positive test = True = 1

$$\text{ACCURACY} = \frac{TP + TN}{TOTAL} = \frac{100 + 50}{165} = 0.91 = 91\% \text{ accurate.}$$

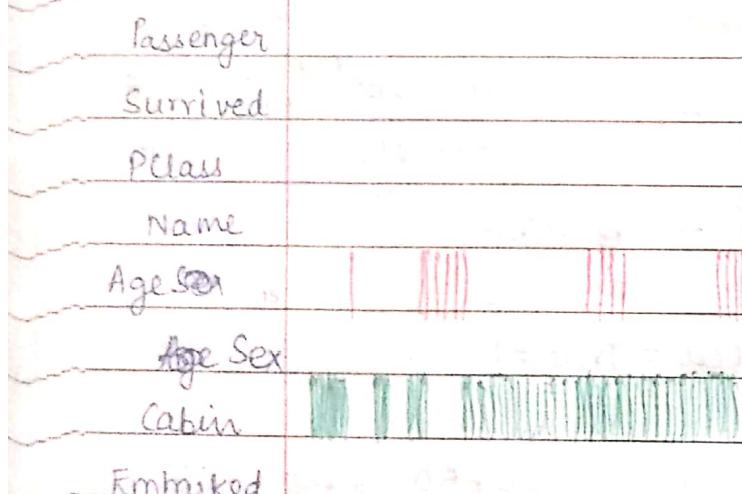
$$\text{MISCLASSIFICATION RATE (ERROR RATE)} = \frac{FP + FN}{TOTAL} = \frac{15}{165} = 0.09 = 9\% \text{ misclassification or error rate}$$

```

import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import seaborn as sns
%matplotlib inline
train = pd.read_csv('titanic_train.csv')
train.head()
train.isnull() → False if it is not null.  

                           True if it is null.
sns.heatmap(train.isnull(), yticklabels=False,
             cbar=False, cmap='viridis')

```



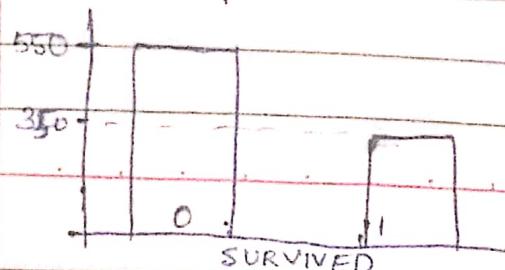
20% age data missing
Red indicates
TRUE, meaning
null values

Age data can be imputed by using the knowledge of other columns to fill in reasonable values for that age column.

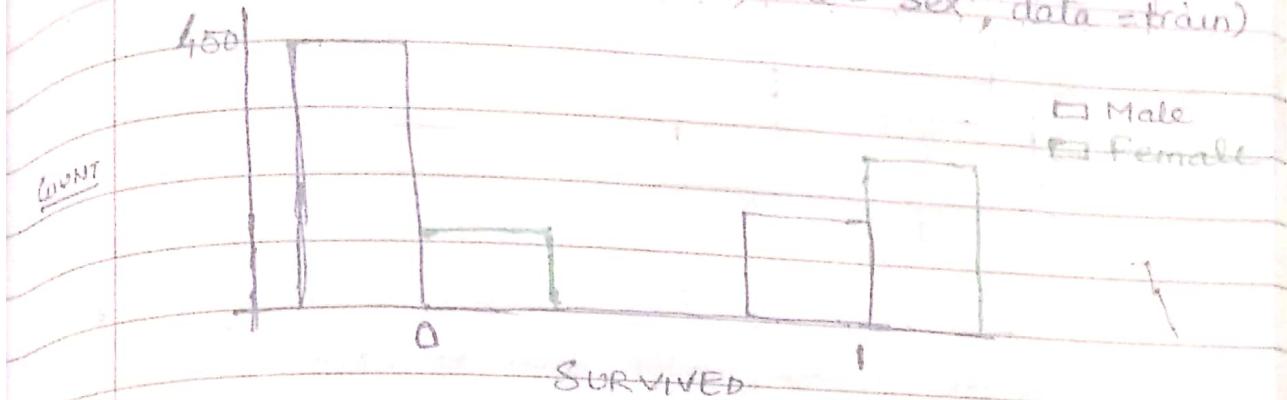
Cabin ~~as~~ column - missing too much of ~~data~~ to do something useful with it. we will probably drop this later or change it to 1 or 0.

sns.set_style('whitegrid')

sns.countplot(x='Survived', data=train)



Camlin Page
Date: _____
sns.countplot(x='Survived', hue='Sex', data=train)



Distribution Plot → Idea of the age of people

sns.distplot(train['Age'].dropna(), kde=False,

bins=30)



or train['Age'].plot.hist(bins=35)

sns.countplot(x='SibSp', data=train)

train['Fare'].hist(bins=40, figsize=(10, 7))

import cufflinks as cf

cf.go_offline()

train['Fare'].iplot(kind='hist', bins=50)

using
interactive
APIs.

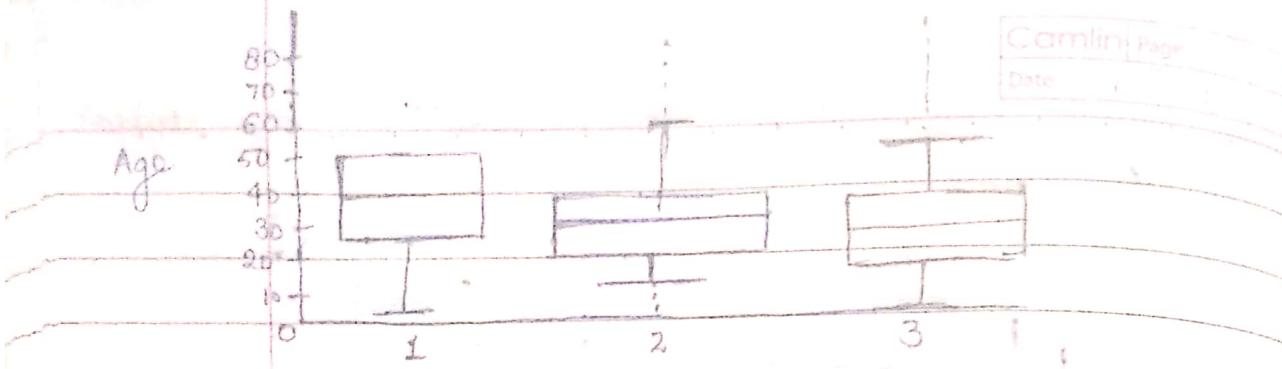
Filling missing data for age column

① Filling in the mean age of all the passengers. This is known as **Imputation**.

② Checking the average age passenger class.

plt.figure(figsize=(10, 7))

sns.boxplot(x='Pclass', y='Age', data=train)



P Class

These average ^{age} values can be used in order to impute the age based off of the past year data.

def impute_age(cols):

Age = cols[0]

Pclass = cols[1]

if pd.isnull(Age):

if Pclass == 1:

return 37

elif Pclass == 2:

return 29

else:

return 24

else:

return Age

```
train['Age'] = train[['Age', 'Pclass']].apply(impute_age, axis=1)
```

sns.heatmap(train.isnull(), yticklabels=False, cbar=False, cmap='viridis')

Passenger

Survived

Pclass

Name

Sex

Age

SibSp

Parch

Ticket

Fare

Cabin

Embarked

In cabin column, there are too many missing points / information. Therefore, drop stat (cabin column).
`train.drop(['Cabin'], axis=1, inplace=True)`

Now, there is only 1 missing value in Embarked.
Dropping any more missing value

`train.dropna(inplace=True)`
we filled few missing values and dropped few
Converting categorical features into dummy variables
using Pandas.

Sex & Embarked → Categorical,

`pd.get_dummies(train['sex'])`

female male:

0	0	1
1	1	0
2	1	0
3	1	0
4	0	1
5	0	1

If ML algorithm is fed both columns, it will know that if it's zero female it can predict perfectly it's going to be 1 male. This is known as multi-collinearity and will mess up the

algorithm. Therefore,

`Sex = pd.get_dummies(train['sex'], drop_first=True)`

`embark = pd.get_dummies(train['Embark'], drop_first=True)`

`train = pd.concat([train, sex, embark], axis=1, inplace=True)`

* Ticket and Name column are not much useful, but we can do feature engineering in features. However, in this scenario we will go ahead and drop these columns.

`train.drop(['Sex', 'Embarked', 'Name', 'Ticket'], axis=1, inplace=True)`

`train.drop(['PassengerId'], axis=1, inplace=True)`

~~X = train.drop(['Survived', axis=1])~~

~~y = train['Survived']~~

~~from sklearn.cross-validation import train_test_split~~

~~X-train,~~

~~X-test,~~ = ~~train-test split(x,y, test_size=0.3, random_state)~~

~~y-train,~~

~~y-test~~

~~from sklearn.linear-model import LogisticRegression~~

~~logmodel = LogisticRegression()~~

~~logmodel.fit(X-train, y-train)~~

~~predictions = logmodel.predict(X-test)~~

~~from sklearn.metrics import classification_report~~

~~print(classification_report(y-test, predictions))~~

	precision	recall	f1-score	Support
0	0.80	0.91 (148/163)	0.85	163
1	0.82	0.65 (68/104)	0.73	104
avg / total	0.81	0.81	0.80	267

~~from sklearn.metrics import confusion_matrix~~

~~confusion_matrix(y-test, predictions)~~

~~array([[148, 15],~~

~~[36, 68]])~~

		PREDICTED		
		0	1	
ACTUALS	0	148	15	163
	1	36	68	104
		184	83	267

$$\text{Accuracy} = \frac{148+68}{267} = \frac{216}{267} = .8089$$

$$\text{precision of } (1) = \frac{68}{83} = .82$$

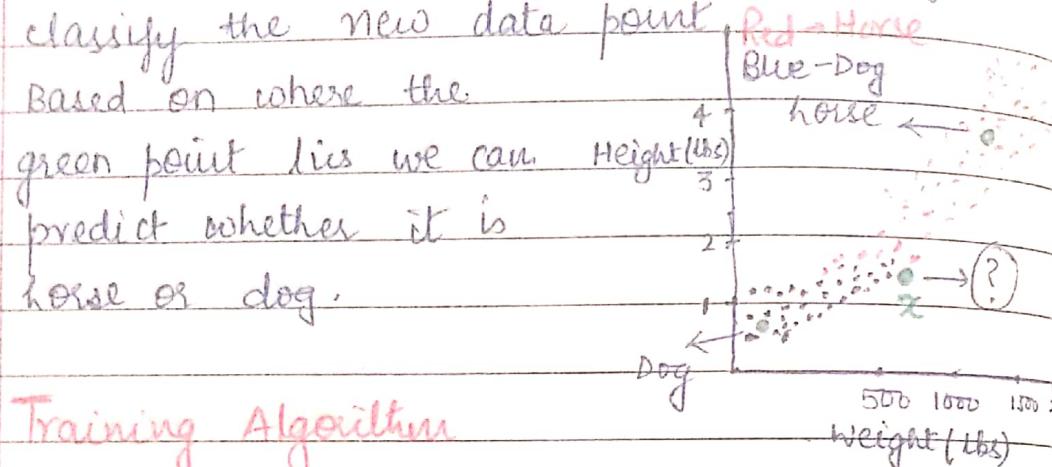
$$(0) = \frac{148}{184} = .80$$

KNN Theory (Chap 4 of Statistical Learning)

→ Classification algorithm.

Ex: We want to know if we get an animal where we only know the height and weight, we want to predict if that new data point is a horse or a dog. We are trying to classify the new data point.

Based on where the green point lies we can predict whether it is horse or dog.



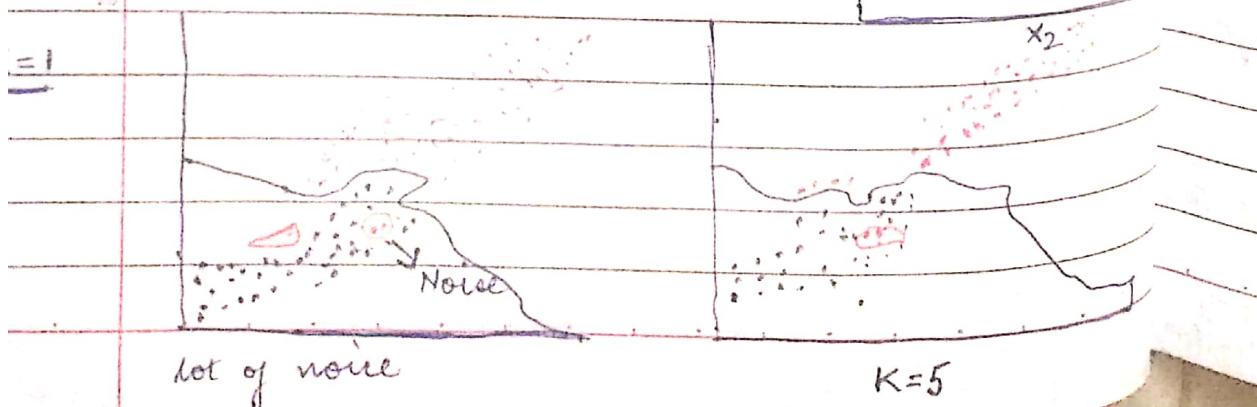
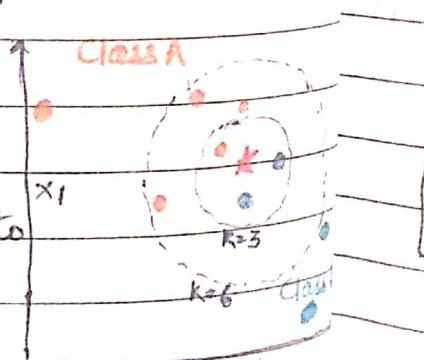
Training Algorithm

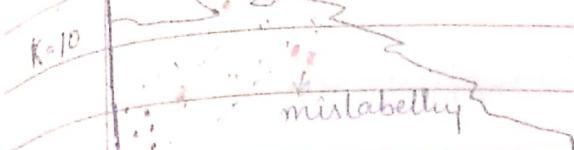
Store all the data.

Prediction Algorithm

1. Calculate distance from x to all points in your data
2. Sort the points in your data by increasing distance from x .
3. Predict the majority label of the "k" closest points.

Choosing k will affect what class a new point is assigned to.





K=50

Camlin Page

Date:

As with increase in K values, we notice it smooth out and create more of a bias in the model or get a clear cut off at the cost of mislabeling some points.

PROS

1. Very simple
2. Training is virtual (sort data)
3. Works with any number of classes
4. Easy to add more data.
5. Few parameters
 - K
 - Distance Metric (Dist. b/w test point and training point)

CONS

1. High Prediction Cost (worse for large data sets)
2. Not good with high dimensional data
3. Categorical features don't work well (as we can't calculate distance)
 - a) as you have to sort entire large data set.
 - b) High prediction cost.
 - c) Testing out what points are closest to data point.
- Throw off ability to measure distances in various dimensions.

```
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns
```

% matplotlib inline

```
df = pd.read_csv('Classified Data', index_col=0)  
df.head()
```

WIT PTI E&T SBI L&E SWG FDJ PIF JHGFI NRI

	WIT	PTI	E&T	SBI	L&E	SWG	FDJ	PIF	JHGFI	NRI	TARGET CLASS
0	1	0	0	0	0	0	0	0	0	0	0
1	0	1	0	0	0	0	0	0	0	0	1
2	0	0	1	0	0	0	0	0	0	0	1

Standardizing everything to same scale

```
from sklearn.preprocessing import StandardScaler
```

```
scaler = StandardScaler()
```

```
scales = scaler.fit(df.drop('TARGET CLASS', axis=1))
```

→ scaled_features = scales.transform(df.drop('TARGET CLASS', axis=1))

→ perform standardization by centering and scaling.

→ gives scaled data

```
df_feat = pd.DataFrame(scaled_features,
```

```
columns = df.columns[:-1]
```

→ standardized data

```
from sklearn.cross_validation import train_test_split
```

```
X_train = df_sealed_df_feat
```

```
y = df['TARGET CLASS']
```

```
X_train, X_test, y_train, y_test
```

```
= train_test_split(X, y, test_size=0.3)
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
Knn = KNeighborsClassifier(n_neighbors=1)
```

```
Knn.fit(X_train, y_train)
```

```
KNeighborsClassifier(algorithm='auto', leaf_size=3,
```

```
metric='minkowski'
```

```
metric_params=None, n_jobs=1,
```

```
n_neighbors=1, p=2, weights='uniform'
```

`pred = knn.predict(x-test)`

evaluation of knn model,

`from sklearn.metrics import classification_report,`

`confusion_matrix`

`print (confusion_matrix(y-test, pred))`

`print (classification_report(y-test, pred))`

`[[151 8]`

`[15 126]]`

	precision	recall	f1-score	support
0	0.91	0.95	0.93	159
1	0.94	0.89	0.92	141
avg/total	0.92	0.92	0.92	300

Elbow method to choose k value

Error rate is an empty list. We will iterate many models using many different values and plot out their error rate and see which one has the lowest error rate.

`error_rate = []`

`for i in range(1, 40):`

`knn = KNeighborsClassifier(n_neighbors=i)`

`knn.fit(x-train, y-train)`

`pred_i = knn.predict(x-test)`

`error_rate.append(np.mean(pred_i != y-test))`

`plt.figure(figsize=(10, 6))`

`plt.plot(range(1, 40), error_rate)`

`color='blue', linestyle='dashed',`

`marker='o', markerfacecolor='red',`

`markersize=10)`

`plt.title('Error Rate vs K Value')`

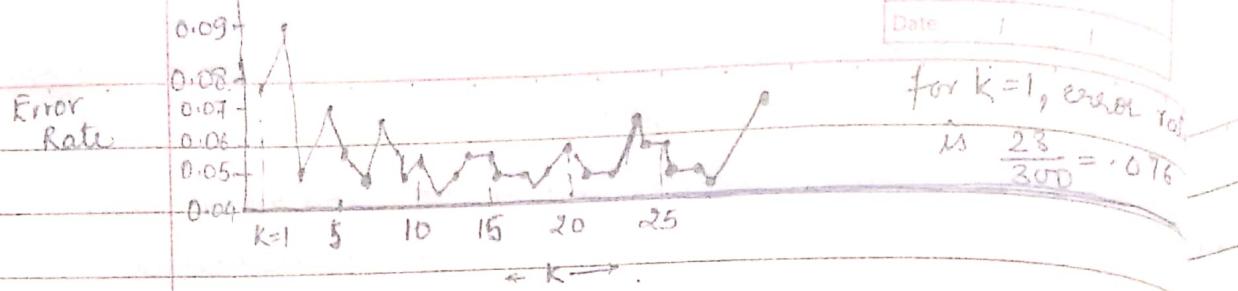
`plt.xlabel('K')`

`plt.ylabel('Error Rate')`

Error Rate vs K Value

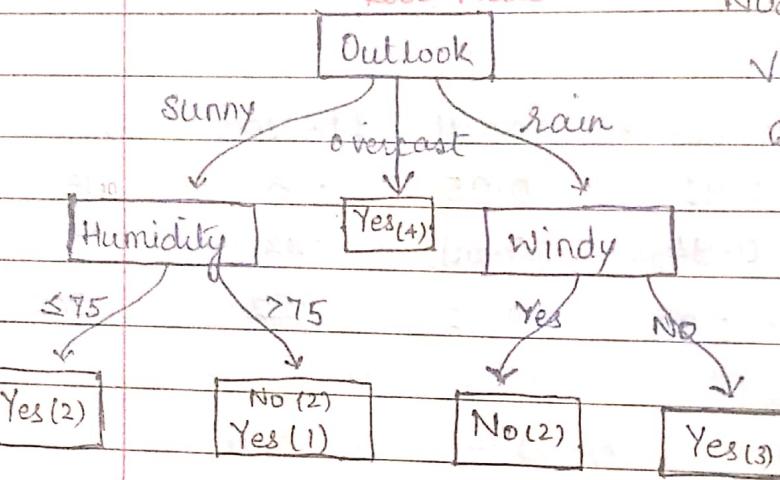
Camlin Page

Date / /



DECISION TREE (Chapter 8)

Root Node



Nodes → Split for the value of a certain attribute.

Root → The node that performs the first split

Leaves → Terminal nodes that

predict the outcome

- Decision trees don't have the best predictive accuracy. This is partially due to the high variance meaning that different splits in the training data can lead to very different trees.
- Bagging is a general purpose procedure for reducing the variance of ML method.

Random Forests → used for very fast classification model
 To improve performance, we can use many trees with a random sample of features chosen as the split.

1. A new random sample of features is chosen for every single tree at every single split.
2. For classification, $m = \sqrt{P}$, where $p = \text{features}$
3. Using 'one very strong feature' in "bagged trees" as top split, will result in an ensemble of similar trees that are highly correlated.
4. Averaging highly correlated quantities does not significantly reduce variance.
5. However, if we leave some features from each split, Random forests "decorrelates" the trees. This will result in reducing the variance of the resulting model by averaging process.

import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

%matplotlib inline

df = pd.read_csv('kyphosis.csv')

df.head()

	Kyphosis	Age	Number/Start
0	absent	71	3 5
1	absent	158	3 14
2	present	128	4 5
3	absent	2	5 1
4	absent	1	4 15

Number - is a number of vertebrae involved in the operation and the

start - Number of the first or top most of vertebrae that was operated on.

`df.info()` gives 81 entries

`sns.pairplot(df, hue = 'Kyphosis')`

from `sklearn.cross_validation import train-test`

`x = df.drop('Kyphosis', axis=1)`

`y = df['Kyphosis']`

`x_train, x-test, y-train, y-test = train-test(x, y, test-size=0.3)`

from `sklearn.tree import DecisionTreeClassifier`

`dtree = DecisionTreeClassifier()`

`dtree.fit(x-train, y-train)`

`predictions = dtree.predict(x-test)`

from `sklearn.metrics import classification_report`
`confusion_matrix`

`print(confusion_matrix(y-test, predictions))`

`print(classification_report(y-test, predictions))`

15 [E17 5] → mistabelling 6 from test set

[1 2]

precision recall f1-score support

	precision	recall	f1-score	support
absent	0.94	0.77	0.85	22
present	0.29	0.67	0.40	3
avg / total	0.67	0.76	0.80	25

	precision	recall	f1-score	support
absent	0.94	0.77	0.85	22
present	0.29	0.67	0.40	3
avg / total	0.67	0.76	0.80	25

	precision	recall	f1-score	support
absent	0.94	0.77	0.85	22
present	0.29	0.67	0.40	3
avg / total	0.67	0.76	0.80	25

Comparing above result with Random forest

from `sklearn.ensemble import RandomForestClassifier`

`rfc = RandomForestClassifier(n_estimators=200)`

`rfc.fit(x-train, y-train)`

`rfc-pred = rfc.predict(x-test)`

`print(confusion_matrix(y-test, rfc-pred))`

`print(classification_report(y-test, rfc-pred))`

[1 2]

precision recall f1-score support

	precision	recall	f1-score	support
absent	0.95	0.86	0.90	22
present	0.40	0.67	0.50	3
avg / total	0.88	0.84	0.86	25

	precision	recall	f1-score	support
absent	0.95	0.86	0.90	22
present	0.40	0.67	0.50	3
avg / total	0.88	0.84	0.86	25

little hard to tell which performed better because it really depends whether we value precision or recall or being present or absent.

It is more important to realize that it is still present vs absent. But it again really depends on the situation and the patient itself and what costs are associated with those decisions.

However, if we compare based on confusion matrix, overall it looks like Random Forest a little better.

Also, as data sets gets larger and larger RF will always do better than single decision tree.

In above ex, dataset was too small therefore RF was not able to outshine decision tree as much.

- df['Kyphosis'].value_counts()

absent 64 // → dataset is unbalanced and
present 17. // that also affect models

TREE VISUALIZATION USING DT

Scikit learn has some built in visualization capabilities for decision trees. However, we they are not used as most often RF are used which are difficult to visualize as there are 100s of trees Vs 1 in DT.

```
from IPython.display import Image  
from sklearn.externals.six import stringIO  
from sklearn.tree import export_graphviz  
import pydot → is not included in anaconda dist.
```

```
features = list(df.columns[1:])
```

```
features
```

```
→ ['Age', 'Number', 'start']
```

additional library that will take a
graph object and create picture

dot_data = StringIO()

export_graphviz(dtrees, out_file=dot_data,

feature_names=feature,

filled=True, round

Go to
www.graphviz.org
to install

graph = pydot.Dot.from_dot_data

(dot_data, get_value())

Image(graph[0].create_png())

↳ require graphviz lib

First choice for very fast classification model
to see what's kind of a baseline accuracy or
precision or recall you can get a model
before using other models or tuning stuff.

Decision tree classifier API

fit(X, y) → build decision tree

predict(X) → predict class for X

Decision tree classifier API

fit(X, y) → build decision tree

predict(X) → predict class for X

Decision tree classifier API

fit(X, y) → build decision tree

predict(X) → predict class for X

Decision tree classifier API

fit(X, y) → build decision tree

predict(X) → predict class for X

Decision tree classifier API

fit(X, y) → build decision tree

predict(X) → predict class for X