

## OHVT

### Unit 1

#### Data

raw or isolated facts from which req info is produced

#### Data science

- art of uncovering insights & trends in data
- uses scientific methods, processes, algorithms and systems to extract knowledge & insights from many structural and unstructured data.

#### problems solved

- Classification : supervised concept, categorizes data into classes
- Regression - method
  - predict continuous outcome ( $y$ ) from one or more  $x$
- Anomaly detection - obs that differs majorly from rest of data
- Recommendations: offering of relevant suggestions
- Clustering : grouping of data pts having similar attributes
- Reinforcement :-

#### Industries

- Banking / Finance - fraud detection, customer segmentation, etc
- Healthcare - disease detection, genetic analysis, survival analysis
- Retail & E-commerce Sales forecasting, profit prediction etc
- Telecommunications network congestion detection etc.
- Energy & Utilities smart grid security, theft detection, failure modeling
- Manufacturing - demand forecasting, quality assurance etc.

## Data Science Lifecycle

→ Business Understanding → ask relevant questions  
→ clear objectives

→ Data Mining → gathering & extracting data

→ Data Cleaning → missing data, duplicate data etc.  
missing data :- Nan data

Duplicacy → duplicate data → row / column

Inconsistent types → numeric as string

Irrelevant

Outliers

Formatting (wrong letter case)

- Data Exploration
- Feature Engineering
- Predictive Modelling
- Data Visualization

### Categories of Data

Quantitative	discrete	Nominal
	cannot be represented into decimals, countable	
	continuous :- only measured	
	can be represented in decimals	
	interval & ratio	
	numeric	
	order	
	difference	
Qualitative	Ordinal	Nominal
	Order data	
	Logical ops	
	$S < M$	
	letter grades	
	(A, B, C)	
How & why	Marital Status	label data
	(Married, Single)	

## Data based on structure

### → Structured

- has a well defined structure
- follows a consistent order
- should be accessed easily

### → Semi-Structured

- has some structure
- not in RDBMS
- lacks rigid schema
  - JSON

### → Unstructured

no structure

text / images / videos

e.g. document

### → Graph data

nodes, edges, properties used to represent & store data  
e.g. linkedin connections

### → Streaming data

continuously generated by diff sources  
req. continuous processing  
youtube, google meet

## → Data Collection

process of collecting, measuring & analyzing diff. types of information using a set of standard validated techniques

### → Primary data collection

data collected from first hand sources is best

#### → interviews

ask que. → respond

high flexibility

#### → observations

evaluate behaviour of people in ob controlled & uncontrolled situations

e.g. observing random people & walking their path to open store

#### → Surveys & Questionnaires

→ broad perspective from large groups

↓ det.

#### → focus groups

similar to interview but in groups

#### → Oral History

collecting opinions & personal experiences of people for a event they were involved in

## → Secondary data collection

data that has already been collected by someone else

#### → Internet

large pool of free & paid resources

just & easy

should only source from authentic sites

#### → Govt archives

→ authentic & verified

→ not always readily available

## → Libraries

- researchers have copies of their research
- imp of authentic info
- also serve as storehouse for business directory, annual reports etc.

## → Data Preprocessing

data mining technique used to transform raw data into useful & efficient format

## → Data Cleaning + Data Transformation + Data Reduction

### → Data Cleaning

- data can have irrelevant and missing parts
- Missing data
  - Ignore the tuples  
only suitable on large dataset, multiple values missing
  - Fill the missing values
    - manually
    - attr. mean
    - most probable value

### → Noisy Data

meaningless data that can't be interpreted  
due to faulty data collection, data entry errors etc

#### → Binning

works on sorted data

divide into segments of equal size → each segment is handled individually

#### → Regression

smooth by fitting in regression model

#### → Clustering

similar data in cluster

outliers outside the cluster

## → Data Transformation

transform into appropriate forms

### → Normalization

In order to scale data values in specified range

$$(-1 - 1) / (0 - 1)$$

### → Attribute Selection

new attr. are constructed from the given set of attr.

### → Discretization

done to replace ~~raw~~ values of numeric attr. by interval levels.

### → Concept Hierarchy Generation

attr. converted from lower level to higher level in hierarchy. eg. city → country

## → Data Reduction

- Analysis is harder when huge amount of data  
aims to inc storage efficiency, reduce data storage & analysis cost

### → Data Cube aggregation

### → Attribute Subset Selection

highly relevant attr. should be used; rest discarded

eg: attr. having p-value > significance level can be discarded

### → Numerosity Reduction

enables to store data model instead of whole data.

eg: Regression model

### → Dimensionality reduction

reduce size of data by encoding

can be lossy / lossless

### → Wavelet transforms

### → PCA (Principal Component Analysis)

## → Data Discretization

Converting attributes values of continuous data into finite set of intervals with min data loss

→ supervised discretization - class data is used

- depending how operation proceeds eg. top down / bottom up

eg. Age :- 1, 5, 9, 4, 7, 11, 13, 18, 17, 14, 19, 31, 33, 37



Attribute Age

1, 5, 4, 11, 14, 17, 31, 33

9 13

After

Discretization - Child      Young      Mature

→ Histogram analysis

→ Binning

→ Cluster Analysis

dividing the vals of n numbers into clusters to isolate a computational feature of n

→ Decision tree Analysis

top down slicing is used

Select attr with least entropy, run recursion on it

→ Correlation analysis

discretizing by Linear regression

get the best neighbouring interval

then large intervals are combined to develop a larger overlap to form final intervals

21/9/22

## DHVT

## Unit - 2

## → Descriptive Statistics

- Statistics :- concerned with describing, interpretation and analyzing of data.
- descriptive → uses analytical methods which provide math to model inferential and predict variation.
- used graphical model to help making numbers visible for communication.

Descriptive → methods of describing data characteristics of data set.

- e.g.: describing weight of product in production line.
- describing, summarizing, organizing + graphical displays

Outlier - data point i.e. significantly,  $>/<$  than other data pts in data set.

- may affect calculation of ds.
- detected using graphing the data
- can be normal or may indicate experimental error or incorrect data. (need to decide if we need to exclude them)

Measures used to describe dataset

## → Measures of position

measure central data tendency (when data is centered)

→ Mean

→ works well when distribution is symmetric & no outliers

→  $\bar{x}$  (sample)

$\mu$  (popn)

### → Median

- half data above it, half below
- reduces effect of outliers
- when data is non-symmetrical
- values should be ordered

### → Mode

- more useful to distinguish b/w unimodal & multimodal distributions (when data has more than one peak)

### → Measures of Spread

- spread : how data deviates from position measure
- gives an indication to amount of variation

### → Range (R)

- highest - lowest values
- can be misleading when data is skewed / presence of outliers
- does not make full use of data

### → Standard Deviation

avg distance of data points from their own mean

- low sd - clustered around mean

- high sd - widely scattered around mean

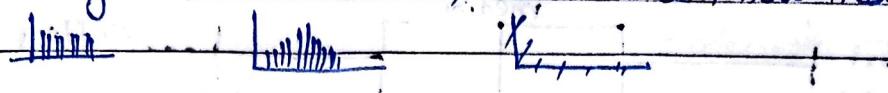
$$s \text{ (sample)} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$\sigma \text{ (popn)}$$

### → Measures of Shape

- plot for distribution
- data will always follow some distribution
- may be symmetrical or non-symmetrical

eg: uniform, normal, camel-back, bow-tied



→ helps us identify which descriptive statistics are more useful in a situation

→ symmetrical? use any (mean = median) →  
skewed? (median)

### → Skewness

- describes if data is distributed symmetrically around the mean.

- -ve → left skewed  
+ve → right skewed

### → Kurtosis

degree of flatness (or peakedness)

if clustered around middle → more peaked

if spread evenly → less more flat

Extra info → IQR (Inter Quartile range)

- divides data into 4 pts. (each - 25%)

- IQR contains middle 50% data

- used when data is not normally distributed.

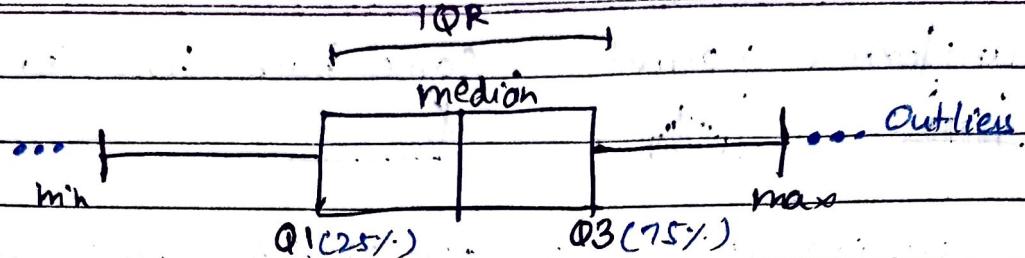
-

### → Box Plot

- chart that shows data from a 5-number summary

- used to show skewness / outliers

- used to find spread.



Minimum -

First Quartile (Q1) :- median of lower half

Median - also considered Q2

Third Quartile (Q3) : median of upper half

Maximum

+ve skewed :- if  $\text{max} - \text{median} > \text{median} - \text{min}$

-ve skewed :- " " " < "

## PIVOT TABLE

df.pivot(index=row, columns=column)

Pivot table is used to summarize and aggregate data inside DataFrame.

→ df.pivot\_table(index=row, columns=col)

(avg) →

aggfunc = "sum"

Pivot Charts  $\xleftarrow{\text{visualisation}}$  Pivot table

summarize large amount of data.

— categories.

— drilling down to result from summary

— filtering, sorting etc to focus on info you want.

— provide, concise & printed reports.

\* Exploring data

\* Change form layout field arrangement

\* change layout of columns, rows & subtotals

\* Change display of blanks & errors

\* change format

## → Heat Map

- graphically representing numerical data where each data pt. is indicated using colors:
- usual → warm-cool
- :: (high data    (low data  
     ::      phs)      phs)
- eg: heat map on website interaction  
     red shows high interaction on that portion.  
     green shows low
- scroll map → based on scrolls
- click map
- mouse tracking map
- useful for cross examining multivariate data
- good for showing variance across multiple variables, generalizing patterns, detecting correlations
- legend is required for it to be successfully read
- categorical data - color coded
- numerical data - requires colour scale
- can be used to show changes over time

## Correlation

- statistical measure that expresses the extent to which two variables are linearly related. common tool for describing simple relationships w/o making a statement of cause of effect.
- -1 to +1
- $r$  (Sample correlation coeff) quantifies strength of relationship
- can't look outside the two variables
- doesn't tell us about cause of effect
- cannot accurately describe non-linear relationships

p-value - we can't conclude that the pop correlation coeff is likely correlation number is  $r$

$r$  - measure of probability used for hypothesis testing

- 0 - weak relationship
- +ve - both variables tend to increase together
- ve - one variable tends to increase while other decrease

## → ANOVA

- Analysis of variance
- make multiple comparisons of several pop means
- instead of pairwise; looks simultaneously
- look for 2 variations - variations b/w sample means; variation within each of our samples
- F-statistic  $\left( \frac{\text{variation b/w samples}}{\text{variation within each sample}} \right)$

## Steps

- sample means for each of our samples as well as mean for all of the sample data
- calculate sum of squares of error  $\sum (x_i - \bar{x})^2$ : (SSE)
- calculate sum of squares of treatment
- square deviation of each sample mean from overall mean.
- multiply it by no. of samples - 1
- SST
- calculate degrees of freedom
  - overall:  $n-1$  (no. of data pts. - 1)
  - treatment:  $m-1$  (no. of samples used - 1)
  - error:  $n-m$

→ calculate mean square of error (MSE)

$$MSE = \frac{SSE}{n-m}$$

→ calculate mean square of treatment (MST)

$$MST = \frac{SST}{m-1}$$

$$F \text{ statistic} = \frac{MST}{MSE}$$

→ Linear Regression

→ supervised machine learning algorithm for predictive analysis  
 → linear approach to modelling the relationship b/w a response (dependent) variable & one or more explanatory variables (independent variables)

Simple Linear regression

→ single independent variable & corresponding target variable

$$\hat{y} = \beta_0 + \beta_1 x$$

dependent variable      | intercept      | slope

→ multiple linear regression

2 or more independent, 1 target variable

can be continuous / categorical

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

assumptions about data

- homogeneity of variance - size of error does not change significantly across values of idv
- independence of observations - idv should not be correlated
- normality - data should follow normal distribution

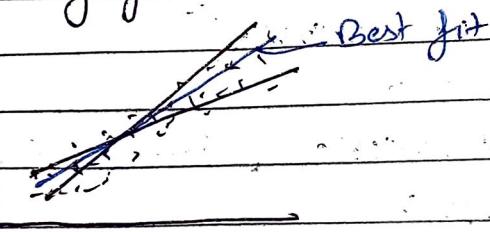
Applications:- predicting crop yields based on rainfall marks scored by student based on no of hours study

Correlation	Linear Regression
→ specify association b/w two var.	defines numerical relation b/w idv & dv.
→ not necessary for variables to be different.	variables are different
→ describes linear relationship	
→ Objective is to find numerical value to express relationship	finds the best fit line on data. Objective is to find random value based on fixed data.

### Simple Linear regression

$$f(n) = \beta_0 + \beta_1 n$$

determination of fit



→ to find coeff of line

↳ use least squares method

$$\text{RSS} - \text{residual sum of squares} \quad e = \text{actual} - \text{predicted}$$

$$= e_1^2 + e_2^2 + \dots + e_n^2$$

Line having min RSS is our best fit for data

### multiple linear regression

→ multiple predictors available

$$f(n_1, n_2, \dots, n_p) = \beta_0 + \beta_1 n_1 + \beta_2 n_2 + \dots + \beta_p n_p$$

Intercept

Coefficients

Predictors

## Qualitative predictors

one-hot encoding to convert qualitative data to continuous data : eg. female = 1 ; female = 0 ; male = 1

eg. Ans

mlr eg: eg of house prices on House age , total rooms

$$\text{eg: } \hat{y} = \beta_0 + \beta_1 \times \text{House age} + \beta_2 \times \text{total rooms}$$

$$\hat{y} = 6989.39 - 72.08 \times \text{House Age} - 0.27 \times \text{total rooms}$$

we get a plane

## Model Estimation

→ after rejecting null hypothesis , we would like to evaluate model done by R statistic  
 Root mean squared error

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

measures proportion of variability in observed response variable

## Assumptions:

- linear relationship b/w predictor & target
- no auto correlation in the value of predictors
- constant variance of errors (homoscedasticity)
- residuals must have normal distribution
- little or absence of collinearity among predictors (no multicollinearity)
- presence of outliers & leverage pts can significantly impact.

## Advantages of LR

- easier to implement and interpret the o/p coeff.
- less complex compared to other algorithms
- over fitting can be avoided using techniques such as dimensionality reduction etc.

## Disadvantages

- sensitive to outliers
- assumptions
- look only on relationship b/w mean of idv & dr.
- mean is not a complete desc of variable.
- LR is r. " " " of relationships among variables

## Model Evaluation using Visualization

why use regression plot?

gives us good estimate of

→ relationship b/w two variables

→ strength of correlation

→ direction of relationship (+ / -)

regression plot :- Scatterplot + fitted line

↳ Residual plot :- measure of how much a regression line vertically misses a data pt

residual values on vertical axis

dr on horizontal axis

→ randomly spread out mean

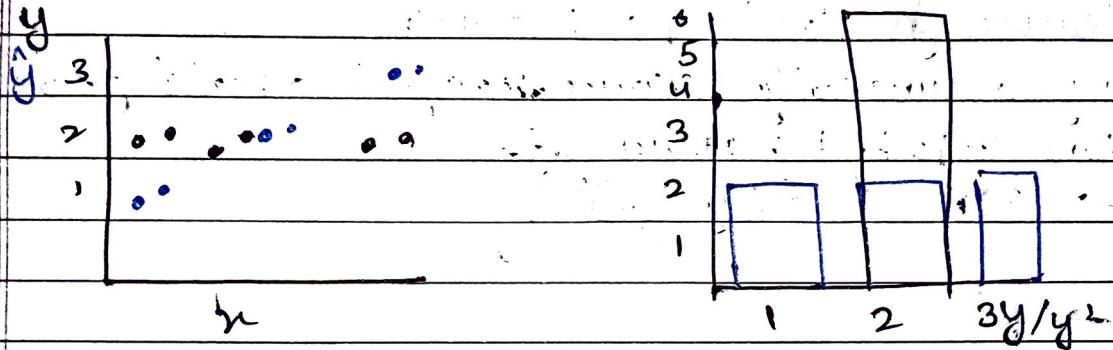
constant on axis

↳ This shows linear plot is off

## Incorrect models

### → Distribution plot

counts predicted value vs actual value



→ fitted values that result from model vs the actual values.

→ from this we can see model follows better than linear.

### → Polynomial regression f pipelines

→ when linear is not best fit

→ transform data into polynomial then use linear

→ describes curvilinear relations

↳ setting higher order terms of predictor variables

### → Quadratic

$$\hat{y} = b_0 + b_1 n_i + b_2 (n_i)^2$$

Cubic

$$\hat{y} = b_0 + b_1 n_i + b_2 (n_i)^2 + b_3 (n_i)^3$$

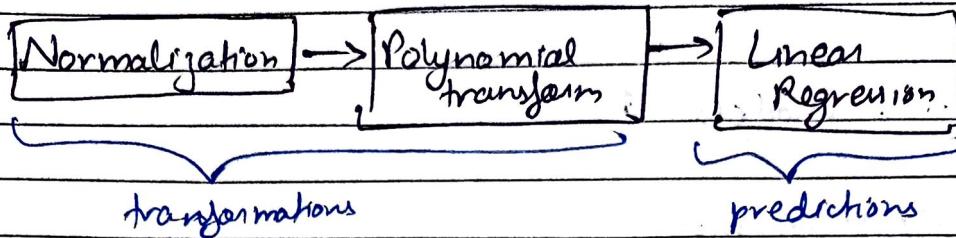
preprocessing

to normalize features when we have multilevel dimensions.

$$b_0 + b_1 n_1 + b_2 n_2 + b_3 n_1 n_2 + b_4 n_1^2 + \dots$$

## Pipelines

there are many steps to getting a prediction



→ Measures for In-Sample Evaluation

→ way to numerically determine how good a model fit on data

→ two important measures

$$\rightarrow \text{MSE} (\frac{\sum (y_i - \hat{y}_i)^2}{n})$$

$$\rightarrow R^2$$

→ coeff of determination

→ method to determine how close data is to a regression line

$$R^2 = 1 - \frac{\text{MSE of regression line}}{\text{MSE of avg. of data}}$$

→ closer to zero → good fit

→ closer to 1 → good fit

→ Prediction and Decision making

How to determine if our model is correct

Look at

→ do predicted values make sense  
→ visualization

→ numerical measures for evaluation

→ comparing b/w diff. models

eg :- Car price based on highway miles per gallon  
price =  $38212.3 - 821 \times \text{highway-mpg}$

→ make sense

↳ we might get negative values for price

→ visualization

Simply visualizing with regression

might show non linear behaviour.

→ using distribution plot

→ numerical measure

→ MSE might be large

→  $R^2$  might be closer to 0.

→ comparing

MSE /  $\text{f}(\text{y})$  rm

Comparing m, b, s, r, p, etc

4/12/22 DHV T

## Unit - 3 - Visualization

- graphical representation of information and data
- accessible way to see & understand trends, outliers & patterns in data.

### Advantages

- quickly see trends / outliers
- easy info sharing
- interactive explore opportunities

### Disadvantages

- easy to make inaccurate assumptions when large data.
- biased or inaccurate
- correlation doesn't always mean causation

Need? interactive, intuitive, easy to share, personalized

General Types? Chart, table, graph, geospatial, infographic, dashboards  
 Specific types? Area Map, Bar chart, Box Plot, Bullet graph, Gantt Chart  
 Heat Map, Histogram, Pie Chart

### 1. Graph Design Principles - for appealing & readable

→ White Space

→ empty space among the elements of graph

→ increase readability, focus the reader's attention

macro - spaces around main content

eg - space behind a figure

micro - space b/w digits, axis & text, b/w bars etc

→ Text

- helps reader understand the content
- eg :- axis labels, title, annotations regarding font size, font family, spacing, paragraph splitting

→ Colour

tool to convey emotions & sensations

characterized by - hue - pure colour, one wavelength  
 brightness - intensity of colour . amount of white light  
 Saturation  
 intensity, amount of shading

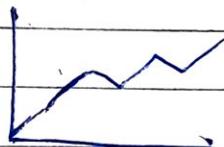
## 2 Graphs in Statistics

→ Bar Graph

- pictorial representation of grouped data , in vertical / horizontal rectangular bars
- length of bar  $\propto$  measure of data

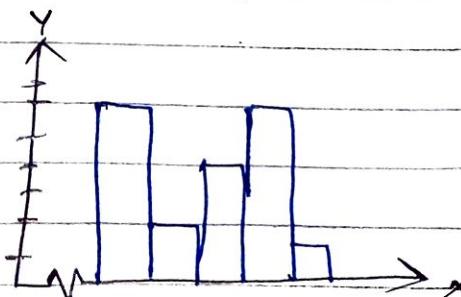
→ Line Graph

- utilizes point of lines to represent change over time
- line shows relation b/w points



→ Histogram

- displays freq of discrete & continuous data in a dataset using connected rectangular bars.
- no. of obs that fall into interval are represented as rectangular bar



### → Pie Chart

- used to represent the numerical proportions of dataset.
- dividing circle into sectors

↳ represents proportion of particular elements as whole.

### → Exponential Graph

representation of exponential functions

$$y = 3^n \quad - \text{Increasing}$$

$$y = 3^{-n} \quad - \text{decreasing}$$

### → Logarithmic graphs

inverse of exponential

$$y = \log_3 n$$

### → Trigonometric graphs

plotted for trigonometric functions



### → freq distribution graph

used to show freq outcomes in a particular sample

### 3 Interpreting data visualizations

- identify what information the chart is meant to convey.
- " info on each axis.
- " range covered by each axis
- look for patterns and trends.
- look for averages / exceptions
- look for bold / highlighted data.
- read specific data
- ensure data is credible

### Common Visualization mistakes

- truncated (shortened) Y axis - broken scale
- misleading cumulative graphs
- ignoring conventions
  - failed calculations
  - visualization type is all wrong
  - displaying too much data
  - try trying too hard to be original
  - making reader work too hard
  - Purposeful bias: deliberate attempt to influence data
    - data omissions / adjustments
- ? - Selective bias : slightly more discreet
- using percentage change in small sample ( $\frac{19}{20} = +5\%$ )
- correlation implying causation
- decoration
- ignoring popn size makes accurate comparisons impossible

#### 4. Graph Selection Matrix

jog google karlo

5. Correlation, bar & paired graph - Excel  
mutual connection b/w two or more sets of data

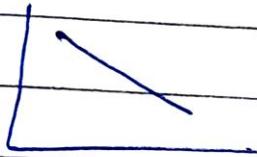
$$r = \frac{\text{covariance}}{SD(x) \times SD(y)} = \frac{n(\sum_{i=1}^n x_i y_i) - (\bar{x})(\bar{y})}{\sqrt{(n \sum_{i=1}^n x_i^2 - (\bar{x})^2)(n \sum_{i=1}^n y_i^2 - (\bar{y})^2)}}$$

= CORREL (array1, array2)

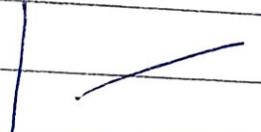
↓  
insert cell range for array (eg A2:A6)

Correlation chart

bivariate scatter plot



-ve, strong



pos, strong



zero

→ select bivariate data x & y in excel sheet

→ Insert → Scatter / Bubble chart

→ add a linear trendline to show correlation

can also change chart types, title, labels etc.

## 6. Integrity

graphical integrity refers to how accurately ~~extreme~~ accurately visual elements represent data.

info can vary widely, desire of tendency to

→ scale data disproportionately so that  
it fits

→ evenly spread values

can result in false impression

### 6 Principles - Edward Tufte

- representation of nos should match their true proportions
- labeling should be clear & detail
- designs should not vary from some ulterior motive.
- well known units for money
- # of dimensions should be same as data
- should not imply unintended context

7

### Visual Integrity

- what's presented accurately represents what's in the data
- no design choices distort/obfuscate the inherent facts analytical findings

→ include info on data provenance  
cite sources

do not affect credibility

→ clearly define data variables

representation makes sense if accompanying text explains what it involves.

prevent misinterpretation

- make sure the data being used is complete
  - eg:- omission of outlier X
  - manipulating axis thresholds X
  - leaving out dependent variables that are correlated to ones included in chart. X
  
- make sure data is consistent
  - present in a way that is correctly interpreted by viewers.
  - eg:- overlaying multiple axes on a single chart if including lines pegged to a different axis. is wrong
  
- be consistent on scale too
  - lie factor :- highlight a finding by scaling the results to make it look more prominent
  
- free from visual noise
  - icons etc
  
- don't filter out data so it can't be viewed
  - sometimes design of visualization, limits data can be viewed

## 7 VISUAL PERCEPTION

Ability to interpret the surrounding environment by processing information that is contained in visible light.  
 Why in visualization? to aid in good decision making  
 drawing insights

→ Visual perception is selective - selectively pay attention,  
 not see everything

→ Our eyes are drawn to familiar patterns:  
 We see what we expect to see.

→ Our working memory is very limited  
we can hold a very limited amount of information in our memory

Visual perception : act of seeing a visual or an image  
handled by ~~cortex~~ visual cortex  
fast & efficient

Cognition :- act of thinking, processing information, making comparisons etc  
in cerebral cortex  
slower & less efficient

→ data vis. shifts the balance b/w perception & cognition  
to use our brain's capabilities to its advantage  
↳ more use of perception, less of cognition

## 8. Memory

how human memory works? - Iconic, Working, Longterm  
interact w/ visualize memory remember

Iconic / Sensory memory

we see a visual, info remains in iconic memory for a tiny period of time.

we process & store info in this time only automatically

↓ preattentive processing

even before paying attention detects

several attributes

## \* preattentive processing

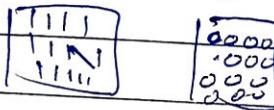
find/odd/containing obj's, etc

automatically occurs in brain prior to conscious awareness  
very fast < 250 ms

## feart feature

features: contrast to surroundings

eg color, orientation, size, motion etc



In graphs? lengths, colors etc

## Working memory / short term memory

We use when we are actually working with the visual info that is of interest to us in processed here.

• stays for about a min.

Capacity: 5-6 similar items

capacity can be increased by chunking

grouping similar items

chunking: when info is done in form of visuals that show patterns, more info can be chunked together

9

## Myths about data visualization

- we visualize because some people are visual learners
- " " for people who have difficulty understanding numbers.
- " " to grab attention with eye catching but inevitable less informative
- best data visualizers are trained in graphic arts
- provide best means of story telling contained in data

## Need of Visual Data

- seeing the big picture
  - ↳ overview
- easily & rapidly comparing values
- seeing patterns among values
- Comparing patterns

## 10 Types of Charts

### a) Change over time

show data over a period of time

use cases → stocks

health stats

chronologies

eg:- Line, bar, stacked bar, candlestick, area, timeline, horizon, waterfall

### b) category comparison

compare data b/w multiple distinct categories

use cases → income across countries

popular venue times

team allocations

eg bar, grouped bar, bubble, parallel coordinates, multi line, bullet

### c) Ranking

show item's posn in ordered list

use cases → election results

performance stats

eg. Ordered bar, ordered column, parallel coordinates

d) part-to-whole

how partial elements add to total

use cases? budgets

revenue

eg. stacked bar, pie, donut, stacked area, treemap,

sunburst

e) correlation

Show correlation

use cases :- income & life expectancy

eg. Scatterplot, bubble, column/line, heatmap

f) distribution

Show how often each value occurs in a dataset

use cases :- pop^n distribution

income distribution

eg Histogram, Box plot, Violin, Density

g) Flow

Show movement of data b/w multiple states

use cases :- fund transfers

vote counts election results

sankey, gantt, chord, network

h) relationship

Show multiple items relate to each other

use cases :- social networks

word clouds

eg network, Venn diagram, chord, sunburst

## Unit 5

### Generalization

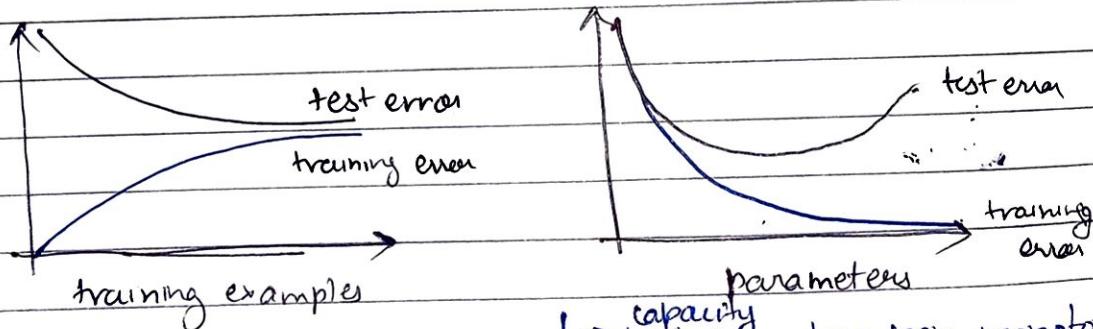
generalization assesses a model's ability to process new data and generate accurate predictions after being trained on a training set.

Overshooting / underfitting will prevent a model from (overfitting / underfitting) | generalizing

### Reasoning

overfitting - model may generalize unintentional / coincidental regularities in training data which are not present in test data.

effects with # of training examples, # of parameters



→ test error decreases as more examples helps in generalization / avoid coincidental regularities

→ training error inc. as small sets are easy to memorize

capacity  
→ test error has non monotonic influence; would like to create powerful enough to learn actual regularities but not powerful enough to exploit accidental regularities

## Measuring

optimising performance of training set is not enough

divide data :-

→ training set

collection of examples on which the network is trained.

→ validation set

fine tune hyperparameters like no. of hidden units & learning rate.

→ test set

evaluate generalization performance

→ Regularization

process of shrinking / regularizing the coeff towards zero.  
used to improve performance by imposing penalty in ML.

$L_1^{\text{norm}}$  → calculated by adding absolute values of vector

$L_2^{\text{norm}}$  → calculated by taking square root of the sum of the squared vector values.

Lasso Regression :- Coeff are penalized to the pt where

→ they reach zero in the Least Absolute Shrinkage & Selection Operator Regression

→ uses  $L_1$ .

→ comes in handy when lot of variables can be used as feature selection

Rd

# Regularization with modified loss functions

↳

## - LASSO (L1)

$$\text{minimize } \left\{ \sum_{i=1}^n \left( y_i - \sum_{j=0}^p w_j x_{0j} n_{ij} \right)^2 + \lambda \sum_{j=0}^p |w_j| \right\}$$

ordinary least squares

regularization term

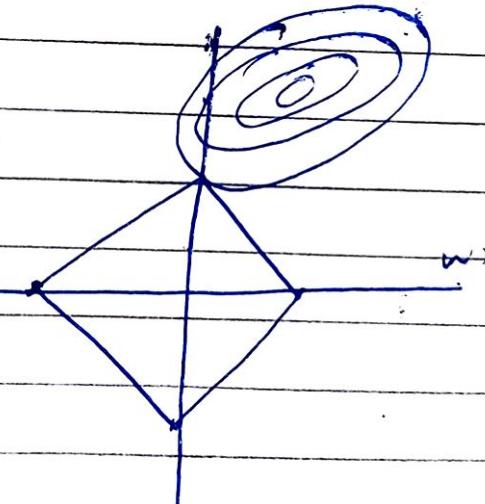
- penalizes regressors by shrinking their weight
- regressors that contribute little are more penalized
- $\lambda$  is weighting factor for regularization to tune output  $\Leftrightarrow$  undraft

→ reduces overfitting

→ eliminate insignificant features

→ selects only significant regressors

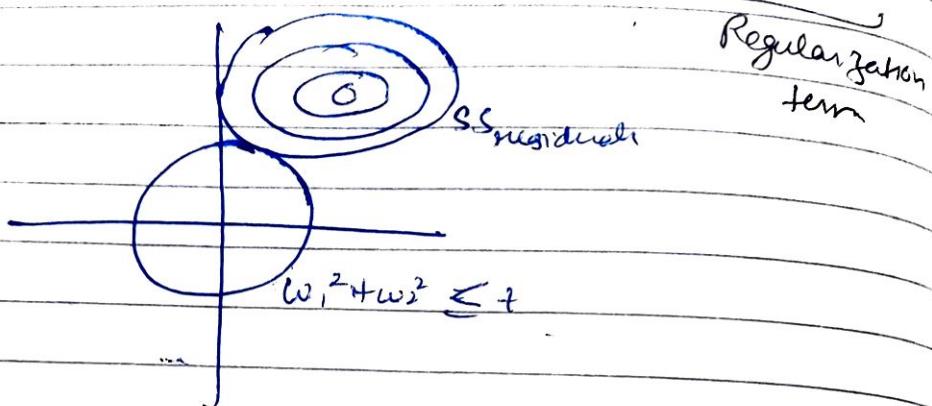
$w_2$



$$(\sum_i (y_i - (w_0 + w_1 x_1 + w_2 x_2))^2 + \lambda \sqrt{w_1^2 + w_2^2})$$

## Ridge (L2)

$$\text{minimize} \left\{ \sum_{i=1}^m \left( y_i - \sum_{j=0}^p w_j x_{i,j} \right)^2 + \lambda \sum_{j=0}^p w_j^2 \right\}$$



- L2 leads to near zero coeffs
- handles multiple correlated regression better

## Elastic Net Regularization

- Combines L1 & L2

each has its own weighting factor

~~Minimize~~ Minimize  $\left\{ \sum_{i=1}^n \left( y_i - \sum_{j=0}^m w_j x_{i,j} \right)^2 + \lambda_1 \sum_{j=0}^m |w_j| + \lambda_2 \sum_{j=0}^m w_j^2 \right\}$

- $\lambda_1$  &  $\lambda_2$  allow
- balance of - attribute elimination
- handling multiple correlated regression
- proper tuning

## Ridge Regression

- When variables in a model are multicollinear
- minimizes no. of inconsequential independent variables but does not eliminate them.
- uses L2.
- $L2 \text{ penalty} = \alpha \text{ square of beta coeffs}$

## Choosing right regularization

### → Ridge

- when all independent vars are included
- collinearity / dependency / codependency b/w variables

### → Lasso

multiple selections and we want feature selection

### → Elastic - net

lot of variables

## → Model Evaluation Metrics

In classification problems, we use 2 types of algos (depending on kind of o/p it creates) :-

→ Class output! - algos like SVM / KNN

probability o/p! - logistic regression, random forest, gradient boosting etc.

## 1. Confusion Matrix

$N \times N$  matrix

$N$  is no. of classes being predicted.

Generally used only with class o/p models

Accuracy:-

Actual Values

		+ve	-ve
predicted values	+ve	TP	FP
	-ve	FN	TN

TP :- True +ve - you predicted +ve if its true

TN :- true -ve - " " " " -ve " "

FP :- false +ve - Type 1 Error

You predicted positive and its false

FN :- false -ve - type 2 Error

You predicted negative but its false

Accuracy :- from all classes, how many were predicted correctly

$$\hookrightarrow \frac{TP+TN}{\text{total}}$$

Precision (true prediction value)

proportion of the cases that were correctly identified

$$\hookrightarrow \frac{TP}{TP+FP}$$

-> We use Sensitivity / Recall :- proportion of actual +ve cases which are correctly identified

$$\hookrightarrow \frac{TP}{TP+FN}$$

Specificity proportion of actual -ve cases which are correctly identified

$$\hookrightarrow \frac{FN}{FP+FN}$$

## 2) F1 Score

when we try to get best precision & recall  
 harmonic mean of " "

$$F_1 = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}}$$

$$F_1 = 2 \cdot \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

why not mean? because HM punishes extreme values more

e.g. precision = 0, recall = 1

$$\text{HM} = 0.5$$

$$\text{AM} = 0 \quad \checkmark$$

→ if we want to give weight to either one of them, we use

$$F_B = \frac{(1 + \beta^2) \text{precision} \times \text{recall}}{(\beta^2 + 1)(\text{precision} + \text{recall})}$$

## 3. Gain and Lift Charts

Concerned to check rank ordering of probabilities

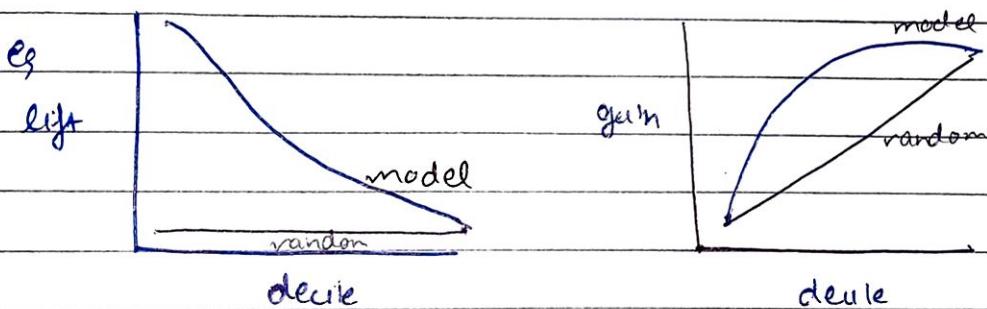
- calculate probability of each obs
- rank these probabilities in decreasing order
- build deciles with each grp having atleast 10% of obs.
- calculate response rate for each deciles - good / bad

- $\frac{\text{prediction}}{\text{without model}}$   
 To measure how much better our model works compared to without the model

- evaluates model performance in a portion of population
  - + decile - simply putting our data into 10 bins
  - we obtain no. of cases (no. of data in decile)  
no. of responses (no. of actual true data in decile)
- Gain - ratio b/w the cumulative number of the no. of responses up to each decile divided by total no. of the obs in data.

Lift - gain ratio percentage to the random percentage at a decile level

measure how much better we can expect to do with the predictive model compared to without it



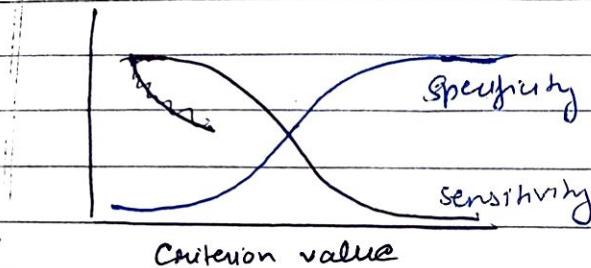
#### 4. Kolmogorov Smirnov chart - KS chart

- measure of degree separation b/w two g-rv distributions
- KS is 100, if scores partition into positives & negatives
- if model cannot diff b/w true g-rv, it is as if it selects randomly so KS = 0.
- KS lies from 0 - 100

## 5. Area Under the ROC curve (AUC - ROC)

ROC → Receiver Operating Characteristic

for a probabilistic model, we get a diff value for each member  
hence, for each sensitivity, we get a diff specificity.



ROC is the plot b/w sensitivity & 1-specificity

↓  
TPR

FPR

some thumb rules

0.9-1 - excellent (A)

$$\frac{TP}{TP+FN}$$

$$\frac{FP}{FP+TN}$$

0.8-0.9 - good (B)

0.7-0.8 - fair (C)

0.6-0.7 - poor (D)

0.5-0.6 - fail (F)

points to remember :-

- for a model which gives class as op, will be represented as single pt in ROC plot. such models cannot be compared with each other as judgement needs to be taken on single metric & not multiple metrics.
- In case of probabilistic model, we get single no. AUC - ROC but we still need to look at the entire curve to make conclusive decisions

## Advantages

Why ROC & not lift?

- lift is dependent of on total response rate of pop<sup>n</sup>, hence if that changes, lift curve will also change
- ROC is independent of total response rate of pop<sup>n</sup>, thus because its numerator & denominator of both n & y will change on a similar scale in case of response rate shift.

## 6. Log Loss

AUC ROC does not take into account the model's capability to predict higher probability that are more likely to be true

$$-\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1-y_i) \cdot \log(1-p(y_i))$$

$p(y_i)$  is the predicted probability of the class

$1 - p(y_i)$  is the " " " " -ve "

$y_i = 1$  for the class & 0 for -ve class (actual values)

- negative mean of log of corrected predicted probabilities
- lower the logloss better the model
- while the AUC is computed with regards to binary classification with varying threshold, log loss actually takes certainty of classification into account

## 7. Gini Coeff

used in classification

$$\rightarrow \text{gini coeff} = 2 \times \text{AUC} - 1$$

its purpose is to normalize AUC so that a random classifier scores 0, and perfect classifier scores 1

$$\rightarrow \text{range: } [-1, 1]$$

## 8. Concordant - Discordant ratio

e.g. 3 students A B C

Likelihood of passing

$$A = 0.9$$

$$B = 0.5$$

$$C = 0.3$$

Total pairs? AB, BC, AC

After year ends, A & C passed, B failed

Choose pairs with one responder, other non-responder

$$\begin{matrix} & | & | \\ AB & & BC \end{matrix}$$

Concordant pair :-  $P(\text{responder}) > P(\text{non-responder})$   
 Non- " " :-  $P(") < "$

50% Concordant ratio

Concordant ratio  $> 60\%$  is good.

Primarily used to check model's predictive power & not for  
 How many customers to target

9

RMSE

Root Mean Squared error

Most popular metric

Assumption  $\rightarrow$  errors are unbiased and follow a normal distribution

$$\sqrt{\frac{\sum_{i=1}^n (\text{predicted}_i - \text{actual}_i)^2}{N}}$$

- "√" empowers this to show large no. of deviations
- "square" displays more robust results (magnitude of error)
- avoids use of absolute error values; highly undesirable
- in more samples, reconstructing error distribution using RMSE is reliable
- highly affected by outliers
- as compared to absolute error, RMSE gives higher weightage & punishes large errors

## 10. RMSLE

used when we don't want to penalize big difference b/w predicted & actual

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (\log(p_i + 1) - \log(a_i + 1))^2}$$

- if  $a$  &  $p$  are small  $\text{RMSE} = \text{RMSLE}$
- if one is big, use RMSLE
- if both are big, use ~~RMSE~~, RMSLE becomes negligible

## 11. R-Squared / Adjusted R-Squared

In RMSE we do not have a benchmark to compare

$$R^2 = 1 - \frac{MSE(\text{model})}{MSE(\text{baseline})} = \frac{MSE \text{ of predictions against actual}}{MSE \text{ of mean prediction against actual}}$$

$$\Rightarrow 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (\bar{y}_i - \hat{y}_i)^2}$$

better the model, higher  $R^2$

## Adjusted R<sup>2</sup>

In adding new features, R<sup>2</sup> either increases or remains the same

R<sup>2</sup> does not penalize for adding features that add no value to model

L

$$\text{adjusted } R^2 = 1 - \left(1 - R^2\right) \left[ \frac{n-1}{n-(k+1)} \right]$$

k : no. of features

n : no of samples

if feature isn't valuable , R<sup>2</sup> remains same & 1 ~~remains the same~~

so overall we subtract a greater value from 1 so adjusted R<sup>2</sup> decreases

from ppt

$R^2$

total variance in actual response

$$SS_{\text{total}} = \sum (y_i - \bar{y})^2$$

error

$$SS_{\text{residual}} = \sum (y_i - \hat{y}_i)^2$$

variance in response given by model

$$SS_{\text{model}} = \sum (\hat{y}_i - \bar{y})^2$$

$$R^2 = \frac{\text{Explained variance}}{\text{Total variance}} = \frac{SS_{\text{reg model}}}{SS_{\text{total}}}$$

$$= 1 - \frac{\text{Unexplained var}}{\text{Total var}} = \frac{SS_{\text{residual}}}{SS_{\text{total}}}$$

issues?

- biased - add regressors, we get better  $R^2$
- may add insignificant predictors
- may start modelling noise
- model overfits

## Adjusted R-Squared ( $\bar{R}^2$ )

Mean Variances

$$MS_{\text{total}} = \frac{\sum (y_i - \bar{y})^2}{n-1} \quad n-1 = \text{dof}$$

$$MS_{\text{model}} = \frac{\sum (\hat{y}_i - \bar{y})^2}{p} \quad \text{dof} = p = \text{no of regression}$$

$$MS_{\text{residual}} = \frac{\sum (y_i - \hat{y}_i)^2}{n-p-1} \quad \text{dof}_{\text{total}} = \text{dof}_{\text{model}} + \text{dof}_{\text{error}}$$

$$\bar{R}^2 = 1 - \frac{MS_{\text{residual}}}{MS_{\text{total}}} = 1 - \frac{SS_{\text{residual}} (n-1)}{SS_{\text{total}} (n-p-1)}$$

$R^2$  vs  $\bar{R}^2$

- |   |  |
|---|--|
| <ul style="list-style-type: none"> <li>→ formula</li> <li>→ increases as more regressors are added</li> <li>→ biased estimate</li> <li>→ does not penalize non-significant terms</li> <li>- Always true</li> <li>- not suitable for statistical test of significance of wts.</li> </ul> | <ul style="list-style-type: none"> <li>increases only if regression is really significant</li> <li>unbiased estimate</li> <li>can be -ve</li> <li>suitable for statistical test of significance of wts.</li> </ul> |
|---|--|

2

## CROSS VALIDATION

train model on one subset of data

validate model on complimentary subset

disadvantage :- we loose a good amount of data from train.

the model ~~can~~

can be high bias

overcome

1

50-50 split

train on first 50, validate on other

train on other 50, validate on first 50.

reduces bias

1

2-fold cross validation

## k-fold Cross-validation

e.g. 7-fold validation

- divide pop<sup>n</sup> into 7 equal samples
- train on 6, validate on 1
- next iteration, do validation on diff sample
- reduces selection bias, reduces variance in prediction power
- avoids overfitting
- once we have 7 models, we take avg of error terms to find which model is best

trade off for k.

Small k :- high selection bias but low variance in performance  
 Large k :- small " " " high "

## → Overfitting & Underfitting in ML

main goal of model is to generalize well

so we need to check overfitting & underfitting

signal :- refers to true underlying pattern of data

noise :- unnecessary and irrelevant data that reduces

bias :- prediction error i.e. introduced in the model | performance  
 (error rate) due to oversimplifying the algorithms.

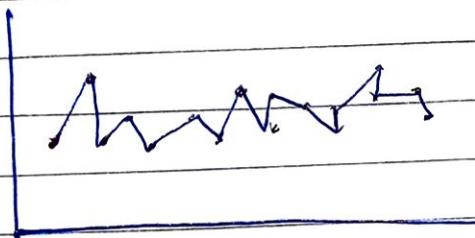
Variance :- perform well with training data, but not with testing data

diff b/w error rate of training & testing data

## Overfitting

- when our model tries to cover all the data pts. more than required data pts.
- model starts catching noise & inaccurate values present in dataset, reduces efficiency & accuracy of model.
- has low bias & high variance

eg:



~~goal~~ goal of model is to get best fit, but here there is no best fit line

Avoid overfitting :- cross validation

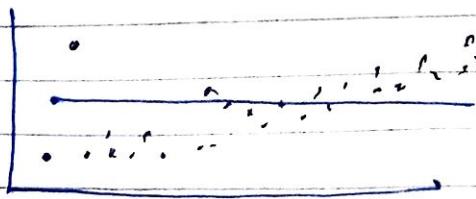
training with more data  
removing features

early stopping the training  
regularization  
ensembling

## Underfitting

- when model is not able to capture underlying trend in data
- may fail to find best fit

eg



Avoid :- increase training time  
in " " no. of features

Page No.		
Date		

↳ Goodness of fit

defines how closely the result of predicted values match true values

good unfitted model — good model — overfitted model

"

( stop at good pt. )

find by → resampling

→ validation dataset