

Business Statistics for Data Science

Population vs Sample

Population Size N

Population mean μ

Population variance σ^2

Sample Size n

Sample mean x

Sample variance s^2

$$\sum_{i=1}^{i=N} x_i$$

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

$$\sum_{i=1}^{i=N} x_i$$

$$\frac{N}{n-1}$$

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

Why $n-1$ for sample variance and not n , because of Bessel's correction

Bessel's correction: you know the value of the n th book or item (use a group of books thickness)

mean,

for measurement of central tendency, we uses median, mode
we use median when mean is skewed by outliers.

When handling missing dataset, if Normal distribution we use mean

if skewed we used median

If missing Categorical data mode.

- If $11 \quad 11 \quad 11 \quad 11$ is too large

find the ratio of categories and apply the ratio to missing data :

Probability

Coin \rightarrow Head

(H)

(T)

$$\rightarrow \frac{1}{2} \quad \text{NU:} = 0.5$$

probability = $\frac{\text{no. of favourable outcome}}{\text{Total no. of possible outcome}}$

Sample space

Event

Mutually Exclusive Events

$$P(A \cup B) = P(A) + P(B)$$

$\downarrow \quad \downarrow \quad \downarrow$
 probability of being a boy or being a girl

Pr of 2 or 5 on a die

$$\Pr(2) + \Pr(5)$$

not mutually exclusive Event

$$P(A \cap B) = P(A) + P(B) - P(A \cap B)$$

$\downarrow \quad \downarrow \quad \downarrow$
 probability of being a boy and a girl

Independent Events

$$P(A \& B) = P(A) * P(B)$$

$$P\left(\frac{A}{B}\right) = P(A)$$

$$P\left(\frac{B}{A}\right) = P(B)$$

rolling a dice and getting heads on a coin

$$S = \{1, 2, 3, 4, 5, 6\}$$

$$A \rightarrow \{2, 4, 6\}$$

$$P(A)$$

cumulative probability: likelihood a value of a random var within a given range.

flip two coins, probability of one or fewer heads

sample space = {HH, HT, TH, TT} $\xrightarrow{x \leq 1}$

$$P(x \leq 1) \rightarrow (x=0)$$

$$P(x \leq 1) = P(x=0) + P(x=1)$$

$$= \frac{1}{4} + \frac{3}{4} = \frac{3}{4}$$

conditional probability

return allowed

$$6 \text{ Red marbles} \xrightarrow{\text{R/B}} \frac{6}{10} \quad \frac{R/B}{4/10} = \frac{16}{100}$$

return not allowed

$$6 \text{ Red Marbles} \rightarrow \frac{6}{10}$$

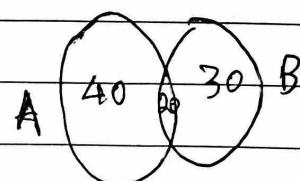
$$4 \text{ black Marbles} \rightarrow \frac{4}{10} : \frac{3}{9}$$

$$P(B|A) = P(A \cap B) / P(A)$$

$$\frac{4}{10} \times \frac{3}{9} = \frac{12}{90}$$

100 people 40 bought sports drinks 30 bought snacks

20 bought sport drinks and snacks



$$\textcircled{1} \quad P(A) = \frac{40}{100} = 0.4$$

Bayes Theorem
Extension of Conditional Probability

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

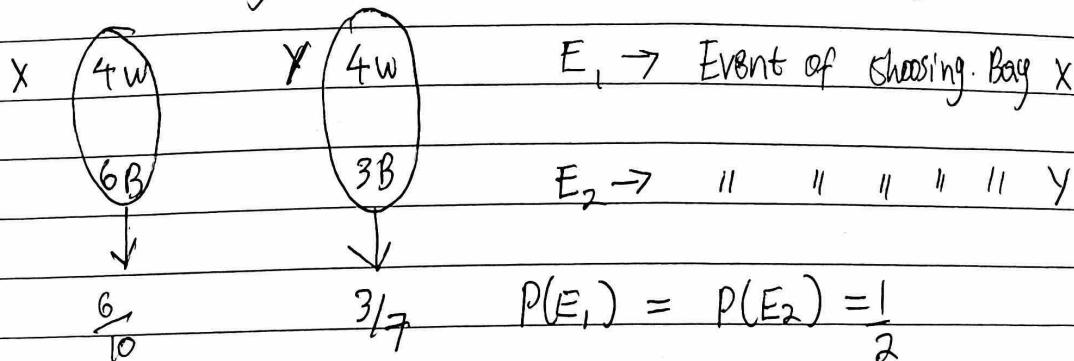
$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

$$P(A \cap B) = P(B|A) \times P(A)$$

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

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

A Bag X containing 4 white & 6 black balls while another Bag Y containing 4 white and 3 black balls. One ball is drawn at random from one of the bags, and it is found to be black. Find the probability that it was drawn from Bag X.



$$P(A|E_1) = P(\text{Black from Bag X}) = \frac{6}{10}$$

$$P(A|E_2) = P(\text{Black from Bag Y}) = \frac{3}{7}$$

$$P(E_1|A) = \frac{P(E_1) \times P(A|E_1)}{P(E_1)P(A|E_1) + P(E_2)P(A|E_2)}$$

(NU:)

~~16:00~~

probability of getting one head

$$P(X=1) = \frac{5}{32} = \frac{5!}{32} = \frac{5!}{1!(5-1)!} \times \frac{1}{32}$$
$$= \frac{5!}{4!} \times \frac{1}{32}$$
$$= \frac{5}{32}$$

$$P(X=2) = \frac{5!}{32} = \frac{5!}{2!(5-2)!} \times \frac{1}{32}$$

$$= \frac{120}{2 \times 6} \times \frac{1}{32} = \frac{10}{32} = \frac{5}{16}$$

$$P(X=3) = \frac{5!}{32} = \frac{5!}{3!(5-3)!} \times \frac{1}{32} = \frac{10}{32}$$

$$P(X=4) = \frac{5!}{32}$$

$$P(X=5) = \frac{1}{3}$$

$$\frac{\frac{1}{2} \times \frac{6}{10}}{\frac{1}{2} \times \frac{6}{10} + \frac{1}{2} + \frac{3}{7}} = \frac{\frac{6}{20}}{\frac{6}{20} + \frac{3}{14}}$$

$$= \frac{6}{20} \div \left(\frac{42}{140} + \frac{30}{140} \right)$$

$$\frac{6}{20} \div \left(\frac{72}{140} \right)$$

$$\frac{6}{20} \times \frac{\cancel{140}}{\cancel{72}} = \frac{84}{900}$$

$$= \frac{42}{450} = \frac{16}{150}$$

$$\frac{7}{12} = \frac{17}{75}$$

probability distribution

uniform Distribution

- discrete uniform : Discrete integers
- continuous uniform : Continuous float

X

Binomial Distribution

$$P(X) = {}_n C_x P^x (1-P)^{n-x}$$

$$2 \times 2 \times 2 \times 2 \times 2 = B_2$$

flip 5 coins probability

Heads = 0

Possible outcome = $2^5 = 32$

$$P(X=0) = \text{TTTT} = \frac{1}{32} = \frac{5!}{32} = \frac{5!}{0! 32!} \times \frac{1}{32} = \frac{1}{32}$$

(nu:)

Calculating Probability with Z-score for Normal Distribution

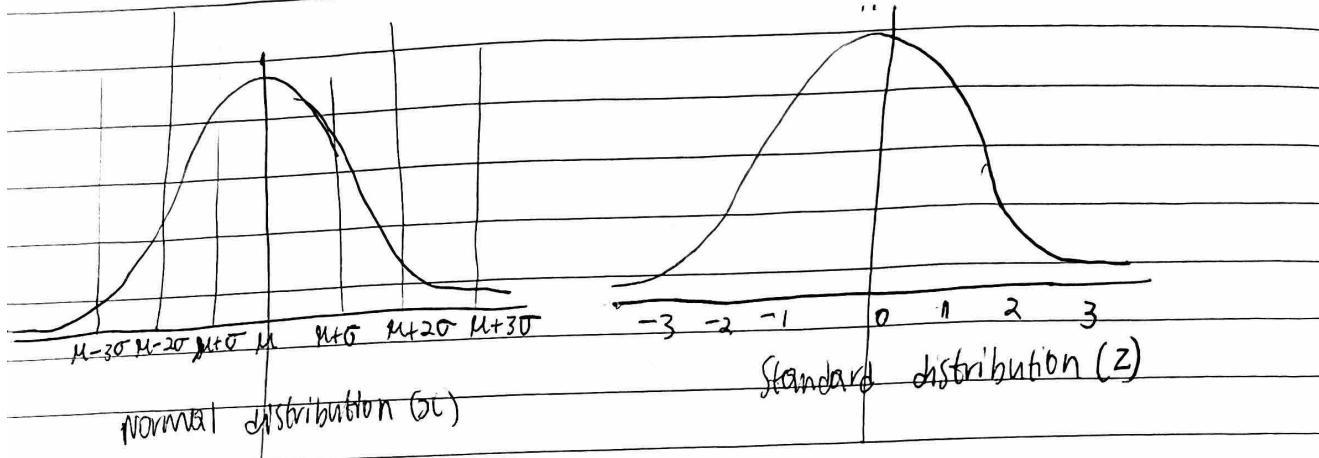
- find out how far the value of x is from μ (in terms of "number of SD")
- The number is called a "Z-score" or "standard score" or "z value"

ie

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

$$x = \mu$$

$$z = \frac{\mu - \mu}{\sigma} = 0$$



using the z table to make predictions

200 students

John scored 700/1000 in a test

how many students did John do better than

Average score is 600

S-P is 150

$$\frac{700 - 600}{150} = 0.67$$

Check z table

$$0.7486 \times 100 = 74.86\% \text{ H8 did better than}$$

Kurtosis

Measures the fourth degree of tailored

Types of Kurtosis

platy kurtosis

Meso kurtosis

Leptokurtic

Poisson Distribution

This is a discrete distribution, occurring over period like an hour

$$P(X=x) = \frac{\lambda^x e^{-\lambda}}{x!} - \cancel{\lambda^{10}} = \cancel{0!} \quad \lambda = 2.71828$$

- Average number of calls for a call center is per hour $\lambda=10$
- probability of having different number of calls in an hour could be calculated as:

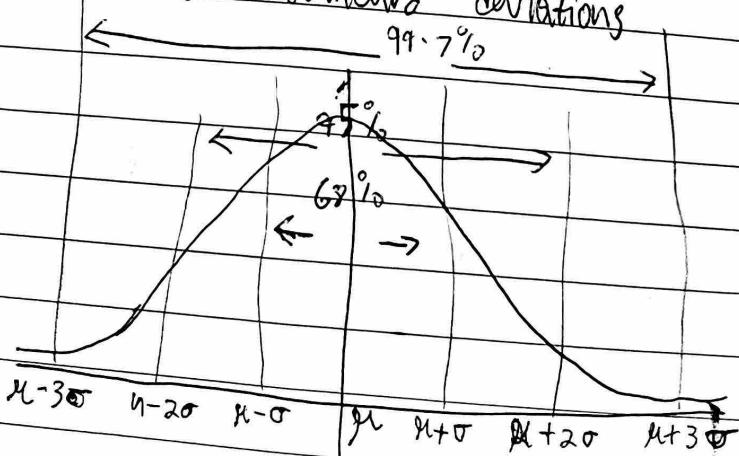
$$P(X=0 \text{ calls}) = \frac{10^0 e^{-10}}{0!} = \frac{1 \cdot 0.000045}{1} = 0.0045\%$$

$$P(X=3 \text{ calls}) = \frac{10^3 e^{-10}}{3!} = \frac{1000 \cdot 0.000045}{6} = 0.75\%$$

Normal distribution (Gaussian distribution)

- Graphically, normal distribution appears as a bell curve

- for all normal distribution, 68.2% of the observations will appear within plus or minus one standard deviation of the mean, 95.4% of observations will fall within +/- two standard deviations; and 99.7% within +/- three standard deviations



- Null Hypothesis: In Inferential Statistics, null hypothesis is a general statement that there is no relationship between two measured phenomena (An assumption)

a company production is $= 50$ unit / per day

- Alternative hypothesis: The alternative hypothesis is the hypothesis used in hypothesis testing that is contrary to the null hypothesis.

a company production is $\neq 50$ unit / day etc

null hypothesis: =, greater than or equal, less than or equal

Alternative hypothesis: $\neq, <, >$

Tailed test

one tailed test: A test of statistical hypothesis, where the region of rejection is on only one side of the sampling distribution

≥ 800

Two tailed test: two sides critical areas

$\neq 800$

p-value: p value or calculated probability, is the probability of finding the observed, or more extreme, results when the null hypothesis of study question is true

Types of test

- (1) t-test
- (2) z-test
- (3) ANNOVA
- (4) Chi-Square Test
- (5) correlation test

under Treatment

NU:

Correlation

Correlation is similar to covariance, also for measuring how two variables are related or moving together

+1 to -1 Range

Direction and Magnitude or useful

$$r_{xy} = \text{Correlation}(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x)} \sqrt{\text{Var}(y)}}$$

$$\frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

$$\text{Cov}(x, y) = \frac{1}{n-1} \sum (x - \bar{x})(y - \bar{y})$$

$$\frac{29.8}{4} = 7.45$$

Standard deviation

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

Sample mean

0 - 0.3 weak linear relationship

0.3 - 0.7 moderate linear relationship

0.7 - 1 strong linear relationship

p-value
finding of study

Hypothesis testing

statistical method that is used in making statistical decisions using experimental data. used for two mutually exclusive events

Types

- (1) t - t
- (2) z - z
- (3) ANOVA
- (4) Chi
- (5) some

nu:

$\frac{75}{100} \times 200 = 150$ students John did better than

2) % of students shorter than 1.5m in a normal distribution
 $\sigma = 0.06\text{ m}$ Mean = 1.6m

$$z = \frac{x - \mu}{\sigma} = \frac{1.5 - 1.6}{0.06} = \frac{-0.1}{0.06} = \frac{100}{6} > \frac{5}{3} = -1.67$$

0.0475

4.75% shorter than
1.5m

4.05 Covariance

Covariance is a statistical term that refers to a relationship between two random variables on how one variable changes when the other one

Covariance can either be positive or negative

$-\infty$ to $+\infty$

Magnitude does not really matter with covariance

number of

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

mean of x mean of y
1st value 2nd value

n = number of data

$$\frac{1}{n-1} \sum$$

When sample size:

1 mark

ANOVA (Analysis of Variance)

It's a statistical technique that is used to check if means of two or more groups are significantly different from each other. It's parametric test for hypothesis testing

H_0 : $\text{mean}_1 = \text{mean}_2 = \text{mean}_3$. There is no significant difference between mean values of each group

H_1 : Means aren't equal

Sum of Squares (SS): Quantifies the variability within the group

Degree of Freedom (DF):

df (between groups): no of groups - 1

df (within groups): no of observations - no of groups

Mean Squared (MS): It is the average variation either between groups or within groups. $MS = SS/df$

F-statistics: It is the test statistic in one way ANOVA

$$F = MS(\text{between groups}) / MS(\text{within groups})$$

If p-value > significance value \rightarrow we accept the null hypothesis
else not

Significance level: 0.05

T-test

Makes use of hypothesis test and alternative test to make decision, if $P > 0.05$, then we follow the null hypothesis

H_0 - Null Hypothesis

H_1 - Alternative Hypothesis

Z-test

A Z-test is used to compare the mean of two given samples and infer whether they are from the same distribution or not.

We do not implement Z-test when the sample size is less than 30.

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

Chi-square Observed

$$\chi^2 = \sum \frac{(O - E)^2}{E} \quad \begin{matrix} \text{Observed} \\ \text{Expected} \end{matrix}$$

H_0 : There is no association between gender & location

H_1 : There is association between gender & location

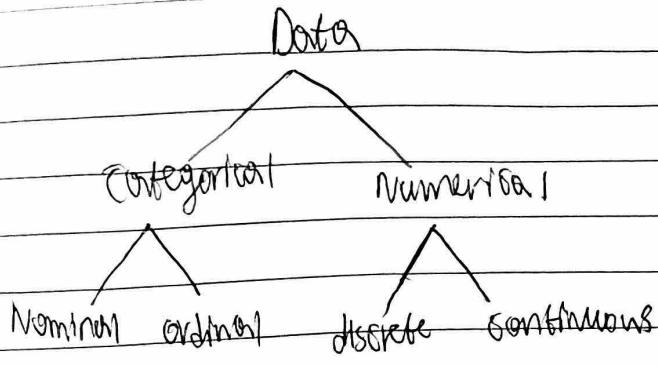
p-value $> 0.05 \rightarrow$ Reject H_0

p-value $\leq 0.05 \rightarrow$ Accept H_0

$$\text{Expected} = (\text{Row total} * \text{Column Total}) / \text{Grand Total}$$

$$df = n_{\text{Row}} - 1 * n_{\text{Column}} - 1$$

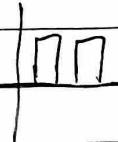
Other Treatment



Type of Data

1 Categorical

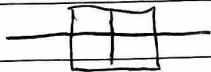
What we see



Type of test

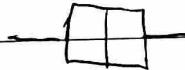
1 sample test

1 Numerical



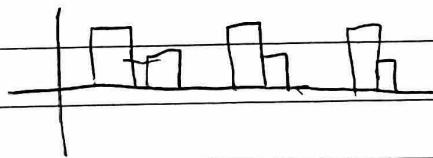
t - test

1 Numerical and 1 categorical



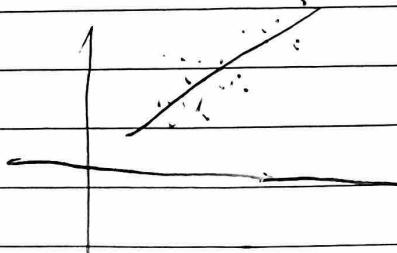
t - test or Anova

2 Categorical



Chi square test

2 Numerical



Correlation

practical cleaning of Skinning model

Row number	int
customer ID	object *
surname	object
credit score	int
Geography	object
Gender	object
Age	int
Tenure	int
Balance	float
Num of products	int
Has credit card	object (binary) *
Is Active member	object (binary) *
Estimated salary	float
Exit	object (binary) *

feature scaling

making large numbers more digestible for our model using scalers

standard scaler (standardization) when you have a normal distribution

ensures that for each feature, the mean is 0 and standard deviation is 1, bringing all features to the same magnitude. Standardization helps scale down features based on the standard normal distribution

min-max scaler (normalization) when your distribution is not normal

Normalization helps you to scale down your features between a range 0 to 1

$$x_{\text{norm}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

Handling missing values.

(1) Delete Rows / columns

delete column if it has 75% missing values and no dependence

2% - 3% of data missing you can input or delete

(2) Replace with (mean) / (median) / (mode)

when there
are outliers

uses for
categorical
data

(3) Algorithm imputation

KNN, Naive Bayes, Random forest are algorithms that have ability to handle missing value.

(4) predicting the missing values

prediction model is one of the advanced method to handle missing values. dataset with no missing values become training set and data set with missing value become the test set and the missing values is treated as target variable.

Practical handling missing values

- filling with mean, median, mode
- forward fill or back fill
- using Time series algorithm (too tedious for this stage)

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.2}{0.4} = 0.5$$

$$\begin{array}{r} 118 \\ 84 \\ \hline 34 \\ 17 \\ \hline 10 \end{array}$$

$$\frac{0.2}{0.6} = 0.33$$

~~or~~ $P(A \cap B) / P(B)$

$$\frac{P(B)P(A \cap B)}{P(B)P(A \cap B) + P(A)P(A \cap B)}$$

$$\frac{0.2(0.4)}{0.2(0.4) + 0.6(0.2)}$$

$$\frac{0.08}{0.08 + 0.12}$$

$$\underline{0.08}$$

$$\underline{0.20}$$

EDA

Data Sampling gathering data from multiple sources as external or internal data collection.

- 1 public data
- 2 private data

Data Cleaning

Clean the data to improve the quality of the data for further data analysis

Handling missing values → Standardization of the data → another treatment

→ Handle invalid values $\textcircled{n.u.}$

Types of analysis

1. Univariate analysis

are Variable analysis
can be used for Numerical and categorical data
mostly categorical data. distribution analysis using box plot
and histogram

2 Bivariate analysis

analysis of two variables, can be done for both numerical
and categorical

two numerical — Scatterplot

2 categorical — Bar chart

Numerical & categorical — box plot, histogram

3. Multivariate Analysis

multiple variable analysis

Numerical Analysis

various analysis of numerical data

1 numerical — Box plot

2 numerical — Scatterplot

for multiple numerical functions you can just use corr() in
python which gives the correlation in form of heatmaps

Outlier Treatment

Outliers are the most extreme values in data
Detect outliers using following methods:

1. Box plot
2. Histogram
3. Scatter plot
4. Z-score
5. Inter quartile range (Values out of 1.5 times of IQR)

Outlier treatment is important for time series forecasting

Handle outlier using following methods

1. Remove the outliers
2. Replace outlier with suitable values by using the following methods
 - Quantile methods
 - Inter quartile range
3. Use ~~more~~ ML models which are not sensitive to outliers. like:- KNN, Decision Tree, SVM, Naive Bayes, Ensemble methods.

Types of Data

Qualitative

Nominal

Gender MF

location

ordinal

Economic status

Grade

Quantitative

Discrete

Whole

Number

continuous

37.7

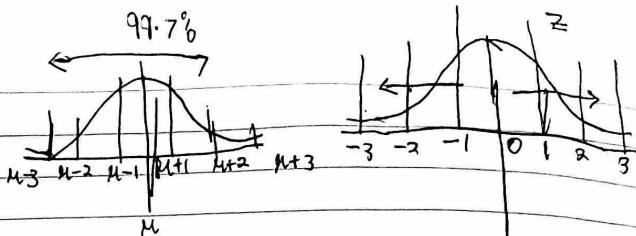
33.2

Example

Standardization

$$\mu = 0$$

Age



$$z - \text{Score} = \frac{x - \mu}{\sigma}$$

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$$

$N-1$ sample

Age	Income (£)	New Value
24	15000	$(15000 - 19000) / 9643.65 = -0.4147$
30	12000	$(12000 - 19000) / 9643.65 = -0.7258$
28	30000	$(30000 - 19000) / 9643.65 = 1.1406$

SKLearn = has libraries to accommodate feature Scaling

Example

Normalization

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Age	Income (£)	New Value
24	15000	$(15000 - 12000) / 18000 = 0.16667$
30	12000	$(12000 - 12000) / 18000 = 0$
28	30000	$(30000 - 12000) / 18000 = 1$

Restart windows

Shift + ESC

Mount -n -o remount, rw/

passwd your-username

reboot

- Model + Dimensionality reduction + visualization

~~Model + latent representation + dimensionality reduction + visualization with DAG~~

- Model + Gradient methods

- Model + dimensionality reduction + LIME

- Model + latent representation + visualization

One-Hot Encoding Example

Dummy
features

Location	loc-m	loc-indo	loc-hn
malaysia	1	0	0
Indonesia	0	1	0
India	0	0	1
India	0	0	1
Mal	1	0	0
man	1	0	0

- Dummy encoding : an extension of one hot encoding, use dummy features. in here one of the features (^{columns}) can be ignored
 $n-1$ = dummy features

Select ~~id~~ ^{id} from Weather
 Select Date from Record Date

Where

Select ID from Weather

I → Where
 + I

derived metrics
create a new variable from the existing variable to get a
insightful information from the data

Age	Age months	derives metrics
25	25 x 12	
28		
29		

- calculated from Data
- From Domain Knowledge
- Feature Binning
- Feature Encoding

Feature Binning

feature binning converts or transforms continuous / numerical variable to categorical variable. It can also identify missing values and outliers.

- Equal width binning (grouping by width size)
- Equal frequency binning (grouping by count)

$$0-30 \rightarrow 2k$$

$$31-45 \rightarrow 2k$$

$$46-50 \rightarrow 2k$$

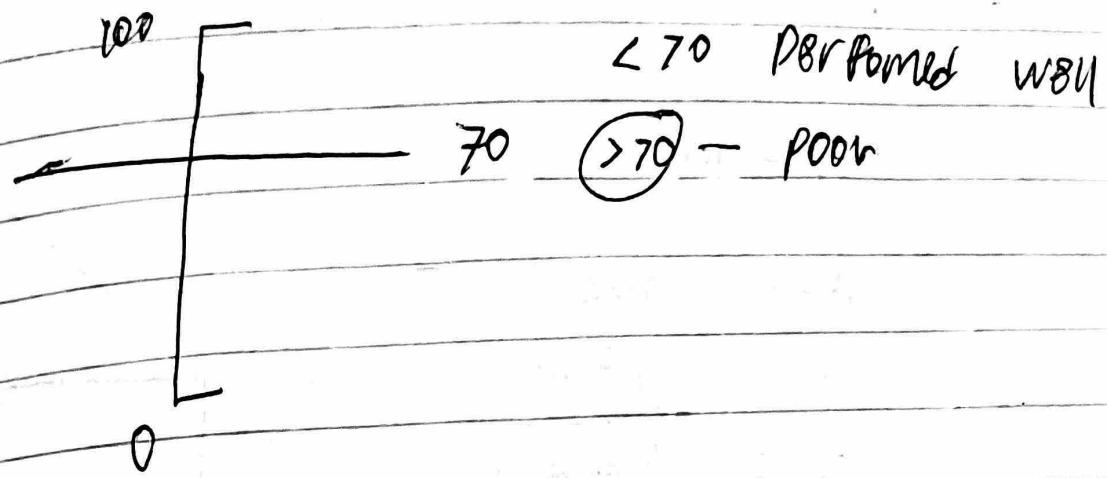
Feature Encoding

- label encoding: label encoding is technique to transform categorical variables by assigning a numerical value to each of the categories. Only used when it's the target / y / dependent variable, becos large number has many
- one-hot encoding: This technique is used when independent variables are nominal. It creates k different columns each for a category and replaces one column with 1 rest of the columns is 0. 0 represents the absence and 1 represents the presence of that category

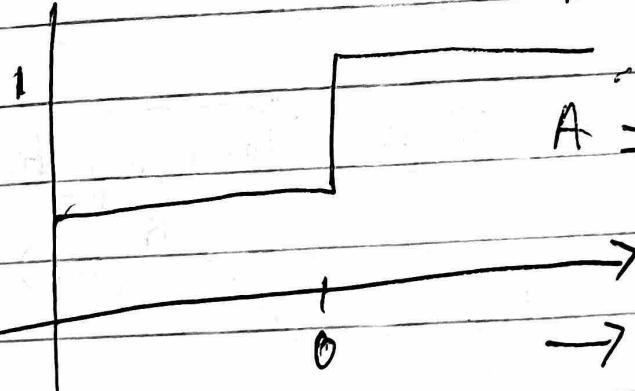
nu:

Activation function: Step functions

Threshold based activation function



Step function

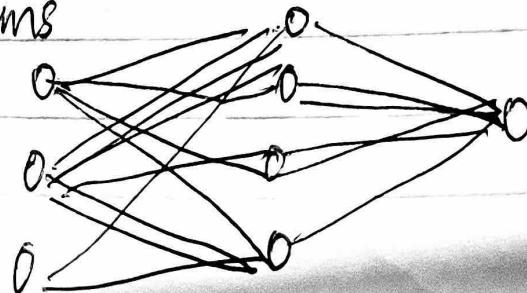


$$A = 1, y > \text{threshold}$$

$$A = 0, \text{otherwise}$$

$$A = 0$$

Step function use Binary classification problems



↑ activated

↓ not activated

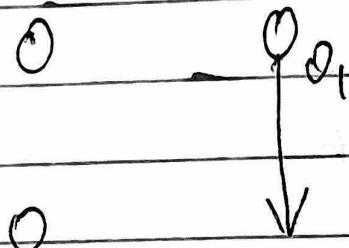
$\phi = 0$ (node is inactive)

$\phi \neq 0$ (node is active)

Activation functions

Part of the hidden layers which decides if neuron

should be fired or not



$$o_i = \phi \sum (w_i x_i) + \text{bias}$$

ϕ ranges from $-\infty$ to ∞

Types of Activation functions

Step functions

Linear functions

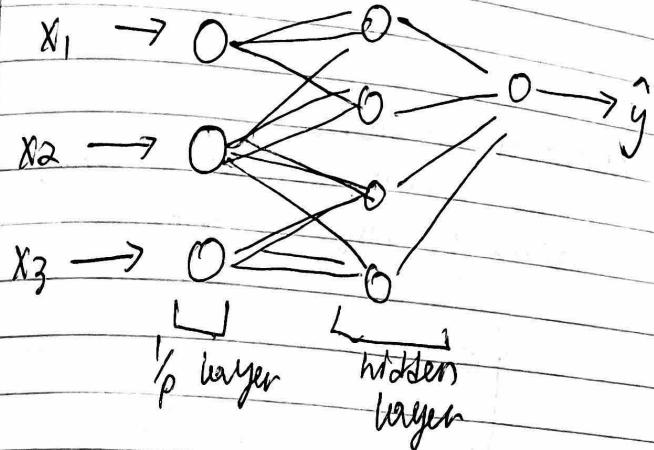
Sigmoid functions

Tanh functions

ReLU ✓

Deep learning

Neuron

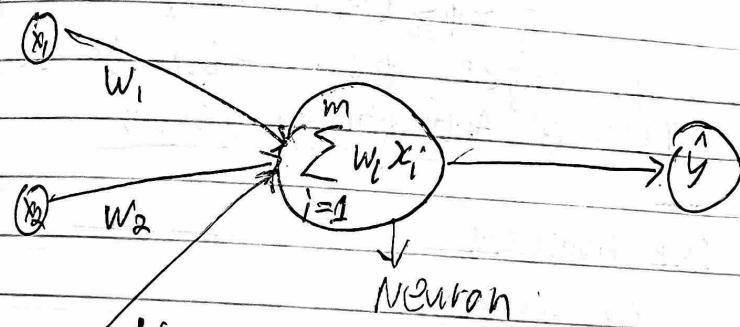


Visualization

Dimensionality
LAE

+ Lime

Visualization



$$w_1 x_1 + w_2 x_2 + w_3 x_3$$

$w_1, w_2, w_3 \rightarrow$ weights (connection strength)

$$\hat{y} = \phi(\sum w_i x_i)$$

→ Activation function

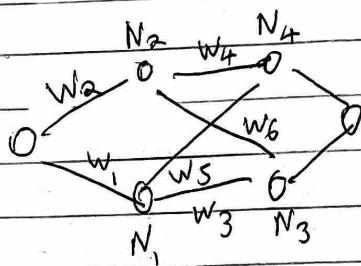
Abstract

Signed

Nov

abstraction

last 10



$$N_1 = w_1 x_1$$

$$N_2 = w_2 x_2$$

$$N_3 = w_3 N_1 + w_4 N_2$$

$$\begin{aligned} N_3 &= w_3 w_1 x_1 + w_4 w_2 x_2 \\ &= m_1 x_1 + m_2 x_2 \end{aligned}$$

Δy_j
to

- Poorly used
- Advises not to be used for multiple hidden layers. (mostly used for first layer)

$$= \sqrt{36 + 9}$$

$$= \sqrt{54}$$

$$\therefore = \sqrt{6 \cdot 9}$$

gradient

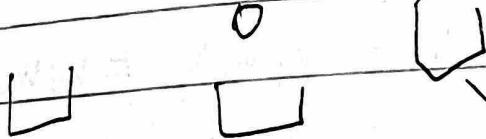
$$y = \frac{\Delta y}{\Delta x} = \frac{OB}{OA} = \sqrt{\frac{(8-8)^2 + (6-3)^2}{(2-8)^2 + (3-3)^2}}$$

$$0 \quad 0 \quad = \frac{3}{6} = 0.5$$

$$0 \quad 0 \quad 6$$

$$0 \quad 0 \quad 0$$

$$0 \quad 0 \quad 0$$



$$y_1 = mx$$

$$y_2 = y_1$$

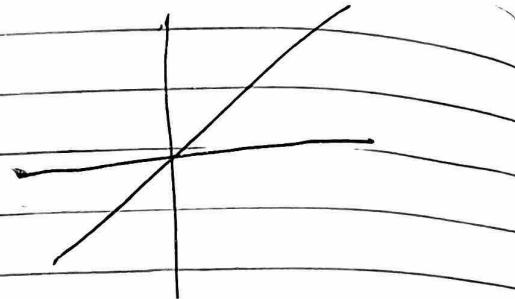
- Step function cannot be used for multiclass problem
- Activation function: Linear function

$$y = mx$$

$$x=0 \quad y=0$$

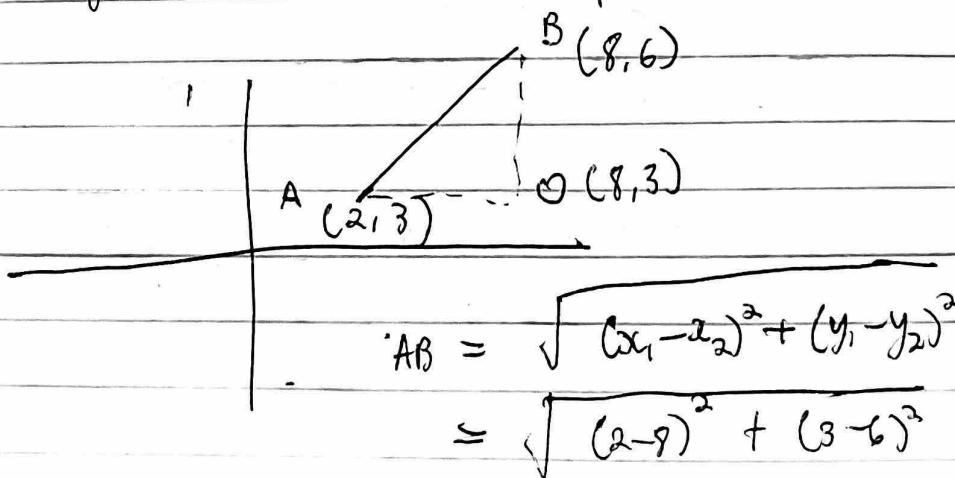
$$x=1 \quad y=m$$

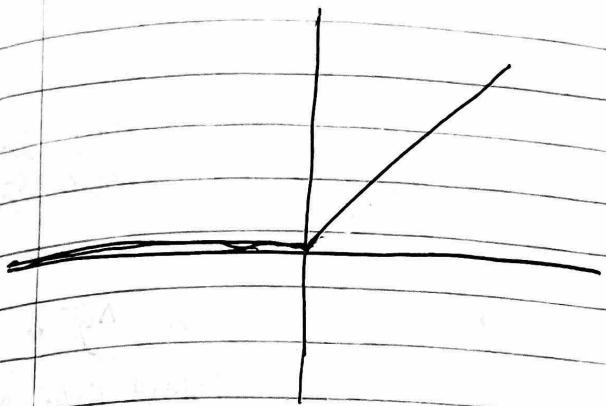
$$x=2 \quad y=2m$$



$$\frac{dy}{dx} = m$$

gradient has no relationship to x





① Non-Linear

$$y = \max(0, x) \text{ range}$$

if $x < 0$ $y = 0$

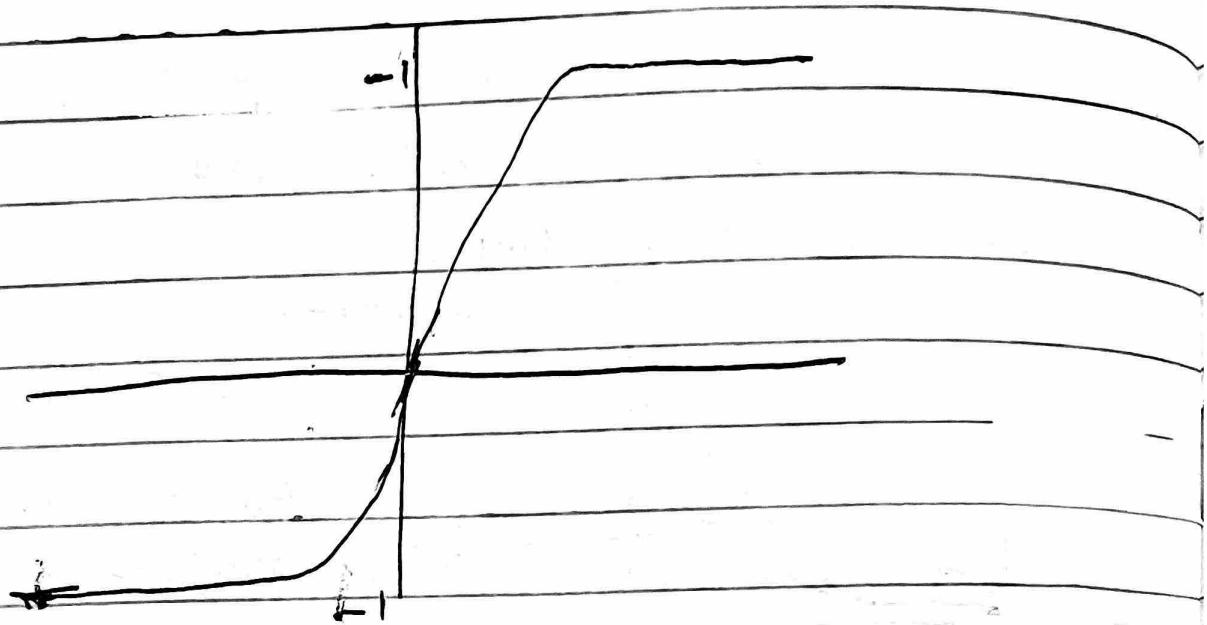
if $x > 0$ $y = x$
 $(0, \infty)$

Layers can be stacked because it's non-linear

linear

	0	0	0	
0	0	0	0	0
0	0	0	0	
0	0	0	0	
1				1
$\frac{1}{L}$				$\frac{1}{L}$

Tanh Function



Similar to Sigmoid

Non-linear

Range '(-1, 1)

~~relationship btw Tanh and Sigmoid~~

$$\tanh(x) = 2 \phi(2x) - 1$$

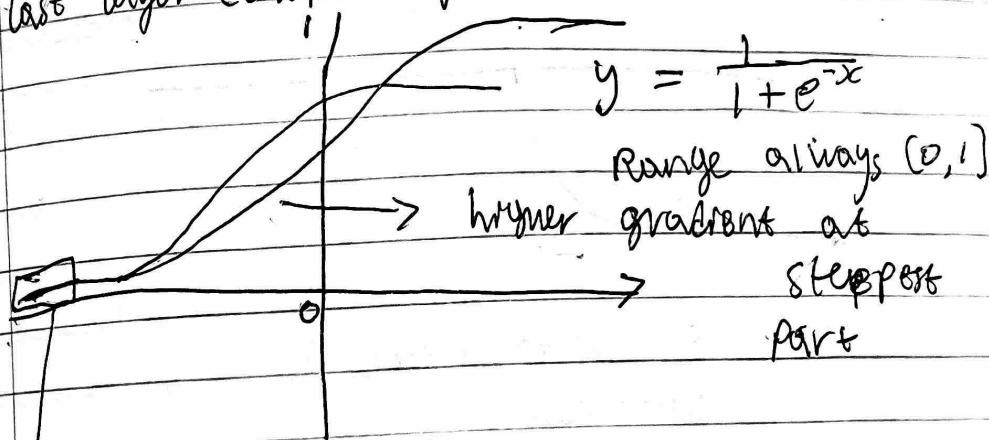
ReLU Function

Rectified linear unit (ReLU)

Activation function Sigmoid function (ϕ)

Sigmoid

Non-linear in nature (one of the best activation function), mostly used in the last layer (output layer)



$\downarrow y \approx 0$ gradient ≈ 0 , here network refuses to learn or drastically becomes low

$$\frac{dy}{dx} = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right) = ?$$

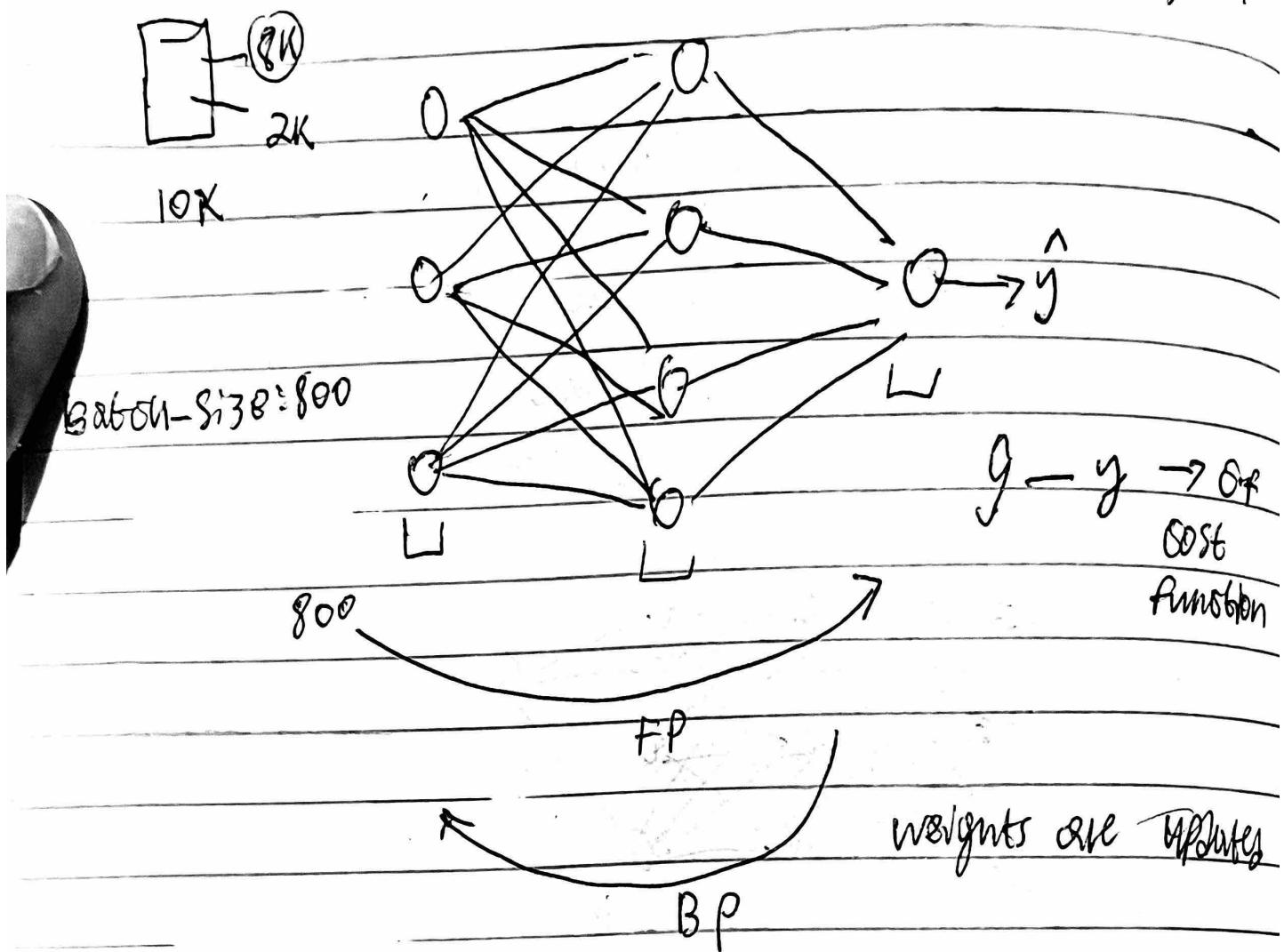
$N_0 N_2$

$w_0 w_2 x_2$

$m_2 x_2$

multiple hidden
layer

$w_1 \ w_2 \ w_3 \ w_4$



Iteration: one complete FP + BP

$$800 \times 10 = 8000 \text{ test}$$

$$\frac{8000}{800} = 10 \text{ cycles}$$

Epoch = the completion of a full training
data points

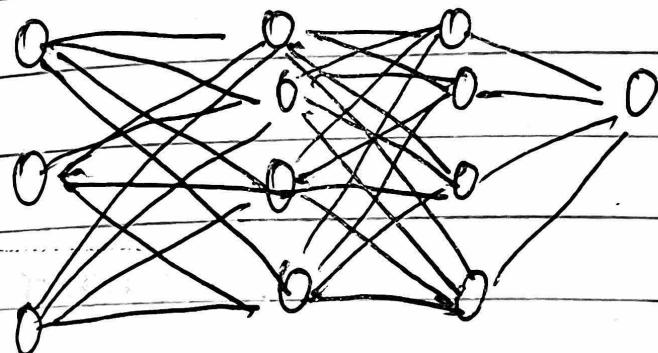
Backpropagation and forward pass

Customer churn prediction

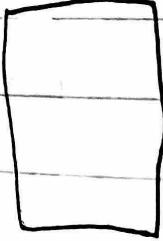
Age

Time

Contract type



10K



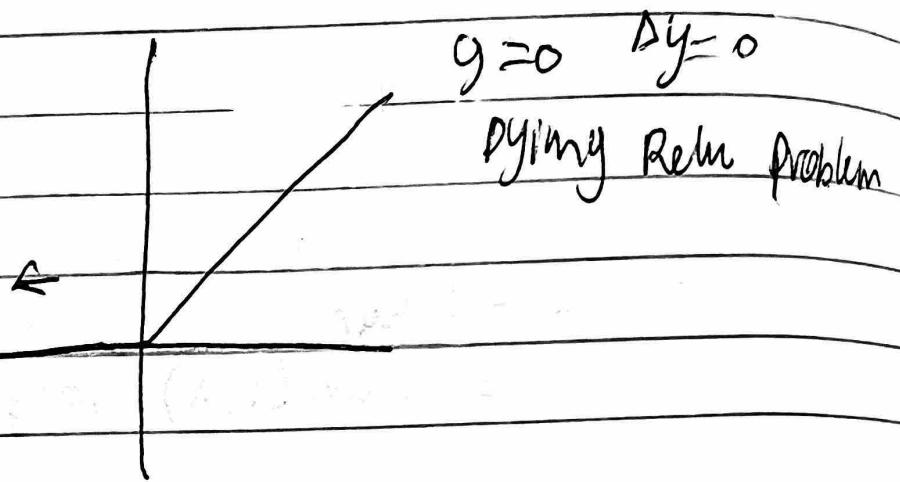
Train

8K

Test

2K

if sigmoid or tanh is used for such an ^{large layer} activation, it will be very slow and computationally expensive. To ensure sparse activation we use ReLu



neurons in this state are not active

To solve dying ReLU problem ~~some neurons~~ neurons are not completely passive to create leaky ReLU

$$y \neq 0 \quad y = 0.001x \text{ leaky ReLU}$$

time

Artificial Neural Networks
features are inputs for input layer

- Traditional machine learning steps
- + Install libraries & importing them
- + Importing dataset
- + Data preprocessing
- + Feature scaling
- + Feature extraction
- + Splitting data into train & test
- + Fitting the classifier to training data predictions
- + Analyzing predictions with actual test data

gradient

minimum of
the steepest descent

y = f(x)

using gradient

In ANN, Step #4 to #6 is replaced by below:

- + Initializing ANN
- + Adding Input, Hidden & Output layers
- + Compiling the ANN
- + Fitting the ANN to the training dataset

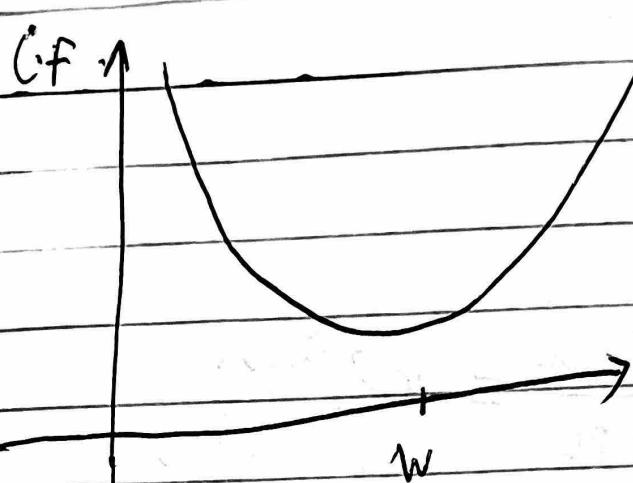
NN training hyperparameters

Epochs is an hyperparameter

Optimizer is an hyperparameter

minimum

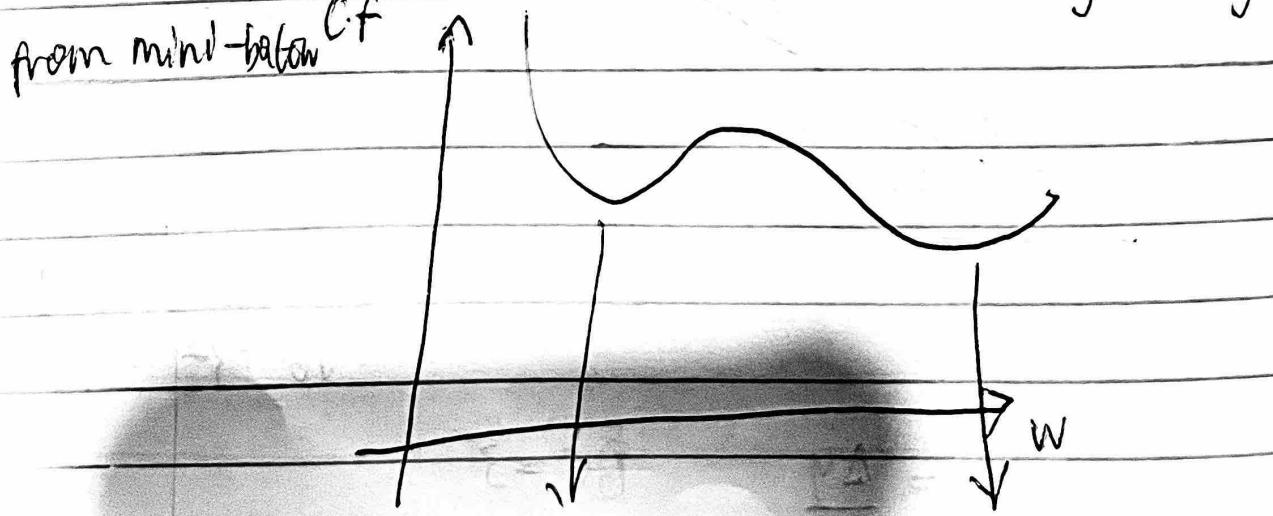
Cost function vs weight over time



Descent

gradient: cost function decreasing, iterative optimization algorithm used to find the minimum of a function by taking steps in the direction of the steepest descent

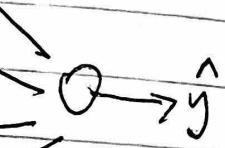
Stochastic gradient descent: continuously finding minimum past the local minimum by using gradient from mini-batch C.F.



local minimum

Global minimum

$w_1 \ w_2 \ w_3 \ w_4 \dots \ w_{16}$



$y - y \rightarrow \text{of}$
cost
function

weights are updated

complete FP + BP

est

etion of a full training

Training 10000

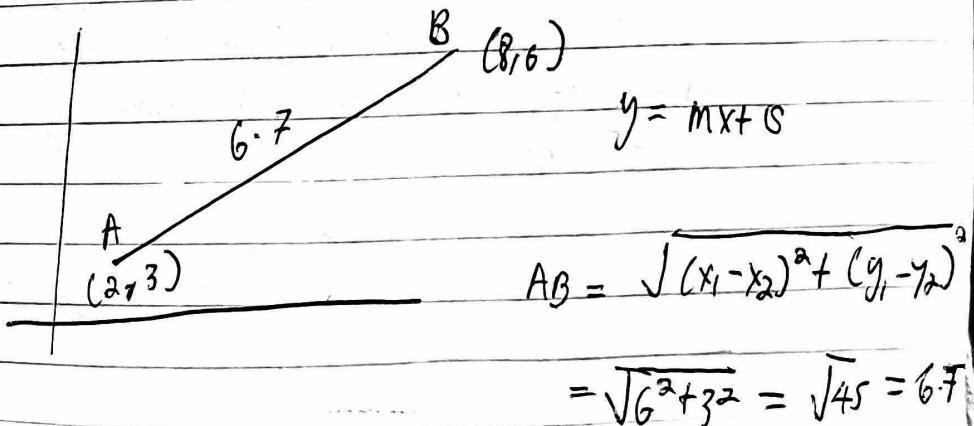
batch-size : 100

epochs : 20

$$\frac{10000}{100} = 100 \xrightarrow{\text{cycles}} 100 \xrightarrow{\text{epochs}}$$

2000 cycles is \rightarrow 20 epochs

Gradient Descent

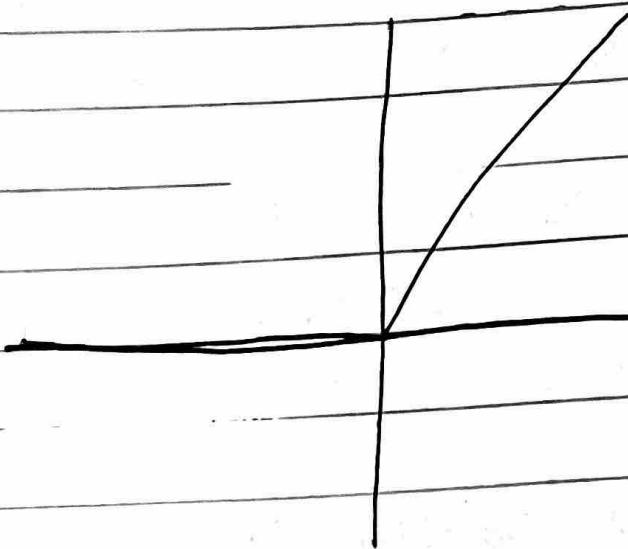


$$g = \frac{\Delta y}{\Delta x} \quad \Delta x = 3$$

$$\Delta y = 6$$

$$g = \frac{3}{6} = 0.5$$

ReLU



$$y = \max(0, x)$$

$$x < 0 ; y = 0$$

$$x \geq 0 , y = x$$

- Pooling: will tell where most of the information is. there are 4 types of pooling in Keras
- max
- Average
- Global max
- Global Average

~~11~~



sov

Hyperparameter Optimization

Manual Search

Random Search = Randomized Search CV

Grid Search CV

Manual

Solving your model again and again
on Adam optimizer and RMS optimizer

Randomized Search CV

defined parameter grid

OPT = ["adam", "RMS OP"]

Epochs = [5, 10, 15, 30, 50]

→ to possible combination

n_iter (iteration)



5 number of iteration,

Grid search

Running each and every combination

Pooling

In pooling we use the pooling methods

0	6	1	0	0
4	1	2	3	1
1	0	2	1	1
0	0	2	3	1
1	0	0	3	1

Max

6	3	1
1	3	1
1	3	1

Average

2.75	1.5	0.25
0.25	2	0.5
0.25	0.75	0.25

Flattening

6	3	1
1	3	1
1	3	1



6
3
1
1
3
1
1
3
1

feature

feature is

x

is 0 with

feature detector goes over our image
to detect features

we uses multiple feature detectors to
develop first $2^n, 3^n$

7×7 image with 3×3 feature
detector, output will be 5×5

* strides are hyperparameters 1 stride is
 5×5 , 2 strides is 3×3 matrix

padding
in p

0	6
4	1
1	0
0	1

PL

ReLU

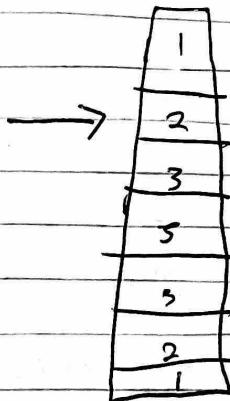
$$x < 0 \quad y = 0$$

$$x \geq 0 \quad y = x$$

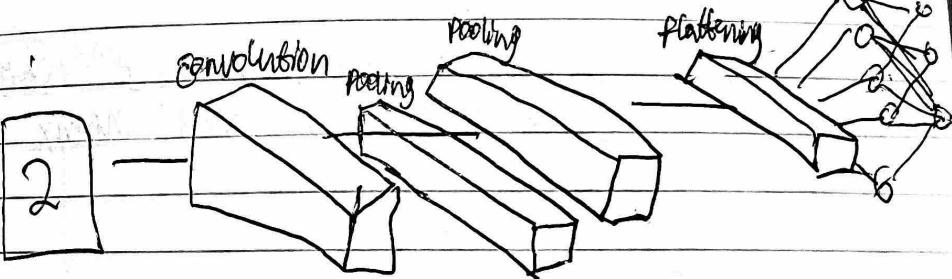
values less than 0 be some 0 with
ReLU

Flattening

1	2	3
5	6	7
3	2	1



fully connect



most of ~~the~~ the
+ types of pooling

Convolution neural network architecture

0 - 255

0 - black

1 - white

(Kerner, filters)

Picture detector: defects features.

Complexity of data determines the epochs

$28 \times 28 \times 1 \rightarrow$ Gray scale

$28 \times 28 \times 3 \rightarrow$ Colors
↓

R, G, B

Convolutional Neural Networks : Batch Size vs

Iterations vs Epochs

Gradient Descent: Optimization algorithm in machine learning used to find the best result

Epochs: When the entire dataset is passed forward and backward through the neural network ~~only once~~. One epoch is with the entire

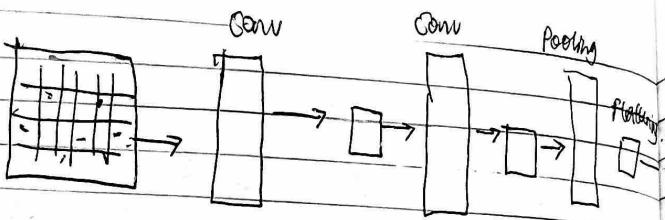
batch: the first flow of data is a batch single forward and backward propagation

iteration is the number of batches required to complete one epoch

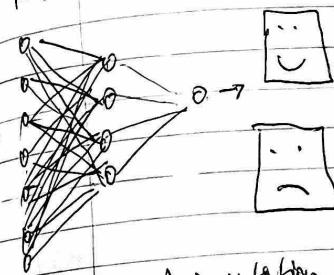
$$\text{training data} = 1000$$

$$\text{batch size} = 100$$

$$\text{iteration} = \frac{\text{training data}}{\text{batch size}} = 10$$



fully connected layer



Imbalanced dataset is the major issue with ML / DL and Data science

Accuracy is not used for imbalanced datasets

We use Classification report, confusion matrix

Augmentation

Used to create new data from existing data

Voice data → amplitude vs time pitch shift

Adding noise

Removing empty space

Adding volume
(+db, -db)

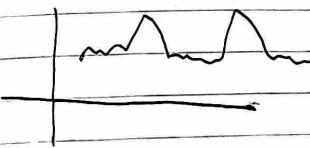


Image → Tiltting (vertical, horizontal)

Adding noise

Blurring

Zooming in/out

