

## - Understanding Machine Learning Algorithms: Handwritten Notes -

### Table of Contents

1. What is Machine Learning?
2. What are the Types of Machine Learning?
3. Supervised Machine Learning
4. Unsupervised Machine Learning
5. Reinforcement Learning
6. Semi-Supervised Learning
7. Steps in ML Project
8. Exploring Step 1 Data Collection
9. Exploring Step 2 Data Preparation
  - Exploratory Data Analysis
  - Data Preprocessing
  - Feature Engineering

**Notes by RaviTeja G**



## 10. Exploring Step 3 - Train Model on Dataset

- Types of Learning
- Under Fitting and OverFitting
- Regularization techniques
- Hyperparameter Tuning

## 11. Exploring Step 4 - Evaluation of a Model

- Evaluation Metrics
- Confusion Matrix
- Recall/Sensitivity
- Precision
- Specificity
- F1 Score
- AUC and ROC Curve
- Analysis of a Model

## 12. Supervised Learning

- Linear Regression
- Regularization Techniques
- Logistic Regression
- Decision Trees
- Ensemble Techniques
- Random Forests
- AdaBoost
- Gradient Boost
- XG Boost
- K-Nearest Neighbours
- Support Vector Machines
- Naive Bayes Classifiers

## 13. Unsupervised Learning

- Clustering Techniques
- K-Means Clustering
- Hierarchical Clustering
- DB Scan Clustering
- Evaluation of Clustering Models
- Curse of Dimensionality
- Principal Component Analysis

## 14. Cheat Sheet of Supervised and Unsupervised Algorithms

# Machine Learning

## Introduction:-

→ At a high-level, machine learning is simply the study or teaching a computer program/algo. how to progressively improve upon a set task that it is given.

So, more practically it is the study of how to build application that exhibit three, iterative improvement.

There are many ways to implement it and largely there are 3 major categories.

- 1) Supervised Learning.
- 2) Unsupervised Learning.
- 3) Reinforcement Learning.

And also there is semi-supervised learning.

Machine learning algorithms are the 'engine' or 'machine learning', meaning it is the algorithms that turn a dataset into a model.

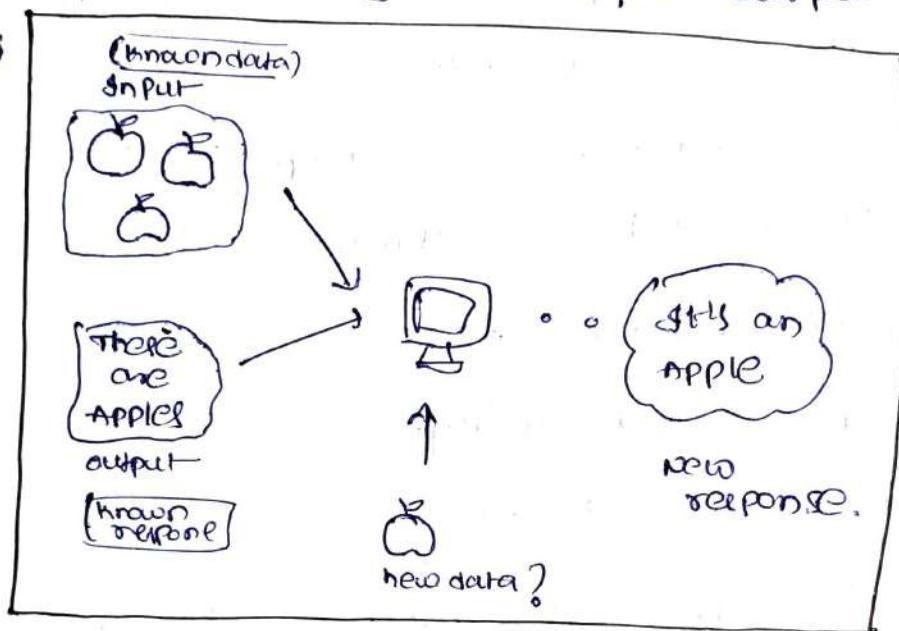
- \* which kind of algorithm works best and what to choose (supervised, unsupervised, classification, regression, etc...) depends on the kind of problem you're solving, the computing resources available & the nature of the data.

## Brickling the types!:-

### ① Supervised learning

Supervised learning is the most popular paradigm for performing machine learning operations. It is widely used for data where there is a precise mapping b/w input-output data.

i.e;



- Given data in the form of example with labels, we can feed a learning algorithm these 'example-label' pairs one by one. overtime, the algorithm will learn to approximate the exact nature of the relationship b/w examples & their labels. (it's called training the data)
- when fully trained, the supervised learning algorithm will be able to observe a new, never-before seen example & predict a good label for it.

→ And so supervised algorithms are called 'task-oriented'. As we provide it with more & more examples, it is able to learn more & more properly so that it can undertake the task & yield us the output more accurately.

It is exhibited in many of following applications  
Face Recognition

Face look app in our phone has been using a supervised learning algorithm that it is trained to recognize your face. Having a system that takes a photo, finds faces & guesses who that is in the photo (suggesting a tag) is a supervised process.

### Spam classification

Spam filter is a supervised learning system. fed emails examples & labels (spam/not-spam) these systems learn how to filter out malicious emails so that the user is not harassed.

### Advertisement popularity

Selecting advertisements that will perform well is often a supervised learning task. many ads you see as you browse are placed there because a learning algorithm said that those were of reasonable popularity (more clicked).

## Types of Supervised Learning

1) Regression

2) Classification

3) Neural networks

4) Boosting techniques.

### Regression

→ Regression algorithms

Predicts a 'continuous value' based on the input variables.

### Regression algorithms

1) Linear Regression

2) Polynomial Regression

3) Exponential Regression

4) Logarithmic Regression

### Applications

① Risk Assessment

② Score prediction etc.

### Classification

→ we use classification algorithms for predicting set of items class or category.

### Classification algorithms

Logistic Regression

1) K-nearest neighbours

2) Decision trees

3) Random forest

4) Support Vector machine (svm)

5) Naive Bayes

6) Ada boost  
7) XG boost

### Applications

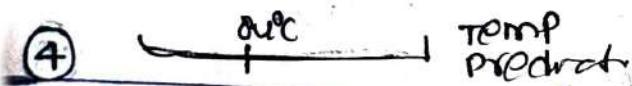
① Fraud detection

② Email spam detection

③ Image classification

④ Diagnosis etc..

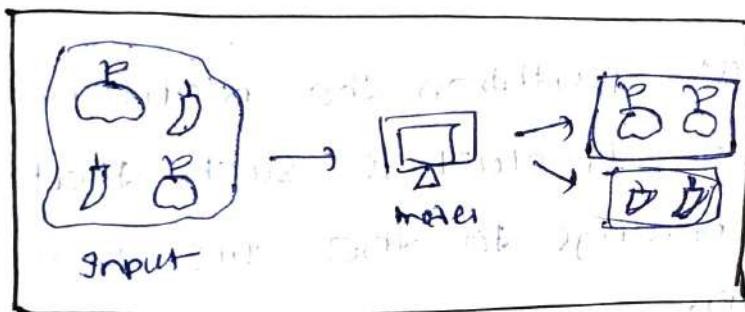
Hot | cold



## ② unsupervised learning

[Based on Hebbian learning)

- To easily understand, in case of unsupervised learning algorithm, the data is not explicitly labeled into diff. classes, that is no labels. The model is able to learn from the data by finding implicit patterns.
- unsupervised learning algorithms identify the data based on their density / structures, similar segments & other similar features.



### Types of unsupervised learning

- 1) Clustering
- 2) Dimensionality reduction & visualization
- 3) Anomaly detection
- 4) Association
- 5) Auto encoders

## 1) clustering

clustering, also known as cluster analysis, is a technique of grouping similar sets of objects in the same group that is diff from the objects in other group.

→ Some of the essential clustering techniques are as follows:

### a) K-means

In this we partition the  $n$  observations in the data into  $K$  clusters such that each observation belongs to the cluster with nearest mean.

### b) DBSCAN

It groups the data based on density. It groups together the points that are given in space & marks the outliers in low density regions.

### c) Hierarchical clustering

A hierarchy of clusters is built.

### d) mean shift

## 2) Dimensionality Reduction & visualization

### a) Principal Component Analysis (PCA):

Reduce data from more dimensions to lower dimensions while attempting to preserve the variance.

→ we reduce the size of data to extract useful information.

### b) t-distributed stochastic neighbor Embedding (t-SNE):

It is used to visualize high-dimensional data in a 2D / 3D space.

## 3) Anomaly detection

Anomaly detection techniques detect outliers in the unlabeled data under an assumption that most of the data example are normal by observing the instance that fit the remainder dataset.

### one-class classification

Train a model on only one-class, if anything lay outside of this class, it may be an anomaly.

Algorithms for doing so include -

- ① one-class k-means
- ② one-class SVM
- ③ Isolation Forest

#### 4) Association

we use association algorithms for allocating co-occurring items / events.

Association algorithms

• Apriori

→ 5) Autoencoder <sup>(Dimensionality Reduction)</sup>

Autoencoders take input data, compress data onto a code, then tries to recreate the input data from learned code.

→ Autoencoder can remove noise from visual data like email, video / medical scans to improve quality.

Note:

→ Autoencoder can also be used to find outliers (anomaly).

### ③ Reinforcement Learning

It doesn't require any data or labeled input or labeled output. All it has is the main data.

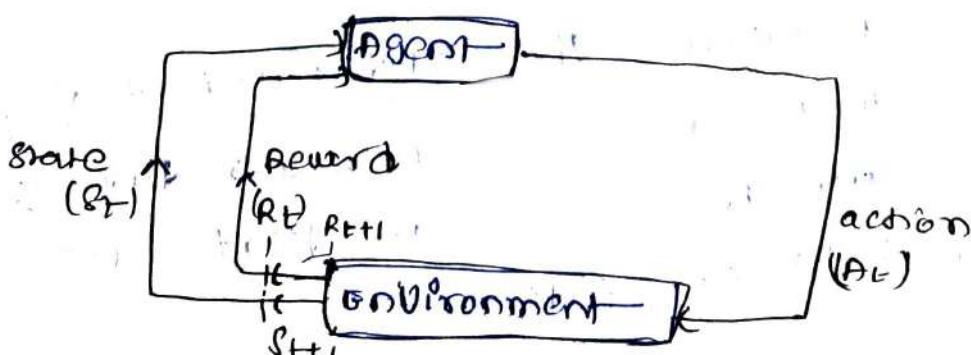
→ Use trial & error method. It finds in a trial then it learns & tries & it learns etc.

→ Games are a real example of RL.

Take a chess game model. All it has is the environment. First AI takes a step & feels the outcome & learns from it and so on.

→ In this kind of RL, AI agents are attempting to find the optimal way to accomplish a particular goal, or improve performance on specific tasks. As the agent takes action that goes toward the goal, it receives a reward.

overall aim → Product of best next step to earn the best reward.

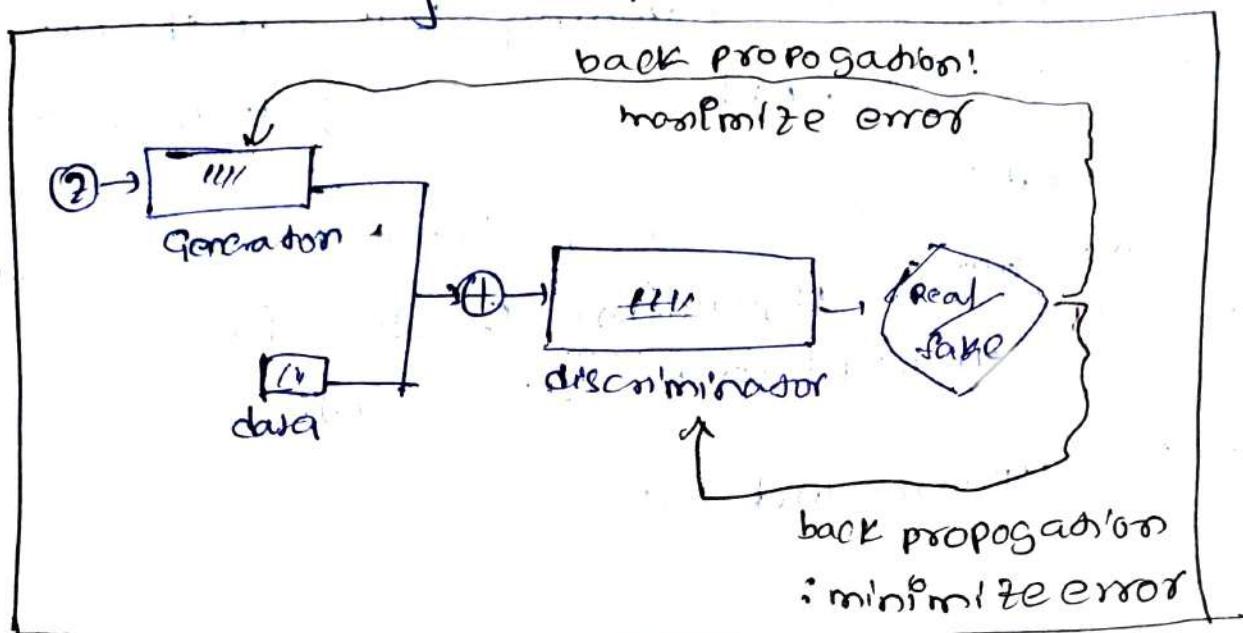


When there is no train. dataset, it learns from experience.

#### ④ Semi-supervised data

- semi-supervised learning is, for the most part, just what it sounds like: a training dataset with both labeled & unlabeled data.
- this method is particularly useful when extracting relevant features from the data is difficult & labelling examples is a time-intensive task for experts.
- A popular training method that starts with a very small set of labelled data is using 'generative adversarial networks' / [GAN's]  
like two deep learning networks in competition, each trying to outsmart the other. That's a GAN.
- one of the networks called generator, tries to create new data points that mimic the training data.
- the other network the discriminator, evaluates whether they are part of training data or fake.

→ the networks improve in a positive feedback loop - as the discriminator gets better at separating the fakes from originals, the generator improves its ability to create convincing fakes.



e.g.

- common situations for this kind of learning are medical images like CT scans (or) MRI's. A trained radiologist can go through & label a small subset of scans for tumors or diseases.
- it would be too time-intensive & costly to manually label all the scans - but the deep learning network can still benefit from the small proportion of labeled data & improve its accuracy compared to fully unsupervised models.

→ It's the basic knowledge about algorithms.

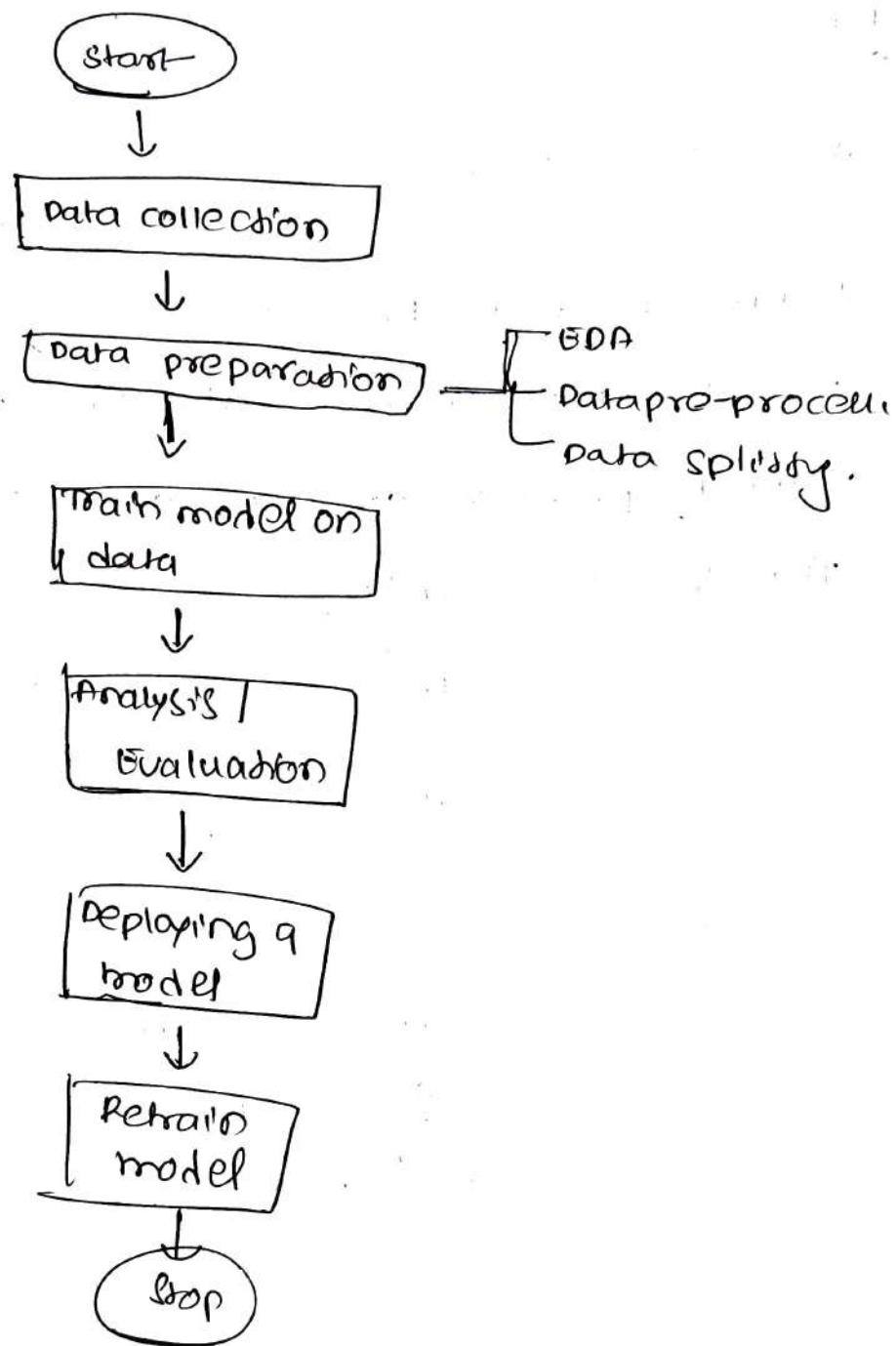
→ Before going deep into algorithms.

Better to know the flow chart of machine learning. I feel ML is like a Christopher Nolan movie, first time you don't understand anything but u should keep watching, At last everything makes sense.

Same works with ML projects, take a small proj on ~~some~~ a algorithm & see it. You will understand the path.

So, let's understand the path.

## ML Road map (Steps in machine learning project)



## \* Understanding the process :-

### Step ①

#### Data collection

##### a) Questions to ask?

- 1) what kind of problem are we trying to solve?
- 2) what data sources already exist?
- 3) what privacy concerns are there?
- 4) is the data public?
- 5) where should we store the data?

##### b) understanding the types of data

###### i) structured data?

Data which appear to be tabulated. Form  
can contain diff types of data.  
like

- i) nominal / categorical.
- ii) numerical
- iii) ordinal
- iv) time series

###### ii) unstructured data?

Data with no rigid structure.

(Image, video, natural language text, speech,

## Step②

### Data preparation

- It has 3 main process.
  - Exploratory data analysis (EDA), understanding data.
  - Data preprocessing, preparing your data to be modelled.
  - Data splitting.
- a) EDA, understanding the data

- what are feature variables (input) and the target variables (output)?

[e.g. For predicting heart diseases, feature variable may be a person's age, weight, heart rate etc. and target variable is whether or not they have heart disease]

- what kind of data do you have?  
structured, unstructured, categorical / numerical.
- Are there missing values? Should you remove them / fill them with feature imputation?
- what are outliers?  
How many of them are there?  
Are they out by much ( $3 + 2\sigma$ )?  
Why are they there?
- Are these questions you could ask a domain expert about the data?

## b) data preprocessing, preparing your data to be modeled

### ① Feature imputation

(filling missing values)

A ML model can't learn on data that isn't there.

i) single imputation: fill with mean / median or

ii) multiple imputation:  
model other missing values  
column

iii) KNN (K-nearest neighbor):  
and fill with what your model finds.

fill data with a value

from another example which is similar.

iv) many more such as random imputation,  
last observation carried forward (for time series),  
moving window, most frequent.

∴ there are ways to handle missing values.

### ② Feature encoding

(turning values into numbers)

A ML model requires all values to be numeric

i) OneHotEncoding:-

Turn all unique values into 1's & 0's where the target value is 1 & the rest are 0's.

e.g. car colors green, red, blue,

a green car's color feature would be [1, 0, 0]

& a red one would be [0, 1, 0]

& a blue one [0, 0, 1]

### ii) Label Encoder :-

Turn labels into distinct numerical values

Eg:-

If your target variables are diff.

animal as dog, cat, bird then 0, 1, 2 respectively

### iii) Embedding encoding :-

Learn a representation amongst all the diff. data points.

Eg:-

A language model is a representation of how diff. words relate to each other.

embeddings are also becoming more widely available for structured data.

## ③ Feature normalization (scaling) or standardization :-

when your numerical variables are on diff scales. (like size of land b/w 1000 - 20000 sq ft).

In such case some ML models don't perform well. Scaling / standardization help to fix this.

### i) Feature scaling (normalization) :-

shifts your values so they always appear b/w 0 & 1.

This is done by subtracting min value & divide with range. (value b/w 0 & 1 but numerical relationship is still there)

### ii) Feature standardization :-

standardized all values so they have a mean of 0 & variance of 1.

This is done by subtracting mean & divide with SD.

values doesn't end up b/w 0-1 (st. var.).

Standardization is more robust to outliers than scaling.

## ④ Feature Engineering

→ Transform data into potentially more meaningful representations by adding domain knowledge.

### i) Decompose

Like, turn a date (such as 2020-06-18 to 20:16: etc.) into hour-dr-day, day-dr-week, day-dr-month, ts-holt etc.

### ii) Discretization

(Turning larger groups to smaller groups)

For numerical variables,

like age you may want to move it into bucket such as over-50 (under-50, 21-30, 31-40 etc.) this process is known as BINNING.

For categorical variables,

such as car color, this may mean combining colors such as ('light-green', 'dark-green', 'lime green') into a single green

### iii) Crossing & Interaction features

(Combining two/more features)

→ depending on situation sometimes adding / subtracting two features could help.

### iv) Indicator features

like using  
is-Paid-traffic  
for Paid

[using other parts of data to indicate something potentially significant]

+ Create a X-missing feature for whenever a column X contains a missing value.

## ⑤ Feature Selection :-

- Selecting the most valuable feature or your dataset to model.
- Potentially reducing overfitting & training time (less overall data & less redundant data to train on) & improving accuracy.

### i) Dimensionality reduction:-

A common dimensionality reduction method, PCA takes a larger no. of dimension (features) & uses linear algebra to reduce them to less dimensions. For example, say you have 10 numerical features. You could run PCA to reduce it to 3.

### ii) Feature Importance :-

Fit a model to a set of data, then inspect which features were most important to the results, remove the least important ones. (can use "TPOT" for this)

## ⑥ Dealing with Imbalance :-

Does your data have 10,000 examples or one day, but only 100 examples the other?

- collect more data (if possible)
- use scikit-learn-contrib imbalanced-learn package
- use SMOTE (synthetic minority over-sampling technique) which creates synthetic samples of your minor class to try & level the playing field.

→ So this is all about <sup>data</sup> Preprocessing.

### c) Data Splitting

1) training set (usually 90-80% of data):

model learns on this.

2) validation set (typically 10-15% of data):

model hyperparameters are tuned on this.

3) test set (usually 10-15%): (Dev set)

model's final performance is evaluated on this.

④ Do not use this dataset to tune the model

#### a. why shouldn't we tune on test set?

A.

To do the regularization we have to take 100 diff. models & try on the test data. And so last we finds (suppose) the best hyperparameter that produces a model with 10% generalization error, say just 5%.

And you launch it to production & found it produces 15% errors.

④ Now, here, is what happened as you trying every hyperparameter on the same test data. Eventually your model learned the test data & fitted to it thus reducing generalization errors but not improving models.

This is called "overfitting".

→ So, to avoid we take validation set / cross-validation

④

### Step③

#### Train model on dataset

It has 3 main process.

a) choose an algorithm

b) overfit the model

c) Reduce overfitting with regularization.

a) Algorithm can be choosed by understanding supervised, unsupervised (Explained previous)

b) Type of learning

i) Batch Learning

All or your data exists in a big static warehouse & you train a model on it.

You may train a new model once per month once you get new data. Learning may take a while & isn't done often. Runs in production without learning (can be retrained).

ii) Online Learning

Your data is constantly being updated. You constantly train new models on it. Each learning step is usually fast & cheap.

Runs in production & learns continuously.

### iii) transfer learning

Take the knowledge one model has learned & use it with your own. Give you the ability to leverage SOTA (state of the art) models for your own problems.  
Helpful if you don't have much data / want compute resources.

### iv) Active learning

Also referred to as "human in the loop" learning. A human expert interacts with a model & provides updated to labels for samples which the model is most uncertain about.

### v) Ensembling

Not really a form of learning,  
more combining algorithms which have  
already learned in some way to get better  
results.

### Underfitting

Happens when your model doesn't perform as well as you'd like on your data.  
try training for a longer & more advance model.

## overfitting

Happens when your validation loss (how your model is performing on the validation dataset, lower is better) starts to increase, or if you don't have a validation set,

\* Happens when the model performs far better on the "training set" than on the "test set".  
(e.g. 99% accuracy on training set, 67% accuracy on test set) → overfitting the training set.

Fix through various regularization techniques.

Regularization → [Techniques to prevent/reduce overfitting.]

i) L1 (Lasso) & L2 (ridge) regularization:

→ L1 regularization sets unneeded feature coefficients to 0.

→ L2 constrains a model's features (won't set them to 0)

ii) Dropout:

Randomly remove parts of your model.

So the network has to become better.

iii) Early stopping:

Stop your model from training before the validation loss starts to increase too much / more generally, any other metric has stopped improving. Early stopping usually implemented in the form of a model callback.

#### iv) Data Augmentation

- manipulate your dataset in  
artificial ways to make it harder to learn.
- For example, if you're dealing with images, randomly rotate, skew, flip & adjust the height of your images.
  - This makes your model have to learn similar patterns across different styles of the same image (harder).
  - Use functions like `ImageDataGenerator` in `keras` or `transforms` in `torchvision`.

#### v) Batch Normalization

standardize inputs as well as adding two parameters (beta, how much to offset the parameter for each layer & epsilon to avoid division by zero) before they go into next layer. This often results in faster training speed since the optimizer has less parameter to update.

maybe a replacement for dropout in some networks.

## Hyperparameter Tuning

→ Run a bunch of experiments with different model settings & see which works best.

### a) Setting a learning rate:

(often the most important hyperparameter)  
Generally

High learning rate = algorithm rapidly adapts to new data

Low learning rate = algorithm adapts slower to new data  
(e.g. for transfer learning)

#### i) Finding the optimal learning rate

Train the model for a few hundred iterations starting with a very low learning rate ( $10^{-6}$ ) & slowly increase it to every large value (e.g. 1). Then plot the loss vs learning rate (using a log-scale for learning rate), you should see a U-shaped curve, the optimal learning rate is about 1-2 notches to the left of the bottom of the curve. (P326 - Handwritten)

ii) Learning rate scheduling involves decreasing the learning rate slowly as model learns more (get closer to convergence). (see Adam Optimizer)

#### iii) cyclic learning rates

Dynamically change the LR up & down b/w some thresholds & potentially speed-up training.

other hyperparameters you can tune

- 1) No. of layers (deep learning networks)
- 2) Batch size [How many es or data your model sees at once.]
- 3) No. of trees (decision tree algos)
- 4) No. of iterations [How many times model goes through the data]  
Depends on algorithm. (can use early-stopping)

## Step ④

### Analysis / Evaluation

to estimate how well a model will generalize to out-of-sample data.

#### a) Evaluation metrics:-

##### i) Classification

ROC & AUC curves

Accuracy

Precision

Recall

F1 score

Confusion matrix

##### ii) Regression

MSE (mean squared error)

MAE (mean absolute error)

R<sup>2</sup> (R-squared)

Task-based metric (create one based on your specific problem).

#### i) Classification accuracy:-

Percentage of correct predictions. ( $\frac{Y_{\text{pred}}}{Y_{\text{true}}}$ )

##### Null accuracy

Accuracy that could be achieved by always predicting the most frequent class.

It's a good way to know the min. we should achieve with our model.

It gives a baseline to compare our

model's accuracy.

Classification accuracy is easy, but it doesn't tell you the underlying distribution or response class.

And it doesn't tell you what type of error your classifier makes. (27)

## a) confusion matrix

FP - False Positive  
FN - False Negative.

Predicted			
		0	1
Actual 0	0	TN	FP
	1	FN	TP

they actually don't have corona. But we predicted they have.

they actually have, but we say they don't

TP% which are actually correct

TN% which are actually wrong

FP% they are actually wrong,

But predicted as true.

FN% they are actually true,  
But predicted as False

Classification

accuracy

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

How often is the classifier correct?

talks about  
correct.

misclassification rate

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

Recall / sensitivity

[True positive rate]

→ How "sensitive" is the classifier to detecting positive instances?

→ when the actual positive is tested, how often it is correct?

Predicted Pos.

TP

$$\text{Recall} = \frac{TP}{TP + FN}$$

→ It depends on FN,

which were actually true, but predicted false.

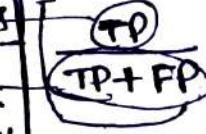
out of total positive ones,  
How many positive I am able to Recall

$$\text{Actual Pos} = TP + FN$$

## Precision & Recall linked with TP

### Precision

Really Pos  
Total Pos.  
actually



$$\frac{TP}{TP+FP} \rightarrow \text{Really Pos}$$

$\frac{TP}{TP+FN}$  → Predicted Pos

→ when a positive value is predicted,  
How often is the prediction correct?

→ How "precise" is the classifier when  
Predicting positive individual.

out of all  
Predicted Pos,  
How many  
really Pos.

$$\text{Precision} = \frac{TP}{TP+FP}$$

→ It depends on FP,  
which were actually  
negative, but  
predicted true.

Q) which metrics to focus on?

Eg ①

1 = COVID 19 (+)

0 = - (Healthy)

Pred		
0	1	
0	TN	FP
1	FN	TP

Healthy predicted as sick.

Sick predicted as healthy.

→ now analyze the cost

→ we know a sick predicted as healthy  
could spread a lot more easily.

So,  $\text{cost of FN} > \text{cost of FP}$

Hence it depends on FN, we go for "Recall".

Eg ②

1 = Spam

0 = Not Spam

Pred		
0	1	
0	TN	FP
1	FN	TP

Not spam predicted as spam.

Spam got into mail box.

→ we know if a not spam (imp. mail) goes to spam,  
that could cause lot of damage.

So,  $\text{cost of FP} > \text{cost of FN}$

Hence, it depends on FP, we go for "Precision".

(29)

Sensitivity P.

$$\boxed{\frac{TN}{TN + FN}}$$

False positive rate

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

FB-Score

when you want to consider both  
 PP & PN like we need both  
 Sensitivity (recall) & precision,  
 In such case we choose F-Beta.

General Formula

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

$$F_B = \frac{(1+\beta^2) \times TP}{(1+\beta^2) \times TP + \beta^2 \cdot FN + FP}$$

In beta, there is possibility for 3 ranges.

- 1) when FP & FN are equally important, like if you can't compromise even a bit both are crucial, then we set beta=1 which is F1-score.

$$\boxed{F1\text{-score} = \frac{2 \times \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}}$$

- 2) when both are worse, but FP is drastic than FN. then we go for beta = 0.5 (usually 0-1)
- 3) when both are worse, but FN is drastic than FP. then we go for beta = 2 (usually 1-10).
- 

### Note

we can also adjust classification threshold to modify performance of classifier.

### Roc curves

can see how specificity & sensitivity are affected by various thresholds, without actually changing the threshold.

### AUC curve & Area under curve

If you randomly choose one + one, AUC represents the likelihood that your classifier will assign a higher predicted probability to the positive observation.

---

## Analysis of model

### 1) Feature Importance

Which feature contributed most to model

Should some be removed?

### 2) Training / Inference time

How long does model take to train?

Is this feasible?

How long does inference take?

Is it suitable for production?

### 3) using what-if tool

what-if rechanges bring in data?

How does this effect the outcome?

### 4) Bias / variance trade-off

→ High bias results in underfitting &  
a lack of generalization to new sample

→ High variance results in overfitting as  
the model finds patterns in the  
noise,

Step 5

## Deploying a model

Put the model into production & see how it goes. Evaluation metrics in notebook are great but until it's production, you won't know how it performs for real.

→ Step 6

## Retrain model

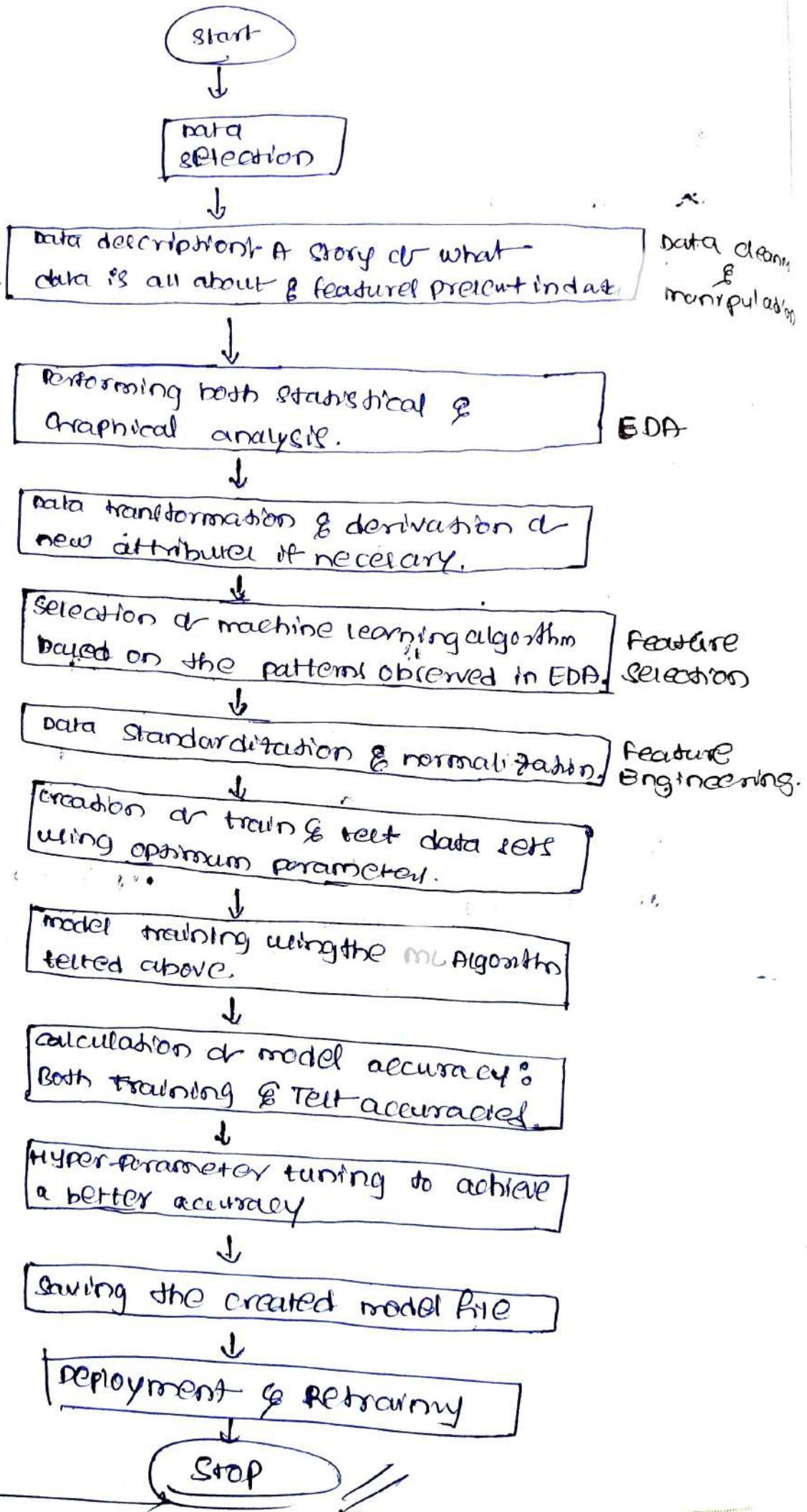
→ See how the model performs after serving based on various evaluation metrics & revisit the above steps as reqd.

→ You will find your model's prediction starts to lag or drift, as in when data source change or upgrade. This is when you will want to retrain it.

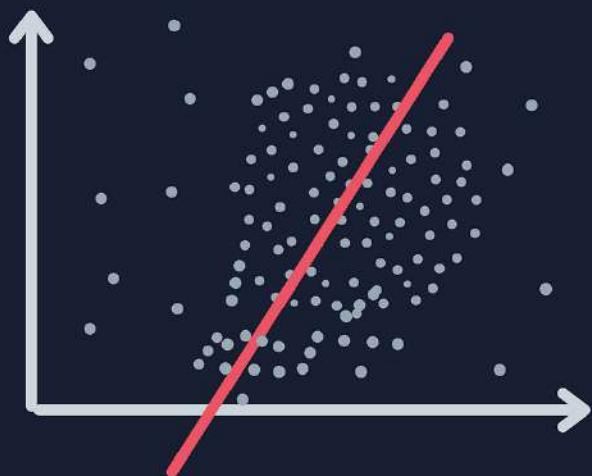
Stop / NO END

DO IT

so, basically the flowchart :-



# Linear Regression



## Table of Contents

1. What is Linear Regression
2. Understanding with an example
3. Evaluating the fitness of the model
4. Understanding Gradient descent
5. Understanding Loss Function
6. Measuring Model Strength
7. Another Approach for LR - OLS

## understanding the algorithms

### ① Simple Linear Regression

Linear model in Regression  
(supervised learning)

Simple linear regression assumes that a linear relationship exists b/w the response variable & the explanatory variable, it models this relationship w/ a linear surface called a "hyperplane".

→ This will be clear with an ex.  
lets take the height & weight sample of few people

Height	Weight
5.6	60
5.7	62
5.8	63
5.9	62
6.1	64
6.2	75

$$\begin{matrix} H \propto W \\ \downarrow \\ W = aH \end{matrix}$$

Now, if I want to know what would be the weight of a person with 6.0 Height?

→ By observation or second like

$$H \propto W$$

→ To remove proportionality we can multiply a constant

$$W = mH$$

→ But, it may also differ by some value,  $l_0, \pm b$

$$W = mH + b$$

→ Now, this is the relation of height & weight.

constants

$w = mH + b$  → now, here if I know the value of 'm' & 'b'. Then I can predict the weight.

Given height

Now, main aim is to find out value of 'm' & 'b'.

→ Let's take a sample, (5.6, 60), (5.7, 62).

$$\begin{aligned} 60 &= m(5.6) + b \quad \text{①} \\ 62 &= m(5.7) + b \quad \text{②} \end{aligned} \quad \left. \begin{array}{l} \text{By this we get} \\ m = 20, b = -52. \end{array} \right.$$

So, our equation  $w = 20H - 52$

lets say with a known height 5.9.

$$w = 20(5.9) - 52 = 66 \text{ kg.}$$

Here we got '66', but we are expecting 62.

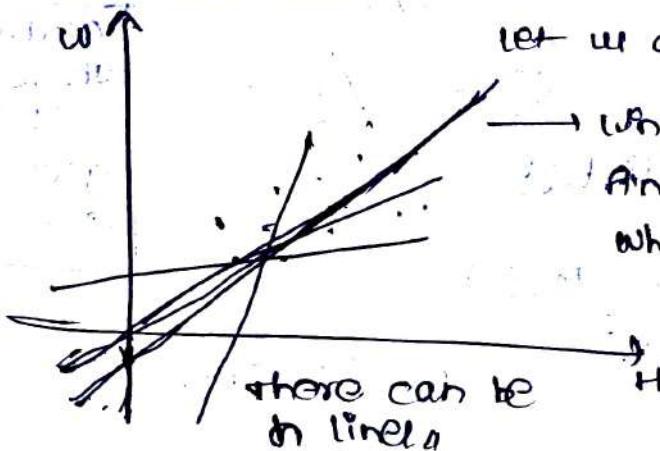
④ It may be or we change ['m' & 'b'] then we could get an exact value (approx.)

Now, Here we get the definition of ML as establishing a mathematical relationship or data. So, that it can predict a new data.

→ So, when we say a model,

④ It is nothing but a mathematical relationship which can predict new data.

Plotting ( $H$  &  $w$ )



Let us assume we have a best fit line.

→ Linear regression always tries to find the linear relation bw variable which is a straight line.

Hence,  $y = mx + c$

→ we could find n diff lines.

But that one line which can predict

all the values approx. is the "Best fit line"

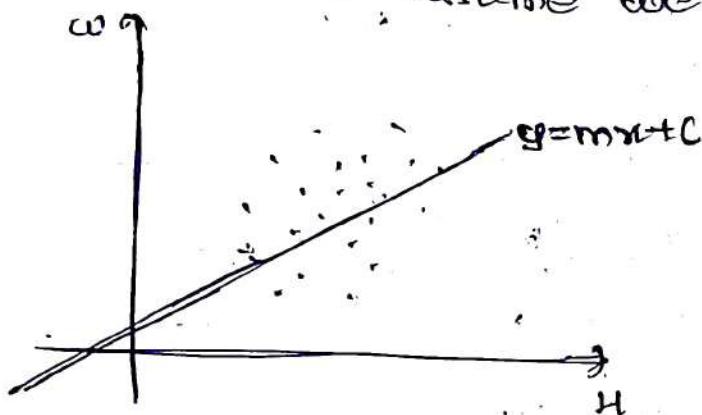
$y = mx + c$  → And this best fit line is called model  
→ we have to find the 'm' & 'c'.

m-slope, c-Intercept

\* Evaluating the fitness of the model with a

cost function

→ First, let us assume we have a best fit line.



But, now we have doubt.

How did we opt this line?  
Why this line?  
What about other lines?

→ Now, if it's the best fit line, then procedure  
if I want the weight or 'w'.

→ It is nothing but the value of y-coordinates  
 $y = mx + c$  like worth of coordinate '6,0'

→ If this line is  $w = 20H - 52$ , and at  $H = 5.9$ ,

I got  $w = 66$ , but actually  $w = 62$ .

So, I can see there is a known error.

$$\text{loss} = |(y - \hat{y})| \quad \text{Predicted}$$

Q) Why did we get this loss?

A) may be best fit isn't being taken.

Always  
there  
will be  
a slight  
error

## Gradient Descent

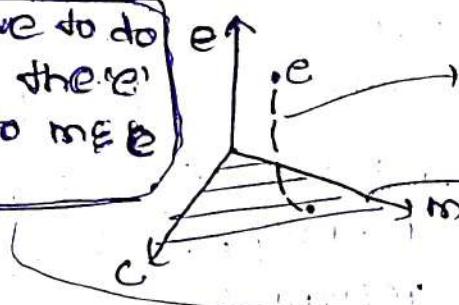
→ Now, to improve my model strength  
I have to reduce the loss.

and this loss depends on line "c" taken.

→ So, I am supposed to shift in  $m$  &  $c$  such that I end up reducing the loss.

We have loss( $e$ ), slope( $m$ ), constant( $c$ ).

All I have to do  
is bring them  
down to  $m$  &  $c$  plane



should bring down &  
find the new coordinate

In this plane we  
have no errors.

]

② Slowly, by examining each input  $m$  &  $c$  gets changed and eventually loss decreased. This process is called learning.

So,

$$m_{\text{new}} = m_{\text{old}} - \eta \Delta m$$

$$c_{\text{new}} = c_{\text{old}} - \eta \Delta c$$

$m_{\text{new}}$  - changed one, new

$m_{\text{old}}$  - old  $m$  value.

$\eta$  - learning rate

$\Delta m$  - Error in  $m$  (de/dm)

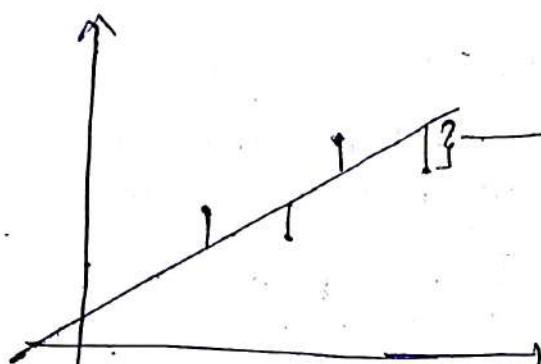
$$\eta \rightarrow (0.0001 - 10) \text{ range.}$$

→ So, here we are taking a new  $m$ .

By subtracting a fraction of the error rate from old  $m$  that fraction is learning rate.

So, now I am supposed to choose the  $m$  &  $c$ , which gives a least possible error. Then, it will be the best fit line.

→ A cost function, also called a loss func., is used to define & measure the error of a model.  
 the diff. b/w the weights predicted by model & the observed weights in training set are called "residuals", or training error / loss.



→ This is the residual / loss which is  $|y_i - \hat{y}_i|$

→ we have to reduce this residual  
 Actually, we have to reduce residual of all data points.

→ we can produce best weight predictor model by minimizing the 'sum of residuals'.

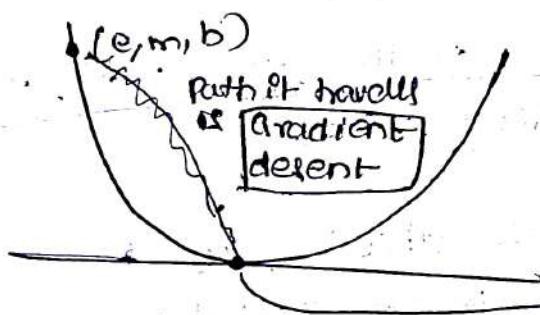
$$(\text{Sum of residuals}) \quad E_n = \sum_{i=1}^n [y_i - (m x_i + c)]$$

so, to cancel the negative value, we take square

$$R(x) = \sum_{i=1}^n y_i^2 = \sum_{i=1}^n [y_i - (m x_i + c)]^2$$

sum of squares  
of residuals,

now, it is a quadrature in 3D.



→ so, we are supposed to shift the  $(c, m, b)$  to the origin  $(0, m, b)$ .

$$\text{where } \frac{dc}{dm} = 0, \frac{dc}{db} = 0.$$

As we can see residual is both a func. of  $m$  &  $b$ ,  
so, differentiating partially w.r.t  $m$  &  $b$  will  
give us:

$$\frac{\partial R}{\partial m} = \sum_{i=0}^n \alpha_i (b + mx_i - y_i)$$

$$\frac{\partial R}{\partial b} = \sum_{i=0}^n \alpha_i (b + mx_i - y_i)$$

so, we know for best line, residual should be min.  
minima or max. occurs where derivative = 0.  
i.e.,

$$\sum_{i=0}^n \alpha_i (c + mx_i - y_i) = 0$$

$$\sum_{i=0}^n \alpha_i (b + mx_i - y_i) = 0$$

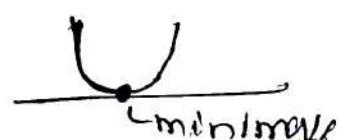
$$\sum_{i=0}^n \alpha_i c + \sum_{i=0}^n \alpha_i m x_i^2 + \sum_{i=0}^n \alpha_i y_i x_i = 0$$

$$\sum_{i=0}^n \alpha_i c + \sum_{i=0}^n \alpha_i m x_i - \sum_{i=0}^n \alpha_i y_i = 0$$

The same eq. can be written in matrix form as:

$$\begin{bmatrix} \sum_{i=0}^n \alpha_i & \sum_{i=0}^n \alpha_i x_i^2 \\ n & \sum_{i=0}^n \alpha_i x_i \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=0}^n \alpha_i y_i \\ \sum_{i=0}^n \alpha_i y_i \end{bmatrix}$$

ideally, if we have an eq. or one  
dependent & one independent variable  
the minima will look like



dependent var  $y = m x + c$  independent var

→ the new values for slope & intercept  
are calculated.

repeat until convergence

{

$$m_{\text{new}} = m_{\text{old}} - \eta \left[ \sum_{i=1}^n (h_0 \cdot x_i - y_i) x_i \right] \xrightarrow{\frac{\partial E}{\partial m}}$$

$$c_{\text{new}} = c_{\text{old}} - \eta \left[ \sum_{i=1}^n (h_0 \cdot x_i - y_i) \right] \xrightarrow{\frac{\partial E}{\partial c}}$$

}

Accuracy this is how LR works.

④ measuring model strength RSE - residual sum of squares RSS - total sum of square

The R-squared ( $R^2$ ) statistic provides a measure of fit.

It takes the form of proportion - (proportion explained).

$$R^2 = 1 - \frac{RSS}{TSS}$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{Product of error changes}$$

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2 \quad \text{mean.}$$

It  $R^2 = 0.75$

It says that our model fits 75% of data points.

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

→ As R<sup>2</sup> is just a linear eq., as we increase the no. of independent variables in our eq., the R<sup>2</sup> increases as well.

But, that doesn't mean that the new independent variables have any correlation with the output variable.

i.e., R<sup>2</sup> will increase, but it is not necessarily model yields better results.

To rectify this, we use adjusted R<sup>2</sup> value which penalises excessive use of such features which do not correlate with the output data.

$$R^2_{adj} = 1 - \frac{(1-R^2)(N-1)}{N-P-1}$$

where

P = no. of predictors

N = total sample size.

$$\text{If } P=0, \quad R^2_{adj} = R^2.$$

Note!

→ Using training data to learn the values of the parameters for simple linear regression that produce the best fitting model is called ordinary least squares (OLS) or linear least square.

Another approach to find m & C.

(to solve eqs)

→ variance is a measure of how far a set of values all spread out.

→ covariance is a measure of how much two variables change together. If the variables increase together, their cov is positive.

o - no relation, -ve → one increase & one decrease

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

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

→ variance of explanatory variable.

→ covariance of the response, explanatory variables.

$$\beta = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

$$y = \beta x + \alpha$$

$$\alpha = \bar{y} - \beta \bar{x}$$

now can see

Note :-

→ the independent variables are uncorrelated with the residual term, also known as "Exogeneity".  
This in layman term generalizes that in no way should the error term be predicted given the value of independent variable.

the error terms have a constant variance.

### Homo<sup>s</sup>cedasticity

- the error terms are normally distributed.
- no multicollinearity, i.e., no independent variables should be correlated with each other or affect another.

### model confidence

(or) is linear regression a low bias / high variance model

(or) a high bias / low variance model?

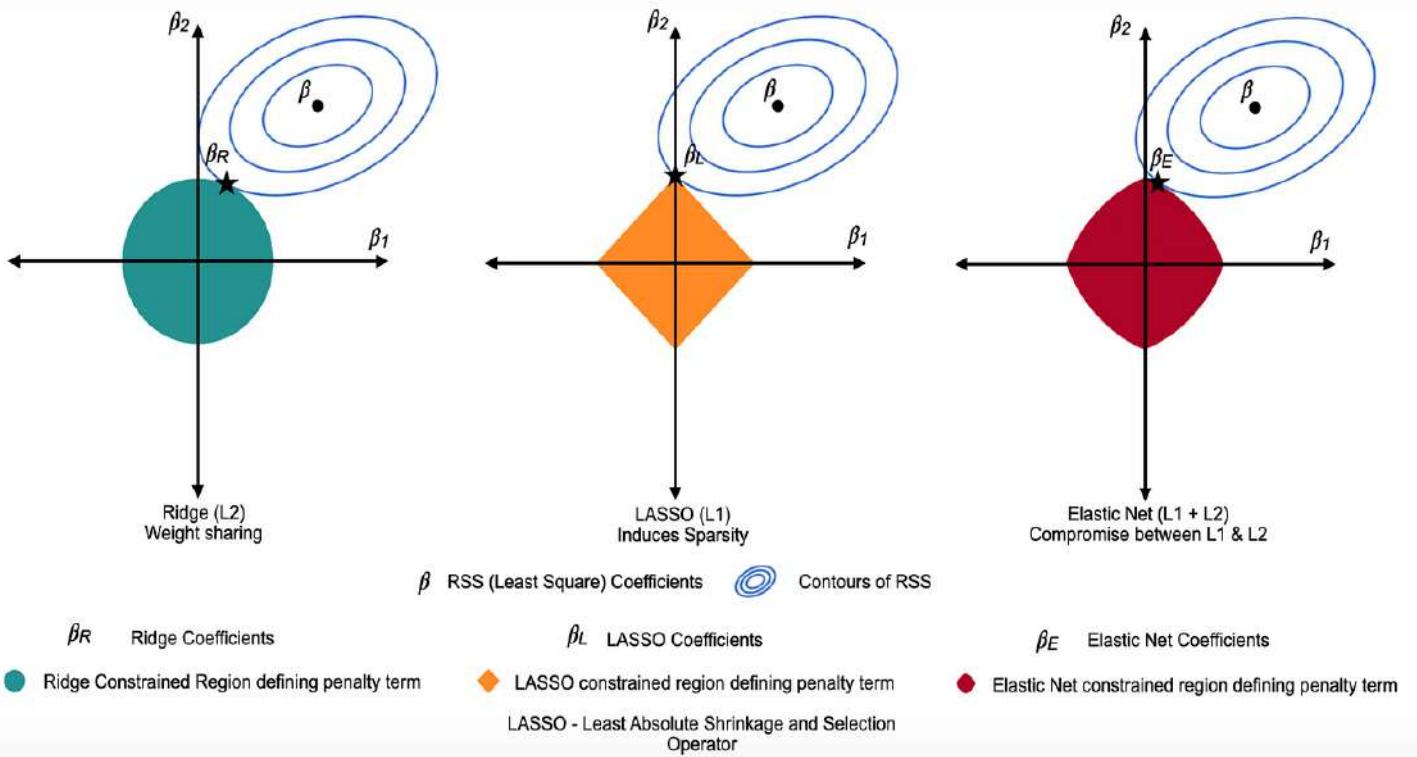
→ It's a "high bias / low variance" model.

→ Even after repeated sampling, the fit line won't stay roughly in the same pos. (low variance)

→ But, the avg'd the models created after repeated sampling won't do a great job in capturing the p'st relationship.

→ Low variance is helpful when we don't have well behaved data.

```
import statsmodels.formula.api as smf  
model_name = smf.ols(formula='y ~ b1, b2')  
model_name.conf_int()
```



## Table Of Contents

1. Understanding Multicollinearity
2. Variance Inflection Factor
3. Regularization
4. Lasso - L1 Form
5. Ridge - L2 Form
6. Elastic Net
7. Difference Between Ridge and Lasso
8. When to use Ridge/Lasso/Elastic Net
9. Polynomial Regression

## Multi-collinearity

- we can define multi-collinearity at the situation where the independent var have strong correlation among themselves.
- the co-eff. in a linear Reg model represent the extent or change in  $y$  when a certain  $x_1(x_1, x_2, x_3, \dots)$  is changed keeping other constant. But if  $x_1 \& x_2$  are dependent then still assumption itself is wrong that (we are) changing one vars keeping others constant all the dependent var will also be changed.

It means model becomes a bit flawed.

$$y_2 = m_1 x_1 + m_2 x_2 + c$$

while,

$$m_2 = a m_1 + d$$

Independent var ( $m_1, m_2$ )  
are related.

more than one linear eqn,

Hence called

multi-collinearity

→ we have redundancy in our model as two variables (or more than two) are trying to convey the same information.

- ⇒ → As the extent of collinearity increases, there is a chance that we might produce an "overfitted model".
- An overfitted model works well with the test data but its accuracy fluctuates when applied to other datasets.

④ can result in a dummy variable trap.

→ By the heatmap we can check collinearity.

Generally a correlation greater than 0.9 (or)

less than -0.9 are to be avoided.

#### ⑤ Variance Inflation Factor (VIF)

Regression on one X var against other X variables.

$$\text{VIF} = \left( \frac{1}{1-R^2} \right)$$

$$\text{VIF} < 5$$

✓ perfect

→ If  $\text{VIF} > 5$ , extreme correlation is avoided.

#### Remedies for multi-collinearity

- ① de-coeffing: if corr not extreme / the var not used the ignore
- ② remove one-var like in dummy var trap (one hot encoding)
- ③ combine correlated vars - combine variables.
- ④ Principle component Analysis (PCA)

## \* Regularization

→ when we use Regr. models to train some data,  
→ there is a good chance that the model will overfit  
the given training dataset.

→ Regularization helps sort this overfitting problem  
by restricting the degree or freedom of given eq.  
P.s,

simply reducing the order degree or a polynomial  
by reducing their corresponding weight.

### Reform for Regularization

→ in a linear eqn, we don't want huge coeff,  
as a small change in coeff can make a large diff.  
so regular constraints weights or such features to  
avoid overfitting.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$$RSS = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j \cdot x_{ij})^2$$

To regularize a model,

shrinkage penalty is added to cost function,

## LASSO (Least Absolute Shrinkage & Selection Operator) L1 form

→ LASSO regression penalizes the model based on the sum or magnitude of the coefficients.

The regularization term is

given by,  $\boxed{\text{regularization} = \lambda * \sum |B_j|}$ .

shrinkage factor

→ Loss after regularization is -

new loss func.

$$\boxed{\text{RSS} + \lambda \cdot \sum_{j=1}^p |B_j| = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j \cdot x_{ij})^2 + \lambda \cdot \sum_{j=1}^p |B_j|}$$

Here, loss is not calculated just by previous values ( $y_i - \hat{y}_i$ ), now it's been split into several parts which decreases loss,

## Ridge Regression (L2 form)

→ Ridge Regression penalizes the model based on the sum of squares of mag of coeff.

$$\boxed{\text{regularization} = \lambda * \sum |B_j|^2}$$

$\lambda$  can be calculated by cross validation

→ loss after regularization,

$$\boxed{\text{RSS} + \lambda \cdot \sum_{j=1}^p B_j^2 = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j \cdot x_{ij})^2 + \lambda \cdot \sum_{j=1}^p B_j^2}$$

Note :- consider  $B_1 \in R_2$  be coeff of LR if  $\lambda = 1$ .

- ① for LASSO,  $\boxed{\beta_1 + \beta_2 \leq S}$  → where 'S' is the max value the eq. can achieve.
- ② for Ridge,  $\boxed{\beta_1^2 + \beta_2^2 \leq S}$

## Elastic net

→ middle ground b/w Ridge & Lasso

→ the regularization term is a simple mix of both Ridge & Lasso, and you can control how much  $\alpha$ .

$$\text{Lenet}(\hat{\beta}) = \frac{\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1})^2}{2n} + \lambda \left[ \left( \frac{1-\alpha}{2} \right) \cdot \sum_{j=1}^m \hat{\beta}_j^2 + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right]$$

→ where  $\alpha$  is the mixing parameter

Ridge ( $\alpha=0$ )

Lasso ( $\alpha=1$ )

## ④ Difference between Ridge & Lasso

→ Ridge Regression shrinks the coefs for those predictors which contribute very less in the model but have huge weights, very close to 0.

$y = 0.0016 - \text{very less}$   
 $R^2 = 0.0016 - \text{very less}$

thus final model will still contain all those predictors, though with less weight.

→ whereas in Lasso, the L1 penalty does reduce some coefs exactly to zero, when we use a sufficiently large tuning parameter  $\lambda$ .

so, in addition to regularizing,

Lasso also performs feature selection.

## → why use regularization?

- it helps to reduce the variance of the model, without a substantial increase in the bias.
- if there's variance in the model that means the model won't fit well for dataset other than training data (which is called overfitting).
- the tuning parameter ( $\lambda$ ) controls this bias & variance trade off, selected via cross-validation.
- when ' $\lambda$ ' is increased to certain point, it reduces the variance without losing any imp prop of data.
- But after certain point model will start lossing imp prop, which will increase bias in the data.

as what should be used plain linear, Ridge, Lasso or Elasticnet?

- Ans
- it is always preferable to have atleast a little bit of regularization, so generally avoid plain linear Reg.
  - Ridge is a good default,
  - But if you prefer that only a few features are actually useful, you should prefer Lasso/EN, since they tend to reduce the useless feature weights down to zero.
  - In general Elastic net is preferred over Lasso since Lasso may behave erratically when the no. of features is greater than no. of training instances or when several features are strongly correlated.

whenever, mean is not 0;  
it's better to do standard scalar

Notes 1

① standardScaler()

$$\frac{x - \bar{x}}{s_x}$$

$\bar{x} = 0$

$s_x \neq 1$

② Two ways to find multicollinearity —

correlation  
VIF

③ 4 ways to remove multicollinearity.

④ Train set, high accuracy; Test set, low accuracy.  
It's overfitting, so regularize it.

→ Even after checking Lasso and Ridge,

it still gives the same accuracy values,  
then you can't tell if it's overfitted.

e.g. just by having huge difference in  
train accuracy & test accuracy than it  
doesn't just mean overfitted. Test before  
judging! In this case model is not problem,  
more data is needed.

⑤

The Linear Regression does not talk  
about the degree or polynomial eqn. in  
terms of dependent var.

Instead, it talks about the degree of coeff.

$$y = a + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

It's not talking about power of  $x_1$ ,  $x_2$ ,

but the power of  $a, b_1, b_2, \dots$  and their powers

Hence it is linear regression

## Polynomial Regression

→ It's a mechanism to predict a dependent var based on the polynomial relationship with the independent variable.

on the eq,  $y = a + bx + cx^2 + \dots + nx^n + \dots$

max. power of  $x$  is called degree of polynomial.

Deg 1 →  $y = a + bx$

Deg 2 →  $y = a + bx + cx^2$

when do we?

④ If data is scattered in a polynomial curve shape (not linear) then keep trying different degrees until the line fits the data.

### Code

```
from sklearn.preprocessing import PolynomialFeatures
```

```
Poly-reg = PolynomialFeatures(degree=2)
```

```
X_Poly = Poly-reg.fit_transform(X)
```

→ fit the quadratic function  
→ keep changing degree until it fits the data

\* Logistic regression is not Regression analysis

## ② Logistic Regression

\* Probability based model

Linear model for classification  
(supervised learning)

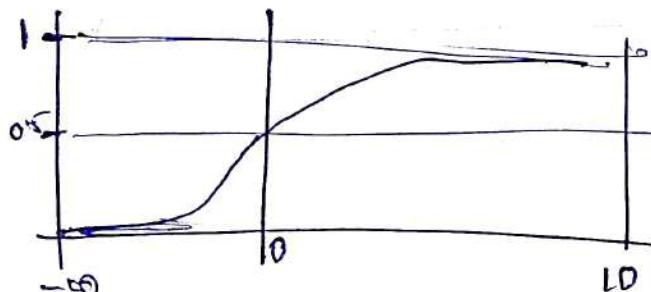
→ Logistic Regression is used for performing  
'classification problems'.

→ It calculates the probability that a given value  
belongs to a specific class. If the probability  
is more than 50%, it assigns the value to  
that particular class,  
else if probability is less than 50%,  
the value is assigned to the other class.

∴ we can say that  
logistic Regr. acts as a binary classifier.

\* Here, As we give input  $x$ , we get  $y$ ,  
where  $y$  is the probability of given variable  
belonging to a certain class.  
Thus, it is obvious that the value of  
 $y$  should be between 0 & 1.

So, we use Sigmoid function as the  
underlying function in Logistic Regression.



$$\sigma(x) = y = \frac{1}{1 + e^{-x}}$$

## so why we sigmoid?

- ① the sigmoid function's range is bounded between 0 & 1, which is what we need.
- ② easy to calculate other functions which is useful during gradient descent calculation.
- ③ or is a simple way to introduce non-linearity to model

the logistic func is given as

$$P(n) = \frac{e^{B_0 + B_1 n}}{1 + e^{B_0 + B_1 n}}$$

we have  $B_0 + B_1 n$

$$P(n) = \frac{e^{B_0 + B_1 n}}{1 + e^{B_0 + B_1 n}} = \frac{e^{B_0 + B_1 n}}{e^{B_0 + B_1 n} + 1} \quad \text{--- ①}$$

$$1 - P(n) = 1 - \frac{e^{h(\theta)}}{1 + e^{h(\theta)}} \quad \text{if } h(\theta) = B_0 + B_1 n$$

$$1 - P(n) = \frac{1}{1 + e^{h(\theta)}} \quad \text{--- ②}$$

$$\frac{\text{from, } \frac{P(n)}{1 - P(n)} = e^{h(\theta)}}{\text{from, } \frac{P(n)}{1 - P(n)} = e^{B_0 + B_1 n}} = e^{B_0 + B_1 n}$$

$$B_0 + B_1 n = \log\left(\frac{P(n)}{1 - P(n)}\right) \rightarrow \text{(logarithmic function)}$$

Now we have,

$$\log\left(\frac{p(m)}{1-p(m)}\right) = B_0 + B_1 x \quad \text{(logit function)}$$

If we input  $x$ , we will get  $p(m)$  (probability)

→ And we get  $p(m)$ , based on the line with  $B_0, B_1$  constants.

This is how logistic regression works

eg, if we have two classes (0,1)

$$y_2 \begin{cases} 0, & \text{if } p(m) \leq 0.5 \\ 1, & \text{if } p(m) \geq 0.5 \end{cases}$$

0.5-threshold

we can change thresholds

→ Based on the scenario

we should change thresholds

cost function / loss function

for a single training instance,

$$\text{cost}(p(m), y) = \begin{cases} -\log(p(m)), & \text{if } y=1 \\ -\log(1-p(m)), & \text{if } y=0 \end{cases}$$

if  $p(m)=1 \& y=1$ ,  $\text{cost}/\log = 0$

if  $p(m)=0 \& y=0$ ,  $\text{cost}/\log = 0$

But,

if  $p(m)=0 \& y=1$ ,  $\text{cost} \rightarrow \infty$

if  $p(m)=1 \& y=0$ ,  $\text{cost} \rightarrow \infty$

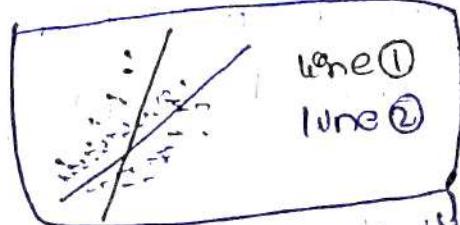
a cost func for whole training set is given as

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y_i \cdot \log(\hat{p}_i) + (1-y_i) \cdot \log(1-\hat{p}_i)]$$

where, m - rows

$y_i$  - expected y

$\hat{p}_i$  - Predicted y

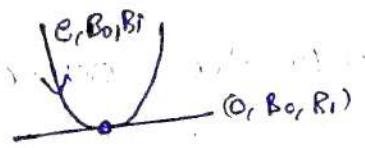


In logit func, RHS we have a line, and we will keep on adjusting the line w.r.t error so that we get correct probability

for from data

→ the values of Parameter ( $\theta$ ) for which the cost func is min is calculated using the gradient descent algorithm (As showed in LR).  
the partial derivative for above cost func is,

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m [\sigma(\theta^T \cdot x_i) - y_i] \cdot x_i$$



### multinomial Logistic Regression

- we can extend Logistic Reg for multi-class classification. The logic is, we train our model for each class and calculate the probability that a specific feature belongs to that class. Once we have trained the model for all the classes, we predict a new value's class by choosing that class for which prob(class) is maximum.
- we rarely use it this way, because we have many other classification models for these scenarios.