

## ML cheat sheet notes

(S1)

① ML → A subfield of AI, we build model which make prediction based on data

② Types → supervised (classification, regression) :- here we have both input and output columns

→ unsupervised (Anomaly detection, clustering, associate Based learning, Dimensionality reduction) :- we have only input column

→ reinforcement learning → no input, output column, self driving cars

→ semi-supervised learning → google photos.

③ steps → ① ~~frame~~ frame the problem, ② gather data ③ data Preprocessing  
④ EDA ⑤ feature selection and engineering ⑥ model development ⑦ ~~training~~ ⑧ optimization, model deployment ⑨ testing ⑩ optimization.

### Regression metrics

① loss function - MSE, MAE, RMSE, Utilization = R<sup>2</sup> score, Adj R<sup>2</sup>  
② MSE =  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ , not robust to outliers  
③ RMSE =  $\sqrt{\text{MSE}}$ , not robust to outliers  
④ MAE =  $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$ , robust to outliers, not differentiable at 0  
⑤ R<sup>2</sup> score → also known as coefficient of determination and goodness of fit, It compare regression model with the mean,  
formula =  $1 - \frac{\text{sum of squared error of model}}{\text{sum of squared error of mean}}$

→ If we add irrelevant column in the value of R<sup>2</sup> score increase

so we use adjusted R<sup>2</sup> score

→ value lie in between 0 to 1

⑥ Adjusted R<sup>2</sup> score

$$F = 1 - \frac{(1-R^2)(n-1)}{(n-1-k)}$$

n = number of rows  
k = number of columns

## ⑦ Classification Matrix

① Recall, Precision, F1 Score, accuracy, confusion matrix

② Accuracy  $\Rightarrow$   $\frac{\text{Total correct prediction}}{\text{Total Prediction}}$

$\Rightarrow$  It does not tell which type of error

$\Rightarrow$  Not accurate for imbalanced data set

$\Rightarrow$  What accuracy is good  $\rightarrow$  Depend on problem

## ③ Confusion Matrix

Actual/Prediction

TP	FN	Type 1 error = FP
FP	TN	Type 2 error = FN

FP = model falsely predicting Positive  $\Rightarrow$  actual = Neg, model = Positive

FN = Model falsely predicting Negative  $\Rightarrow$  actual = Pos, model = Negative

Type 1 = Precision = Email spam classifier

Type 2 = Recall = Cancer detection

Type 1 + Type 2 = F1 Score

Multiclass  $\rightarrow$  like cat, dog, goat

		Prediction		
		cat	dog	goat
cat	cat	1	3	5
	dog	2	4	6
goat	goat	7	8	9

## ④ Precision

Out of all positive prediction, How many truly belong their

$P = \frac{TP}{TP + FP}$ , variant  $\rightarrow$  weighted and macro (for multiclass system)

## ⑤ Recall

Out of all positive cases How many get predicted correctly

$R = \frac{TP}{TP + FN}$   $\rightarrow$  variant  $\rightarrow$  weighted, micro

⑥ F1 Score  $\rightarrow$  When Both T1 and T2 error important, variant  $\rightarrow$  weighted macro.

## ⑧ Linear Regression

① Types - single, multiple, polynomial

② In Sklearn, it is implemented in 2 ways

(i) least square Ordinary least square  $\rightarrow$  linear regression

(ii) Gradient descent  $\rightarrow$  SGDRegressor()

③ Assumption  $\rightarrow$  linear relation b/w input and output

$\rightarrow$  No multicollinearity

$\rightarrow$  residual normally distributed

Formula

For class  $\rightarrow 1, 2, 3$

macro precision =  $(P_1 + P_2 + P_3)/3$

weighted precision = \_\_\_\_\_

④ Gradient Descent → Optimization technique, if you give a differentiable function to it, it will return minima of it.

⑤ Bias → The inability of ML Algo to capture the pattern of data.

Low captured pattern will lead to more Bias,

⑥ ~~Variance~~ Variance → Difference in fit in different datasets

⑦ Bias variance trade off

→ generally you have either High Bias with low variance, or

Low Bias with high variance

Low Bias + High variance = Overfitting

High Bias + Low variance = Underfitting

3 ways to deal with it → Regularization, Boosting, Bagging

⑧ Regularization

→ Adding a term in loss function to reduce overfitting

→ L<sup>2</sup>(Ridge), L<sup>1</sup>(Lasso), Elastic Net

→ Ridge regularization (L<sup>2</sup>) → (i) Apply when you know all columns are important, term =  $\lambda \sum_{i=1}^m w_i^2$  m = 10k,  $\lambda = \text{HP}$  (0, ∞)

→ Lasso regularization (L<sup>1</sup>) → (i) Apply when you know all columns are not important, term =  $\lambda \sum_{i=1}^m |w_i|$

→ Elastic Net regularization → Apply when you don't know are these columns important or not, combine lasso + ridge,  $\alpha(w)^2 + \beta|w|$ , 0, 1, ∞

⑨ Logistic Regression → Here Polynomial works and softmax Perceptron trick, Stream implementation us Sigmoid function

⑩ Decision Tree → Big structure of nested if-else condition, use hyperplanes which are parallel to any one of the axes to cut coordinate system into hyper rectangles, used for regression problem called → CART (Classification and Regression Tree), overfitting

⑪ Hyperparameter tuning → 2 ways GridSearchCV, RandomSearchCV,

⑫ K mean clustering →

Algo → Find number of clusters by elbow method (a)

→ Assign n random centroid

→ Make clusters based on centroid, And point closer to the centroids.

→ Assign centroid again in those clusters

→ If those centroid same as previous then finish, else

elbow method

① From 2 to 10 or 100 we calculate WCSS, and plot the WCSS over graph, and the elbow point is the number of clusters

② with in cluster sum of squared distance = WCSS

(P1)

(P2)

(P3)

(P4)

(P5)

(P6)

$$WCSS = d_1^2 + d_2^2 + d_3^2$$

$$W_1 = 100$$

$$W_2 = 50$$

$$W_3 = 10$$

## (g) Ensemble Learning

(i) wisdom of the crowd, (ii) types (iii) ~~accuracy, low bias, variance~~

Voting, Bagging, AdaBoost, RF

Robust

Boosting Stacking

ORF

accuracy, low bias, variance

AdaBoosting, RF, xgboost

(ii) Voting → Assumption → each model should have min accuracy → 0.51,

regression, All model should be independent and dis-similar.

classification → multiple different model train on a dataset and most frequent (~~highest~~ voting), highest probability (soft voting), mean will get n.

(iii) Bagging → Boosting + Aggregation. Boosting (use random set of data from a data set with replacement), use only 1 type of algo, at prediction majority win and mean. • Types → Bagging, Pasting, Random subspace, Random Patches, Bagging gives better result than Pasting, use Random Patches and subspaces should be used while dealing with dimensional data.

(iv) Random Forest → Difference b/w Bagging and Boosting →

① Bagging support all different algo, Random Forest only use decision tree

② Bagging has model level column sampling (For each model there will be randomly columns get selected) whereas Random Forest has node level column sampling (For each node of every decision tree, randomly columns get selected)

(v) Bagging vs Boosting → Bagging → ① types of models used: - models with low bias and high variance ② parallel learning possible ③ base model weights is equal. Boosting → ① types of models used: - model with high bias and low variance ② sequential learning ③ base model weight is not equal

(vi) AdaBoost → stage wise additive method, build by adding multiple weak learners, all the weak learners are added stage by stage, weak learners are model who has very less accuracy (mostly use decision stumps as a group of weak learners (Decision Tree with maxheight=1)) note → In adaboost classification problem → we use +1 and -1 instead of 1 and 0

(vii) Gradient Boosting → sequential stage wise addition, GB uses the gradient (loss) of model as input to the next model and it's goes on.

(X) Adaboost vs GB → ① Adaboost → we use Decision Stump which have max depth=1, we assign a different weight to each model

② GB → Here we use decision trees which have max depth in bw (8-32), we assign a single learning rate for each model.

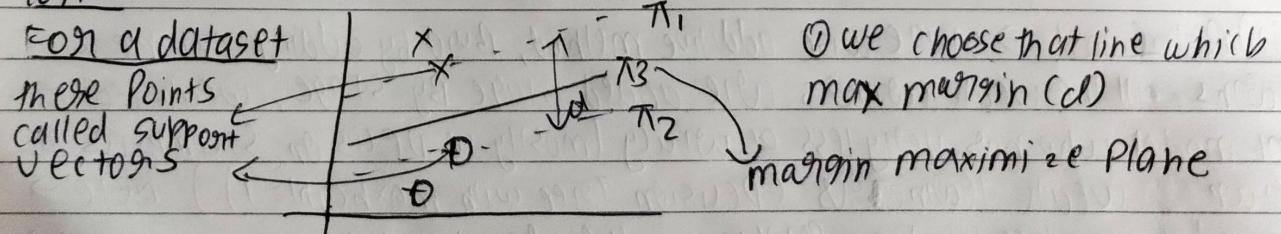
5.12 Stacking and Blending Ensembles :- First we train multiple models (called base models) on data, then we give base models output to a new model (called meta model) as input with data output columns. There are 2 way to do this → ① Blending / Hold out method  
② Stacking / K-fold - method

5.13 Hierarchical clustering → Type ① Agglomerative clustering ② Divisive clustering,  
mostly used, use dendrogram method to find the number of clusters.

5.14 K-Nearst Neighbors → KNN doesn't really "learn" from training data, it simply stores all the points and waits until it needs to make a prediction (lazy learner), when a new data point comes for prediction KNN looks at  $K$  nearest neighbors for prediction (majority win for classification and mean for regression) uses Euclidean distances, Dis → slow for large data, outliers can effect, fail for high dimensional data. To decide the number of clusters make multiple model with grid search CV and choose the best one.

## 5.15 SVM

core Intuition = Maximize the margin B/w dataset, Base = logistic regression logic



working → First we draw a line like  $T_3$ , parallel to which we draw 2 lines  $T_1, T_2$  which is touching the support vectors, then we find out the  $d$ , max  $d$  line will be the answer

Benefits → ① Robust to outliers ② able to work for non-linear dataset with the help of kernels,

## Kernel Trick

- used to classify data which is non-linear
- Here we use a maths functions and applied it on data so that now that data set become separable.
- it convert low dimensional feature space into higher dimensional feature space