

Machine learning

1. Linear Regression

Assumptions in regression analysis:

* Linear relationship: There should be a linear and additive relationship between dependent (response) and independent (predictor) variable.

To check: Residual vs fitted values.

1. If there is any pattern, then it is non-linear.
2. If a funnel shape is there, then its non-constant variance

Multicollinearity:

It affects target variable. Standard errors tend to increase. Confidence become wider and less precise.

check : Scatter plot. $VIF < 4$ means no multicollinearity.
 $VIF > 4$ Serious.

Homoscedasticity or equal variance:

Non-constant variance arises in presence of outliers or extreme leverage value.

To check: Residual vs fitted plot. The plot exhibits funnel shape patterns.

Normal distribution of error terms

If not, confidence intervals may become wide or narrow. Few unusual data points

Check: Q-Q plot, Shapiro-Wilk test

The q-q quantile-quantile is a scatter plot which helps us validate the assumption of normal distribution.

The straight line shows normality

No Autocorrelation

When the residuals are dependent on each other, there is autocorrelation.

See the x-y chart, you get to know it

No. of observation > No. of predictors

The standard errors are uncorrelated

Evaluation metrics of Regression model

Mean Absolute Error (MAE) :

It is a simple metric which calculates the absolute difference between actual and predicted values divided by total no. of points in our dataset.

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

Advantage when to use : you want to know how close the predictions are to the actual on average.

Advantages : The MAE is in the same unit as the output variable.

* It is most robust to outliers

Disadvantages : The graph of MAE is not differentiable so have to apply various optimizers like GD which is differentiable

Mean Squared error (MSE)

It states the finding of squared difference between actual and predicted value.

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

Advantages:

The graph of MSE is differentiable, so we can use loss function.

Disadvantages:

- * It penalizes the outlier
- * It is not in same unit

Root Mean Square error (RMSE):

It is a simple square root of mean squared error:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

Advantages:

The output is in same unit.

It is insensitive to outliers.

When to use:

If you are concerned about large error, RMSE is good.

R²

It is how well the model performance is

$$R^2 \text{ squared} = 1 - \frac{SS_r}{SS_m}$$

~~R²~~ squared means the proportion of variance explained by the model

Disadvantage:

When adding new features, R² starts increasing or remains constant or increasing, it never decreases. because it assumes that while adding more data variance of the data increases.

Adjusted R²

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

where

n = no. of observations

k = no. of independent variables

2. Logistic Regression

Assumptions for logistic regression:

The dependent variable must be categorical in nature.

The independent variable should not have multi-collinearity.

Logistic:

Instead of fitting a regression line, we fit S shaped logistic function.

Logistic regression is the ratio of correctly idea

$$e^z = \frac{P(x)}{1-P(x)}$$

After applying natural log,

$$P(x; b, \omega) = \frac{e^{\omega x + b}}{1 + e^{\omega x + b}} = \frac{1}{1 + e^{-\omega x - b}}$$

If we use the MSE of linear regression, then it will give a non-convex graph with local minima.

$$\text{Log loss} = \frac{1}{N} \sum_{i=1}^N -(y_i * \log(\hat{y}_i) + (1-y_i) * \log(1-\hat{y}_i))$$

$$\text{Cost function} = \begin{cases} -\log(h(x)) & \text{if } y=1 \\ -\log(1-h(x)) & \text{if } y=0 \end{cases}$$

If learning rate is a big number then
we miss the minimum point

Metrics to evaluate classification model

1. Precision

It explains how the predicted positives cases actually turned out to be positive.

It is helpful where false positives is a higher concern than false negatives.

Eg: Importance of precision in recommendation leads to customer churn.

$$\text{Precision} = \frac{\text{True positive}}{\text{True positive} + \text{False positive}}$$

2. Recall

It explains how many of the actual positive cases we were able to predict.

It is helpful metric where False negative is of higher concern than False positive.

Eg: Medical field

$$\text{Recall} = \frac{\text{True positive}}{\text{True positive} + \text{False negative}}$$

Precision can be seen as a measure of quality and recall as a measure of quantity.

Higher precision means that an algorithm returns more relevant results than irrelevant ones and high recall means that an algorithm returns most of the relevant results.

3. Recall (Sensitivity) F1 Score

It is a harmonic mean of precision and recall.

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

It punishes extreme values.

It is an effective metric in following cases:

when FP & FN are equally costly

when more data doesn't effectively change the outcome.

True Negative is high.

Precision can be seen as a measure of quality and recall as a measure of quantity.

Higher precision means that an algorithm returns more relevant results than irrelevant ones and high recall means that an algorithm returns most of the relevant results.

3. Recall (Sensitivity) F1 Score

It is a harmonic mean of precision and recall.

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

It punishes extreme values.

It is an effective metric in following cases:

when FP & FN are equally costly

when more data doesn't effectively change the outcome.

True Negative is high.

4. AU-ROC

The Receiver operator characteristic (ROC) is a probability curve that plots the TPR (True Positive Rate) against FPR (False positive Rate) at various threshold values and separates the signal from noise.

AUC (Area under curve) is the measure of ability of a classifier to distinguish between classes.

X axis \rightarrow False positive rate

Y axis \rightarrow True positive rate

Higher the value of X means means high False Positive.

The choice of threshold depends on ability to balance between FP and FN.

5. Log loss

$$\text{log loss} = y \log(p) + (1-y) \log(1-p)$$

For imbalanced data, ROC helps.

3. Naïve Bayes

It is used to predict models for binary or multi-classification labels. It operates on conditional probabilities, which estimate the likelihood of a classification based on the combined factors while assuming independence between them.

Assumption of Naïve Bayes

It assumes that all features in the input data ~~and some other~~ are independent of each other. Used in sentimental analysis, spam filtering and text classification.

'Bayes' theorem

It is used to determine the probability of a hypothesis with prior knowledge.

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

Process:

1. Convert the given dataset into frequency tables.
2. Generate likelihood table by finding the probabilities of given features.
3. Use Bayes theorem to calculate the posterior probability.

Eager learner - Naïve Bayes.

Types of Naïve Bayes

Gaussian - It follows normal distribution. It assumes that the values are sampled from the Gaussian distribution.

Multinomial - It assumes it is taken from multinomial data.

Bernoulli - Assumes the predictor variables are independent Bernoulli variables. used for document classification.

It (GaussNB) has segregated the datapoints with fine boundary.

* It's a eager learner

5. SVM

The shortest distance between the observations and the threshold is called the margin.

Maximal margin classifier They are super sensitive to the data.

It involves misclassification for variance-bias tradeoff.

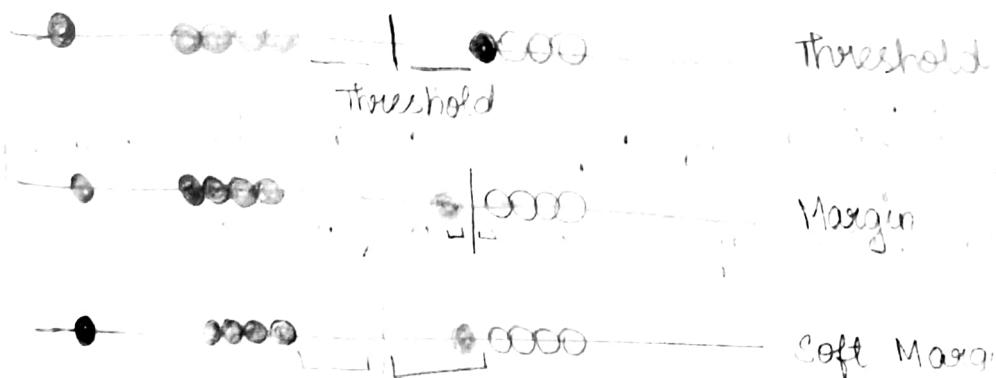
Soft Margin :

When we allow misclassifications, the distance between the observations and threshold is called soft margin

uses CV to determine how many misclassification and observations to fall inside of the soft margin.

Soft Margin classifier (support vector) Machine

It comes from that observations on the edge and within soft margin are called support vectors.



When 3-dimensional, it is plane. If 4 or more, it is subspace, it is hyperplanes.

Support Vector Machine

- Idea is, start with relatively one and move to higher dimension
- Find support vector classifier to get line

How to decide the transformation?

- Kernel function decides.
when degree = 1,
it computes the relationship between each pair of observations in 1 dimension and used to find the SVC.

RBF:

It uses infinite dimension. It is like weighted KNN.

SVM: They ^{behave} as if they were in the higher dimensions.
They don't actually do the transformation.

Polynomial Kernel :

$$(a \times b + r)^d$$

r is coefficient, d is the degree, a & b are observation

We get a dot product.

We need to calculate the high dimensional relationship is to calculate the Dot products between each pair of points.

Radial Basis Kernel :

$$e^{-\gamma(a-b)^2}$$

It uses infinite dimensions. It behaves like kNN
 γ is determined by cross validation.

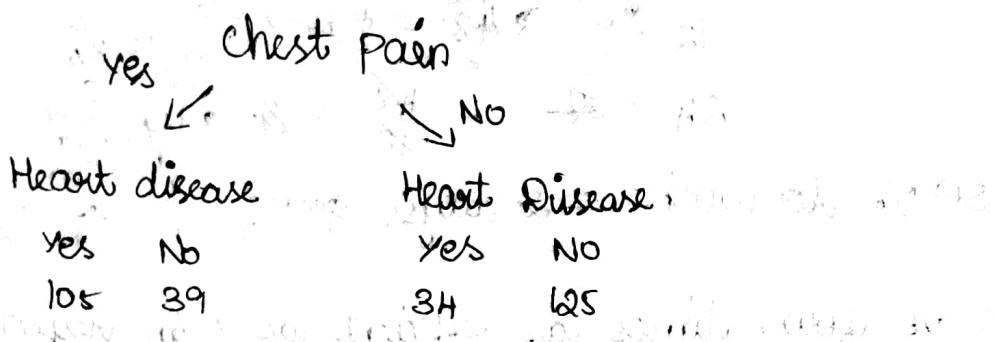
The less the value, further they are.

Taylor series

Turn $f(x)$ into infinite sum.

6. DECISION TREE

Step 1: Take a column and target and build a tree. Repeat those for all column



Because none of the leaf nodes are 100% Yes or No. They all are Impure.

To determine which separation is best, we need a way to measure and compare impurity.

For a leaf,

$$\text{The Gini impurity} = 1 - (\text{Probability of Yes})^2 - (\text{Probability of No})^2$$

$$\begin{aligned} \text{Gini Impurity for a Node} &= \text{Weighted average of} \\ &\quad \text{Gini impurities for the leaf nodes.} \\ &= \left(\frac{144}{144+159} \right) 0.305 + \left(\frac{159}{144+159} \right) 0.336 \end{aligned}$$

The one

Step 2 : The one with the lowest Impurity is made the root node.

Step 3 : Take another column for left and select lowest branch entropy.

Step 4 : Do the same thing for right side.

Ans

REGRESSION TREES :

Find the mean of the sort adjacent values.

Find Gini's Impurity for each of them.

STATISTICS

CENTRAL LIMIT THEOREM:

In order to apply Central limit theorem, there are four conditions to be met.

1. Randomization

2. Variables should be independent and identical

3. When sample is drawn without replacement, the sample size should be no longer than 10% of the population

4. Sample size needs to be sufficiently large

Statement :

Suppose that a large sample of observation is obtained, each observation being randomly produced in a way does not depend on other, the average of observed values is computed. If the procedure is performed many times, resulting in an collection of

Observed averages, the CLT says if the sample size was large enough, the probability distribution of these averages will closely approximate a normal distribution.

Law of large numbers

It is a theorem that describes the average of the results obtained from a large number of trials should be close to the expected value and tends to become closer to expected value as more trials are performed.

The LLN applies to the average. Therefore, while

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

Standard deviation vs standard error

Standard error: The standard deviation of the means of the mean is called standard error.

S.D: It quantifies the variation within a set of measurements.

$$\text{S.E} = \frac{\sigma}{\sqrt{n}}$$

Confidence Intervals

It is interval that covers the percentage of means of means sample means. It is calculated after Bootstrapping.

Population parameters

Why not dividing by N underestimates the variance?

When we divide by n , the value around the sample mean is always less than the value around population mean.

We need squares to find the derivative and finding the minimum value will be harder.

Approximated after calculation

Likelihood.

$$L(\text{mean} = 32 \text{ g}, S.D = 2.5 \mid \text{mouse weighs } 34 \text{ grams}) \\ = 0.12$$

while probability is $P(\text{mouse} > 34 \text{ grams} \mid m = 32, S.D = 2.5)$

probabilities are the areas under a fixed distribution. Likelihood are the y-axis values for fixed data points with distributions that can be moved.

Maximum likelihood Estimation.

MLE is a statistical method used to estimate the parameters of a probability distribution that best explains a given set of data.

It aims to find the values of these parameters that make the observation most probable.

Conditional probability

$$P(E_i|A) = \frac{P(E_i)P(A|E_i)}{\sum_{j=1}^n P(E_j)P(A|E_j)}$$

is the Bayes theorem

According to conditional probability,

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

1. A bag contains 4 white and 6 black balls while bag I contains 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 it is from bag I

$$P(I|B) = \frac{\frac{1}{2} \times \frac{6}{10}}{\frac{1}{2} \times \frac{6}{10} + \frac{1}{2} \times \frac{3}{7}} = \frac{\frac{6}{20}}{\frac{6}{20} + \frac{3}{14}}$$

2. A man is known to speak the truth 2 out of 3 times. He throws a die and reports that the no. obtained is a four. Find the probability that no. obtained is actually a four.

$$P(E_i | A) = \frac{P(E_i) P(A|E_i)}{\sum_{i=1}^n P(E_i) P(A|E_i)}$$

$$\therefore P(E_i | A) = \frac{1}{6} \times \frac{2}{3} = \frac{1}{9}$$

$$(four | true) = \frac{\frac{1}{6} \times \frac{2}{3}}{\frac{2}{3} \times \frac{1}{6} + \frac{1}{3} \times \frac{5}{6}}$$

which shows 3 cases of situation given by

which makes 3 cases of situation. & prob. obtained

$$= \frac{2+5}{18} = \frac{7}{18}$$

which shows 3 cases of situation given by

$$= \frac{3}{7}$$

$$\frac{3}{7}$$

$$\frac{1}{9} \times \frac{1}{2}$$

$$= \frac{1}{18}$$

$$= \frac{1}{6} + \frac{2}{6} = \frac{1}{3}$$