

Machine Learning Algorithm

Evaluating performance for classification

1. Accuracy = $\frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$

→ Accuracy useful when target classes are balanced

2. Recall = $\frac{\text{Number of True Positive}}{\text{Number of True Positive} + \text{Number of False Negative}}$

$$\hat{P}(R-P) = \frac{1}{1+e^{-\frac{1}{R}}}$$

→ Ability of a model to find all the relevant cases within a dataset

3. Precision = $\frac{\text{Number of True Positive}}{\text{Number of True Positive} + \text{Number of False Positive}}$

→ Ability of a model to identify only the relevant point

4. F1-score = $2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$

→ harmonic mean of precision and recall

→ optimal blend of precision and recall

$$\hat{P}(R-P) = \frac{1}{1+e^{-\frac{1}{R}}}$$

Evaluating performance for Regression

- Evaluation metrics like accuracy or recall aren't useful for regression problems
- need metric designed for continuous values.

1. Mean Absolute Error (MAE)

- mean of the absolute value of errors

- easy to understand

$$\frac{1}{n} \sum_i^n |y_i - \hat{y}_i|$$

simply said it's the average of the difference between our prediction and true value.

- However, MAE won't punish large errors, so we have MSE.

2. Mean Squared Error (MSE)

- mean of the squared errors

- larger errors are noted more than with MAE (TO PUNISH OUTLIERS)

- more popular

$$\frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2$$

3. Root mean squared error (RMSE)

- root of MSE.

- so that the unit of the metrics is same as Y

- most popular

$$\sqrt{\frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2}$$

Machine Learning

1. Simple Linear Regression

Predicting a response using a single feature

→ a method to predict dependent variable (Y) based on values of independent variable (X).

→ it is assumed that the 2 variables are linearly related.

∴ hence, we try to find a linear function that predicts the response value (y)

as accurately as possible as a function of the feature of (x)

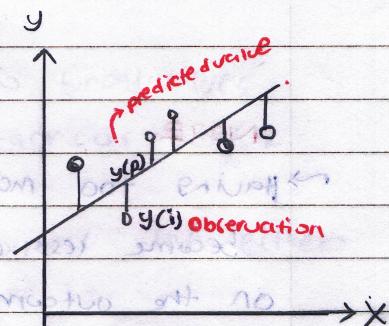
How to find best fit line

→ we are trying to minimize the errors in prediction by finding "line of best fit"

→ trying to minimize the length between the observed value (y_i) and the predicted value (\hat{y}_i) from our model (y_p)

+ finding equation and establishing relationship → slope

$$\text{Dependent Variable } \leftarrow Y = b_0 + b_1 X, \quad \begin{matrix} \downarrow \text{slope} \\ \text{Intercept} \end{matrix} \quad \begin{matrix} \downarrow \\ \text{Independent Variable} \end{matrix}$$



$$\min \{ \sum (y_i - y_p)^2 \}$$

2. Multiple Linear Regression

→ similar to simple linear regression just that multiple linear regression attempts to model the relationship between 2 or more features

$$\text{Dependent Variable } \leftarrow Y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n \rightarrow \text{multiple independent Variable}$$

For a successful regression analysis it is essential to validate these assumptions.

1. Linearity \rightarrow The relationship between dependent and independent variables should be linear
2. Homoscedasticity \rightarrow (constant variance) of the errors should be maintained
3. Multivariate normality \rightarrow Multiple regression assumes that the residuals are normally distributed

(+) If there is lack of multicollinearity \rightarrow it is assumed that there is little or no bias (\hat{Y}) due to multicollinearity in the data. (-) Multicollinearity occurs when the features (independent variables) are not independent of each other. $\text{ad} = V \rightarrow$ High bias

NOTE

\rightarrow Having too many variables could potentially cause our model to become less accurate especially if certain variables have no effect on the outcome or have a significant effect on other variables. There are various methods to select the appropriate variable

1. Forward selection
2. Backward Elimination

3. Bi-Directional Comparison

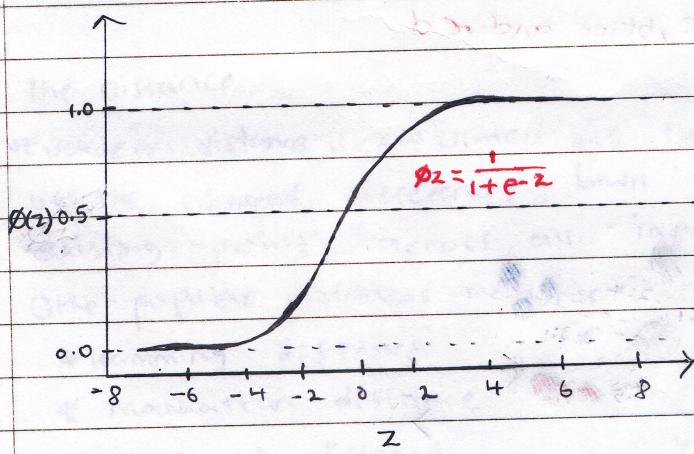
Bi-directional comparison is a method of selecting variables from a set of variables.

3. Logistic Regression

- used for classification problems
 - to predict the group to which the current object or observation belongs to
 - gives a discrete binary outcome between 0 to 1
 - ex. to predict a person will vote or not in upcoming elections
- How does it work?
- Logistic regression measures the relationship btwn the dependent variable (our label, what we want to predict) and the one or more independent variables (our features) by estimating probability using its underlying logistic function.

Making predictions

- These probabilities must then be transformed into binary values
 - in order to actually make a prediction. This is task of logistic function or sigmoid function
 - This values btwn 0 to 1 will be transformed into either 0 or 1 using threshold classifier
- Sigmoid function
- S-shaped curve that can take any real-valued number and map it into a value btwn range of 0 and 1, but never exactly at those limits.



Type of Logistic Regression

1. Binary Logistic regression

→ The categorical response has only 2 possible outcomes.

ex: spam or not

2. Multinomial Logistic Regression

→ 3 or more categories without ordering.

ex: predicting which food is preferred more (veg, non-veg, veg+nonveg)

3. Ordinal Logistic Regression

→ 3 or more categories with ordering.

ex: movie rating from 1 to 5

4. K-Nearest Neighbours - KNN classification algorithm

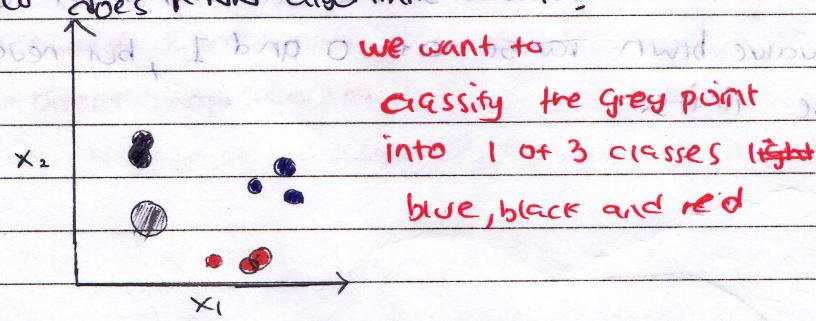
→ Simple & yet most used classification algorithm

→ KNN is a non-parametric means that it does not

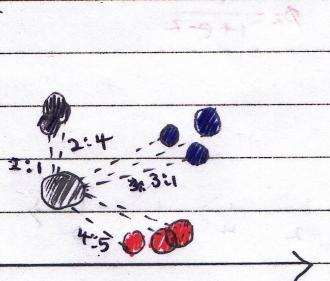
make any assumptions on the underlying data distribution

instance-based (means that our algorithm doesn't explicitly learn a model. Instead it chooses to memorize the training instances and used in a supervised learning setting)

How does K-NN algorithm work?



Start by calculating
the distance between
the grey point and
k-nearest points



Making predictions

→ To

- When used for classification - the output is a class membership (predicts a class - a discrete value). No decision
- there are 3 key elements of this approach: a set of labelled objects, distances and no to a neighbor + an e.g., a sets of stored records, a distance between objects, and the value of K , the number of nearest neighbors.

Making predictions.

- To classify an unlabeled object, the distance of this object to be labelled is computed, its k -nearest neighbors are identified, and the class label of the majority of nearest neighbors is then used to determine the class label of the object.

For real-valued input variables, the most popular distance measures is Euclidean distance.

- 2.1 → 1st NN
- 2.4 → 2nd NN
- 3.1 → 3rd NN
- 4.5 → 4th NN

- # of votes for each class
- CLASS 1 wins the vote
- Point 1 is therefore predicted to be of class 1

The Distance

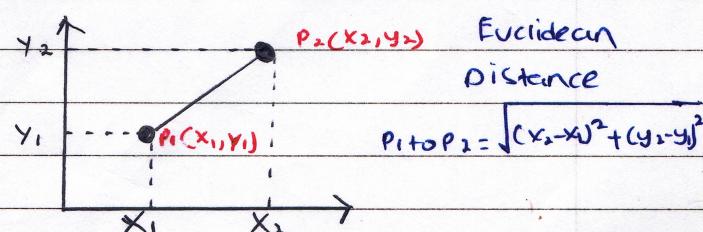
- Euclidean distance is calculated as the square root of the sum of the squared differences between a new point and an existing point across all input attributes.

Other popular distance measures include:

* Hamming distance

* Manhattan distance

* Minkowski distance



Value of K

- Finding value of K is not easy
- A small value of K means that noise will have a higher influence on the result and a large value make it computationally expensive.
- It depends a lot on your individual cases, sometimes it's best to just run through each possible value of K and decide for yourself.

18.11.2021

5. Support Vector Machines

- What is SVM?
- SVM is a supervised machine learning algorithm which can be used for both classification or regression, but mostly in classification problems.

→ In this algorithm, we plot each data item as a point in n-dimensional space (where n is the number of features) with the value of each feature being the value of a particular coordinate.

How is data classified?

- We perform classification by finding the hyperplane that differentiates the 2 classes very well.
- In other words the algorithm outputs an optimal hyperplane which categorizes new examples.

What is optimal Hyperplane?

- For SVM, it's the one that maximizes the margins from both tags.
- In other words, the hyperplane whose distance to the nearest element of each tag is the largest.

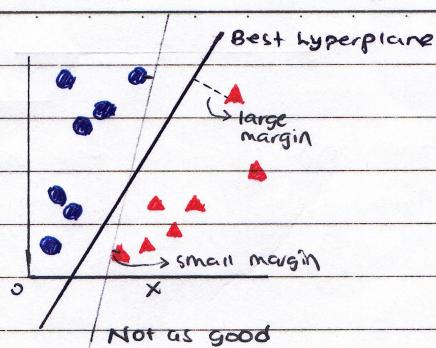
$$(w_0 + w_1 x_1 + w_2 x_2) = 0$$

$$(x_1, x_2)$$

Distance from origin *

Distance from both tags *

Margin *



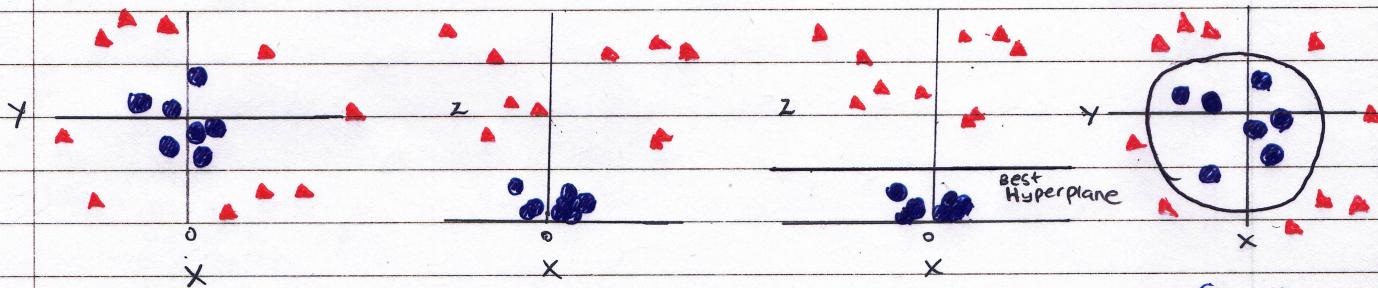
* labelled data

* In 2D, the best hyper-plane is simply a line

* Not all hyper planes are created equal

* The best hyperplane is with a large margin

Nonlinear Data



(1) In above case we cannot draw a linear boundary. We will now add a 3rd dimension. We create a new z dimension, and we rule that it can be calculated a certain way that is convenient for us: $z = x^2 + y^2$ (center for a circle)

(2) This will give us a 3 dimensional space. From a diff perspective the data is now in linearly separated groups. All values for z will be always positive because z is squared sum of both x and y .

(3) since we are in 3 dimensions now the hyper-plane is a plane parallel to the x -axis at a certain z . we change the hyperplane, one that maximize the margins from both of the classes.

(4) Now we map back to 2D. Our decision boundary is a circumference which splits both tags using sum. We set a circle as hyper plane.

Tuning parameters

1. Kernel
2. Gamma
3. Regularization
4. Margin.