

Classification

Monday, June 21, 2021 6:21 PM

1. KNN

Algorithm :

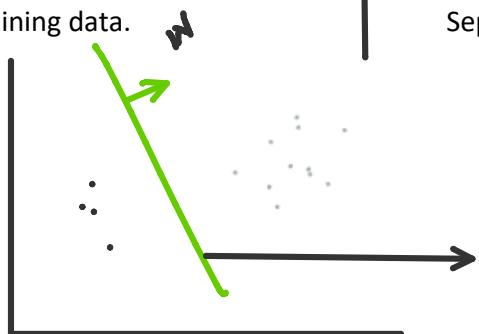
1. Initialize the List of Distances = [].
2. For each of x_j in the dataset:
 Compute the distance $d(x_i, x_q) = d_i$
 Store them in list as (x_i, y_i, d_i)
3. Sort the list on the basis of distances.
4. Take the smallest 'k' distances and store them in list 'KNN' .
5. For each x_i in KNN list:
 if y_i is true:
 count_pos = +1
 else:
 count_neg = +1
6. If (count_pos > count_neg) :
 print('Positive class')
 else:
 print('Negative class')

There is no training in KNN algorithm.

1. Logistic Regression

Linear Regression

- Supervised Learning
- Regression Task
- $D_n = \{(x_i, y_i)\}_{i=1}^n | y_i \in R\}$
- Task : To find the line that best fits The training data.



Logistic Regression

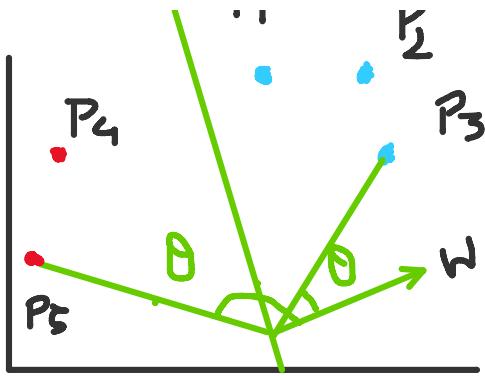
- Supervised Learning
Classification Task
 $D_n = \{(x_i, y_i)\}_{i=1}^n | y_i \in \{-1, 1\}\}$
Task : Find the line that BEST Separates Line between +ve & -ve

W = Norm of the Hyperplane/ plane
i.e. It always perp. to the line

$$y = mx + c \cong w^T x + w_0$$



$$\vec{w} \cdot p_1 = +ve$$



$$\begin{aligned}\vec{w} \cdot p_1 &= +ve \\ \vec{w} \cdot p_2 &= +ve \\ \vec{w} \cdot p_3 &= +ve \\ \vec{w} \cdot p_4 &= -ve \\ \vec{w} \cdot p_5 &= -ve\end{aligned}$$

We have assumed that the data is separable by line/plane.

So the formula is,

$$\text{Signed Distance} = \sum_i^n (y_{actual} * Y_{predicted})$$

$$\text{Where } Y_{predicted} : |w| \cdot |p| \cos\theta \text{ or } mx + c$$

The vectors which gives +ve values from above formula are considered to be correctly classified whereas vectors which gives overall -ve values are considered to be misclassified.

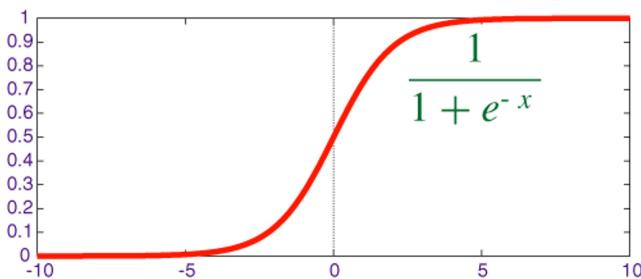
Best Separator will gives less misclassifications or more correct classification

Optimum value of m & c :

$$m^*, c^* = \arg_{m,c} \max \{ \sum_i^n y_i * (mx_i + c) \}$$

(Maximizing the signed distance)

After removing outliers, still there are some outliers which were not detected by EDA, Hence we have to convert the Signed distance to sigmoid function



It basically more robust to outliers. It neglects the effect of outliers. It shows the tapering effect. We are doing some sort of scaling on signing distance. Its value between 0 to 1.

So,

$$m^*, c^* = \arg_{m,c} \max \{ \sum_i^n y_i * (mx_i + c) \}$$

$$= \arg_{m,c} \max \{ \sigma [\sum_i^n y_i * (mx_i + c)] \}$$

$$= \arg_{m,c} \max \{ \sum_i^n \frac{1}{1 + \exp\{-\chi_i * (mx_i + c)\}} \}$$

Optimization theory :

$$\begin{aligned}\max f &\approx \min 1/f \\ \max f &\approx -\min 1/f\end{aligned}$$

$$= \arg_{m,c} \min \{ \sum_i^n 1 + \exp\{-\chi_i * (mx_i + c)\} \}$$

Decision Tree :

A Decision tree is a flowchart like tree structure, where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (terminal node) holds a class label.

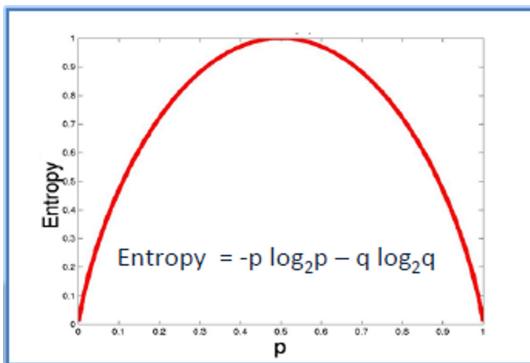
Entropy ($H(y)$) : Measure of Randomness. [y is o/p]

It captures the randomness/ uncertainty of variable. Its values always lies between 0 and 1.

1. It is 0 when all points in variable are same. If we have 50% +ve and 50% -ve, then Entropy is 1.

$$0 \leq H \leq 1$$

$$H(y) = - \sum_i^k p(y_i) \cdot \log_2(p(y_i))$$



$$\text{Entropy} = -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$$

Information Gain

The information gain is based on the decrease in entropy after a dataset is split on an attribute. Constructing a decision tree is all about finding attribute that returns the highest information gain .

$$\text{IG}(Y, f) = H(Y) - \sum_i^k \left\{ \frac{|D_i|}{|D|} * H_{D_i}(Y) \right\}$$

where $\sum_i^k \left\{ \frac{|D_i|}{|D|} * H_{D_i}(Y) \right\}$: Weighted Entropy

Step 1: Calculate entropy of the target.

Step 2: The dataset is then split on the different attributes. The entropy for each branch is calculated. Then it is added proportionally, to get total entropy for the split. The resulting entropy

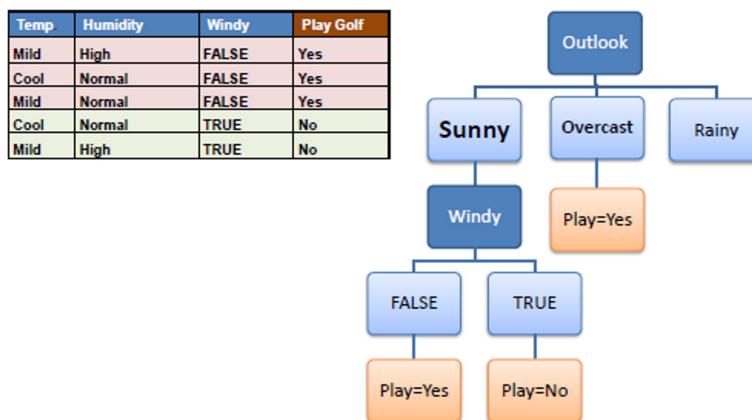
is subtracted from the entropy before the split. The result is the Information Gain, or decrease in entropy.

Step 3: Choose attribute with the largest information gain as the decision node, divide the dataset by its branches and repeat the same process on every branch.

Step 4 : those attributes who have more than 0 Entropy needs further split while those have 0 doesn't need.

The diagram shows a decision tree structure. The root node is 'Outlook'. It branches into three nodes: 'Sunny', 'Overcast', and 'Rainy'. The 'Sunny' node has a child node 'Windy'. The 'Overcast' and 'Rainy' nodes do not have further children. To the right of the tree are three tables representing the data for each branch:

- Sunny:** Contains 5 rows of data with columns: Outlook, Temp, Humidity, Windy, Play Golf. Rows: (Sunny, Mild, High, FALSE, Yes), (Sunny, Cool, Normal, FALSE, Yes), (Sunny, Cool, Normal, TRUE, No), (Sunny, Mild, Normal, FALSE, Yes), (Sunny, Mild, High, TRUE, No).
- Overcast:** Contains 4 rows of data with columns: Outlook, Temp, Humidity, Windy, Play Golf. Rows: (Overcast, Hot, High, FALSE, Yes), (Overcast, Cool, Normal, TRUE, Yes), (Overcast, Mild, High, TRUE, Yes), (Overcast, Hot, Normal, FALSE, Yes).
- Rainy:** Contains 5 rows of data with columns: Outlook, Temp, Humidity, Windy, Play Golf. Rows: (Rainy, Hot, High, FALSE, No), (Rainy, Hot, High, TRUE, No), (Rainy, Mild, High, FALSE, No), (Rainy, Cool, Normal, FALSE, Yes), (Rainy, Mild, Normal, TRUE, Yes).



Advantages & Disadvantages of Decision Tree :

- We don't require to convert categorical values to one hot encoding because Decision Tree never compares the relationship between categorical values. We need to do just label Encoding and it also avoids to increase Dimensionality.
- We don't need to do Standardization/ Normalization to data because we don't need to keep all data in similar range like in other algorithms(because they have to calculate some distance in algorithm).
 - Overfitting : (High Variance and low Bias)
 - Real Valued Column : If model has real valued column, then for every it check entropy less than and entropy greater than (for Ex. $H < 170$ and $H > 170 \forall$ point) from the whole column and if data is huge it will increases the computation time for algorithm.

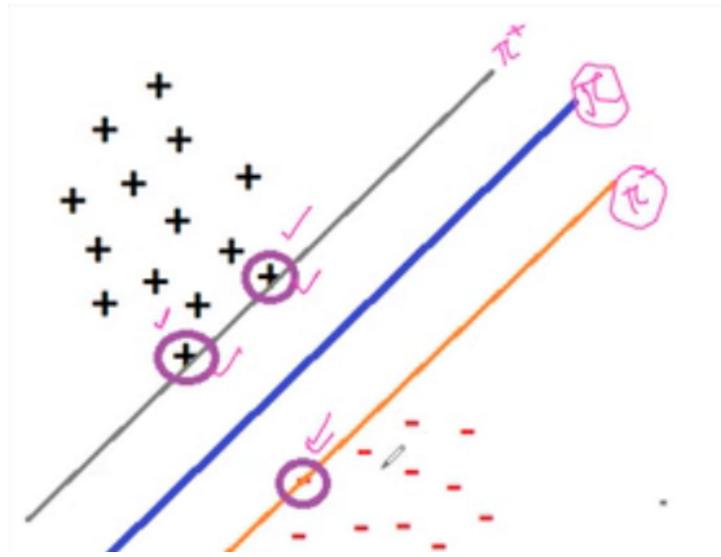
Gini Impurity : measure of randomness similar to entropy. It require less computational power than Entropy.

$$I_g = 1 - \sum_i^k p(y_i)$$

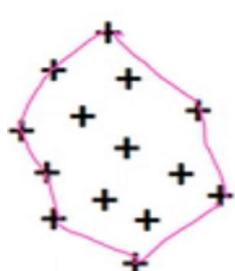
$$-0.5 \leq g \leq 0.5$$

SVM (Support Vector Machine) :

SVM tries to find the hyperplane that best separate the +ve and -ve as widely as possible. Also known as MARGIN MAXIMIZING hyperplane.

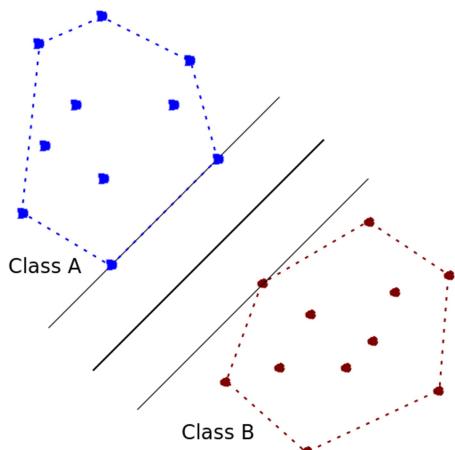


Convex Hull : A figure in which if we draw a line between two points in figure, the line always present inside the figure.

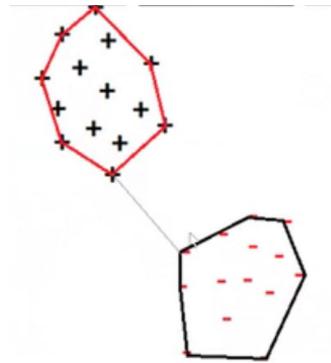


Steps for SVM :

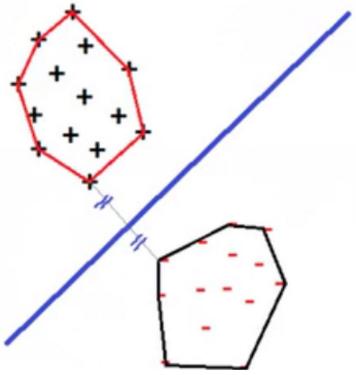
1. Draw the convex Hull for +ve and -ve points.



2. Finding the shortest line between the classes.



3.Bisect the line into two parts.



Kernel trick :

A Kernel Trick is a simple method where a Non Linear data is projected onto a higher dimension space so as to make it easier to classify the data where it could be linearly divided by a plane.

In Logistic Regression, we have to find the best line that separates the +ve and -ve .

$y_{actual} * y_{pred} \geq 0$ i.e signed dist. is +ve.

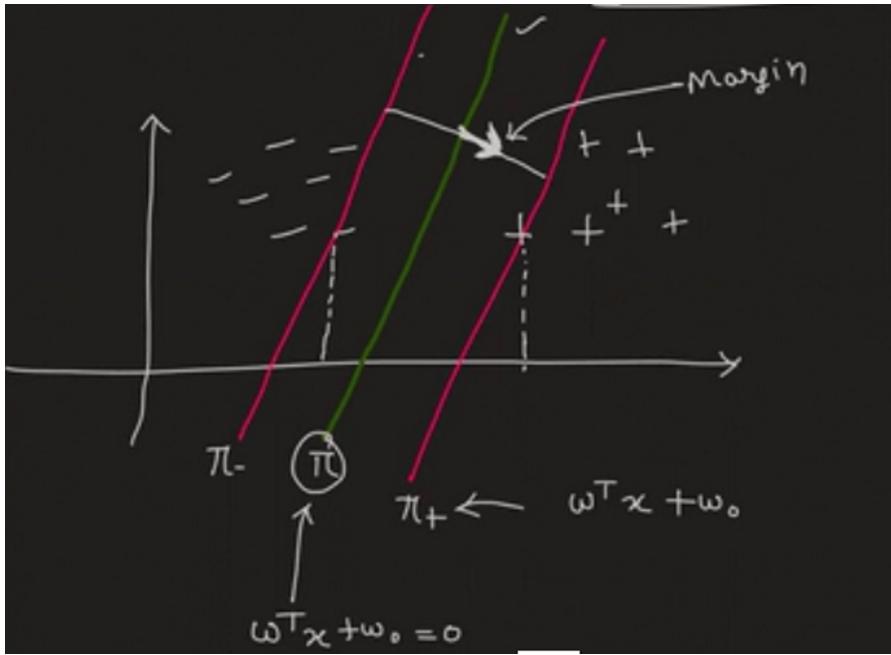
In SVM, we have to maximizing the margin distance(π_+ , π_-)

Hard Margin SVM Formulation :

Maximum Margin such that every point in the training data is correctly classified.

Max. Margin such that

$$y_i * (mx_i + c) \geq 0 \quad \forall i$$



Line equation for $\pi = mx + c$ or $w^T x + w_o$
 Similarly for π_+ , $\pi_- = mx + c$ or $w^T x + w_o (+/-)k$

Only difference between them that they have different intercept w_o .
 So,

$$\begin{aligned}\pi_+ : w^T x + w_o &= 1 && \text{(Take 1 for assumption)} \\ \pi_- : w^T x + w_o &= -1\end{aligned}$$

Then,

$$\begin{aligned}w^T x_{pos} + w_o &= 1 \\ -w^T x_{neg} + w_o &= -1 \\ \therefore w^T x_{pos} - w^T x_{neg} &= 2 \\ \therefore w^T(x_{pos} - x_{neg}) &= 2 \\ \therefore \frac{w^T(x_{pos} - x_{neg})}{|w|} &= \frac{2}{|w|}\end{aligned}$$

This thing is called projection of distance between margin on \rightarrow_w

Max. $\left[\frac{2k}{|w|} \right]$
 such that

$$y_i * (w^T x + w_o) \geq 0 \quad \forall i$$

This is called Hard Margin SVM Formulation.

Dual Form :

Lagrange Multiplier helps to achieve the Dual Form

$$\max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^\top x_j$$

s.t. $0 \leq \alpha_i$ & $\sum_{i=1}^n \alpha_i y_i = 0$

Kernel Trick \rightarrow

$$\max \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

Linear SVM \rightarrow

$$K(x_i, x_j) = (x_i^\top x_j) = x_i \cdot x_j$$

They use quadratic equation to achieve higher dimension and to project data to higher dimension because Linear SVM is not able to classify all types of data.

More higher polynomial equations are used to classify the data.

Quadratic Kernel SVM :

$$K(x_i, x_j) = x_i^\top x_j^2$$

Polynomial Kernel SVM :

$$K(x_i, x_j) = x_i^\top x_j^d$$

RBF Kernel(Radial basis function) :

$$K(x_i, x_j) = \exp \left\{ \frac{-||x_i - x_j||}{2\sigma^2} \right\}$$

The best thing about RBF is gives results most of the time but you have to do experiments.

Kernel Trick has huge training time complexity, that's why SVM has higher T.C.

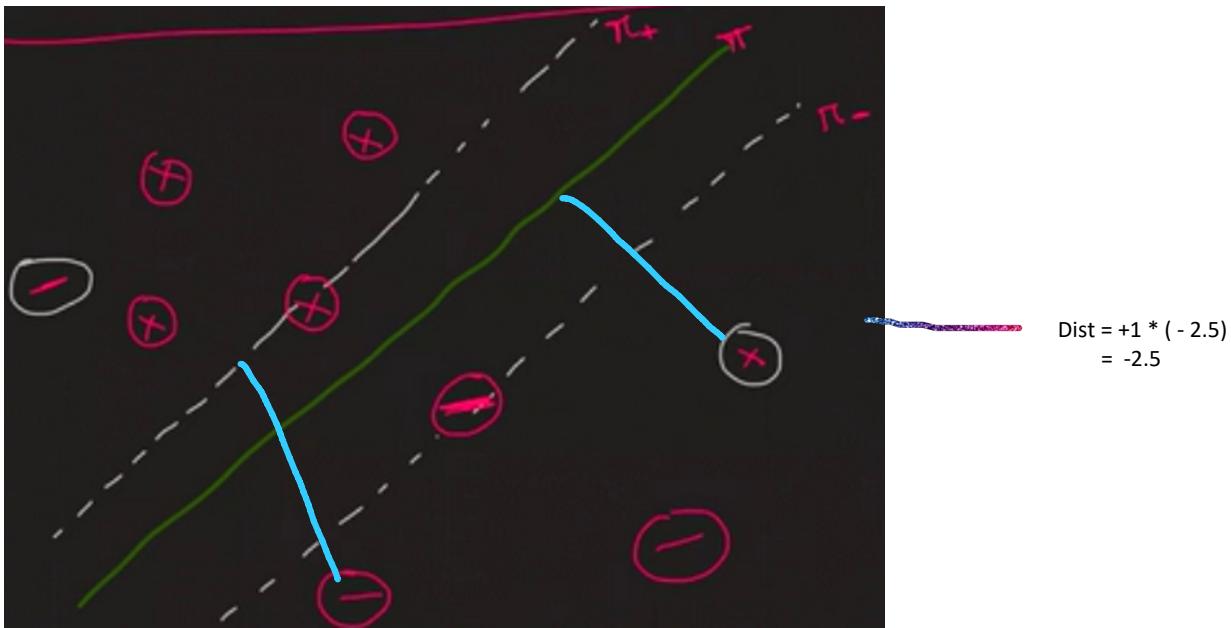
Soft Margin SVM :

$$\max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

s.t. $0 \leq \alpha_i \leq C$ & $\sum_{i=1}^n \alpha_i y_i = 0$

Dual form & Soft Margin SVM

Only difference in $0 \leq \alpha_i \leq C$ part.



$$\text{Dist} = -1 * 2.5 \\ = -2.5$$

We have to use slack variable means we have to subtract it with 1 without changing its value.

$(1 - (+3.5)) = -2.5$ where 3.5 value is called slack variable denoted by ξ_i (zeta).

So, in general

if $\xi_i \geq 0$, point is misclassified according to π_+, π_-
& $\xi_i \leq 0$, point is classified according to π_+, π_-

$\therefore \xi_i \uparrow$, then the point goes to incorrect direction(π_+, π_-)

$$\boxed{y_i (\omega^\top x_i + \omega_0) > 1} \rightarrow \boxed{\xi_i < 0}$$

$$\boxed{y_i (\omega^\top x_i + \omega_0) \leq 1} \rightarrow \boxed{\xi_i \geq 0}$$

According to Optimization theory,

$\max \text{ margin}$ $\min \frac{1}{2} \text{margin}$ $\xrightarrow{\text{Hyperparameter}}$ $\max f = \min (-f)$
 $\boxed{\min \left\{ \frac{1}{2} * \|\omega\| + C \sum_{i=1}^n \xi_i \right\}}$

s.t. $y_i (\omega^\top x_i + \omega_0) \geq 1 - \xi_i \quad \forall i$
 $\xi_i \geq 0$

This is Soft Margin SVM formulation for more and more noise.

Where $\xi_i \geq 0$ means data contains more noise.

$C * \frac{1}{n} \sum_{i=1}^n \xi_i$ is called Hinge Loss.

$\max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$
 s.t. $\boxed{0 \leq \alpha \leq C} \leftarrow \quad \& \quad \sum_{i=1}^n \alpha_i y_i = 0$

Dual form of SOFT MARGIN SVM

Only difference in $0 \leq \alpha \leq C$ part of Hard form

KNN > SVM > LogReg in terms of time complexity.

Naïve Bayes :

Conditional Probability : $P(E1 | E2) = \frac{P(E1 \cap E2)}{P(E2)}$

Independent Events : $P(E1 | E2) = P(E1)$ $\because E2$ and $E1$ are independent.

Suppose $P(E1 | E2) = \frac{P(E1 \cap E2)}{P(E2)}$ and $P(E2 | E1) = \frac{P(E2 \cap E1)}{P(E1)}$

Then, According to Set Theory $P(E1 \cap E2) = P(E2 \cap E1)$

$$\Rightarrow P(E2 | E1) = \frac{P(E1 \cap E2)}{P(E1)}$$

$$\Rightarrow P(E1 \cap E2) = P(E2 | E1) * P(E1) \quad \text{----- (Equation : 2)}$$

Put Equation 2 in $P(E1 | E2)$

$$\Rightarrow P(E1 | E2) = \frac{P(E2 | E1) * P(E1)}{P(E2)}$$

This is called Bayes Theorem.

In this Algorithm we have to calculate $P(C_k | x_q)$

Where C_k - Class Label and x_q = Query or datapoints given

We have to find the label which has highest $P(C_k | x_q)$ and Take the label as Output.

For Example in Iris Dataset :

$$P(\text{Species} = \text{Setosa} | x_q) = ?$$

$$P(\text{Species} = \text{Virginica} | x_q) = ?$$

Select the label which has maximum probability.

$$\begin{aligned} P(C_k | x_q) &= \frac{P(x_q | C_k) P(C_k)}{P(x_q)} \text{ According to Bayes theorem} \\ &= \frac{P(x_q \cup C_k)}{P(x_q)} \text{ or } \frac{P(x_q, C_k)}{P(x_q)} \\ &= \frac{P(x_1, x_2, x_3, \dots, x_d, C_k)}{P(x_q)} \\ &= \frac{P(x_1 | x_2, x_3, \dots, x_d, C_k) \cdot P(x_2 | x_3, \dots, x_d, C_k) \cdot P(x_3 | \dots, x_d, C_k) \cdot \dots \cdot P(x_d | C_k) \cdot P(C_k)}{P(x_q)} \end{aligned}$$

This is called Chain Rule of Probability

$$= \frac{P(x_1 | C_k) P(x_2 | C_k) P(x_3 | C_k) \dots P(x_d | C_k) P(C_k)}{P(x_q)}$$

Naïve Assumption : Input variables are independent of each other.

$$= \frac{P(C_k)}{P(x_q)} \prod_i^d P(x_i | C_k)$$

$$\hat{y} = arg_k \max \left\{ \frac{P(C_k)}{P(x_q)} \prod_i^d P(x_i | C_k) \right\}$$

Different types of NB :

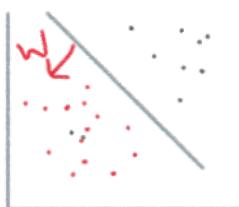
1. Multinomial NB : Text Data
2. Gaussian NB : Real Value
3. Bernoulli NB : Binary Feature

Maximum Likelihood / probabilistic approach for Logistic Regression

X_q - query point

$w^T x$ - dot product

Here $y = (0,1)$



Given a X_q , what is its probability that it belongs to class 1

$$P(y = 1 | X_q) = \sigma(w^T x)$$

$$P(y = 0 \mid Xq) = 1 - \sigma(w^T x)$$

$$\begin{aligned} P(Y_i \mid Xq) &= (\sigma(w^T x))^{Y_i} + (1 - \sigma(w^T x))^{1-Y_i} \text{ also called likelihood of } Xq \\ \text{If } Y_i = 1 &= (\sigma(w^T x))^{Y_i} \\ Y_i = 0 &= (1 - \sigma(w^T x))^{Y_i} \end{aligned}$$

We need to maximize the likelihood for each training data point.

$$\max \prod_{i=1}^n P(Y_i \mid X_i) = \max(P(Y_1 \mid X_1) + P(Y_2 \mid X_2) + \dots + P(Y_n \mid X_n))$$

According to optimization theory,

$\max f(x) = \log(\max f(x))$, it is monotonically increasing function

$$\begin{aligned} \max \prod_{i=1}^n P(Y_i \mid X_i) &= \max(\log(P(Y_1 \mid X_1)) + \log(P(Y_2 \mid X_2)) + \dots + \log(P(Y_n \mid X_n))) \\ &= \max \sum \log((P(Y_i \mid X_i))) \\ &= \max \sum \log((\sigma(w^T x))^{Y_i} + (1 - \sigma(w^T x))^{1-Y_i}) \\ &= \max \sum Y_i \cdot \log((\sigma(w^T x))^{Y_i} + (1 - Y_i) \log(1 - \sigma(w^T x))^{1-Y_i}) \\ &\quad \text{since } \log b^a = a \log b \\ &= \max \sum Y_i \cdot \log((\sigma(w^T x))^{Y_i} + (1 - Y_i) \cdot \log(1 - \sigma(w^T x))^{1-Y_i}) \\ &= -\min \sum Y_i \cdot \log((\sigma(w^T x))^{Y_i} + (1 - Y_i) \log(1 - \sigma(w^T x))^{1-Y_i}) \\ m^*, c^* &= -\min \sum Y_i \cdot \log(P_i) + (1 - Y_i) \log(1 - P_i) \end{aligned}$$