

Logistic Regression

Logistic Regression is a classification algorithm. It is used to predict a binary outcome (1 / 0, Yes / No, True / False) given a set of independent variables.

Logistic regression as a special case of Linear regression when the outcome variable is categorical, where we are using log of odds as dependent variable. In simple words, it predicts the probability of occurrence of an event by fitting data to a logit function.

$$\underset{\omega}{\operatorname{argmax}} \sum_{i=1}^m y_i w^T x_i$$

signed distance

$w^T x_i$ } dist from x_i to π (w is a unit vector)

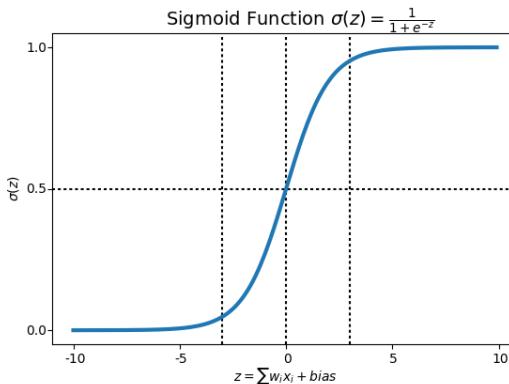
$y_i w^T x_i$: +ve $\Rightarrow \pi$ as defined by w correctly classifies x_i

\hookrightarrow -ve \Rightarrow incorrectly classifies x_i

Here we will always trying to maximize signed-distance

Sigmoid Function

- Used to map predicted values to probabilities.
- The function maps any real value into another value between 0 and 1.
- In machine learning, we use sigmoid to map predictions to probabilities.



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

Monotonic Function-

If x increases then $f(x)$ is also increases.

Eg $\log(x)$

$$\begin{aligned} \text{if } g(x) \text{ is a monotonic fn.} \\ \arg \min_x f(x) = \arg \min_x g(f(x)) \\ \arg \max_x f(x) = \arg \max_x g(f(x)) \end{aligned}$$

Mathematical formulation of Objective Function

$$\begin{aligned} \text{geometric} \quad w^* = \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) \\ \text{Prob. method} \quad w^* = \arg \min_w \sum_{i=1}^n y_i \log p_i - (1-y_i) \log(1-p_i) \\ p_i = \sigma(w^T x_i) \end{aligned}$$

Weight Vector

Optimal w^* of above equation is called Weight Vector.

Case 1. Given a feature f_i , if the weight vector component corresponding to feature f_i is positive, then for any query point x_q , if the feature value corresponding to the i th feature increases, its probability of belonging to +ve class increases.

Case 2. Given a feature f_i , if the weight vector component corresponding to feature f_i is negative, then for any query point x_q , if the feature value corresponding to the i th feature increases, its probability of belonging to -ve class increases.

Need of Regularization -

We all know Machine learning is about training a model with relevant data and using the model to predict unknown data. By the word unknown, it means the data which the model has not seen yet. We have trained the model, and are getting good scores while using training data. But during the

process of prediction, we found that the model is underperforming when compared to the training part. Now, this may be a case of **over-fitting** which is causing incorrect prediction by the model. Regularizing the model, in this case, can help resolve the problem.

Regularization -

Regularizations are techniques used to reduce the error by fitting a function appropriately on the given training set and avoid overfitting.

It is a technique which make slight modification to the learning model/algo such that model generalizes better. This in turn improves the model's performance on the unseen data as well.

L₂ regularization : Overfitting vs Underfitting:

$$\hat{w}^t = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

let $z_i = y_i w^T x_i$

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-z_i))$$

$w = \underset{i=1}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-z_i)) \geq 0$, So minimal value of w is 0.

Ques- When does this zero occur?

Ans- If z_i is positive and tends to + infinity then $\exp(-z_i)$ tends to become 0.

1. $z_i = +ve$, It means x_i is correctly classified by w
2. z_i tends to infinity, then to reach to minima(reach to least value ie 0)

If we pick w such that- As y_i and x_i are fixed so we can modify only w

- All training point are correctly classified
- z_i tends to +ve infinity. Then such a w is Best " w ".

With best w if there are outlier then it cause overfitting. Or $w_i = +infinity$ or $-infinity$ To get rid of +infinity and -infinity we use Regularization.

L1 Regularization -

L1 is also called Lasso Regression.

L1 is differ from L2 only in penalizing the high coefficients. It uses $|W_j|$ (modulus)instead of squares of W , as its penalty. In statistics, this is known as the L1 norm.

In the case of the L1, penalty has the effect of forcing some of the coefficient estimates to be exactly equal to zero when the tuning parameter λ is sufficiently large. Therefore, the lasso method also performs variable selection and is said to yield sparse models.

$$\begin{aligned}
 & \|\omega\|_2^2 \text{ for reg.} & w_i \rightarrow +\infty \\
 & \downarrow & \text{or } w_i \rightarrow -\infty \\
 & \|\omega\|_1 \text{ for reg.} & \|\omega\|_1 = \sum_{i=1}^d |w_i| \\
 & w^* = \underset{\omega}{\operatorname{argmin}} \left(\underbrace{\text{logistic loss}_{\text{for training data}}}_{\text{will avoid } w_i \rightarrow +\infty \text{ or } -\infty} + \lambda \|\omega\|_1 \right) & \xrightarrow{\text{hyperparam}} \lambda \|\omega\|_1 \\
 & & \xrightarrow{\text{L1-reg}}
 \end{aligned}$$

- One benefit we get from L1 is Sparsity.

L2 Regularization -

$$\begin{aligned}
 & \text{regularization} \\
 & \omega^* = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i \omega^T x_i)) + \lambda \omega^T \omega
 \end{aligned}$$

- First term is called Loss-term.
- Second part of eqn is called Regularization-term
- Second part will ensure that w_i don't tends to $+\infty$ and $-\infty$

$$\begin{aligned}
 & \min \left(\text{loss-fn over training data} + \text{reg.} \right) \\
 & \text{cross-validation} \\
 & \lambda = 0 \Rightarrow \text{overfit} \rightarrow \text{high variance} \\
 & \lambda = \text{large} \Rightarrow \text{underfit} \rightarrow \text{high bias}
 \end{aligned}$$

Here, λ is the tuning parameter that decides how much we want to penalize the flexibility of our model. The increase in flexibility of a model is represented by increase in its coefficients, and if we want to minimize the above function, then these coefficients need to be small.

L1 Regularization is also called Lasso Regression.

L2 Regularization is also called Ridge Regression.

What does Regularization achieve?

A standard least squares model tends to have some variance in it, i.e. this model won't generalize well for a data set different than its training data.

Regularization, significantly reduces the variance of the model, without substantial increase in its bias. So the tuning parameter λ , used in the regularization techniques described above, controls the impact on bias and variance. As the value of λ rises, it reduces the value of coefficients and thus reducing the variance.

Till a point, this increase in λ is beneficial as it is only reducing the variance(hence avoiding overfitting), without loosing any important properties in the data. But after certain value, the model starts loosing important properties, giving rise to bias in the model and thus underfitting. Therefore, the value of λ should be carefully selected.

Regularization - <https://towardsdatascience.com/regularization-in-machine-learning-76441ddcf99a>
[\(https://towardsdatascience.com/regularization-in-machine-learning-76441ddcf99a\)](https://towardsdatascience.com/regularization-in-machine-learning-76441ddcf99a)

Sparsity -

$$w = \langle w_1, w_2, w_3, \dots, w_d \rangle$$

- If w vector is sparse then solution of logistic Regression is said to be Sparse.
- **If we use L1 reg then all less-important feature become zero.**
- **If we use L2 reg then all less-important feature become small but not necessarily zero.**

$$f_1, f_2, \dots, f_i, \dots, f_d$$

- Above are Features, f_i is less important feature
- For f_i corresponding weight vector is w_i .

$$w = \langle w_1, w_2, \dots, w_i, \dots, w_d \rangle$$

- For L1 reg. w_i become zero.
- For L2 reg. w_i become small value but not necessarily zero.

Linear regression is the standard algorithm for regression that assumes a linear relationship between inputs and the target variable. An extension to linear regression involves adding penalties to the loss function during training that encourage simpler models that have smaller coefficient values. These extensions are referred to as Regularized Linear Regression or Penalized Linear Regression.

Elastic-net -

Elastic net is a popular type of regularized linear regression that combines two popular penalties, specifically the L1 and L2 penalty functions.

$$\omega^* = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-z_i)) + \lambda_1 \|\omega\|_1 + \lambda_2 \|\omega\|_2^2$$

Brief Summary of Regularization -

A regularization technique is in simple terms a penalty mechanism which applies shrinkage (driving them closer to zero) of coefficient to build a more robust. Few common ways to regularize models are -

1. L1 Regularization or Lasso Regularization -

This add regularization terms in the model which are function of absolute value of the coefficients of parameters. The coefficient of the parameters can be driven to zero as well during the regularization process. Hence this technique can be used for feature selection and generating more parsimonious model

2. L2 Regularization or Ridge Regularization -

This add regularization terms in the model which are function of square of coefficients of parameters. Coefficient of parameters can approach to zero but never become zero.

3. Combination of the above two such as Elastic Nets -

This add regularization terms in the model which are combination of both L1 and L2 regularization.

2. Probabilistic Interpretation Gaussian Naive Bayes

Probabilistic interpretation of Logistic regression

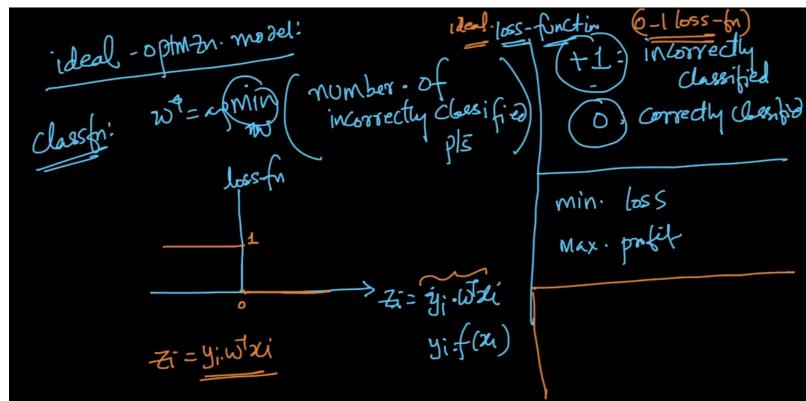
Logistic regression : ① Geometry & simple algebra
 ↳ ② Probability
 ↳ ③ loss minimization

Prob: $\omega^* = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n -y_i \log p_i - (1-y_i) \log(1-p_i) + \text{reg}$

where $p_i = \sigma(\omega^T x_i)$

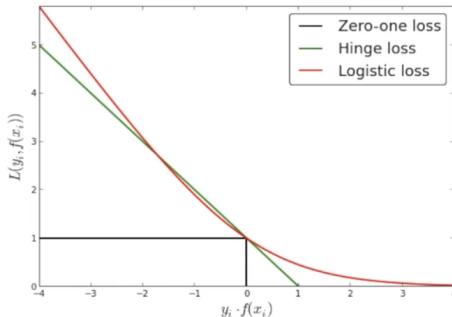
+1 or 0

Loss Minimization interpretation of LR



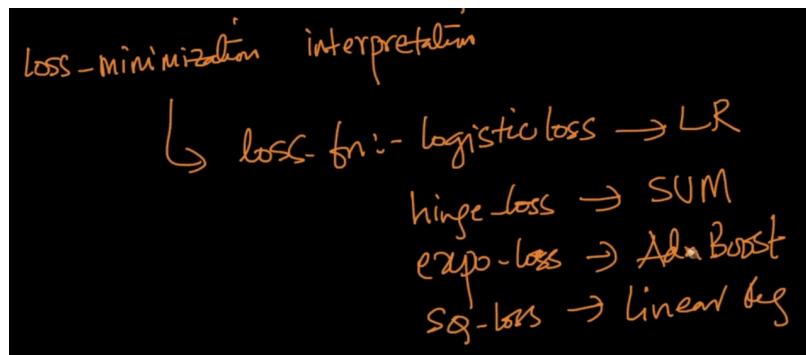
$$\overset{\text{ideal}}{=} \left[w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n 0-1 \text{ loss}(x_i, y_i, w) \right]$$

$$0-1 \text{ loss}(z_i) = \begin{cases} 1 & \text{if } z_i < 0 \\ 0 & \text{if } z_i \geq 0 \end{cases}$$



- To minimize the loss function to make the optimize model by using differentiation in calculus.
- $f(x)$ is differentiable when function is continuous, if its not continuous then its not differentiable.
- Only smooth curves are differentiable.
- Here Hinge-loss is continuous but not differentiable.
- As loss func is not continuous so its not differentiable. So we use Logistic-loss

Loss minimization Interpretation are-

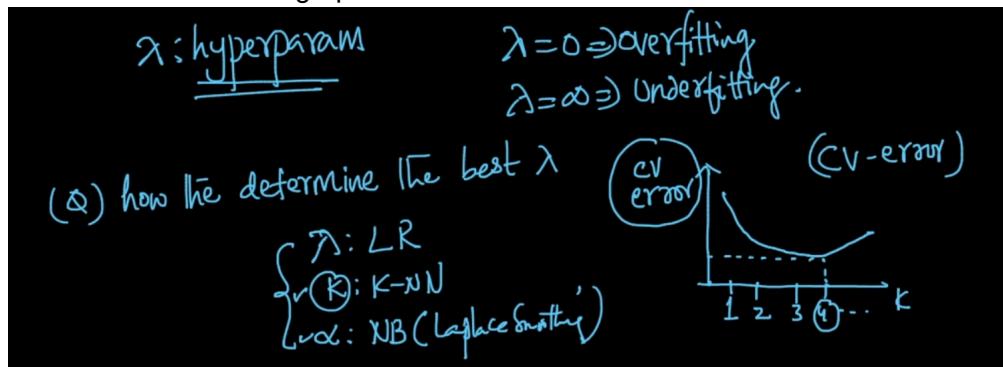


Hyperparameter Tuning -

The hyperparameter is the default parameter that works in all situations. They are termed as an important part of a model. It is not necessary that you can only use the default parameters you can make tweaks if the situation demands. It is important to have three sets in which data is divided like training, testing, validation set whenever you make tweaks in the default parameter to get the necessary accuracy so as to stop data leaks.

- Hyper parameters control the behavior of the algorithm that is used for modeling.
- Hyper parameters are passed in the arguments during the initialization of the algorithm. (Ex. Specifying the criterion for decision tree building)
- If you want to check about the hyperparameters for an algorithm you can make use of the function `get_params()`.

Eg- For k in K-NN we use below graph CV error VS K value.



- K in K-NN is an integer value.
- λ in Logistic Regression is real no.
- To find λ we use Grid Search.

During tuning of the hyper parameters the data should always be divided into three parts that are training, validation, and testing so as to stop data leak.

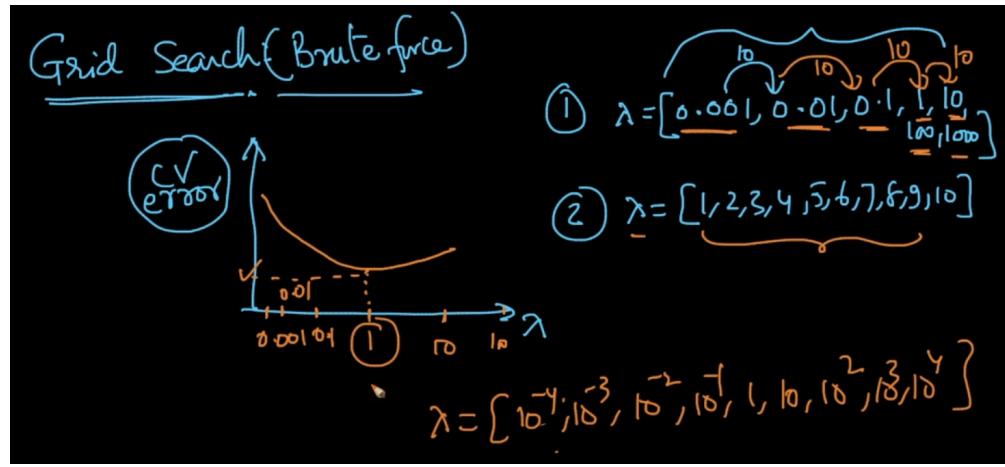
Approach to find Hyper-parameter -

1. **Grid Search CV** - It takes all the parameters combination.
2. **Randomized Seach CV** - It can sample a given number of possibilities with certain distribution from a parameter space.

Grid Search

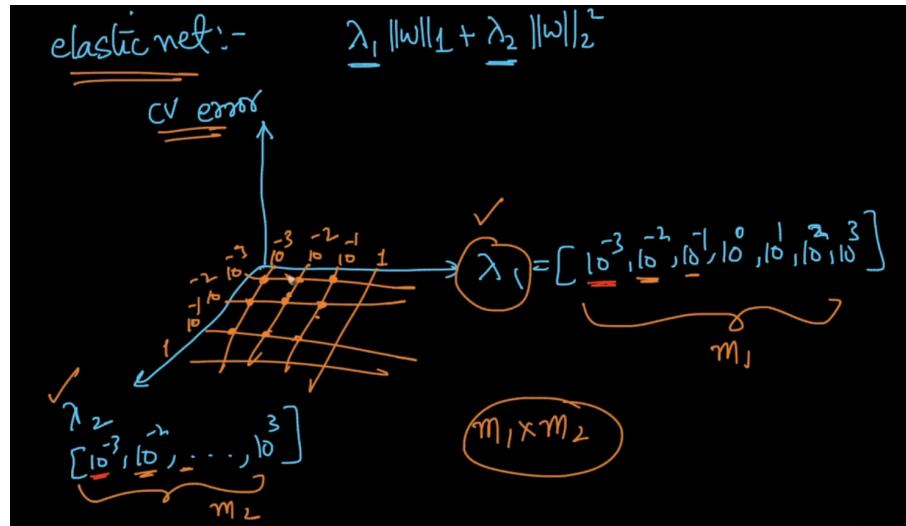
- Grid-searching is the process of scanning the data to configure **optimal hyper-parameter** for a given model.
- It iterates through every parameter combination and stores a model for each combination.
- Every model that is built is validated and ranked.
- The best performing model having best hyperparameters values are taken.

We can take any range of value for λ .



Grid Search in Elastic-net

- Here we have λ_1 and λ_2 .



In General,

- As hyperparameter increases, the no. of times model needs to be trained increases exponentially.
- Hence Grid Search is not good when there are a lot of parameters. So we use **Random Search**.

$$\lambda_1 : 1 \text{ hyperparam.} - (m_1) \text{ times train}$$

$$\lambda_1, \lambda_2 : 2 \text{ hyperparam.} - (m_1 \times m_2)$$

$$\lambda_1, \lambda_2, \lambda_3 : (m_1 \times m_2 \times m_3)$$

The **benefit** of grid search is that it is **guaranteed** to find the optimal combination of parameters supplied.

The **drawback** is that it can be very time consuming and computationally expensive.

Random Search-

- In Random Search, a random combinations of the hyperparameters are used in a given interval, to find the best solution for the built model.
- It searches the specified subset of hyperparameters randomly instead of exhaustively.
- The major benefit being decreased processing time.
- Its has proven to yield better results comparatively to Grid Search.
- The **drawback** of random search is that it yields **high variance**.

We are **not guaranteed** to find the optimal combination of hyperparameters.

Feature/Column Standardization -

Standardization is another scaling technique where the values are centered around the mean with a unit standard deviation. This means that the mean of the attribute becomes zero and the resultant distribution has a unit standard deviation.

Data Standardization = Feature/Column Standardization = Mean Centering = Scaling all are same.

Column / feature standardization

$$x'_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j} : \text{standardization}$$

- Here x_{ij} is datapoint,
- μ_j is Mean of Column
- σ_j is Standard Deviation of Column

Feature Importance

- Only valid for independent feature.
- If features are independent only then we can use $|w_j|$ value of weight vector or value as a Feature importance value.

(K-NN) :- feature - imp \rightarrow forward feature Seln.
 (NB) :- $P(x_i | y=+1)$ \rightarrow features which are important
 (LR) :- $w_j's \rightarrow$ to determine f.I

- Here w_j is weight corresponding to feature, FI is Feature Importance.

Eg. Given a set of datapoint to distinguish b/w male and female, here we give more weight to hair, hair is considered as one of the feature.

Collinearity or Multicollinearity

- Here when we get another feature by some mathematical operation.
- Below is an eg. of collinear, when we have more than 2 feature which are co-related then its Multicollinearity.

Collinearity: f_i, f_j
 s.t. if, $f_i = \alpha f_j + \beta$
 then f_i & f_j are collinear - Q

Ques- Why does weight of feature is not useful as Feature Importance if features are Collinear?

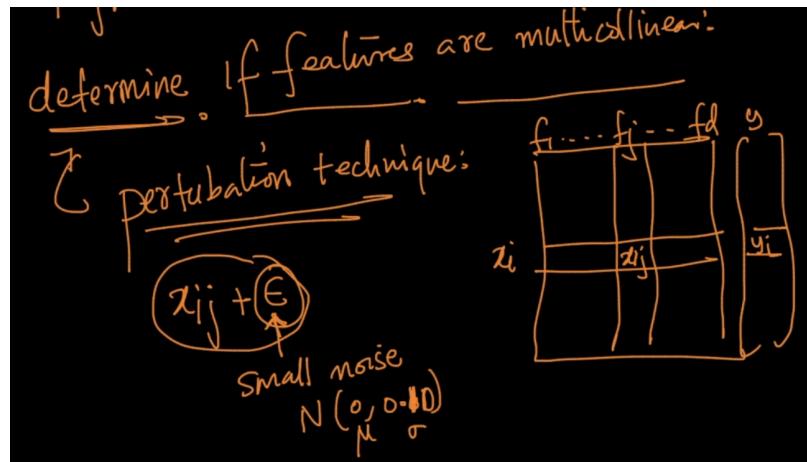
Ans- Different weight of feature will result same classifier.

- If features are collinear then weight vector can change arbitrarily which imply w_j or weight vector can not be used as Feature Importance.

Determine Feature are Multicollinear

Perturbation Technique -

- Before Perturbation, we will adding noise we will computer weight vector, $w = \langle w_1, w_2, \dots, w_j, \dots, w_d \rangle$
- Here we add small noise to every x_{ij} .
- After perturbation, We will again trained Logistic Regression and compute $\tilde{w} = \langle \tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_j, \dots, \tilde{w}_d \rangle$.
- If w and \tilde{w} differ drastically then we can conclude that Features are Collinear. It means we can not use $|w_j|$ as Feature Importance value.
- If w and \tilde{w} differ slightly then we conclude that $|w_j|$ as Feature Importance.



Conclusion - We should always check Collinearity before start using $|w_j|$ as Feature Importance value.

- Forward Feature Selection technique can be used in any models

Train & Runtime Space and Time Complexity

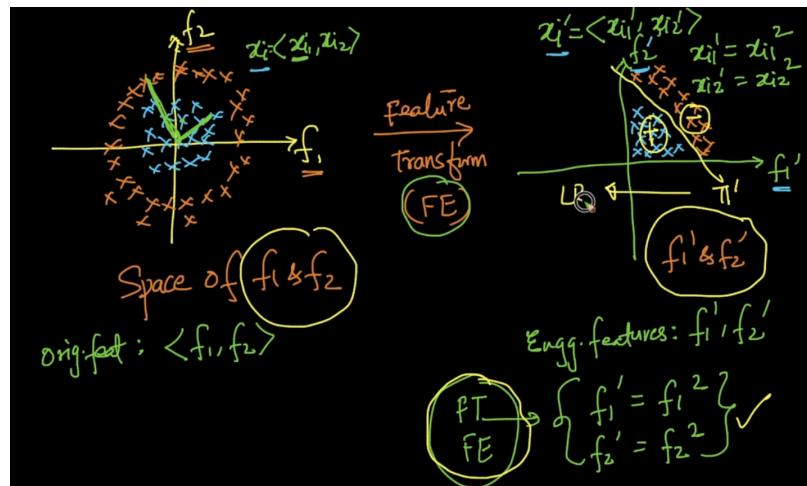
n = No. of points in DTrain, d = Dimension

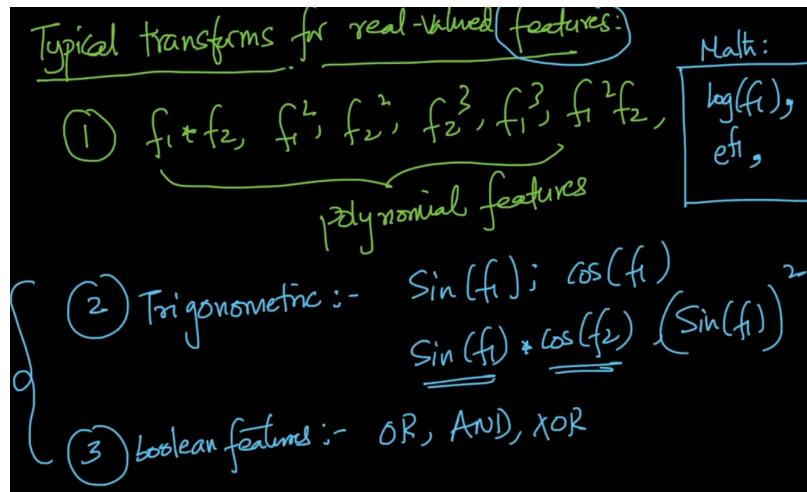
1. Train
 - Time Complexity - $O(nd)$
2. Runtime
 - Time Complexity - $O(d)$
 - Space Complexity - $O(d)$

Non-Linearly separable data & Feature Engineering/ Transform

Ques- For a given datapoints having two concentric circle, Can we distinguish using plane?

Ans- Yes by doing Feature transform.





Fundamental Techniques of Feature Engineering-

<https://towardsdatascience.com/feature-engineering-for-machine-learning-3a5e293a5114>
[\(https://towardsdatascience.com/feature-engineering-for-machine-learning-3a5e293a5114\)](https://towardsdatascience.com/feature-engineering-for-machine-learning-3a5e293a5114)

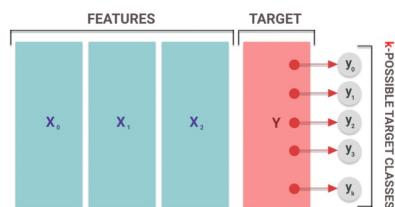
Generalized Linear Model

Log-reg = GNB+ Bernoulli

1. Multinomial LR
2. Linear Regression
3. Poisson Regression

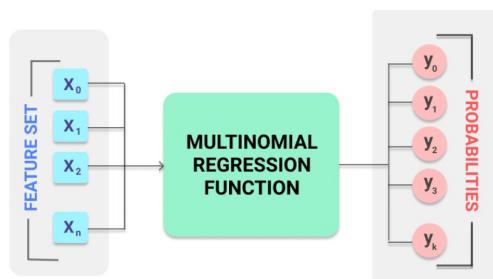
1. Multinomial LR -

MLR can predict one out of k-possible outcomes, where k can be any arbitrary positive integer.



1.1 How does MLR work?

The multinomial regression function is a statistical classification algorithm. What this means is that once we feed the function a set of features, the model performs a series of mathematical operations to normalize the input values into a vector of values that follows a probability distribution.



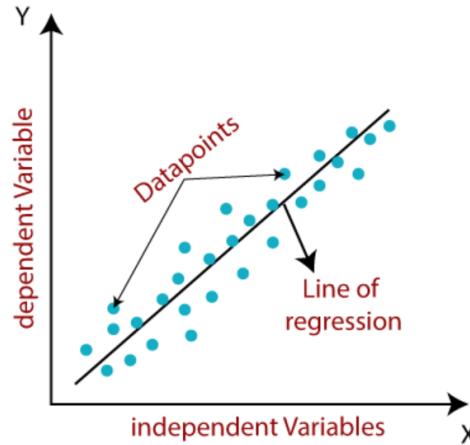
- The input that we give to the model is a feature vector, X, containing features $x_1, x_2, x_3 \dots, x_n$.
- The output we get is a probability vector Y, containing probabilities $y_1, y_2, y_3 \dots, y_k$ for the k target classes.
- Here, $y_1 + y_2 + y_3 \dots + y_k = 1$, since the total probability of all the possible events in a system is always 1.

Finally, the outcome with the highest probability will be the predicted outcome for the given feature set.

2. Linear Regression -

Linear Regression is a type of regression analysis where the number of independent variables is one and there is a linear relationship between the independent(x) and dependent(y) variable.

Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable.



3. Poisson Regression -

Poisson regression is used to model count data, assuming that the label has a Poisson distribution. For example, you might use it to predict the number of calls to a customer support center on a particular day.

Refer - <https://towardsdatascience.com/generalized-linear-models-8738ae0fb97d>
[\(https://towardsdatascience.com/generalized-linear-models-8738ae0fb97d\)](https://towardsdatascience.com/generalized-linear-models-8738ae0fb97d)