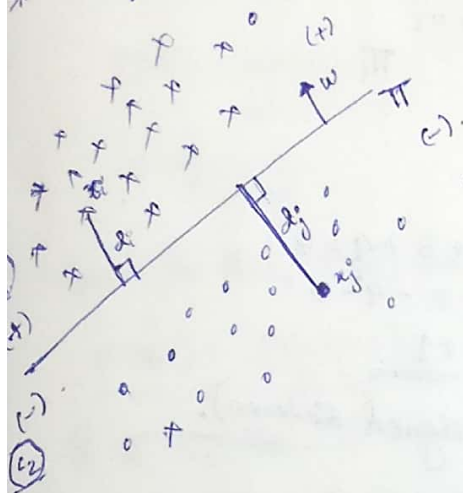


# Logistic Regression - (classification technique).

## Geometric Intuition

\* Assumption of LR:- classes are almost/perfectly linearly separable.



Given:-  $D_n = \{+ve, -ve\}$  pts.

task:- find :-  $w$  &  $b$

$$\Pi: w^T x + b$$

find a plane  $\Pi$  that best separates +ve pts from negative pts.

Say  $y_i = +1$  for class 1  
&  $y_i = -1$  for class 2.

$$y_i \in \{-1, 1\}.$$

$$d_i = \frac{w^T x_i}{\|w\|} \quad \text{If } w \text{ is unit vector, } \|w\| = 1.$$

$$\text{So, } d_i = w^T x_i.$$

Say  $\vec{w}$  points towards the class  $c_1$ .

$\vec{w}$  &  $\vec{x}_i$  in same dir<sup>n</sup>. Hence  $d_i = +ve$ .

$\vec{w}$  &  $\vec{x}_j$  in opposite dir<sup>n</sup>. Hence  $d_j = -ve$ .

classifier  $\Rightarrow$  if  $w^T x > 0$   
then  $y_i = +1$   
else  $y_i = -1$

Assume plane passes through origin for simplifying.

If we consider,  $y_i w^T x_i > 0$  for correctly classified pts.

&  $y_i w^T x_i < 0$  for mis-classified pts.

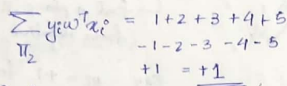
For classifier to be good,  $\rightarrow$  min # misclassifications.  
(or)  
max # correctly classified pts.

Find a plane  $\Pi$  s.t.  $y_i w^T x_i > 0$  for as many pts as possible.  
( $w, b$ )

mathematical optimization  
fun<sup>n</sup> that we need to solve

$$\max_w \sum_{i=1}^n y_i w^T x_i \quad (\text{Or}) \quad w^* = \operatorname{argmax}_w \left( \sum_{i=1}^n y_i w^T x_i \right)$$

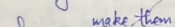
no  $\phi$  — buffer.



So,  $\pi_2$  is better classifier than  $\pi_1$ .  
terrible classifier.

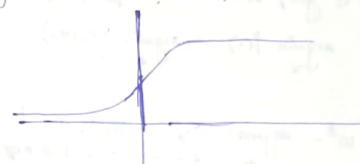
- this difference is due to one single extreme outlier point is impacting the model.

4 signed distance range -  
make them small. tapering off.



The graph shows a curve on a coordinate system. The vertical axis is labeled 'signed distance' and the horizontal axis is labeled 'range'. The curve starts at a negative value on the y-axis, crosses the x-axis, and then levels off towards a positive value. An arrow points to the leveling-off part of the curve with the label 'tapering off.'

- optimizing for

$$\underset{w}{\operatorname{argmax}} \sum_{i=1}^n y_i w^T x_i \quad \rightarrow \quad \underset{w}{\operatorname{argmax}} \sum_{i=1}^n f(y_i w^T x_i)$$
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$T(x) \quad \begin{array}{l} \max: -1 \\ \min: 0 \end{array}$$

$$\sigma(0) = 0.5,$$

$$\sigma(z_i) = P(y_i = 1) \quad z_i = \omega^T x_i$$

So, the points lies on the decision boundary. So,  $P(y_i=1)$  should be 50% that's exactly what  $\sigma(-\chi_i)$  gives.

There would be many other functions to do the tapering off. Why  $\sigma$ ?

## Optimization Problem

maximization problem

$$w^* = \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \left( \frac{1}{1 + \exp(-y_i w^T x_i)} \right) \rightarrow \text{optimization problems.}$$

$\varepsilon x \rightarrow \log(x)$  is monotonically increasing for  $x \geq 0$ .

Let  $g(x) = \log(x)$ .

$$x^* = \underset{x}{\operatorname{argmin}} f(x).$$

$$x' = \operatorname{argmin}_x g(f(x))$$

claim  $x^* = x'$  as  $g(x)$  is monotonic function

So, if  $g(x)$  is monotonic function,

$$\arg\min_x f(x) = \arg\min_x g(f(x))$$

$$\text{So, } w^* = \arg\max_w \sum \log \left\{ \frac{1}{1 + \exp(-y_i w^T x_i)} \right\}$$

$$w^* = \arg\max_w \sum -\log(1 + \exp(-y_i w^T x_i))$$

$$w^* = \arg\min_w \sum \log(1 + \exp(y_i w^T x_i)) \quad (+1 \times -1)$$

Alternative  
formulation  
using probabilistic  
approach

$$w^* = \arg\min_w \sum_{i=1}^n -y_i \log p_i - (1 - y_i) \log(1 - p_i)$$

$$p_i = \sigma(w^T x_i) \quad \text{Both are the same.}$$

Weight Vector

$$w^* = \arg\min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

weight-vector.  $w = \langle w_1, w_2, \dots, w_d \rangle$   
 $f_1 \quad f_2 \quad \quad \quad f_d$

decision:  $x_i \rightarrow y_i$  if  $w^T x_i > 0$  then  $y_i = +1$   
else  $y_i = -1$

Case I:  $(w_i = +ve)$   
If  $f_i \uparrow$ ,  $w_i x_{gi} \uparrow \Rightarrow w^T x_i \uparrow \Rightarrow \sigma(w^T x_i) \uparrow \Rightarrow P(y_i = 1) \uparrow$

Case II:  $(w_i = -ve)$   
If  $f_i \uparrow$ ,  $w_i x_{gi} \downarrow \Rightarrow w^T x_i \downarrow \Rightarrow \sigma(w^T x_i) \downarrow \Rightarrow P(y_i = 1) \downarrow$   
 $\downarrow$   
 $P(y_i = -1) \uparrow$

So, probability of class increased on increasing the  $i$ th feature value in any data point if  $w_i$  is +ve & vice versa.

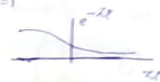
Le regularization: Overfitting vs Underfitting

$$w^* = \arg\min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

$$\text{Let } z_i = y_i w^T x_i$$

$$\text{Then } w^* = \arg\min_w \sum_{i=1}^n \log(1 + \exp(-z_i))$$

$$\times \exp(-z_i) \geq 0$$



$$\text{So, } \log(1 + \exp(-z_i)) \geq \log(1) \quad \{ \text{monotonic property} \}$$

$$\log(1 + \exp(-z_i)) \geq 0$$

Minimal value of  $\log(1 + \exp(-z_i))$  is '0'.

It occurs when  $z_i \rightarrow +\infty$ .

$$\text{So, } \sum_{i=1}^n \log(1 + \exp(-z_i)) \rightarrow 0 \text{ when } z_i \rightarrow +\infty \text{ for all } i \in \{1, 2, \dots, n\}$$

So, we need to adjust  $w_i$ 's such that all  $z_i$ 's  $\rightarrow \infty$ .

such a  $w$  is the best  $w$ .

But the problem with this is overfitting.

We may overfit to the given dataset while trying to achieve  $z_i \rightarrow \infty$   $\forall i \in \{1, 2, \dots, n\}$

\* For perfectly fitting each and every point, we may end up having all  $w_i$ 's as  $+\infty$  or  $-\infty$

\* So, one way we can limit the model from overfitting to given data is by kind of limiting its magnitude.

That's what we do in regularization.

Regularization:

$$w^* = \arg\min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda w^T w$$

loss term regularization term

Case I:  $w_i$ 's  $\rightarrow +\infty / -\infty$

loss term  $\rightarrow 0$

regularization term  $\rightarrow \infty$

Case II:  $w_i$ 's  $\rightarrow$  small

loss term  $\rightarrow$  large

regularization term  $\rightarrow 0$



So, optimising for both will end up giving us model best fit.

$\lambda$  - hyper-parameter. — found via CV.

Small  $\rightarrow$  overfitting.  
Large  $\rightarrow$  underfitting.

Many optimisation functions in ML are of form  
 $\min \left( \text{loss function} + \text{regularisation} \right)$   
of training data

L1 regularisation and sparsity

$$\|w\|_2^2 = w^T w$$

$$w^* = \arg \min_w \underbrace{\sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))}_{\text{logistic loss}} + \underbrace{\lambda \|w\|_2^2}_{\text{L2-regularisation}}$$

\* Now, instead of L2-regularisation, we can also use L1-regularisation.

$$\|w\|_1 = \sum_{i=1}^d |w_i| \quad \text{— L1-norm}$$

$$w^* = \arg \min_w \left( \text{logistic loss for training data} \right) + \lambda \|w\|_1$$

Advantage of L1-reg  $\rightarrow$  causes sparsity (in  $w$ -vector).

\*  $w_i = 0$  for all  $w_i$  corresponding to less important features.  
\* In L2-reg  $w_i$ 's becomes small for less imp. features, but not '0'.

Why L1 causes sparsity as compared to L2?

$\rightarrow$  see in optimisation, intuitive soln.

Elastic-Net :- 
$$w^* = \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2$$

Probabilistic Interpretation

For derivation, refer pdf - [cs.cmu.edu/~tom/mlbook/NBayerLogReg.pdf](https://cs.cmu.edu/~tom/mlbook/NBayerLogReg.pdf).

From geometry we get

$$w^* = \arg \max_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \text{reg.}$$

$+1 \text{ or } -1$

From probability, after derivatn, we'll get

$$w^* = \arg \min_w \sum_{i=1}^n -y_i \log p_i - (1 - y_i) \log(1 - p_i) + \text{reg.}$$

$+1 \text{ or } 0$

where  $p_i = \sigma(w^T x_i)$ .

We'll see that both are the same only. Ignore the reg. term for the time being.

Case I:  $y_i = +ve$ .

$$\begin{aligned} \text{geom} &:- y_i = +1 \\ \text{prob} &:- y_i = +1. \end{aligned}$$

$$\text{geom}:- \log(1 + \exp(-w^T x_i)) \quad \checkmark$$

$$\text{prob}:- -1 * \log \frac{1}{1 + \exp(-w^T x_i)} = \log(1 + \exp(-w^T x_i)) \quad \checkmark$$

Case II:  $y_i = -ve$ .

$$\begin{aligned} \text{geom} &:- y_i = -1. \\ \text{prob} &:- y_i = 0. \end{aligned}$$

$$\text{geom}:- \log(1 + \exp(w^T x_i))$$

$$\text{prob}:- -(1-0) \log \left[ 1 - \frac{1}{1 + \exp(-w^T x_i)} \right]$$

$$= -\log \left[ \frac{1 + \exp(-w^T x_i) - 1}{1 + \exp(-w^T x_i)} \right]$$

$$= \log \left[ \frac{1 + \exp(-w^T x_i)}{\exp(-w^T x_i)} \right]$$

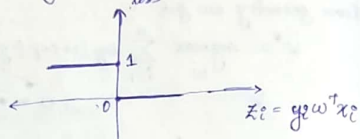
$$= \log \left[ 1 + \frac{1}{\exp(-w^T x_i)} \right] = \log[1 + \exp(w^T x_i)]$$

## Loss minimization interpretation of LR

An ideal loss function would be something like following—

→ +1: incorrectly classified pts.  
→ 0: correctly classified pts.

Such a loss function is called 0-1 loss function.



$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n 0-1 \text{ loss}(x_i, y_i; w) \quad 0-1 \text{ loss}(x_i) = \begin{cases} 1 & \text{if } x_i < 0 \\ 0 & \text{if } x_i \geq 0. \end{cases}$$

But a big issue with selecting this loss fun<sup>n</sup>:- non-differentiable.

\* Since it's not differentiable, we try to find a loss function which is close to 0-1 loss but differentiable.

\* One such approximation is logistic loss. — gives L.R. (Logistic Regression)

Another approximation would be hinge loss. — gives SVM.

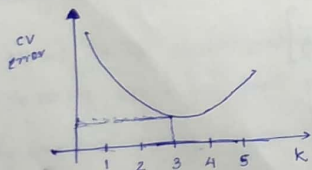
Square loss — Linear regression

exp. loss — AdaBoost.

Refer:- Loss functions for classification — Wikipedia for graphs.

## Hyperparameter Search/Optimization

In KNN, we used CV error to determine ~~error~~ hyperparameter values.



But k in KNN is integer.

whereas  $\lambda$  in L.R. is a real no.

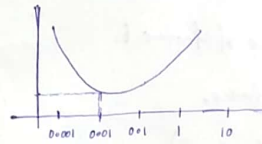
It can have infinitely many possibilities.

So, we have different techniques for finding values of  $\lambda$ .

## 1) Grid Search

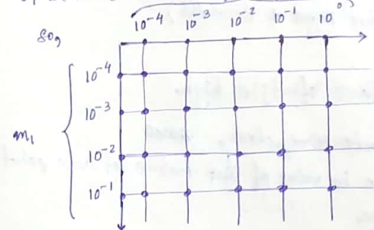
$$\lambda = [0.001, 0.01, 0.1, 1, 10, 100].$$

By doing the 10x jump, we can search a much larger space.



Grid search is a brute force technique.

If we have elastic-net,  $m_2$   $\lambda_1 \|w\|_1 + \lambda_2 \|w\|_2$ , we have 2 parameters so,



for 2 hyper-parameters,  $(m_1 \times m_2)$  values we have to search if we are doing grid search.

So, as #hyperparameters  $\uparrow$ , #times the model needs to be trained increases exponentially.

\* In deep learning, we have 10s & 100s of hyperparameters.

So, we have an alternate approach.

## 2) Random Search

$\lambda \in [10^{-4}, 10^4]$  — randomly pick values in the given interval.

Scientifically proven that random search is as good as grid search and also at times faster than grid search.

Both search available in sci-kit.

Column/feature standardization (Mean centering & scaling).

Just like in K-NN, here also in Logistic Regression, we are dealing with distances from the separating hyperplane.

So, different scales of different features can impact the model.

\* So, it's mandatory to perform feature standardization before training on Train.

### Feature Importance and Model Interpretability

In KNN, we found feature imp. by forward feature select<sup>n</sup>.

In NB, using  $P(X_i | Y=+1)$  we found importance of feature  $i$ .

In LR, we can use  $w_j$ 's to get feature importance.

Assumption:- All features are independent. (Naive Bayes).

Actually, LR is (Gaussian Naive Bayes + Bernoulli)

If  $|w_j|$  is large, then the importance of  $f_j$  is high.

If  $w_j$  is +ve and large, it impacts +ve class. ~~more~~  
with increase in value of that feature of data point  
& vice-versa.

Ex predict gender:- male or female  
(+1) (-1)

(h1) hair\_length:  $|w_{HL}|$  is large &  $w_{HL}$  is +ve.

$|w_{HL}| \uparrow$ ;  $P(Y_2 = -1) \uparrow$  as probability of being female inc<sup>s</sup>  
with increase in hair length.

(h2) height:  $|w_H|$  is large and  $w_H$  is +ve.

$|w_H| \uparrow$ ;  $P(Y_2 = +1) \uparrow$

Model Interpretability → Get the most important features & analyse them to give reasoning of the model behaviour.

we can handpick top 10 features or so and do the above.

### Co-linearity and Multi-collinearity

Co-linearity:-  $f_i$  &  $f_j$  are co-linear if

$$f_i = \alpha f_j + \beta$$

Multi-collinearity:- if  $f_1, f_2, f_3, f_4$  s.t.

$$f_1 = \alpha_1 f_2 + \alpha_2 f_3 + \alpha_3 f_4 + \alpha_4$$

then  $f_1, f_2, f_3$  &  $f_4$  are said to be multi-collinear.

Q: why does  $|w_j|$  not be useful as F.O. if features are co-linear / multi-collinear?

→ say  $w^* = \langle 1, 2, 3 \rangle$        $x_q = \langle x_{q1}, x_{q2}, x_{q3} \rangle$ .

Now say we have  $f_2 = 1.5 f_1$ .

$$\begin{aligned} \text{Then, } w^{*T} x_q &= 1 \cdot x_{q1} + 2 \cdot x_{q2} + 3 \cdot x_{q3} \\ &= x_{q1} + 2(1.5 x_{q1}) + 3 x_{q3} \\ &= 4 x_{q1} + 3 x_{q3} = \langle 4, 0, 3 \rangle. \end{aligned}$$

So,  $\langle 1, 2, 3 \rangle$  &  $\langle 4, 0, 3 \rangle$  - both feature vectors give the same result. ~~same~~

Thus, due to the linear relat<sup>n</sup> among features, we can change the weight vector without affecting the result at all yield.

So, we can't decide the feature importance uniquely.

So, before using  $w_j$  as feature importance metrics, we must determine whether features are multi-collinear or not.

### Perturbation technique →

(shake a little). For all datapoints of  $D_{train}$ , take each  $x_{ij}$  & add a

little noise to each

$$x_{ij} + \epsilon \quad \text{--- small noise } N(\mu, \sigma) \quad (0, 0.10).$$

Train model before & after perturbations



before  $w = \langle w_1, w_2, \dots, w_d \rangle$   
 after  $\tilde{w} = \langle \tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_d \rangle$

If  $w_i$  &  $\tilde{w}_i$  differ significantly, then features are colinear.  
 $\Downarrow$

can't use  $|w_i|$  as feature importance.

\* we always have forward selection technique for feature importance.

### Train & Runtime Space & Time Complexity

Train LR: It is time taken in solving the optimizat<sup>n</sup> equation

$$w^* = \underset{w}{\text{argmin}} (\text{logistic loss} + \text{reg}^n).$$

We'll see details in optimization chapter using SGD (stochastic gradient descent).

It takes  $O(nd)$  roughly.

$n$  = # of training pts.

$d$  = dimensionality.

### Runtime LR:

Space complexity  $\rightarrow O(d)$

Time complexity  $\rightarrow O(d)$ .

So, if  $d$  is small, LR is apt. for low latency applications.

Also very memory efficient.

In case of high dimensionality, we can make use of  $L_1$  regularization to create sparsity in  $w$  so as to reduce calc<sup>n</sup>s during runtime.

We can adjust  $\lambda$  to create proper amount of sparsity.

#  $\lambda \uparrow$ , bias  $\uparrow$ .

$\lambda \downarrow$ , latency  $\uparrow$  & variance  $\uparrow$

So, we have to decide based on our business model & requirements

### Real world cases

Decision surface  $\rightarrow$  Linear / hyperplane.

Assumption  $\rightarrow$  data is linearly separable or almost linearly separable.

Imbalanced data  $\rightarrow$  upsampling / downsampling.

Outliers  $\rightarrow$  less impact due to  $\sigma(x)$  function tapering off nature. But not completely avoided.

In order to handle them, we use following filtering technique:-

①  $D_{\text{train}} \rightarrow w^*$

②  $x_i \rightarrow w^* x_i$  - distance of  $x_i$  from  $\Pi$  to  $x_i$  to  $D_{\text{train}}$  give  $D'_{\text{train}}$ .

③ remove pts which are very far away from  $D_{\text{train}}$  give  $D'_{\text{train}}$ .

④  $D'_{\text{train}} \rightarrow \tilde{w}^* \rightarrow$  final sol<sup>n</sup>.

Missing values  $\rightarrow$  standard mean / median imputation techniques.

Multi-class: One Vs Rest  $\rightarrow$  Typically

extensions  $\left\{ \begin{array}{l} \text{Softmax} \\ \text{Multinomial LR} \end{array} \right. \rightarrow$  deep learning

Similarity Matrix  $\rightarrow$  can be used with an extension LR  $\rightarrow$  kernel LR

### Best and Worst cases

\* almost linearly separable.

\* low latency requirement.

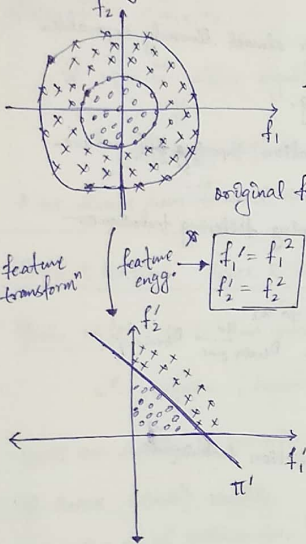
\* very fast to train

large  $d \rightarrow$  chances of linearly separable high.

$\downarrow$   
 for low latency  $\rightarrow$   $L_1$  reg used.

\* LR works decent even upto 1000 dimensions.

## Non-linearly separable data



Can we separate the classes using LR?  
→ can't separate the two classes by drawing a plane.

original features  $\langle f_1, f_2 \rangle$

In this transformed plane, we can find a hyperplane  $\Pi'$  that separates the two classes.

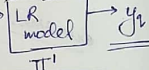
Similarly, when we get a query point  $x_q$

$$x_q = \langle x_{q1}, x_{q2} \rangle$$

↓ FT or FE

$$x'_q = \langle x'_{q1}, x'_{q2} \rangle$$

$$\langle x'^2_{q1}, x'^2_{q2} \rangle$$



But how to come up with feature transforms?

→ comes with practice.

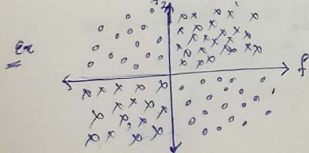
Given data → in form of concentric circles

Eq<sup>n</sup> of circle →  $f_1^2 + f_2^2 = r$

Eq<sup>n</sup> not linear in  $f_1$  &  $f_2$ .  
but linear in  $f_1^2$  &  $f_2^2$ .

Most important aspect of Applied ML.

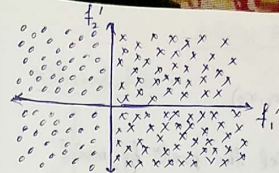
- ① feature engg. (creativity reqd., Art part of ML).
- ② Bias-variance tradeoff.
- ③ Data Analysis.



how to transform to use LR?

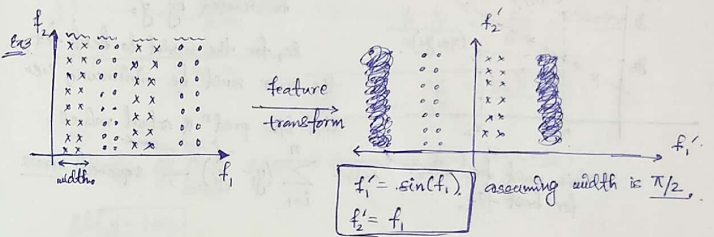
→ the dataset is called XOR dataset.

XOR	0	1
0	0	1
1	1	0



$$\begin{matrix} f'_1 = f_1 \cdot f_2 \\ f'_2 = f_2 \end{matrix}$$

Now linearly separable.



Typical transforms for real value features.

- ①  $f_1 \cdot f_2, f_1^2, f_1^3, f_1^2 f_2$ , etc. — polynomial features.
- ②  $\sin(f_1), \cos(f_1), \sin(f_1) \cdot \cos(f_2)$  — trigonometric features.
- ③ boolean features → OR, AND, XOR.
- ④  $\log f_i, \epsilon^{f_i}$ .

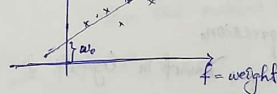
## Linear Regression

↳ does actual regression.

ex → predict height ∈ R given weight, gender, ethnicity, hair color, etc.

consider only one feature → weight for simplicity.

linear regression → find a line that best fits the given data.



$$\text{height} = w_1 \times \text{weight} + w_0$$

$$y = mx + c$$

In general, we have to predict a hyperplane that best fits the given data.