

# Logistic Regression

1-

## I. Probability LR

Logistic Regression is an approach to learning function of the form  $f: X \rightarrow Y$ , or  $P(Y|X)$  in the case where  $Y$  is discrete valued and  $X = \langle x_1, \dots, x_n \rangle$  is any vector containing discrete or continuous variables. LR = Gaussian Naive Bayes + Bernoulli

Parametric model assumed by Logistic Regression in case  $Y$  is boolean is :

$$P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)} \quad \text{--- (1)}$$

and

$$P(Y=0|X) = \frac{\exp(w_0 + \sum_{i=1}^n w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)} \quad \text{--- (2)}$$

Sum of probabilities (1) & (2) must be 1.

To classify any given  $X$  we generally want to assign the value  $y_k$  that maximizes  $P(Y=y_k|X)$ .

Example: large label  $Y=0$  is assigned  $Y=0$  if following condition holds

$$1 < \frac{P(Y=0|X)}{P(Y=1|X)} \quad \text{--- (3)}$$

$$P(Y=1|X) < P(Y=0|X)$$

Substituting (1) & (2) in (3)

$$1 < \exp(w_0 + \sum_{i=1}^n w_i x_i) \quad \text{--- (4)}$$

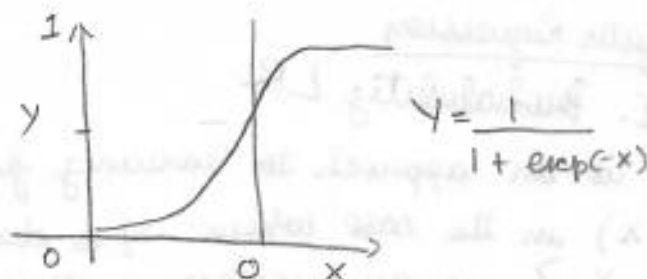


Figure: logistic Function

$$\ln(1) < \ln(\exp(w_0 + \sum_{i=1}^n w_i x_i))$$

$$0 < w_0 + \sum_{i=1}^n w_i x_i$$

### Derivation

We now derive the parametric form of  $P(Y|X)$

$$P(Y=1|X) = \frac{P(Y=1) P(X|Y=1)}{P(Y=1) P(X|Y=1) + P(Y=0) P(X|Y=0)}$$

Dividing both the numerator and denominator by the numerator yields:

$$P(Y=1|X) = \frac{1}{1 + \frac{P(Y=0) P(X|Y=0)}{P(Y=1) P(X|Y=1)}}$$

or equivalently

$$P(Y=1|X) = \frac{1}{1 + \exp\left(\ln \frac{P(Y=0) P(X|Y=0)}{P(Y=1) P(X|Y=1)}\right)}$$

because of our conditional independence assumption we can write this

$$\begin{aligned} P(Y=1|X) &= \frac{1}{1 + \exp\left(\ln \frac{P(Y=0)}{P(Y=1)} + \sum_i \ln \frac{P(x_i|Y=0)}{P(x_i|Y=1)}\right)} \\ &= \frac{1}{1 + \exp\left(\ln \frac{1-\pi}{\pi} + \sum_i \ln \frac{P(x_i|Y=0)}{P(x_i|Y=1)}\right)} \quad (5) \end{aligned}$$

Note the final step expresses  $P(Y=0)$  and  $P(Y=1)$  in terms of the binomial parameter  $\pi$ .

Now consider just the summation in denominator of equation (5). Given our assumption that  $P(X_i|Y=y_k)$  is gaussian, we can expand this term as follows:

$$\begin{aligned}
 \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)} &= \sum_i \ln \frac{\frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(X_i - \mu_{i0})^2}{2\sigma_i^2}\right)}{\frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(X_i - \mu_{i1})^2}{2\sigma_i^2}\right)} \\
 &= \sum_i \ln \exp\left(\frac{(X_i - \mu_{i1})^2 - (X_i - \mu_{i0})^2}{2\sigma_i^2}\right) \\
 &= \sum_i \left(\frac{(X_i - \mu_{i1})^2 - (X_i - \mu_{i0})^2}{2\sigma_i^2}\right) \\
 &= \sum_i \left(\frac{(X_i^2 - 2X_i\mu_{i1} + \mu_{i1}^2) - (X_i^2 - 2X_i\mu_{i0} + \mu_{i0}^2)}{2\sigma_i^2}\right) \\
 &= \sum_i \left(\frac{2X_i(\mu_{i0} - \mu_{i1}) + \mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}\right) \\
 &= \sum_i \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}\right) \quad \text{---(6)}
 \end{aligned}$$

Subs. (6) in (5)

$$P(Y=1|X) = \frac{1}{1 + \exp\left(\ln \frac{1-\pi}{\pi} + \sum_i \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}\right)\right)}$$

Or equivalently,

$$P(Y=1|X) = \frac{1}{1 + \exp\left(w_0 + \sum_{i=1}^n w_i X_i\right)}$$

where

$$w_i = \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2}$$

and where

$$w_0 = \ln \frac{1-\pi}{\pi} + \sum_i \frac{H_{i1}^2 - H_{i0}^2}{2\sigma_i^2}$$

$$P(y=0|x) = 1 - P(y=1|x) = \frac{\exp(w_0 + \sum_{i=1}^n w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$$

One reasonable approach to training logistic Regression is to choose parameter values that maximizes the conditional data likelihood.

$$w \leftarrow \arg \max_w \prod_e P(y^e | x^e, w)$$

$y^e$  denotes the observed value of  $y$  in the  $e^{th}$  training example

$x^e$  denotes the observed value of  $x$  in the  $e^{th}$  training example

This conditional data likelihood  $\rightarrow$  log likelihood.

$$\ell(w) = \sum_e y^e \ln P(y^e=1|x^e, w) + (1-y^e) \ln P(y^e=0|x^e, w)$$

Note that  $y$  can take values 0 or 1.

$$= \sum_e y^e \ln \frac{P(y^e=1|x^e, w)}{P(y^e=0|x^e, w)} + \ln P(y^e=0|x^e, w)$$

$$= \sum_e y^e \ln (w_0 + \sum_i w_i x_i^e) - \ln (1 + \exp(w_0 + \sum_i w_i x_i^e))$$

Gradient ascent,

$$\frac{\partial \ell(w)}{\partial w_i} = \sum_e x_i^e (y^e - \hat{p}(y^e=1|x^e, w))$$

Optimize weights  $w$ .

$$w_i \leftarrow w_i + \eta \sum_e x_i^e (y^e - \hat{p}(y^e=1|x^e, w))$$

$\downarrow$   
step size / learning rate

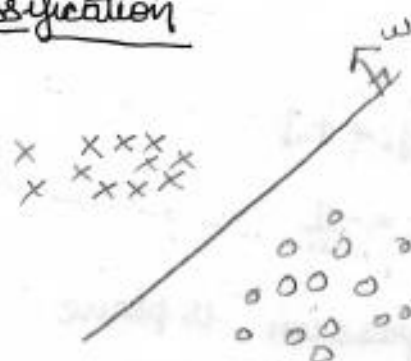
## Regularization in Logistic Regression

Overfitting the training data is a problem that can arise in logistic regression especially when data is very high dimensional and training data is sparse. One approach to reduce overfitting is regularization

$$w \leftarrow \arg \max_w \sum_e \ln P(y^e | x^e, w) - \frac{\lambda}{2} \|w\|^2$$

## ☆ II Logistic regression (Geometric Intuition)

Used for classification



x → +ve class pt  
o → -ve class pt

if my data is (linearly separable)  
almost linearly separable

2D: line } linear  
nD: hyperplane } surface

2D: line  $y = mx + c$

plane  $w^T x + b = 0$

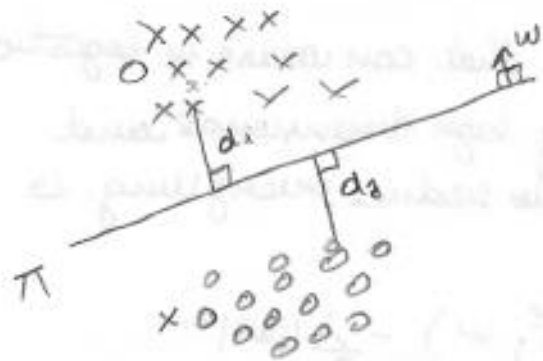
Assumption of LR: classes are almost / perfectly linearly separable



given:  $x_n = \{+ve, -ve\}$

Task: find:  $(w \& b)$

$\Pi$  that best separates +ve pts from -ve pts



$y_i = +1$  : +ve pt  
 $-1$  : -ve pts

$$d_i = \frac{w^T x_i}{\|w\|} ; w \text{ is the normal to the plane.}$$

$\|w\| = 1 \Rightarrow \text{unit vector}$

$$d_i = w^T x_i > 0$$

$$d_j = w^T x_j < 0$$

classifier.

if  $w^T x_i > 0$  then  $y_i = +1$

if  $w^T x_i < 0$  then  $y_i = -1$

decision surface logistic Regression is plane

Case 1:  $y_i * w^T x_i > 0$

$\rightarrow y_i = +1$

⊙ is correctly classified pt

Case 2:  $y_i = -1$  : -ve pt

$$y_i * w^T x_i > 0$$

$$y_i = -1 \rightarrow -1$$

is correctly classifying pt

Case 3:  $y_i = +1$  (+ve pt)

$w^T x_i < 0 \Rightarrow$  LR is saying  $x_i$  is -ve class

$$\boxed{y_i w^T x_i < 0} \Rightarrow y_i = +1$$

$$LR = -1$$

misclassified pt.

Case 4:

$y_i = -1$  (-ve pt)

$$w^T x_i > 0$$

$$y_i w^T x_i < 0 \Rightarrow y_i = -1$$

misclassified pt.

Objective  $\rightarrow$  minimum number of misclassification  
maximum # correctly classified

$$w^* = \underset{w}{\operatorname{argmax}} \left( \sum_{i=1}^n y_i w^T x_i \right)$$

Sigmoid function



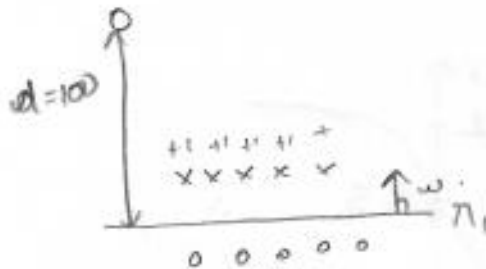
$$w^* = \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \overbrace{y_i w^T x_i}^{\text{signed distance}}$$

$w^T x_i$ : dist. from  $x_i$  to  $\pi$

signed distance

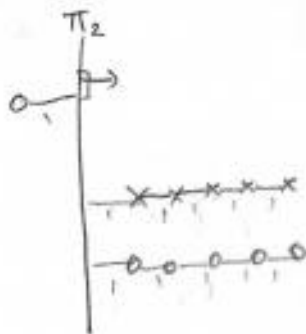
$y_i w^T x_i \rightarrow +ve$  correctly classified  
 $y_i w^T x_i \rightarrow -ve$  incorrectly classified

Case 1:



$$\underset{w}{\operatorname{argmax}} \sum_{i=1}^n y_i w^T x_i = 1+1+1+1+1 + 1+1+1+1+1 - 100 = -90$$

Case 2



$$\underset{w}{\operatorname{argmax}} \sum_{i=1}^n y_i w^T x_i = 1+2+3+4+5 - 1-2-3-4-5 + 1 \rightarrow \text{total} = 1$$

Between Case 1 and Case 2. Case 2 will be chosen as our objective is to find  $\underset{w}{\operatorname{argmax}} \sum_{i=1}^n y_i w^T x_i$  classifier. But it is terrible classifier.  $\pi_2$  is better



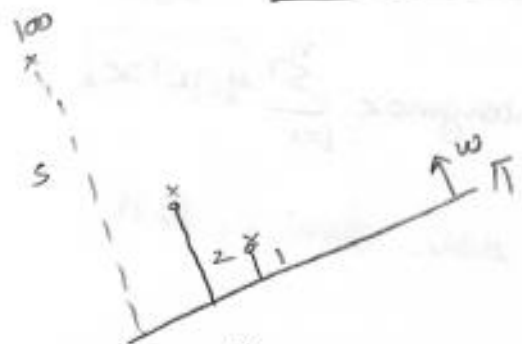
Because of single outlier point hyperplane changed.

So we introduce sigmoid function

idea - instead of using signed distance

if signed distance is small :- use it.

is large : make it small



$$\omega \sum_{i=1}^n y_i \omega^T x_i$$

$$\arg \max_w \sum_{i=1}^n \underbrace{f(y_i w^T x_i)}_{\text{signed dist}}$$

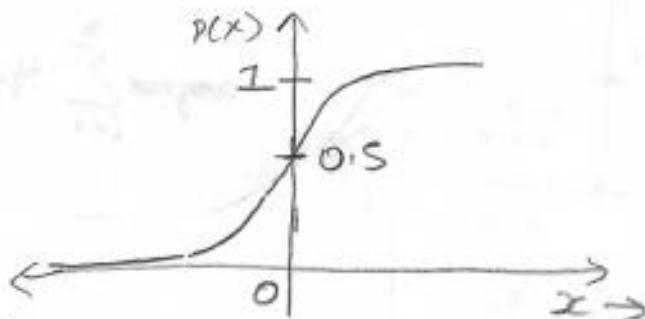
Sigmoid function  $\sigma(x)$

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

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

$$\mu_{\text{max}} : 1$$

min : 0



$$\text{weight} \times \sum_{i=1}^n \sigma(y_i w^T x_i)$$

Reason to choose sigma

1) Nice probabilistic interpretation

point lie on hyperplane  $P(y_i=1)=0.5$



max. sum of signed dist  $\rightarrow$  outliers

$\downarrow$   
 $\sigma(x) \rightarrow$  sigmoid

$$w^* = \arg \max_w \sum_{i=1}^n \sigma(y_i w^T x_i)$$

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

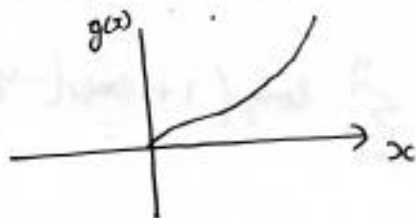
$\rightarrow$  less impacted by outliers

### Mathematical formulation of Objective function

$\rightarrow$  monotonic functions  $\div g(x)$

$x \uparrow \quad g(x) \uparrow \rightarrow$  monotonically increasing fn.

if  $x_1 > x_2$  then  $g(x_1) > g(x_2)$



$\log(x) \rightarrow$  monotonic function

Theorem 1:  $\rightarrow$  if  $g(x)$  is a monotonic fn.

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

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

Using Theorem 1

$$w^* = \operatorname{argmax} \sum_{i=1}^n \frac{1}{1 + \exp(-y_i w^T x_i)}$$

$g(x) : \log(x) : \text{monotonic fcn.}$

$$w^* = \operatorname{argmax} \sum_{i=1}^n \log(\sigma(y_i w^T x_i))$$

$$w^* = \operatorname{argmax} \sum_{i=1}^n \log \left[ \frac{1}{1 + \exp(-y_i w^T x_i)} \right]$$

Using  $\Rightarrow \log \frac{1}{x} = -\log(x)$

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

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

$\Leftarrow \operatorname{argmax} f(y) = \operatorname{argmin} \bar{f}(y)$

◇ Weight-Vector

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

↓  
weight vector

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

↓  
 $x_1, x_2, x_3, \dots, x_d \rightarrow \text{features}$

Decision  $x_q \rightarrow y_q$

if  $w^T x_q > 0$  then  $y_q = +1$   
if  $w^T x_q < 0$  then  $y_q = -1$

### Interpretation of $w$ :

$$\begin{aligned} \textcircled{1} \text{ If } w_i = +ve, \quad x_{qi} \uparrow &\rightarrow (w_i x_{qi}) \uparrow \\ &\Rightarrow \sigma(w^T x_q) \uparrow \\ &\Rightarrow P(y_q = +1) \uparrow \end{aligned}$$

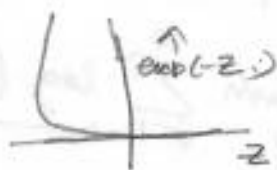
$$\begin{aligned} \textcircled{2} \text{ If } w_i = -ve, \quad x_{qi} \uparrow &\rightarrow (w_i x_{qi}) \downarrow \\ &\Rightarrow \left( \sum_{i=1}^n w_i x_{qi} \right) \downarrow \\ &\Rightarrow \sigma(w^T x_q) \downarrow \\ &\Rightarrow P(y_q = +1) \downarrow \\ &\Rightarrow P(y_q = -1) \uparrow \end{aligned}$$

### ☆ L2 Regularization: Overfitting and Underfitting

$$w^* = \underset{w}{\operatorname{argmin}} \sum \log(1 + \exp(-y_i w^T x_i))$$

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

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



is always +ve

$$\boxed{\exp(-z_i) \geq 0}$$

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

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-z_i)) \geq 0$$

it occurs when  $z_i \rightarrow \infty \quad \forall i$

$$z_i \rightarrow \infty$$

$\hookrightarrow$  (w) modify my  $w$  in such a way that  $z_i \rightarrow \infty$

①  $z_i = +\infty$ ;  $x_i$  is correctly classified by  $w$ .

$$z_i = +\infty;$$

if I pick  $w$  st

a) all training pts are correctly classified

b)  $z_i \rightarrow \infty$

then we get best  $w$

↓  
overfitting

(perfect job on training data  
which does not tell perfect job  
on test data)

$$\textcircled{2} \textcircled{w_i} \rightarrow +\infty$$

$$\text{or } \textcircled{w_i} \rightarrow -\infty$$

one key aspect

$w_i$  is a normal to Hyperplane

$$w^T w = 1$$

To get rid of this problem. regularization is used

$$w^* = \underset{w}{\operatorname{argmin}} \underbrace{\sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))}_{\text{loss term}} + \underbrace{\lambda w^T w}_{\text{L1: Regularization term}}$$

$$w_i \rightarrow \infty$$

$$w_i \rightarrow -\infty$$

} → regularizer will not allow it

case 1. →

$$w^* = \underset{w}{\operatorname{argmin}} \underbrace{\sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))}_0 + \underbrace{\lambda w^T w}_{w \text{ large}}$$

$\lambda \rightarrow$  hyper parameter

$\rightarrow -$

$\lambda = 0 \Rightarrow$  overfit or high variance

$\lambda = \infty \Rightarrow$  underfitting as influence of regularizer increase.  
or high bias

### L1 Regularization and sparsity

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \underbrace{\lambda \|w\|_2^2}_{L2 \text{ regularizer}}$$

$$L1: w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \underbrace{\lambda \|w\|_1}_{L1 \text{ reg}}$$

$\rightarrow$  hyperparameter  
will avoid  $w \rightarrow +\infty$   
 $- \infty$

~~disadvantages~~  $\rightarrow$

Sparse  $\rightarrow$   $w = \langle w_1, w_2, \dots, w_d \rangle$   
Solution to LR is said to be sparse if many  $w_i$ 's are zero.

If we use  $L1$  reg in LR, all the less important features becomes zero.

$f_1, f_2, f_3, \dots, f_i, \dots, f_d$   
 $\rightarrow$  less impor

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

we use  $L1$  regularization  $\rightarrow$  creates sparsity.

# Comparing Geometric and Probability L.R

geom:  $w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \text{reg}$

$\swarrow \quad \searrow$   
 $+1 \quad -1$

probability:  $w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n -y_i \log p_i - (1-y_i) \log(1-p_i)$

$\swarrow \quad \searrow$   
 $+1 \quad 0$

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

Case 1  $y_i: +ve$

geom  $y_i = +1$

prob:  $y_i = +1$

geom:  $w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-w^T x_i)) \quad \text{--- (1)}$

prob:  $w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n -1 \log\left(\frac{1}{1 + e^{-w^T x_i}}\right) =$

$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + e^{-w^T x_i}) \quad \text{--- (2)}$

(1) & (2) are same

Case 2:  $y_i: -ve$

geom:  $y_i = -1$

prob:  $y_i = 0$

geom  $w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(w^T x_i)) \quad \text{--- (3)}$

prob  $w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n -(\log 1 - \frac{1}{1 + e^{w^T x_i}})$

$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n -\log\left(\frac{1 + e^{w^T x_i}}{1 + e^{w^T x_i}} - 1\right) \quad \text{--- (4)}$

$$w^* = \arg\min \sum -\log \frac{e^{w^T x_i}}{1 + e^{w^T x_i}}$$

$$= \arg\min \sum + \log \frac{1 + e^{-w^T x_i}}{e^{w^T x_i}}$$

$$\log(-x) = \frac{1}{\log x}$$

$$= \arg\min \sum \log \frac{(1 + e^{-w^T x_i})}{e^{-w^T x_i}}$$

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

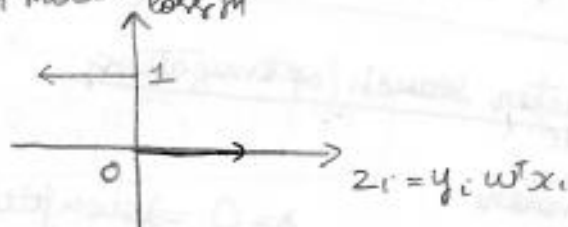
$$= \sum \log (1 + \exp^{w^T x_i})$$

(3) and (4) are same

### III. Loss minimization interpretation of LR.

Ideal optimization model

0-1 loss function



$z_i \rightarrow +ve$  loss fn should be 0 (for correct classification)  
 $z_i \rightarrow -ve$  (misclassification) loss fn should be 1

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

Solve optimization problem in ML

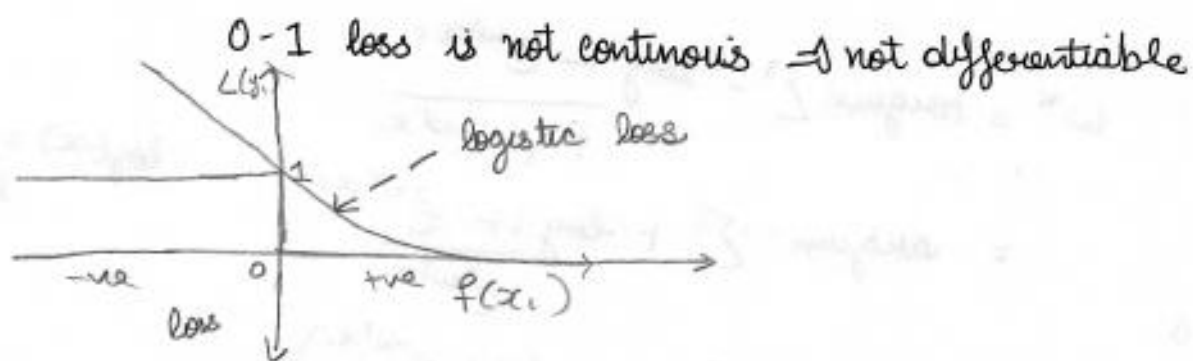
↳ differentiation in Calculus

Column / feature standardization

$x_1, x_2, x_3, \dots, x_d$

$x_i \in \mathbb{R}^d$





☆ logistic regression is approximation of 0-1 loss.

- 1) In 0-1 loss  $z_i > 0$  value is 0 but in logistic loss it is not zero but tends to become 0.
- 2) On the other hand negative side 0-1 gives value of 1 but logistic loss gives value greater than 1

### Loss-minimization interpretation

1) loss for logistic loss  $\rightarrow$  LR

2) hinge loss  $\rightarrow$  SVM.

3) SQ loss  $\rightarrow$  Linear loss

### Hyperparameter search/optimization

$\lambda$ : hyperparameter

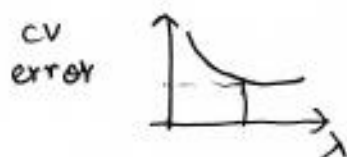
$\lambda = 0 \Rightarrow$  overfitting

$\lambda = \infty \Rightarrow$  underfitting

Q How to determine the best  $\lambda$ ?

( $\lambda$ ) in LR is a real number.

1) Grid Search (Brute force)



①  $\lambda = [0.001, 0.01, 0.1, 1, 10, 100, 1000]$

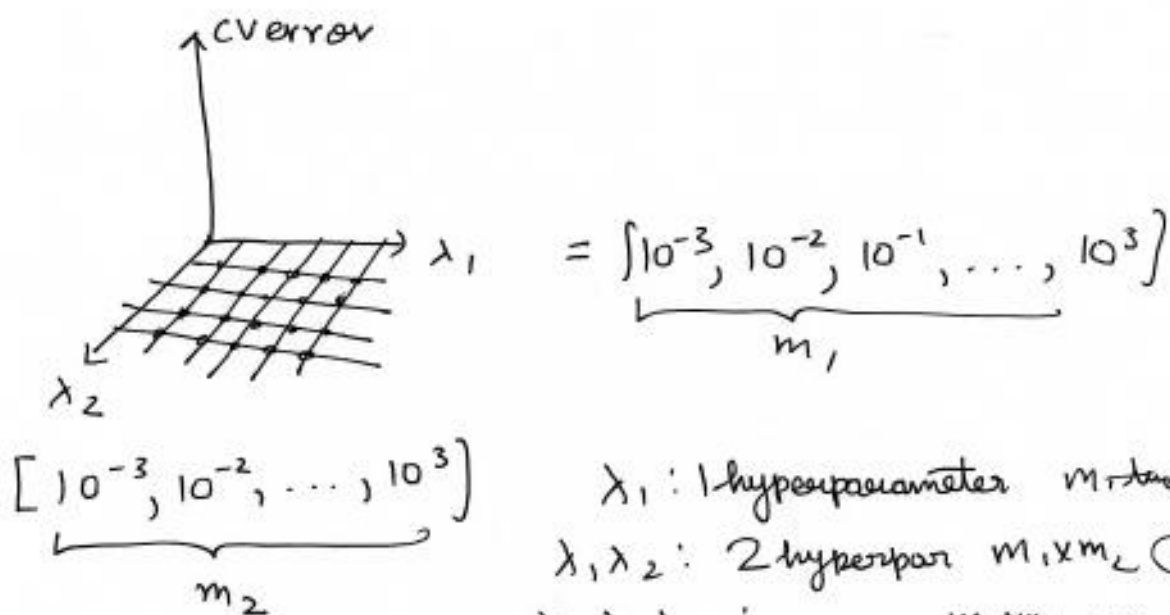
②  $\lambda = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]$

elastic net :-

-9-

$$\lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^2$$

$\downarrow$   $\downarrow$   
 $L_1$  regularization  $L_2$  regularization

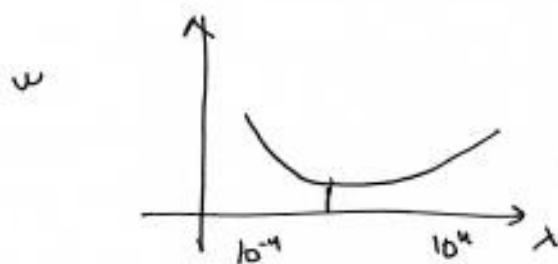


$\lambda_1$  : 1 hyperparameter  $m_1$  times  $(m_1)$   
 $\lambda_1, \lambda_2$  : 2 hyperpar  $m_1 \times m_2$   $(m_1^2)$   
 $\lambda_1, \lambda_2, \lambda_3$  :  $m_1 \times m_2 \times m_3$   $(m_1^3)$

Grid search not good when we have many hyper parameter

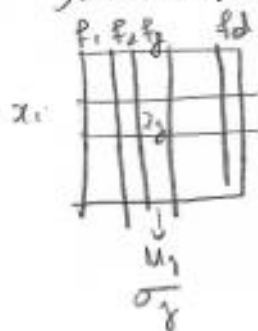
2) Random Search :-

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



$\lambda \rightarrow$  hyperparameter search/optm  
 $\left\{ \begin{array}{l} \rightarrow \text{Grid search} \\ \rightarrow \text{Random search} \end{array} \right.$

Column / feature Standardization



$$x_i \in \mathbb{R}^d$$

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