# Logistic Regression

## 24.1 Geometric intuition of Logistic Regression.

→ classification Technique

→ Simple & Elegant model

    NB : Probabilistic model/Tech
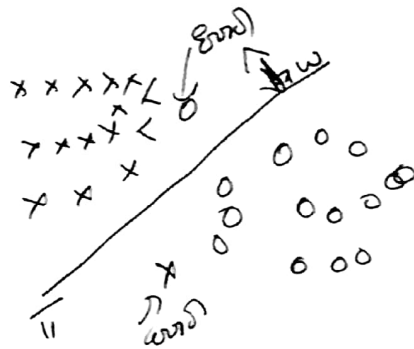
    LR : geometric intuition.

Logistic Regression can be interpreted using below Technique
    ↳ Geometry
    ↳ Probability
    ↳ Loss function.



$x \to$ +ve

$o \to$ -ve

2D : line   } linear
nD : hyper plane } surface.

If my data is linear separable

If the $\pi$ passes through Origin = b=0

$$w^T x + b = 0 \quad \Rightarrow \quad w^T x = 0$$

Assumption of Log Reg is class are almost/ perfectly linearly separable

$$\pi = w^T x + b$$

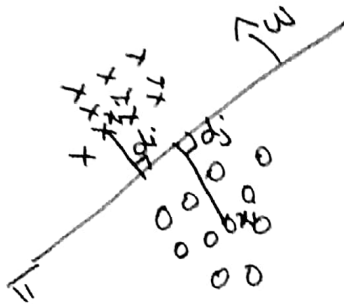$$y_n = \{+ve, -ve\} \to \text{given to us}$$

Task to
Find ÷ w & b

such that the line separates both the +ve & -ve bts

## Assumptions

NB: conditional independence of features

K-nn: Neighbourhood.



$y_i = +1$ : +ve pts

$-1$ : -ve pts

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

$d_i$ = distance of point from plane.

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

$$\|\omega\| = 1 \Rightarrow \text{unit vector.}$$

$d_j = \omega^T x_j$

Since $\omega$ & $x_i$ are on the same side $d_i = \omega^T x_i > 0$

Since $\omega$ & $x_j$ are not on the same side $d_j = \omega^T x_j < 0$

### classifier says is

If $\omega^T x_i > 0$ then $y_i = +1$

$\omega^T x_j < 0$ then $y_i = -1$ } line passes through origin.

$\Rightarrow$ Decision surface in LR is a plane

$\rightarrow$ Classifier to be V.good

↳ min # misclassified

𝔰

↳ max # correctly classified pts

⟵

as many pts as possible to have $y_i * \omega^T x_i > 0$

$\max \sum_{i=1}^{n} y_i \omega^T x_i$    optimal $\omega^*$ means best hyperplane.

$\omega^*_{optimal} = \underset{(\omega)}{argmax} \left. \sum_{i=1}^{n} y_i \omega^T x_i \right\} $ optimization problem

## 2.2 Sigmoid function & Squashing:

$$argmax \sum_{i=1}^{n} \underbrace{y_i \omega^T x_i}_{\hookrightarrow \text{Signed distance.}}$$

$\omega^T x_i$ distance from $x_i$ to $\pi$   ($\omega$ is a unit vector)

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

$\quad \hookrightarrow$ : -ve $\Rightarrow$ incorrectly classifies $x_i$



5 ✓ (correct)
5 ✗ (correct)
1 ✗ (misclassified)

Case 1: $\pi_1$ is my separator, $\sum_i y_i \omega^T x_i = 1+1+1+1+1+1+1+1+1+ 1 - 100$

$$\boxed{\cdots -90}$$

Case 2:



$\sum_{i=1}^{n} y_i' \omega_2^T x_i' = 1+2+3+4+5 -1 - 2 -3 -4 -5$
$\qquad\qquad\qquad\qquad\qquad +1 \text{ (outlier)}$

$\qquad\qquad\qquad = +1$

5 ✓
5 ✗
1 ✓

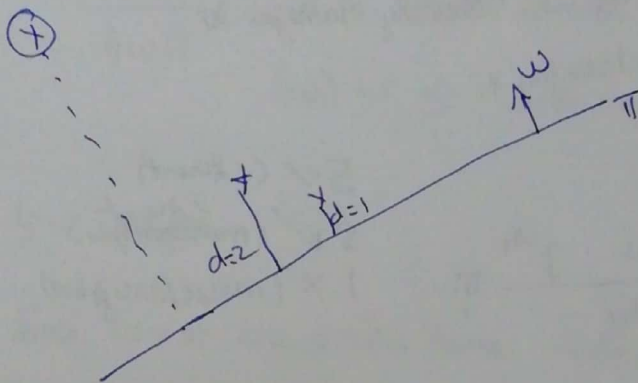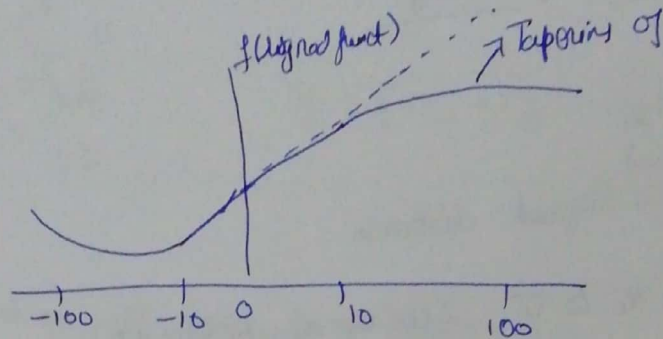One single extreme/outlier pt is changing my model (hyperplane) in Case 1 which is very bad.

max. sum of signed distances not outlier prone.

## Squashing

idea: Instead of loving signed distance.

If signed distance is small → use it as is

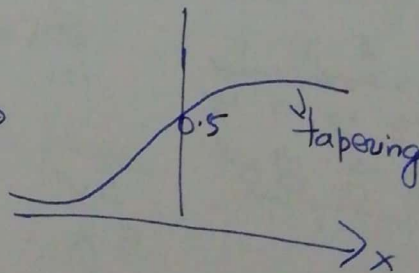" " large → make it smaller as possible.



$$\underset{w}{argmax} \; \sum_{i=1}^{n} f(y_i w^T x_i)$$

$$\underbrace{\qquad}_{\text{Signed distance.}}$$

## Sigmoid function

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



max: 1
min: 0
$\sigma(0) = 0.5$

we will change our problem to

$$\underset{w}{argmax} \; \sum_{i=1}^{n} \sigma(y_i w^T x_i)$$

(f) →word is very large → $P(y=1) = 0.9999$ ⎫
⎬ Problistic Interpretation.
→ $P(y=1) = 0.5$ ⎭

$* \to \omega^T x_i = 0$

Max. Sum of Signed dist → outlier problem.
↓

$\sigma(x) \to$ Sigmoid
↳ tapperou linear
↳ problistic model.

max. Sum of transformed signed interpretation.

$$\omega^* = \underset{(\omega)}{argmax} \sum_{i=1}^{n} \sigma(y_i \omega^T x_i)$$

↓

$$\omega^* = \underset{(\omega)}{arg\ max} \sum_{i=1}^{n} \frac{1}{1+exp(-y_i \omega^T x_i)} \quad \leftarrow less\ impacted\ by\ outlier.$$

distax : $(-\infty, \infty)$

⎰ (squashing
⎱ using σ function
$d\ (0\ to\ 1)$

why sigmoid function?

→ Easy to differentiate
→ problistic interpretation.

## 24.3  Mathematical formulation of objective function

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

→ monotonic functions : $g(x)$
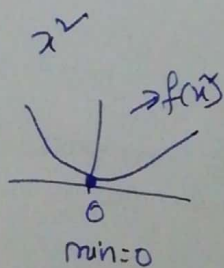
$x\uparrow$ ; $g(x)\uparrow$  monotonically inGrease fn.

If $x_1 > x_2$ then $g(x_1) > g(x_2)$ then it is called monotonial functn.

eg $\log(x) > 0$ ; should be $> 0$

### Optimation problem :-

$x = \underset{x}{argmin}(x^2) = 0$

best$(x)$

$x^2$



min : 0

$x^2$ is mono inGreasi⁀ when $x > 0$

$x^2$ is " deGrea when $x < 0$

---

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

$x^* = argmin(f(x))$ ; $f(x) = x^2$

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

$x' = argmin \log(x^2)$

### claim

$x^* = x'$

---

If $g(x)$ is a monotonic function

$\underset{x}{argmin} f(x) = \underset{x}{argmin} g(f(x))$

$\underset{x}{argmax} f(x) = argmax \, g(f(x))$

$x\uparrow$  $g(x)\uparrow$

$x\uparrow$  $g(x)\downarrow$

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

$g(x) = \log(x)$ : monotonic fn.

$$\omega^* = \underset{\omega}{\arg\max} \sum_{i=1}^{n} \log\left(\frac{1}{1+\exp(-y_i \omega^T x_i)}\right)$$

$$\log(1/x) = -\log(x)$$

$$\omega^* = \underset{\omega}{\arg\max} \sum_{i=1}^{n} -\log(1+\exp(-y_i \omega^T x_i)) \quad \Big\} \ \text{geometry}$$

$$\boxed{\arg\max f(x) = \arg\min \overline{f(x)}}$$

$$\omega^* = \underset{\omega}{\arg\min} \sum_{i=1}^{n} \log(1+\exp(\overset{\nearrow \ \xi -1 \ \partial + 1\xi}{-y_i \omega^T x_i}))$$
$$\underset{\smile}{}$$
Singed dist

$$\boxed{\log(e^x) = x}$$

$$\underset{(\omega)}{\arg\min} \sum_{i=1}^{n} \log(1+\exp(-y_i \omega^T x_i))$$

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

Probability method
$$\omega^* = \underset{(\omega)}{\arg\min} \sum_{i=1}^{n} -y_i \log P_i - (1-y_i)\log(1-P_i)$$

$$\boxed{P_i = \sigma(\omega^T x_i)}$$

Scanned by CamScanner

## 24.4 Weight Vector

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

↓

weight vector $(w) = <w_1, w_2, w_3, w_4 \cdots w_d>$

$\in R^d$ → d features of weight vector $w$

$$w = \quad <w_1, w_2, w_3 \cdots w_d>$$
$$<f_1, f_2, f_3 \cdots f_d>$$

### decision  If $x_q \to y_q$

If $w^T x_q > 0$  Then $y_q = +1$

$\quad w^T x_q < 0$  Then $y_q = -1$

### Problistic functa

$$\sigma(w^T x_q) = P(y_q = +1)$$
$$(0 \text{ to } 1)$$

### Interpretation of $w$ ÷

If $w_i = +ve$, $x_{q,i} \uparrow \Rightarrow (w_i x_{q,i}) \uparrow$
↑
$(f_i)$

$$\sum_{i=1}^{d} (w_i x_{q,i}) \uparrow$$

$$\sigma(w^T x_q) \uparrow$$

$$P(y_q = +1) \uparrow$$

If $w_i = -ve$, $x_{q,i} \uparrow \Rightarrow (w_i x_{q,i}) \downarrow$

$$= \sum_{i=1}^{d} w_i x_{q,i} \downarrow$$

$$= \sigma(w^T x_q) \downarrow$$
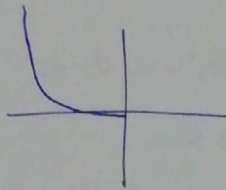
$$= P(y_q = +1) \downarrow \quad P(y_q = -1) \uparrow$$

24.5   L2 Regularization: Overfitting & underfitting

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

let $z_i = y_i w^T x_i \rightarrow$ Signed distance.

$\sim \underset{w}{\arg\min} \sum_{i=1}^{n} \log(1 + \exp(-z_i))$

Plot $(\exp(-z))$ is always $> 0$



$\sum_{i=1}^{n} \log(1 + \exp(-z_i)) \geq 0$     $\Rightarrow \log(1) = 0$

$\log(2) > \log(1)$

$\log(1+\delta) > \log(1)$

$\delta \geq 0$

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

minimal value of $\sum_{i=1}^{n} \log(1 + \exp(-z))$ is zero

If $z_i = +ve$ , $z_i \rightarrow +\infty$

Then $\exp(-z_i) \rightarrow 0$

$\log(1 + \exp(-z_i)) = 0$     Since $\log(1) = 0$

-If I pick my $w$ such that

(a) all training points are correctly classified
(b) $z_i \rightarrow \infty$
                                    ↓
                              overfit

Then to that is called best $w$.

If we make $w_i \rightarrow \infty$ or $-\infty$ we will reach minima $= 0$

Regularization (To avoid $\omega^T$ to become $\to \infty$ or $+\infty$)

$$\omega^+ = \underset{(\omega)}{\arg\min} \sum_{i=1}^{n} \log(1 + \exp(-y_i \omega^T x_i)) + \lambda \omega^T \omega = \lambda \sum_{j=1}^{d} \omega_j^2$$

loss term → ← Regularization term

$\lambda \|\omega\|_2^2$

$\boxed{\omega^T \omega = 1}$

$\lambda \omega^T \omega$
$= \lambda \sum_{j=1}^{d} \omega_j^2$

When $\lambda = 0$, overfit → high variance, by making $z_i \to \infty$ & $-\infty$

$\lambda =$ v. large; underfit → high bias, we are ignoring loss term

we have to find the right $\lambda$ using CV

## 24.6 L1 regularization and Sparsity

$z_i \to +\infty$

$$\omega^+ = \underset{(\omega)}{\arg\min} \underbrace{\sum_{i=1}^{n} \log(\underbrace{1 + \exp(-y_i \omega^T x_i)}_{z_i})}_{\text{Logistic class}} + \underbrace{\lambda \|\omega\|_2^2}_{L_2\text{-reg}}$$

Alternative to L2 reg is L1

$\|\omega\|_2^2$ for reg
$\downarrow$
$\|\omega\|_1$ for reg: $\|\omega\|_1 = \sum_{i=1}^{d} |\omega_i|$

$$\omega^+ = \underset{(\omega)}{\arg\min} \left(\begin{array}{c}\text{Logistic loss for} \\ \text{training data}\end{array}\right) + \underbrace{\lambda \|\omega\|_1}_{L_1\text{-reg}}$$

→ hyper parameter

will avoid $\omega_i \to +\infty$
$\omega_i \to -\infty$

SP

Scanned by CamScanner

## Sparsity:

$$W = \langle w_1, w_2, \ldots w_d \rangle$$

Solution to LR is said to be sparse if many $w_i$'s are zero

If we use $L_1$ reg in LR, all the unimportant (or) less important become zero

$$f_1, f_2 \ldots \overset{\rightarrow \text{Less important}}{\textcircled{f_i}} \ldots f_d$$

$$W = \langle w_1, w_2 \ldots w_i \ldots w_d \rangle$$

$$\downarrow$$

zero if $L_1$ is used.

If $L_2$ reg is used; $w_i$ becomes a small value but not necessarily zero.

(Q) why does $L_1$ reg create sparsity in $w$ as compared to $L_2$ reg

ppl generally use $L_1$ than $L_2$

→ Elastic net:— Either $L_1$ & $L_2$

$$w^+ = \underset{w}{argmin} \; \sum_{i=1}^{n} log(1 + exp(-z_i)) + \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^{\vee}$$

we have to find two hyper parameters $\lambda_1$ & $\lambda_2$

## 24.7 Probabilistic Interpretation: Gaussian Naive Bayes

cxs.cmu.edu/~tom/mlbook/NBayesLogReg.pdf

$$LR \Rightarrow GNB + Bernoulli$$
$$\downarrow \qquad\qquad \downarrow$$
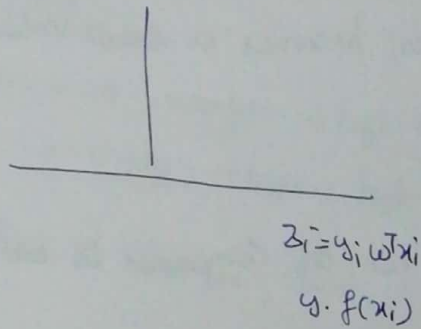$$p(x_i | y_i) \qquad y \sim Bernoulli$$

## 24.8 Loss minimization Interpretation

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

$$z_i = y_i w^T x_i = y_i f(x_i)$$

If we build a ideal optimization model.

$$w^* = \underset{w}{\text{argmin}} \ (\text{num. of incorrectly classified pts})$$



$$z_i = y_i w^T x_i$$
$$y_i \cdot f(x_i)$$

funct<sup>n</sup>

+1: incorrectly classified

0: correctly classified

min: loss

max: profit

24.9    hyperparameter & random Search

     $\lambda$ = hyper parameter
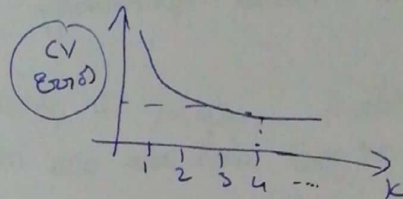
     $\lambda = 0 \Rightarrow$ overfitting

     $\lambda = \infty \Rightarrow$ underfitting

(a) How to find the best $\underline{\lambda}$

         $\lambda = KR$

         $K = KNN$

         $\alpha = NB$ (Laplace Smoothing)



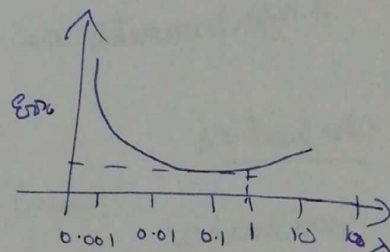K in KNN is an integer, which takes value $\{1, 2, 3 \ldots\ n\}$

$\lambda$ in LR is a real number.

     $\lambda \in \mathbb{R} \quad \begin{cases} \lambda = 0.1234 \\ \lambda = 0.2386 \end{cases}$

$\Rightarrow$ One technique to find $\lambda$ is GRID Search.

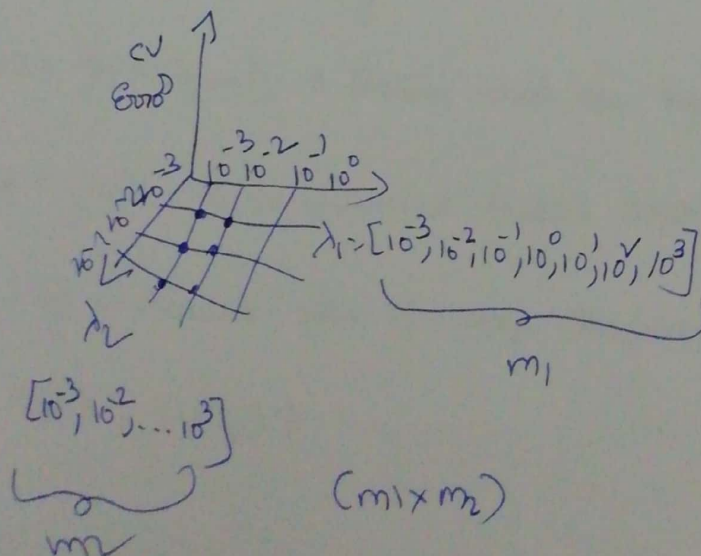Case 1 :-   $\lambda = [0.001, 0.01, 0.1, 1, 10, 100, 1000]$

Case2 :-   $\lambda = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]$

$\rightarrow$ Generally ppl select a large window.

     $\lambda = [10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10, 10^{1}, 10^{2}, 10^{3}, 10^{4}]$

Elasticnet :- $\lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^{\vee}$



     $\lambda_1 = [10^{-3}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{\vee}, 10^{3}]$

                 $m_1$

     $[10^{-3}, 10^{-2}, \ldots 10^{3}]$

        $m_2$        $(m_1 \times m_2)$

## Grid search

$\lambda : 1$ hyper parameter + $(m_1)$

$\lambda_1 \lambda_2 : 2$ " $\div m_1 \times m_2$

$\lambda_1 \lambda_2 \lambda_3 : 3$ " $\div m_1 \times m_2 \times m_3$

$\{$ as # hyper parameters increase, the # times model needs to be trained increases Exponentially

## Grid Search

is not good when there are more hyper parameters

To Overcome this issue we have another technique called

## Random Search

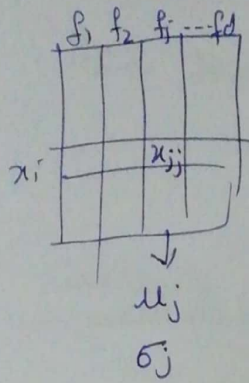$\lambda = [10^{-4}, 10^4] \leftarrow$ randomly pick values in the given interval

→ Random Search is almost as good as grid search Especially when # hyper parameters are large.

## Other functions

GridSearchCV
RandomizedSearchCV

24.10 **Column Standardization (z-score)**



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

$x_{ij}' = \dfrac{x_{ij} - u_j}{\sigma_j}$ : Standardization

Even in Logistic Regression its mandatory to perform feature standardization. before training

mean-Centering
&
Scaling
} Standardization.

24.11 **Feature importance & model interpretability**

$\quad f_1 \quad f_2 \quad f_j \quad fd$

$w \rightarrow \quad w_1 \quad w_2 \quad w_j \quad wd.$

assume if all features are independent (Naive Bayes)
feature importance can be acheived based on the weights

In k-non ÷ feature imp → forward feature Selection.
⤷ we cannot get directly.

NB ÷ $P(x_i | y = +1)$ → features which are important

LR ÷ $w_j's$ → to determine feature importance.

$|w_j|$ = absolute value of weigh corresponding to $f_j$

$|w_j| \uparrow \quad ; \quad (w^T x_j) \uparrow$

<u>Case 1</u> $\quad w_j = +ve \;\&\; large$ ; $\sum\limits_{j=1}^{d} w_j \cdot x_{qj} \Rightarrow w^T x_q$

$\qquad L P(y_q = +1) \uparrow$

<u>Case 2</u> : $w_j : -ve \;\&\; large$ ; $\sum\limits_{j=1}^{d} w_j x_{qj} \Rightarrow P(y_q = -1) \uparrow$

We can determine the important features in LR based on the weights

E.g. Predict the gender : male & female
$\qquad\qquad\qquad\qquad (+1) \quad (-1)$

<u>1 feature</u>: hair_length $= |w_{hl}|$ is large

$\qquad\qquad\qquad \phi_b \; w_{hl} : -ve$

$\qquad\qquad w_{hl}\uparrow \; ; \; P(y_q = -1) \uparrow$

<u>2 featur</u> : height $\uparrow$ ; $P(y_q = +1) \uparrow$
$\qquad\qquad\qquad\qquad\qquad\uparrow$
$\qquad\qquad\qquad\qquad\quad male$

$\quad w_h = +ve.$

<u>model interpretability</u>

$x_q \longrightarrow \boxed{y_q = +1} \rightarrow$ Top features
$\qquad\;\; \searrow \boxed{y_q = -1}$ $\qquad \downarrow$ hair_length, height

24.12 **Collinearity of features**

feature Importance : features are independant

$$|w_j| \text{ as F.I values}$$

**Collinearity (Ov) multicollinearity**

collinearity :- $f_i, f_j$

s.t if $f_i = \alpha f_j + \beta$

Then $f_i \& f_j$ are collinear.

**multicollinearity**

If $f_1, f_2, f_3 \& f_4$ Such that

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

Then $f_1, f_2, f_3 \& f_4$ are said to be multicollinearity

(Q) why does $|w_j|$ not be ureful as f.I if features are collinear?

$$D = \langle x_i, y_i \rangle_{i=1}^{n}$$

$$w^* = \langle 1, 2, 3 \rangle \quad ; \quad x_q = \langle x_{q_1}, x_{q_2}, x_{q_3} \rangle$$
$$\quad\quad f_1, f_2\ f_3$$

$$w^{*T} x_q = x_{q_1} + 2 x_{q_2} + 3 x_{q_3}$$

If $f_2 = 1.5 f_1 \Rightarrow f_1 \& f_2$ are collinear.

$$w^T x_q = x_{q_1} + 3 x_{q_2} + 3 x_{q_3} = 4 x_{q_1} + 3 x_{q_3}$$

$$\langle 4, 0, 3 \rangle$$
$$\uparrow \quad \uparrow \quad \uparrow$$
$$x_{q_1} \ x_{q_2} \ x_{q_3}$$

$$w^* = \langle 1, 2, 3 \rangle \quad {}^{\to f_3\ h\ imp}$$

$$\tilde{w} = \langle 4, 0, 3 \rangle \quad \therefore \text{ assumptions are completely changing}$$
$$\quad {}_{\uparrow f_1\ is\ imp} \quad \text{if features are collinear} \Rightarrow \text{weight vector can can}$$

change arbitrarily $\Rightarrow |w_j|$ can be used for feature importance.

$|w_j|$ as F.I

determine if features are multicollinear

↑ Perturbation technique

↳ means to change the values a little by adding a $\epsilon$



$\rightarrow x_{ij} + \epsilon$

↑ Standardized ↑ Small noise

$$N(0, 0.01)$$

Before Perturbation :- $w = <w_1, w_2, \ldots w_j \ldots w_d>$

after    "    :- $\tilde{w} = <\tilde{w}_1, \tilde{w}_2 \ldots \tilde{w}_j \ldots \tilde{w}_d>$

If $w_i$ & $\tilde{w}_i$ differ significantly than your features are collinear

$|w_j|$'s as F.I cannot be used

## 24.13 Test/Run time space & time Complexity

Train LR :- Solving Logistic Regression problem.

Training time of LR is $O(nd)$

After this we get $w^* = <w_1, w_2, w_3 \ldots w_d>$

$w^{*T} x_q > 0 \rightarrow +ve$

$< 0 \rightarrow -ve.$

At runtime we have to store $w^T$

Space $= O(d)$   → memory Efficient as well

Time $= O(d)$

If d is small than LR is vv good for low latency application

$x_q \rightarrow (1ms) \rightarrow y_q$

If d is large ∵ d ≈ 1000

$w^T x$ ÷ 1000 multi & addition.

→ L1 reg ÷ Sparsity ($w_j$'s corresponding to less important features = 0)

λ ÷ reasonabley

λ↑ ; Sparsity ↑

more of $w_j$ = 0 $\begin{cases} 50 \text{ mult } \& 50 \text{ additions} \\ \downarrow \text{ latency} \end{cases}$

Bias vs Latency

λ↑ ; Bias ↑, latency ↓

## 24.14   Real world cases

Decision surface ÷ Linear / hyperplane.   { +
                                            -

assumption ÷ data is linearly separable & almost linearly separable.

Imbalanced data : upsampling & down sampling.

Outliers : less impact ∵ of σ(x)

→ $D_{train}$ → $w^*$

→ $x_i$ → $w^{*T} x_i$ :: distance from π to point $x_i$ (plane)

→ remove points which are very far away from π

from $D_{train}$ → $D'_{train}$

→ $D'_{train}$ → $\tilde{w}^*$
                    ↳ final solution.

missing ÷ Standard imputation.

multiclass : one Vs Rest ← typically

$\begin{cases} \text{max Ent model} \leftarrow \text{Extension to LR} \\ \text{softmax classifier} \leftarrow \text{deeplearni} \\ \text{multinomial LR.} \end{cases}$

Similarity matrix : Extension to LR → Kernal LR

## Best & worst cases

→ almost lr separable
→ low-latency requirement ($L_1$ reg)
→ very fast to train.

## Large dimensionality

→ d is large, chance data is linearly separable is high
 ↓
low latency → $L_1$ regularization.

24.15 Non-linearly separable data & feature Engineering



(a) Can we use LR to separate the classes

feature
Transform/Engineering.
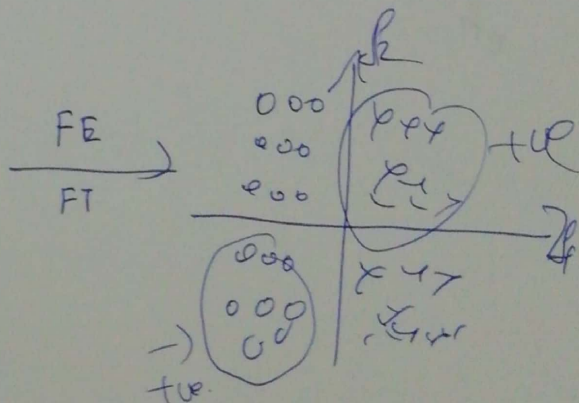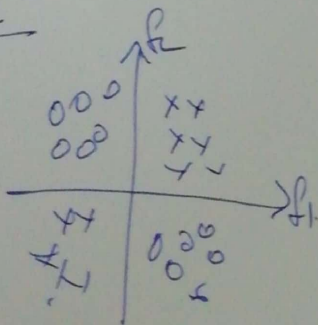(FE)

$\{f_1' = f_1^v ; f_2' = f_2^v\} \rightarrow$ FE

$x_q = \langle x_{q1}, x_{q2}\rangle$
$\downarrow$ FT of FE

$x_q' = \langle x_{q1}', x_{q2}'\rangle \rightarrow \boxed{LR} \rightarrow y_q$

(2) how to know which transform to apply
   → By Superver.

Case 2



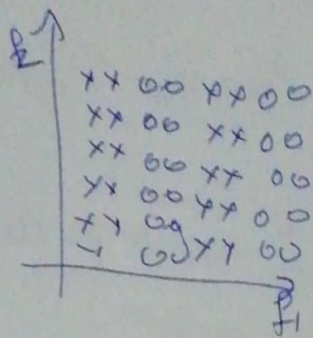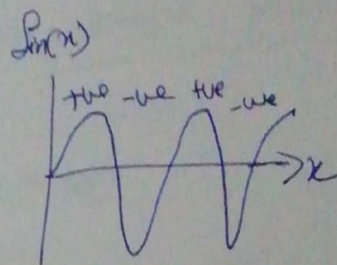$\xrightarrow[\text{FT}]{\text{FE}}$


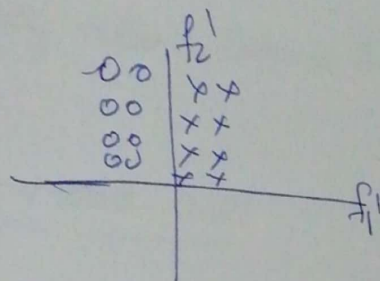+ve

XOR



$f_1' = f_1 \times f_2$

$f_2' = f_2$

③



$$f_1' = \sin(f_1)$$
$$f_2' = f_2$$



Typical transform for real value features:

① $f_1 * f_2$, $f_1^\vee$, $f_2^\vee$, $f_2^3$, $f_1^3$ } polynomial features

② Trignometric features
$$\sin(f_1) \; ; \cos(f_1)$$
$$\sin(f_1) * \cos(f_2)$$
$$\sin(f_1^\vee)$$

③ boolean features :- OR, AND, XOR

④ other
$$\log(f_1)$$
$$e^{f_1}$$