

Quadratic Surfaces - Circles, Parabola

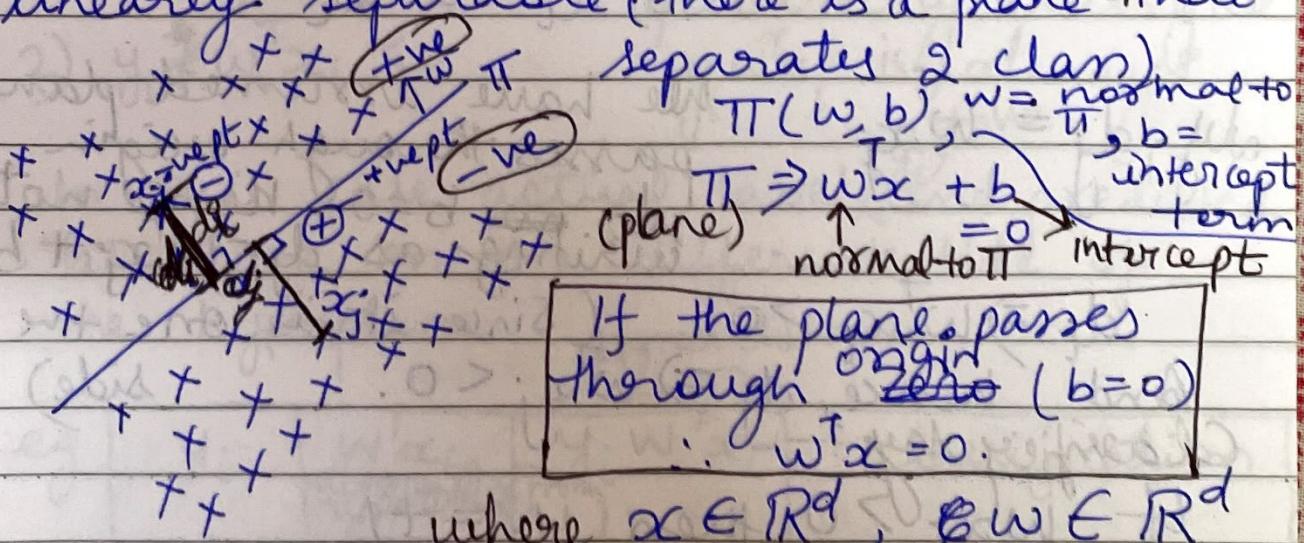
Ch 8 - Logistic Regression

VI: Geometric Intuition of Logistic Regression

Not regression. It is actually a classification technique.

Geometric Intuition:

Assume that classes are almost perfectly linearly separable (there is a plane that



where $x \in \mathbb{R}^d$, $w \in \mathbb{R}^d$
 b = scalar, x = vector, w = vector

Assume, classes are almost perfectly linearly separable.

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

Find: w & b . Find a plane Π that best separates the pts from negative points.

Say $y_i = +1$ for class +ve

$y_i = -1$ for class -ve.

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

$$d_i = \frac{\underline{w^T x_i}}{\|w\|} \quad \begin{cases} \text{we assumed:} \\ \hookrightarrow w \text{ is normal to the plane} \\ w = \text{unit vector.} \end{cases}$$

$$\therefore \|w\| = 1.$$

So, $d_i = w^T x_i$

Also, $d_j = w^T x_j$

We have assumed plane passes through origin
so $b=0$. Hence not writing as $d_j = w^T x_j + b$.

$$d_i \Rightarrow w^T x_i > 0 \quad (\text{since it is one +ve side})$$

And hence, $d_j \Rightarrow w^T x_j < 0$.

Classifier says:-

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

$$w^T x_j < 0 \text{ then } y_j = -1$$

Case 1:

~~so $y_i * w^T x_i > 0$ is correctly classifying the point.~~

Case 2:

~~& $y_i * w^T x_j < 0$ is correctly classifying points.~~

Case 3:

~~$y_i = +ve$ point.~~

~~if x_i is in -ve pt's direction.
but x_i is actually +ve. Misclassified~~

NOTE:

$$1) y_i * w^T x_i > 0$$

In this case actual label (y_i) and predicted label \hat{y}_i are on the same side of the hyperplane

| eg | y_i | $w^T x_i$ <small>(signed distance from hyperplane)</small> | $y_i \cdot w^T x_i$ | classification |
|----|-------|---|---------------------|----------------------|
| | -ve | -ve | +ve | Correctly classified |
| | +ve | +ve | +ve | " " |

$$2) y_i * w^T x_i < 0$$

In this case actual label y_i and predicted label \hat{y}_i are on the different side of the hyperplane

| eg | y_i | $w^T x_i$ <small>(signed distance from hyperplane)</small> | $y_i \cdot w^T x_i$ | classification |
|----|-------|---|---------------------|------------------------|
| | +ve | -ve | -ve | Incorrectly classified |
| | -ve | +ve | -ve | " " |

If incorrectly classified $y_i \cdot w^T x_i$ is always -ve.

Logistic Regression will predict sigmoid($w^T x_i$) < 0.5 as negative and sigmoid($w^T x_i$) ≥ 0.5 as positive, here threshold is 0.5

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

optimization function

V2: Sigmoid Function and Squashing

w^* For every plane there
are unique perpendiculars
optimal w
(best hyperplane's
perpendicular)

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

optimal w
not max of $y_i w^T x_i$

where,

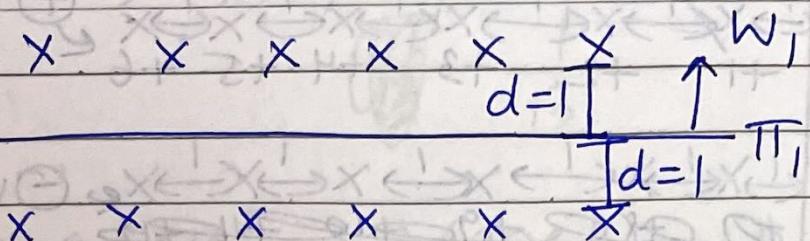
$w^T x_i$ = distance from x_i to π
 w is unit vector.

if signed distance is +ve,
plane (π) defined by w ~~is~~
correctly classifies x_i

if signed distance is -ve,
plane (π) defined by w
incorrectly classifies x_i

\ominus neg pt

$$d = 100 \text{ units}$$



Case 1: Π_1 is the class separating hyperplane

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

For +ve points, $y_i = +1$, $w^T x_i = +1$, $y_i w^T x_i = +1$

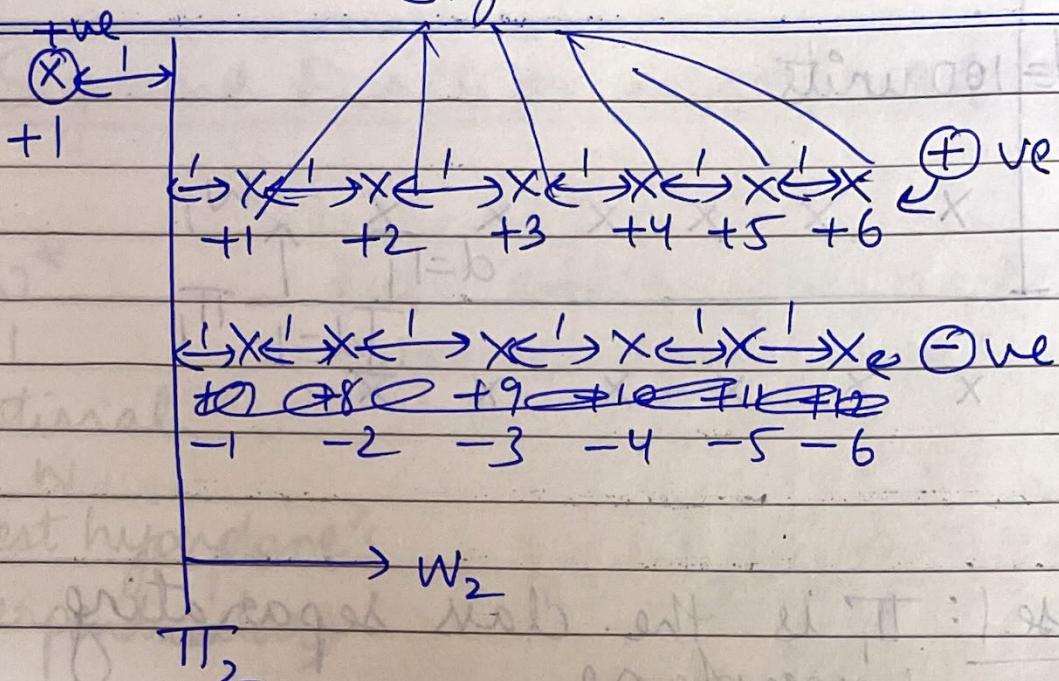
For -ve points, $y_i = -1$, $w^T x_i = -1$, $y_i w^T x_i = +1$

So for all 12 points $\sum_{i=1}^n y_i w^T x_i = 1+1+1+1+1+1 (= +ve)$

$$-100 \quad (\text{for -ve in +ve point area})$$

Case 2: Π_2 is the class separating hyperplane.

$y_i w_2^T x_i$
signed distance



$$\sum_{i=1}^n y_i w_2^T x_i$$

Sum all case 2 signed distance \Rightarrow

$$\Rightarrow 6 - 6 + 1$$

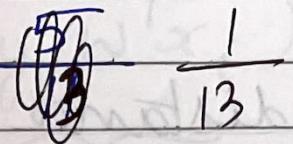
$$\Rightarrow +1$$

But our objective was to find w^* to find maximum signed distance.

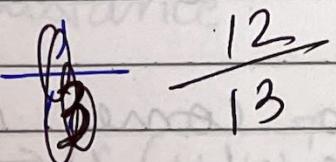
$$w_{\pi_1^*} = -88 < w_{\pi_2^*}^* = +1$$

So π_2^* is my classifier.
but it's doing a terrible job.

Accuracy in case 2:-



Accuracy in case 1:-



So π_1 is better than π_2 .

It is happening: we have an extreme outlier point which is preferring π_2 as the hyperplane.

So, maximizing sum of distance is not outlier prone (outlier can impact)

So, we need to modify the optimization function w.r.t. w*

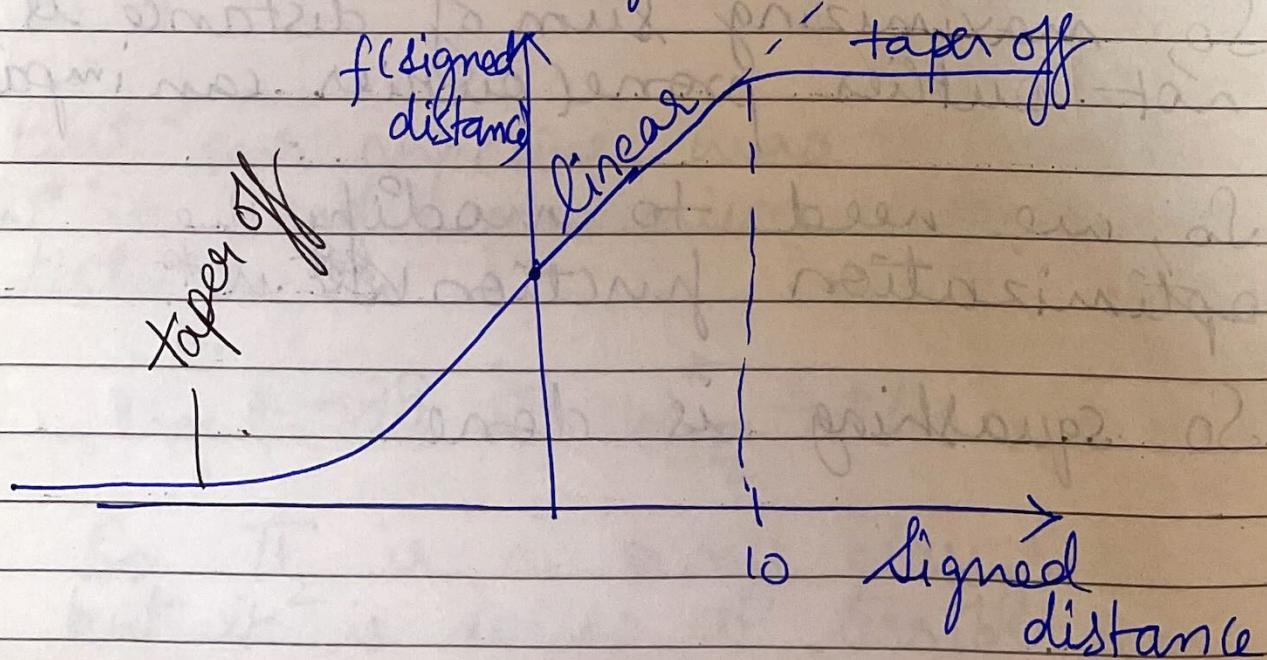
So squashing is done.

Instead of using signed distance,
we use:

If signed distance is small ~~do~~
use it as it is, but if large, make
it a smaller value.

I want to come up with a
function of signed distance
 $f(\text{signed distance})$ which we
use when signed distance \uparrow beyond
a threshold say 10 the ~~the~~ $f(\text{signed-}\text{distance})$ should taper off.

Below 10(threshold), $f(\text{signed-dist})$
should \uparrow with signed distance.



So we should now solve this:

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

signed
distance

This function $\rightarrow f(y_i w^T x_i) = \sigma(x)$
This $\sigma(x)$ is sigmoid function.

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

Sigmoid function.

$$0 \leq \sigma(x) \leq 1$$

$$\sigma(0) = 0.5, \sigma(0) \text{ means } \sigma(\text{dist}=0) = 0.5$$

So this becomes

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

* We chose σ as it has probabilistic interpretation. $\text{dist}=0$ means pt. lies on line
~~off~~ So prob = 0.5 only. Other it's off and prob = 0.5.

σ func. transforms signed distance

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

- comes from e^-

Sigma func. is easy to differentiation and probabilistic interpretation.

It is squashing as it is squashing value $\text{dist}(-\infty, \infty)$ to $\text{dist}(0, 1)$.

Any argmax or argmin func. needs to be differentiable to solve the optimization problem.

V3: Mathematical Formulation of Objective Function writing e^{-x} and ~~*Actual~~ $\exp(-x)$ is same.

Optimization problem -

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

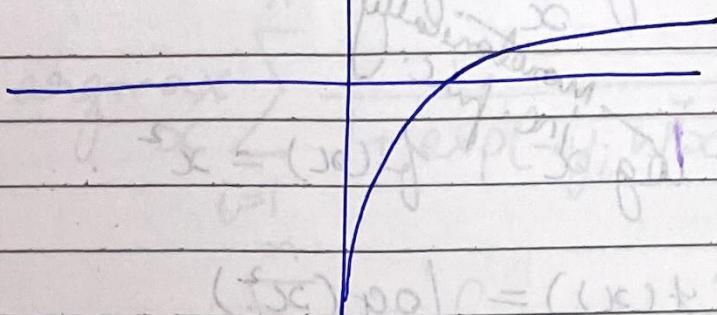
Monotonic function Examples

As $x \uparrow$ if $g(x) \uparrow$, such a func. is called monotonically increasing function.

If $x_1 > x_2$ then $g(x_1) > g(x_2)$ then $g(x)$ is said to be monotonically increasing function for $\Delta x > 0$.

~~for $x > 0$ metu pihorom filosiotaram
for $x < 0$ " preserens "~~

$\log(x)$



$\log(x)$ is only defined for only $x > 0$, not for negative values.

As $x \uparrow$, $\log(x) \uparrow$.
 $x \downarrow$, $\log(x) \downarrow$

e.g Optimization problem :

$$x^* = \underset{x}{\operatorname{argmin}} x^2 \quad | \quad f(x) = x^2$$

Best x

$x^* = 0$ as this is the min. value.

x^2 is not monotonically increasing function.

~~x^2 has 2 parts :-
 x^2 is monotoni~~

monotonically increasing when $x > 0$ | $x^*, g(x)$
 " decreasing " $x < 0$ | $x^*, g(x)$

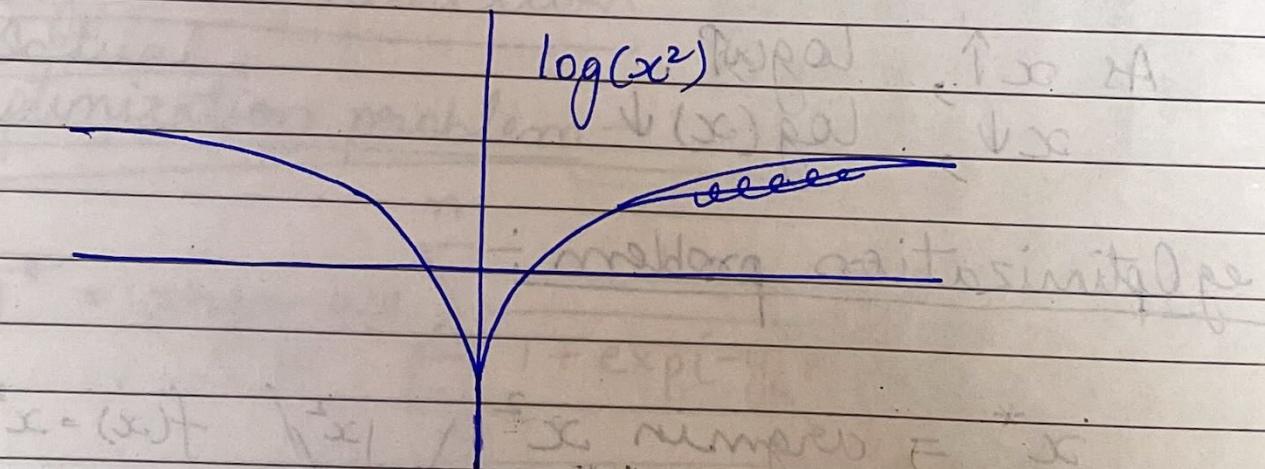
$$x^* = \underset{x}{\operatorname{argmin}} f(x)$$

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

$$g(x) = \log x \quad \text{monotonically inc. func.} \quad f(x) = x^2.$$

$$\text{So } g(f(x)) = \log(x^2)$$

Claim $\Rightarrow x^* = x'$ as $g(x)$ is monotonically increasing function.



Property of monotonic function :-

If $g(x)$ is a monotonic function:-

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

OR

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

contd. Actual optimization problem

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

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

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

So,

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

Rule in Optimization:

Maximizing $f(x) = \operatorname{argmin}_x -f(x)$ Actual optimization problem

So, $\operatorname{argmax}_{w^*} -f(x) = \operatorname{argmin}_x f(x)$

$$L.R. \quad w^* = \operatorname{argmin}_{w \in \mathbb{R}^n} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

geometrically.

So in this formulation, we assumed y_i to be either +1 or -1 but not 0. Not much outlier problem.

Probabilistically,

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

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

Here $y_i = 1 \text{ or } 0$:
not +1 or -1.

V4: Weight Vector

The optimal w^* is the weight vector.

This w^* is also a d-dim vector.

So

$$w^* \in \mathbb{R}^d \text{ as } x_i \in \mathbb{R}^d$$

$$w^* = \langle w_1, w_2, \dots, w_d \rangle$$

$$\text{Also } \langle f_1, f_2, \dots, f_d \rangle$$

If $w^T x_q > 0$ then $y_q = +1$ For every

why $w^T x_q < 0$ then $y_q = -1$ feature there
is a weight associated with

To find prob. $\sigma(y_i w^T x_i)$ it.

gives prob whether
($y_i = +1$) or (-1)

Interpretation of w^*

Cases :-

1) If w_i is +ve] corresponding to f_i ,
as if $x_{q_i} \uparrow$, then $w_i x_{q_i} \uparrow$

$$\text{So } \sum_{i=1}^d w_i x_{q_i} \uparrow$$

$$\text{Then } \sigma(w^T x_q) \uparrow$$

$$\text{Then } P(y_q = +1) \uparrow$$

(Case 2 :- If w_i is -ve]

$$\text{If } x_{q_i} \uparrow, w_i x_{q_i} \downarrow$$

$$\text{Then } \sum_{i=1}^d w_i x_{q_i} \downarrow$$

$$\text{Then } \sigma(w^T x_q) \downarrow$$

$$\text{Then } P(y_q = +1) \downarrow$$

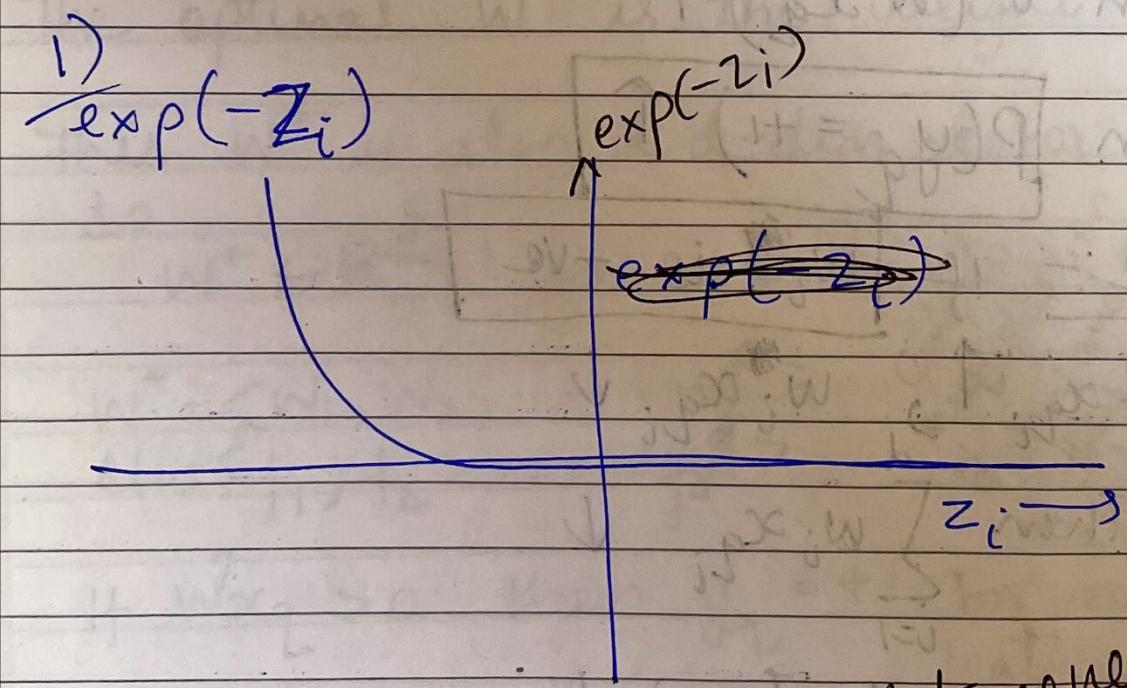
So $P(y_q = -1) \uparrow$ as $P(y_q = +1) = 1 - P(y_q = -1)$

V5: L2 Regularization: Overfitting and Underfitting

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

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

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



$\exp(-z_i) \geq 0$: plot never crosses x -axis and comes below.

As $\log(1) = 0$
 So adding something to $\log(1 + \exp(-z_i))$ must be greater than or equal to 0.

$$\underline{2} \quad \log(1 + \exp(-z_i)) \geq 0$$

If each z_i is ≥ 0 , then the whole term.

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

Minimal value of $\sum_{i=1}^n \log(1 + \exp(-z_i))$

Now, is 0.

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

Case 1: If z_i is +ve

As $z_i \rightarrow +\infty$, $\exp(-z_i) \rightarrow 0$

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

If z_i is +ve, all points are correctly classified.

If I pick W s.t.

- a) All training points are correctly classified by W
- b) $z_i \rightarrow +\infty$

then best W is found.

Perfect fitting on training data means overfitting.

But what if some points are outliers?

- ① Then the model is overfitting.
- ② $w_i \rightarrow +\infty$ or $-\infty$
ie, w_i are getting very large.

then minima value becomes 0.

We missed one key aspect while constructing the optimization problem. that W is a normal. which means that $W^T W = I$ but we didn't use it anywhere. So we use regularization in optimization.

$$w^* = \underset{w}{\operatorname{argmin}} \left[\sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda w^T w \right]^{1/2}$$

where $w^T w = \sum_{j=1}^d w_j^2$ $\xrightarrow{\text{regularization term}}$

which can also be written as

$$\|w\|_2^2$$

\hookrightarrow square of L₂ norm of w.

We minimize both loss + regularizer.

Regularization term is ~~making~~ avoiding
w_i term to go to either +∞ or -∞

There is a tug of war b/w loss & regularizer

λ = hyperparameter (not given to us) term.
regularizer
no so,

If $\lambda = 0$, overfit in training data
high variance.

If $\lambda \rightarrow \infty$ large Underfit in training
high bias no influence of data.

are

V6: L1 regularization and sparsity

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

Now, instead of L2 regularization, we can also use L1 regularization.

$$\|W\|_1 = \sum_{i=1}^d |w_i| \quad \begin{array}{l} \text{— L1 norm} \\ \text{absolute means only +ve value.} \end{array}$$

Even if $w_i \rightarrow +\infty$ or $-\infty$
absolute makes it +ve.

So,

$$W^* = \underset{W}{\operatorname{argmin}} \left(\underset{\text{training data}}{\text{logistic loss for}} \right) + \lambda \|W\|_1$$

Dots same work as $\sqrt{2}/1$

Advantage of using L1 or L2 regularization

Sparsity means, $w = \langle w_1, w_2, \dots, w_d \rangle$

Solution to LR is said to be sparse if many w_i 's are ~~not zero~~ ~~zero~~

Then solution to Log Reg is said to be sparse.

Advantage of L1 regularization

1) On use of L1 regularization, $w_i = 0$ for all w_i corresponding to less important features.

Elastic Net

Using both L1 & L2 regularization.

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

2 hyperparameters :-

λ_1 and λ_2 .

Zeros

- 3 ways to derive logistic regression :-
- 1) Geometry & simple algebra.
 - 2) Probability
 - 3) Loss minimization.

V7: Probabilistic Interpretation: Gaussian

Naive Bayes + Bernoulli distribution = LR

(Sec. 3.1 Tom Mitchell - CMU.edu)

From geometry we get,

$$W^* = \underset{W}{\operatorname{argmax}} \sum_{i=1}^n \log(1 + \exp(-y_i W^T x_i)) + \text{reg.}$$

y is either +1 or -1

(In Naive Bayes, if features are real valued, then conditional probabilities are Gaussian distributed and $y_i = +1$ or 0 are Bernoulli random variables)

From probabilistic, we get : (like in columns)

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

Logistic Reg =

Gaussian Naive

reg. Bayes + Bernoulli

random variable

where $y_i = +1$ or 0.

$$\text{and } p_i = \sigma(W^T x_i)$$

Case 1: $y_i = +ve$ geom : $y_i = +1$

prob : $y_i = +1$

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

$$\text{prob : } -1 + \log \frac{1}{1 + \exp(-W^T x_i)}$$

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

- of log with
reciprocal gets
inverted

$$= \log(1 + \exp(-W^T x_i)) \quad \begin{cases} \text{Second term is} \\ 1 - y_i = 1 - 1 = 0 \\ \text{so second term is 0} \end{cases}$$

Case 2: $y_i = -ve$

$$\text{geom} = y_i = -1$$

$$\text{prob} = y_i = 0$$

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

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

Using $-\log(x) = \log(\frac{1}{x})$, we get:

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

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

Using $\frac{1}{e^{-x}} = e^x$, we get.

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

V8: Loss Minimization Interpretation

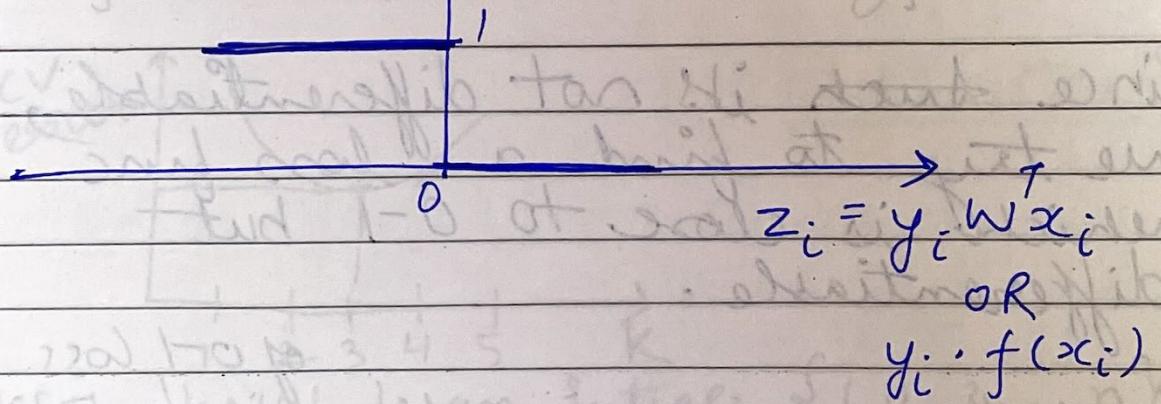
Ideal loss func:-

+1: Incorrectly classified point

0: correctly classified point

$$w^* = \underset{w}{\operatorname{argmin}} \left(\text{no. of incorrectly classified points} \right)$$

Loss func.



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

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

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

To solve optimization in ML we use differentiation in calculus.

But a big issue with selecting this loss func. is it is non-differentiable.

For a func. to be differentiable it must be continuous.

There is a discontinuity at $z_i = 0$.

But it is not since $0-1 \text{ loss} = 1$ when $z_i < 0$ and 0 when $z_i > 0$.

Since stuck its not differentiable, we try to find a loss func. which is close to 0-1 but differentiable.

One such approximation is logistic loss for Logistic Regression.

Another approximation of 0-1 loss would be

* Hinge loss for SVM

* Square loss for Linear Regression

* Exponential loss for AdaBoost.

By just changing the loss func slightly you get many different ML models.

V9: Hyperparameter Search: Grid Search and Random Search

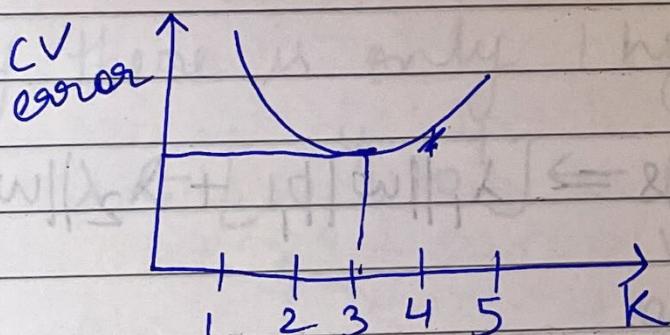
λ = hyperparameter (when doing regularization)
If $\lambda=0$, overfitting, $\lambda=\infty$ underfitting.

Q. How to determine the best λ .

$$\lambda = LR$$

$$K = KNN$$

$$\alpha = NB \text{ (Laplace Smoothing)}$$



K in KNN is an integer $\{1, 2, 3, \dots, n\}$

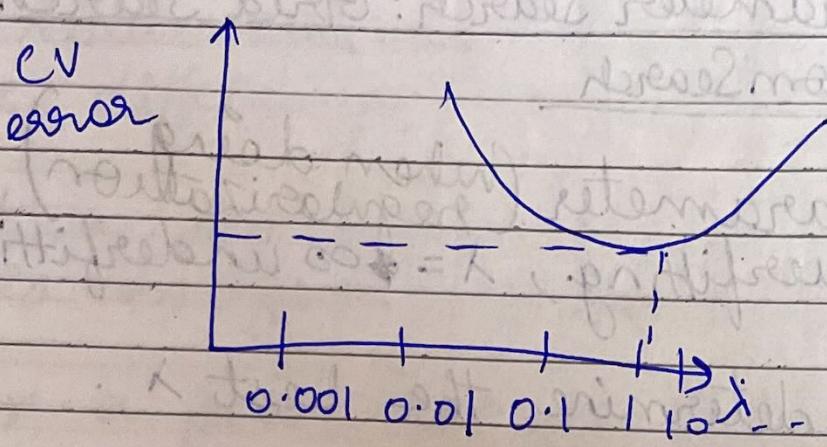
λ in LR is a real number

$$\lambda \in \mathbb{R}, \lambda = 0.1234, \text{etc.}$$

Grid Search, it is somewhat
brute force technique.

eg - ~~λ~~ = [0.001, 0.01, 0.1, 1, 10, 100, 10k]
OR

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

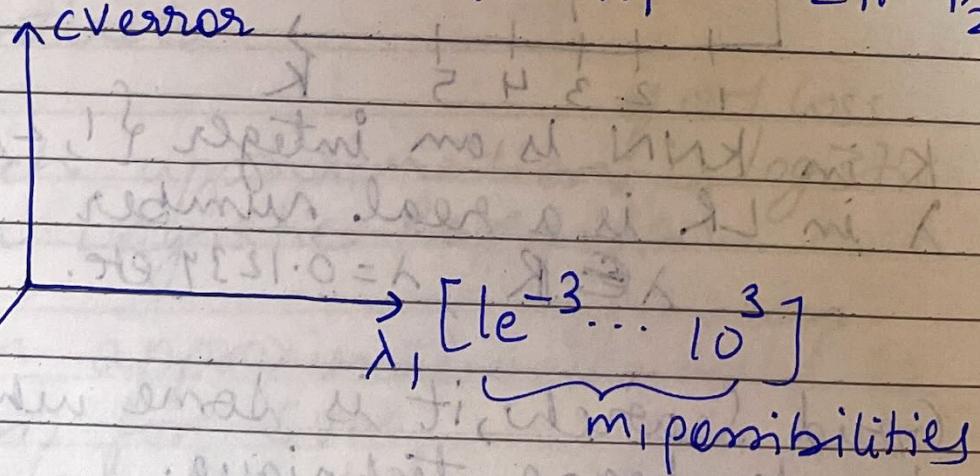


$$\lambda = [1e^{-4} \dots 10^4]$$

$\frac{1}{10^{-4}}$

Elastic Net

$$\text{hyperparameters} = \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^2$$



$$\lambda_1 [1e^{-3} \dots 10^{-3}]$$

m_1 possibilities

Total possibilities =

$$[1e^{-3} \dots 10^{-3}]$$

$\underbrace{\quad\quad\quad}_{m_2 \text{ possibilities}}$

$$M_1 * M_2$$

Probabilistically, Logistic Regression is nothing but Gaussian Naive Bayes + Bernoulli

Problem with Grid Search:-

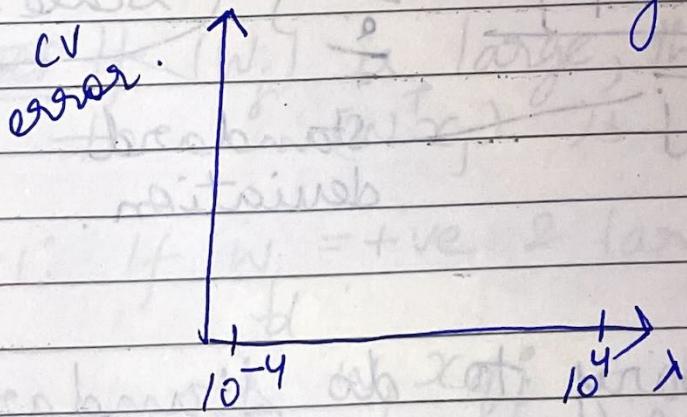
We have to train and compute m_i times where i is the no. of hyperparameter.

So grid search m_i hyperparameter search technique increases exponentially. If you have k hyperparameters, You will have to train m^k hyperparameters.

Random Search

If there is only 1 hyperparameter:-

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



This technique is as good as grid search.

V10: Column Standardization

| x_i | f_1 | f_2 | f_j | f_d |
|-------|-------|-------|-------|-------|
| x_i | | | | |

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

modified point

$$x'_{ij} = x_{ij} - \mu_j$$

$$\alpha_j$$

First we compute μ_j and σ_j of that column and then we do column standardization. For each column we do standardization as different column could be measured on a different

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

centering scale.

standard deviation.

It is compulsory to do standardization in Logistic Regression.

Feature/Column Standardization is also called Mean centering and scaling.

Mean centering = $x_{ij} - \mu_j$ sub
Scaling = σ_j division

Probabilistically, Logistic Regression is nothing but Gaussian Naive Bayes + Bernoulli

VII: Feature Importance and Model Interpretability

general feature

$$f_1, f_2, f_3, \dots, f_j, \dots, f_d \downarrow \\ w \rightarrow w_1, w_2, w_3, \dots, w_j, \dots, w_d \text{ (feature weights)}$$

Assumption: All features are independent
If this assumption is true,

We can get feature importance using w_j .

We can take

LR: $\propto w_j$'s to determine feature importance.

$|w_j|$ = absolute value of weight corresponding to feature f_j .

If (w_j) is large, then contribution to $(w^T x_j)$ is large.

Case 1: If w_j = +ve & large.

$$\sum_{j=1}^d w_j x_{q,j} = w^T x_q$$

It impacts Prob. of query point
 $P(Y_q = H)$

fraction

Case 2: w_j : -ve and large

$$\sum_{j=1}^d w_j x_{q,j}, P(y_q = -1) \text{ is large.}$$

When something is large (+ve or -ve) ^{its impact} is large.
So, we can determine the important feature in LR through feature weights

e.g. → Predict gender: male & female
 $(+1)$ (-1)

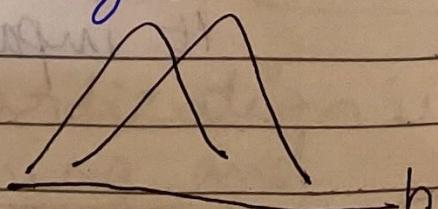
*① Hair length $= |w_{hL}|$ is large. Then, $w_{hL} = -ve$
As $|w_{hL}| \uparrow$; $P(y_q = -1) \uparrow$

So hL will have a large negative weight value.

*② Height $\uparrow = |w_h|$ is positive

As $|w_h| \uparrow$; $P(y_q = +1) \uparrow$

w_h will get a medium positive weight.



Model Interpretability = Give reasoning when absolute weight of feature is large.

x_2

$$y_2 = +1 \text{ or } y_2 = -1$$

V12% Collinearity of features or Multicollinearity

Feature Importance can be calculated using absolute value of weights of independent features only.

But collinearity/multicollinearity is calculated for non independent features..

$$f_i > f_j$$

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

then f_i & f_j are collinear

Multicollinearity

Collinearity among multiple features if f_1, f_2, f_3 & f_4 s.t.

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

then f_1, f_2, f_3, f_4 are multicollinear.
(all α are constant)

Q Why does $|w_j|$ not be useful as f. I. if features are collinear?

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

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

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

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

$$\tilde{W} \Rightarrow W^T x_q = x_{q_1} + 3x_{q_1} + 3x_{q_3} = 4x_{q_1} + 3x_{q_3}$$

Now weight vector is $= \langle 4, 0, 3 \rangle$

$$W^* \Rightarrow \text{Original feature vector} = \langle 1, 2, 3 \rangle \\ \tilde{W} = \langle 4, 0, 3 \rangle$$

Both weight vector give same classifier
If features are collinear, then weight vector can change arbitrarily.

If

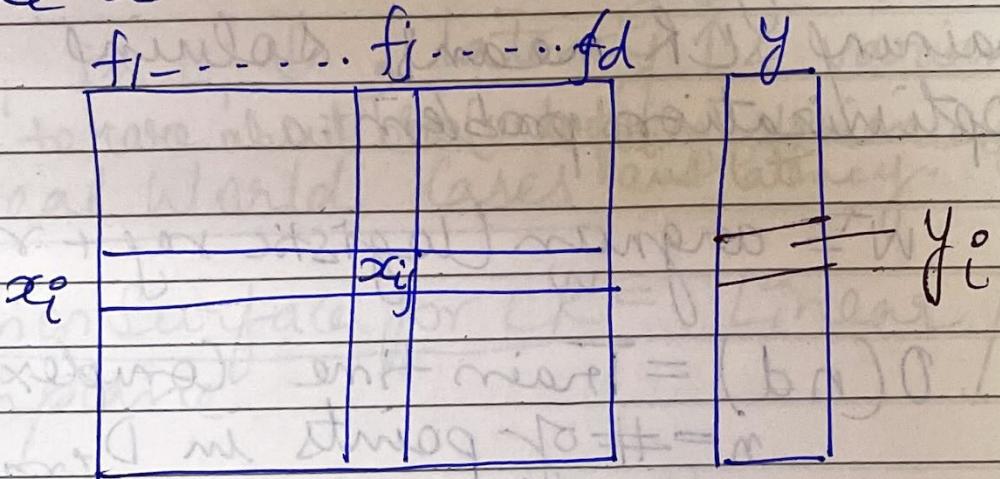
$|w_j|$ cannot be used as f. I.
Earlier with $\langle 1, 2, 3 \rangle$, we felt $f_1 > f_2 > f_3$.
But now, $\langle 4, 0, 3 \rangle$ $f_1 > f_3 > f_2$ in terms of F. I.

Determine if features are multicollinear or not?

Ans → Perturbation Technique.

Perturbation technique

Shake a little



~~x_{ij}~~ $x_{ij} + \epsilon$ Small noise/error from $N(\mu, \sigma^2)$

Before adding noise calculate weight vector and after adding noise again calculate weight vector.

If w_p and \tilde{w}_p differ significantly then data (features) are collinear.

Then you can't use $|w_j|'$'s as feature importance.

IMP. Forward feature Importance
can be used for any model not just KNN.

V13: Train & Runtime & Space Complexity

Training LR means solving optimization problem.

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

$O(nd)$ = Train time complexity.

n = # of points in D_{train} .

d = features in D_{train} .

We just have to store weight vector $w^* = [w_1, w_2, \dots, w_d]$ in memory.

Runtime Complexity

Time = $O(d)$

d = # of features.

Space = Od

d = # of features.

If d is small, the LR is v. good for low latency applications.

$w_i \neq 0 \rightarrow \text{imp. features} = 0.$

In case of high dimensionality, we can make use of L1 reg. To create sparsity in W so as to reduce calculations during runtime.

As λ_1 in $L_1 \uparrow \rightarrow$ bias \uparrow Underfitting
 \therefore latency \downarrow

λ_1 in $L_1 \downarrow \rightarrow$ latency \uparrow & variance \uparrow

You have to care about bias (underfitting), variance (overfitting) and latency.

V14: Real World Cases

~~Decision Surface for LR = Linear / Hyperplane~~

~~Assumption = Data is linearly or almost linearly separable.~~

~~Feature Imp & Interpretability - we can use two. If given features are not multicollinear.~~
~~Imbalanced data = Upsampling/Downsampling~~

Outliers = Less impact due to outliers due to sigmoid function.

but not completely avoided.

(1) So, first calculate W^* from D_{train} .

(2) Calculate W^*x_i on each x_i using W^{*T} . You get distance from plane Π to x_i .

(3) Remove points which are very few far away from Π from

D_{Train} and you get D'_{Train}

(4) calculate \tilde{w}^* using D'_{Train}

So, \tilde{w}^* is the final solution.

Missing values = imputation method like mean, median, mode.

Multiclass = One vs. Rest \leftarrow Typically,

(i) Maximum Entropy models \rightarrow Extension to LR.

(2) Softmax classifier in Deep learning

(3) Multinomial LR.

Similarity Matrix = Regular LR does not work on similarity matrix so we use kernel & Logistic Regression (Research paper)

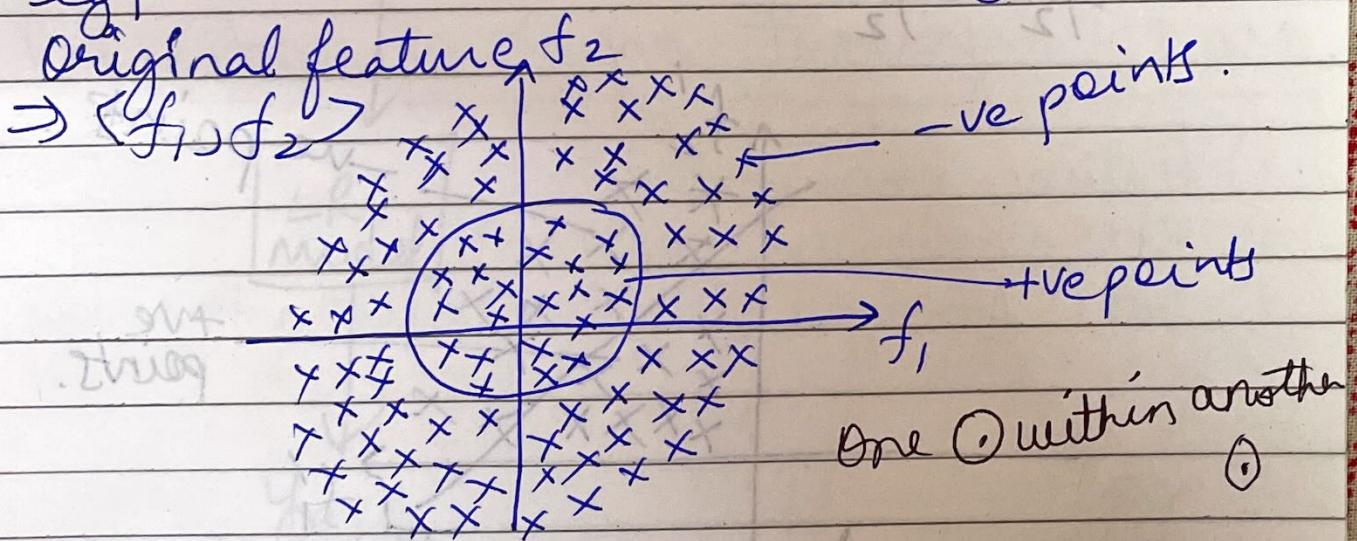
Best Case = If data is almost linearly separable and (LI Regions) then the LR works best.

Large Dimensionality = chances of linearly separable is high.

But if you want a low latency system^{exp}, when d is large then use L1 regularization.

V15: Non-linearly Separable data and feature engineering

e.g.



Can we use Logistic Regression to separate the classes.

→ Drawing a plane in the space of f_1 and f_2 .

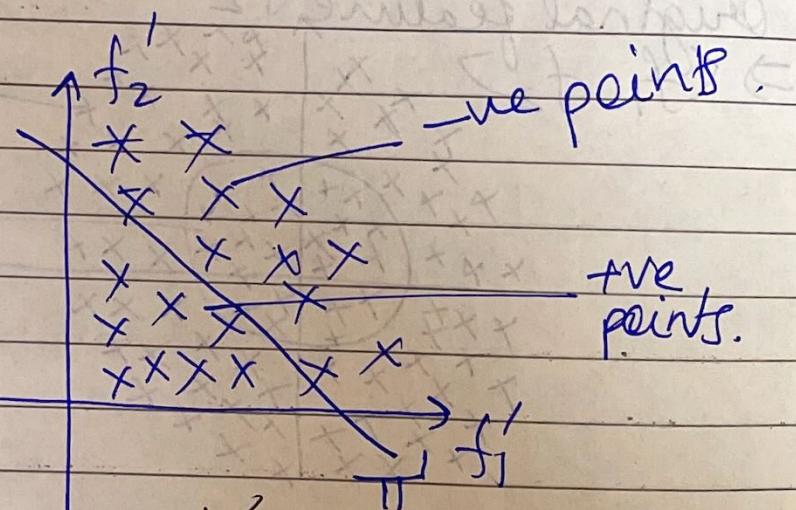
→ Cannot separate +ve and -ve pts by drawing a plane.

So we do feature transformation or feature engineering and we modify the features to a completely different space.

Engineered features: $f'_1, f'_2 \dots$

$$f'_1 = f_1^2$$

$$f'_2 = f_2^2$$



We have engineered features, also we will do same engineering of query points.

$$x'_i = \langle x'_{i1}, x'_{i2} \rangle$$

$$x'_{i1} = x_{i1}^2$$

$$x'_{i2} = x_{i2}^2$$

$$x_q = \langle x_{q_1}, x_{q_2} \rangle$$

↓ FE/FT

$$x'_q = \langle x'_{q_1}, x'_{q_2} \rangle$$
$$\begin{matrix} x'_{q_1} \\ x'_{q_2} \end{matrix}$$

LR
model

y_q

Q How did you know which transform to apply?

Comes with practice.

These were 2 concentric circles.
So we use equation of a circle using geometry.

quadratic

$$\boxed{f_1^2 + f_2^2 = r} \quad (\text{not linear in } f_1 \text{ & } f_2)$$

But linear in $f_1^2 + f_2^2$.

So $\boxed{f'_1 + f'_2 = r}$ linear

Equation of line: $ax + by + c = 0$.

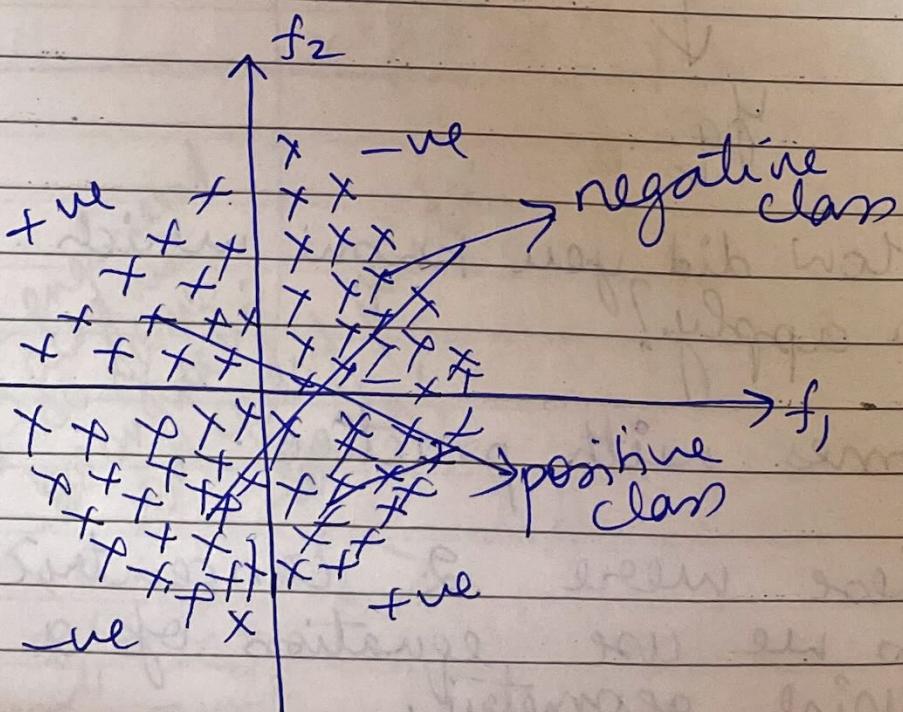
Most important aspect of applied AI

- 1) Feature Engineering
- 2) Bias-variance tradeoff
- 3) Data Analysis & Visualization technique.

ML+

~~eg²~~

~~x₀ &
data~~



| XOR | 0 | 1 |
|-----|---|---|
| 0 | 0 | 1 |
| 1 | 1 | 0 |

FE/FT

f_2'



f_1'

the side of f_1'

$$f_1' = f_1 * f_2 \quad \text{FE/FT.}$$

$$f_2' = f_2$$

the side of f_1'

They are now
linearly separable.

eg 3

f_2

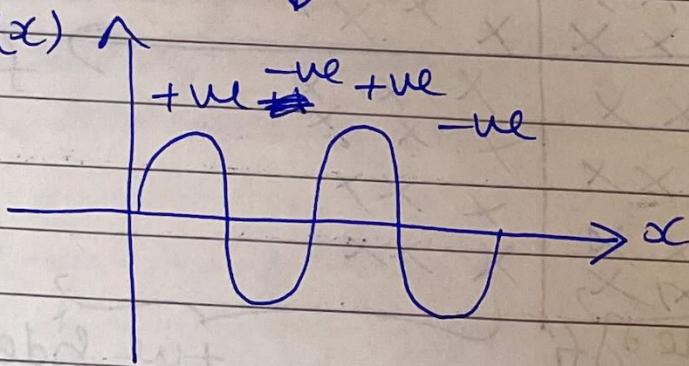
| +ve | -ve | +ve | -ve | +ve | -ve |
|-----|-----|-----|-----|-----|-----|
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |
| xx | xx | xx | xx | xx | xx |

not linearly
separable

f_1

↓ FE/FT

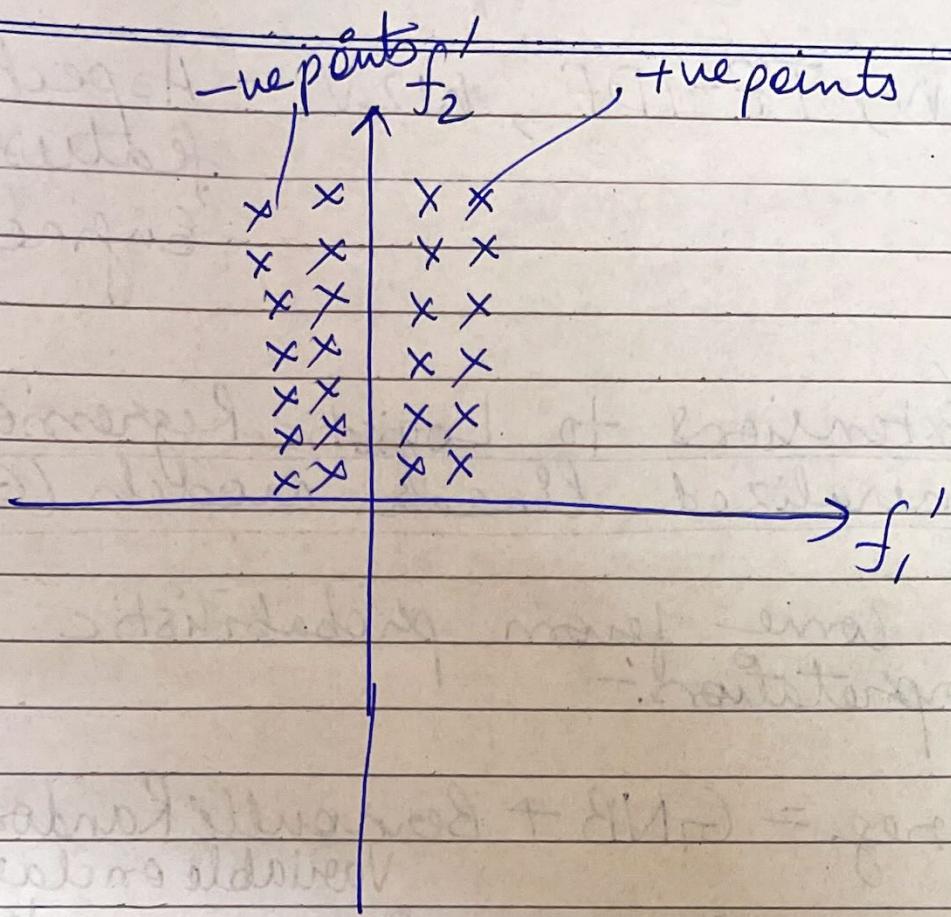
$\sin(x)$



$$f'_1 = \sin(f_1)$$

$$f'_2 = f_2$$

FE/FT



Typical transforms for real valued features:

polynomial features

- 1) Product of features ($f_1 * f_2$)
- 2) Square of " (f_1^2, f_2^2)
- 3) Cube of " (f_1^3, f_2^3)
- 4) Trigonometric functions
 $\sin(f_1), \cos(f_1)$
 ~~$\sin(f_2), \cos(f_2)$~~
- 5) Polynomial trigonometric functions:
 $\sin(f_1) + \cos(f_2)$
 $\sin(f_1)^2$
- 6) Boolean features : OR, AND, XOR
- 7) $\log(f_1), e^{f_1}, \text{etc.}$

In Sklearn Logistic Regression $C = \frac{1}{\lambda}$

BOW, TF-IDF, W2V - Aspects of feature engineering.

V17: Extensions to Logistic Regression - Generalized linear models (GLM)

They come from probabilistic interpretation:-

log reg = GNB + Bernoulli Random Variable on class labels
change it slightly :-

① Multinomial LR :-
Distribution changes from Bernoulli to Multinomial dist.

Useful when you have multiclass classification.

② Linear Regression (Actual Reg. tech.)

$$P(y|x) \sim N(\mu, \sigma^2) \quad y_i \in \mathbb{R}$$

Gaussian assumption

③ Poisson Regression :-
for count prediction for Poisson distributed.
If you want to predict counts say predict how many times a machine fails.