

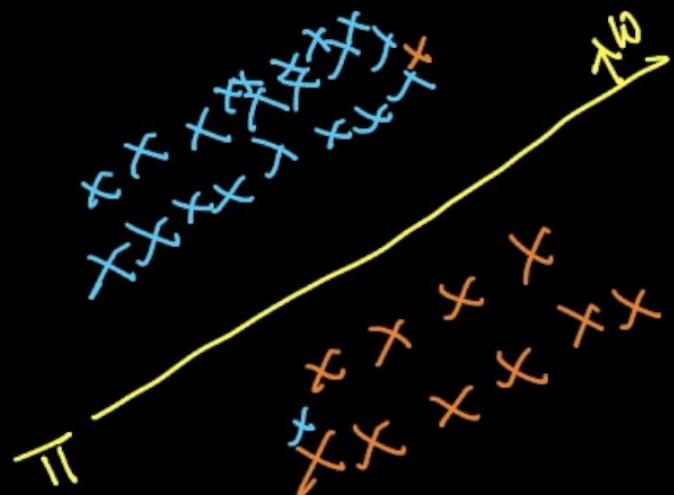
if the TS passes through origin:- $b=0$
 $\underline{w}^T \underline{x} = 0$

$$\text{TS: } \underline{w}^T \underline{x} + b = 0$$

vec \leftarrow ~~$\underline{x} \in \mathbb{R}^d$~~ ; $w \in \mathbb{R}^d$ \rightarrow vec
scalar \leftarrow $b \in \mathbb{R}^1$



Assumption of LR: classes are almost/perfectly }
linearily separable }
} X

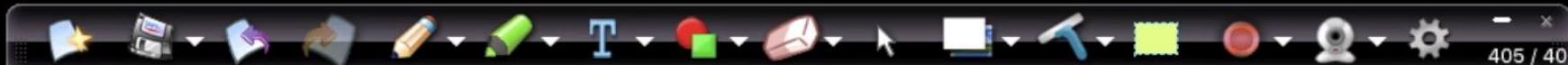


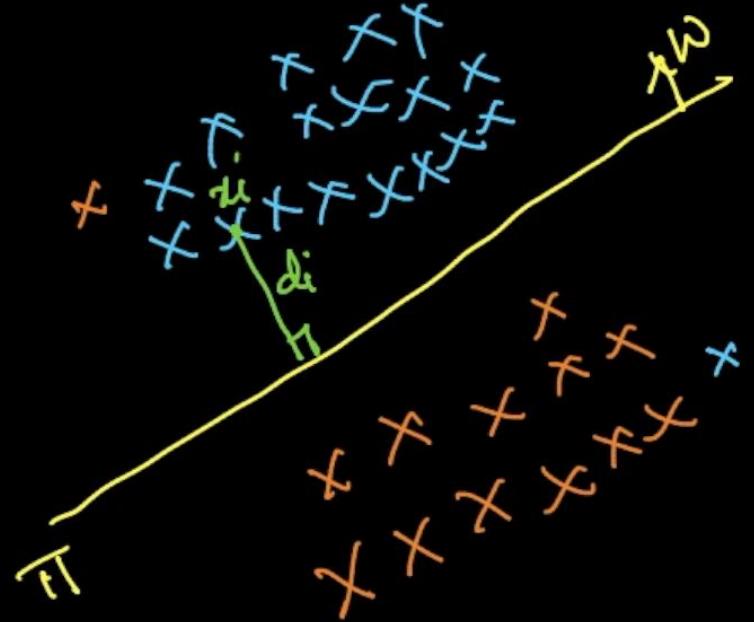
$$\pi: \mathbf{w}^T \mathbf{x} + b$$

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

Task: find: w & b

{ NB: Condition indep of features | Task: π that ^{best} separates +ve pls
K-NN: Neighborhood from -ve pls }





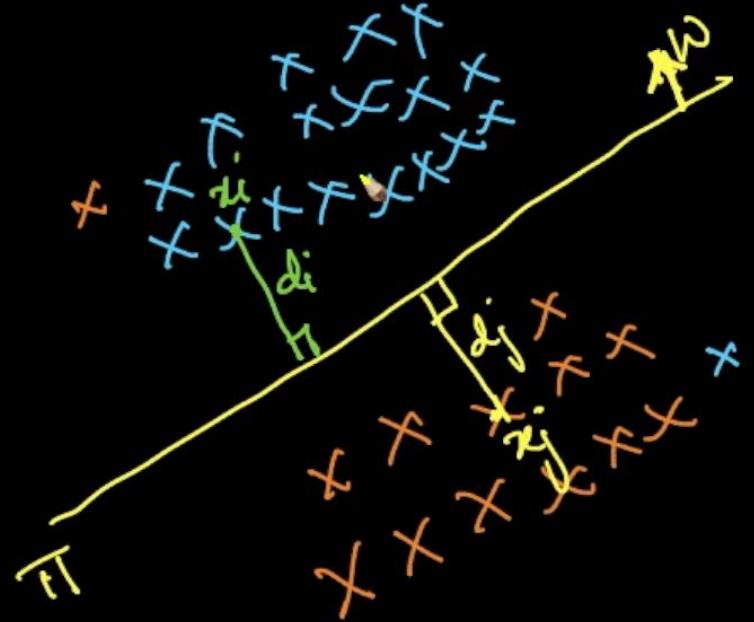
$y_i = \begin{cases} +1 & \text{+ve pls} \\ -1 & \text{-ve pls} \end{cases}$

$\begin{cases} 1: \text{+ve pls} \\ 0: \text{-ve pls} \end{cases}$

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

$$d_i = \frac{\vec{w}^T \vec{x}_i}{\|\vec{w}\|}; \quad \vec{w} \text{ is the normal to the plane}$$

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



$$y_i = \begin{cases} +1 & \text{+ve pls} \\ -1 & \text{-ve pls} \end{cases} \checkmark$$

$$\begin{cases} 1: \text{+ve pls} \\ 0: \text{-ve pls} \end{cases} \times$$

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

$$d_i = \frac{\vec{w}^T \vec{x}_i}{\|\vec{w}\|} ; \quad \vec{w} \text{ is the normal to the plane}$$

$$d_j = \vec{w}^T \vec{x}_j \quad \|\vec{w}\| = 1 \Rightarrow \text{unit vec } \vec{w}$$



$$d_i = \frac{w^T x_i}{\|w\|}$$

w is the normal to the plane
 $\|w\| = 1 \Rightarrow$ unit vector

$$d_j = w^T x_j$$

$$y_i = \begin{cases} +1 & \text{+ve pls} \\ -1 & \text{-ve pls} \end{cases} \checkmark$$

$$\begin{cases} 1: \text{+ve pls} \\ 0: \text{-ve pls} \end{cases} \times$$

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

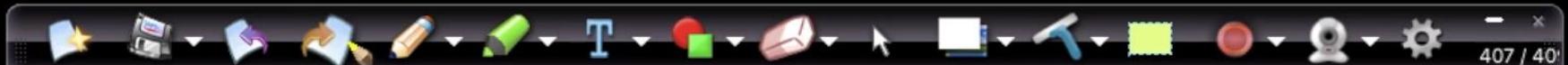
$$d_i = \vec{w}^T \vec{x}_i > 0$$

$$d_j = \vec{w}^T \vec{x}_j < 0$$

classifies

$$\left\{ \begin{array}{l} \text{if } \vec{w}^T \vec{x}_i > 0 \text{ then } y_i = +1 \\ \text{if } \vec{w}^T \vec{x}_i < 0 \text{ then } y_i = -1 \end{array} \right.$$

decision surface in LR
; line / plane



Case 1: (+ve pt) $y_i * \underbrace{w^T x_i}_{y_i = +1} > 0 \rightarrow$ if $y_i \neq +1 \rightarrow$ +ve pt
 $w^T x_i > 0 \Rightarrow$ classifier is saying its +ve pt
 \hat{w} is correctly classifying the pt

Case 2: (-ve pt) $y_i = -1 : -ve pt$
 $w^T x_i < 0 \Rightarrow$ LR concluding that x_i is a -ve pt
 $y_i w^T x_i > 0$



both +ve & -ve pts

$$y_i \underline{w^T x_i} > 0$$

LR model is correctly classifying the pt x_i

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

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

$$\boxed{y_i w^T x_i < 0}$$

$$\begin{array}{l} y_i = +1 \\ \text{LR: } (-1) \end{array}$$

misclassified



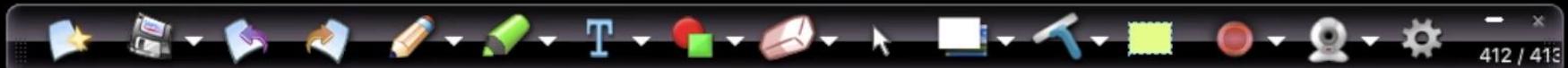
Case 4: $\begin{cases} y_i = -1 \Rightarrow \text{ve class} \\ w^T x_i > 0 \Rightarrow LR \text{ is saying } x_i \text{ is ve pt} \end{cases}$

↳ misclassified

$y_i w^T x_i < 0$

$$\max \sum_{i=1}^n y_i w^T \vec{x}_i$$

n datapls



$$\max_{\mathbf{w}} \sum_{i=1}^n (y_i \mathbf{w}^\top \mathbf{x}_i)$$

variable

n datapts
Train $\rightarrow D_n$

Change / Vary \mathbf{w}

$$\left\{ \begin{array}{l} \text{optimal } \mathbf{w}^* = \arg \max_{\mathbf{w}} \left(\sum_{i=1}^n y_i \mathbf{w}^\top \mathbf{x}_i \right) \\ \text{variable} \end{array} \right.$$

Math:
optimization
problem

$$\max_{\mathbf{w}} \sum_{i=1}^n (y_i \mathbf{w}^\top \mathbf{x}_i)$$

variable

$\frac{n \text{ datapts}}{\text{Train}} \rightarrow D_n$

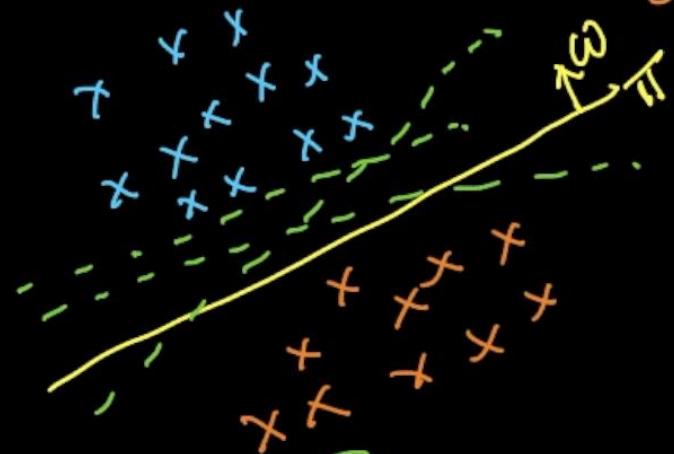
Change / Vary \mathbf{w}

$$\left\{ \begin{array}{l} \text{optimal } \mathbf{w}^* = \arg \max_{\mathbf{w}} \left(\sum_{i=1}^n y_i \mathbf{w}^\top \mathbf{x}_i \right) \\ \text{variable} \end{array} \right.$$

Math. optimization problem

$\left\{ \begin{array}{l} +ve: \text{correctly classified} \\ -ve: \text{misclassified} \end{array} \right.$

Squashing & Sigmoid function:



$$\pi_i: \underline{\omega_i}$$

✓ $\underline{\underline{\omega}}^* = \underset{\omega}{\operatorname{arg\max}}$

optimal ω

best/optimal $\omega \rightarrow \omega^*$

optimization problem

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

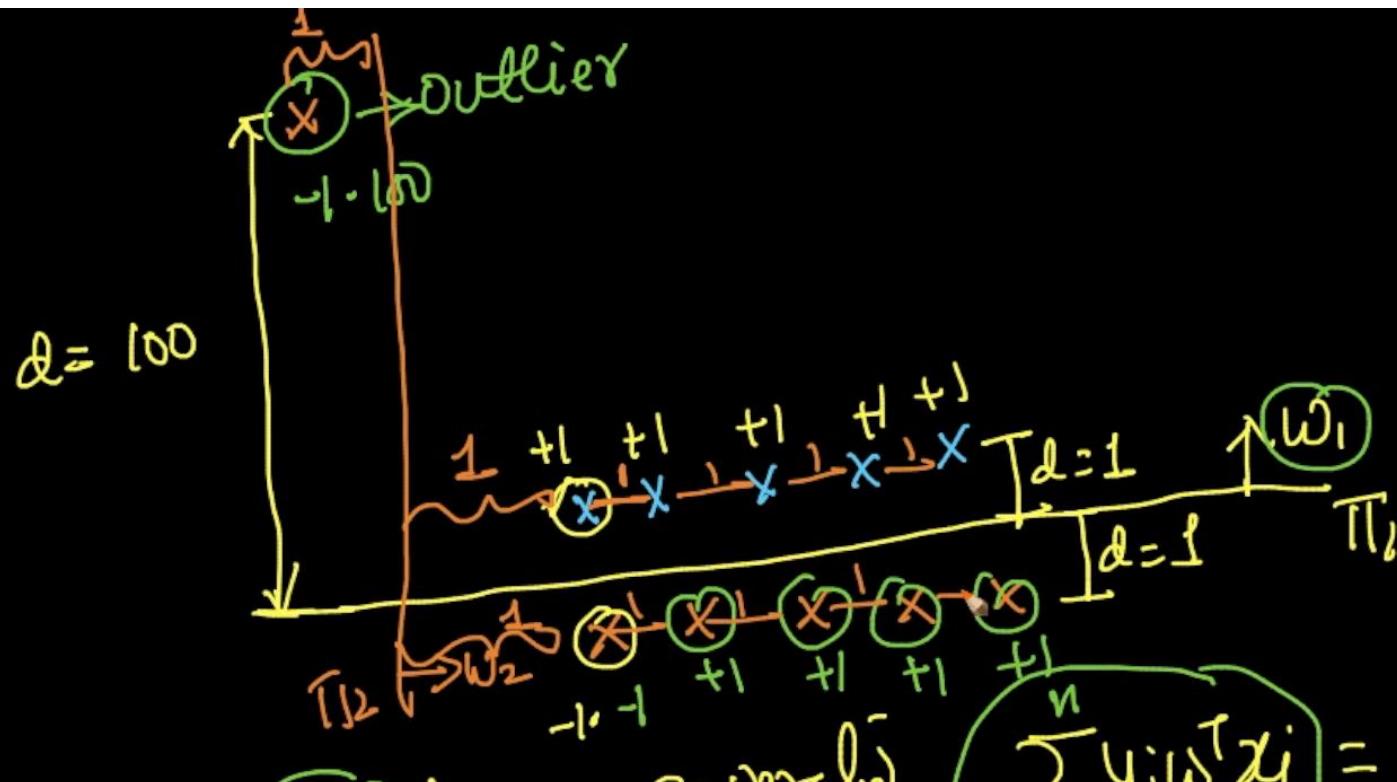
$$\arg \max_w \sum_{i=1}^n y_i w^T x_i$$

↑ signed distance
+1 $\frac{n-1}{n}$

$w^T x_i$: dist from x_i to Π (w is a unit vector)

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



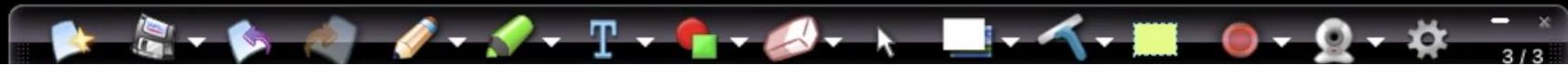


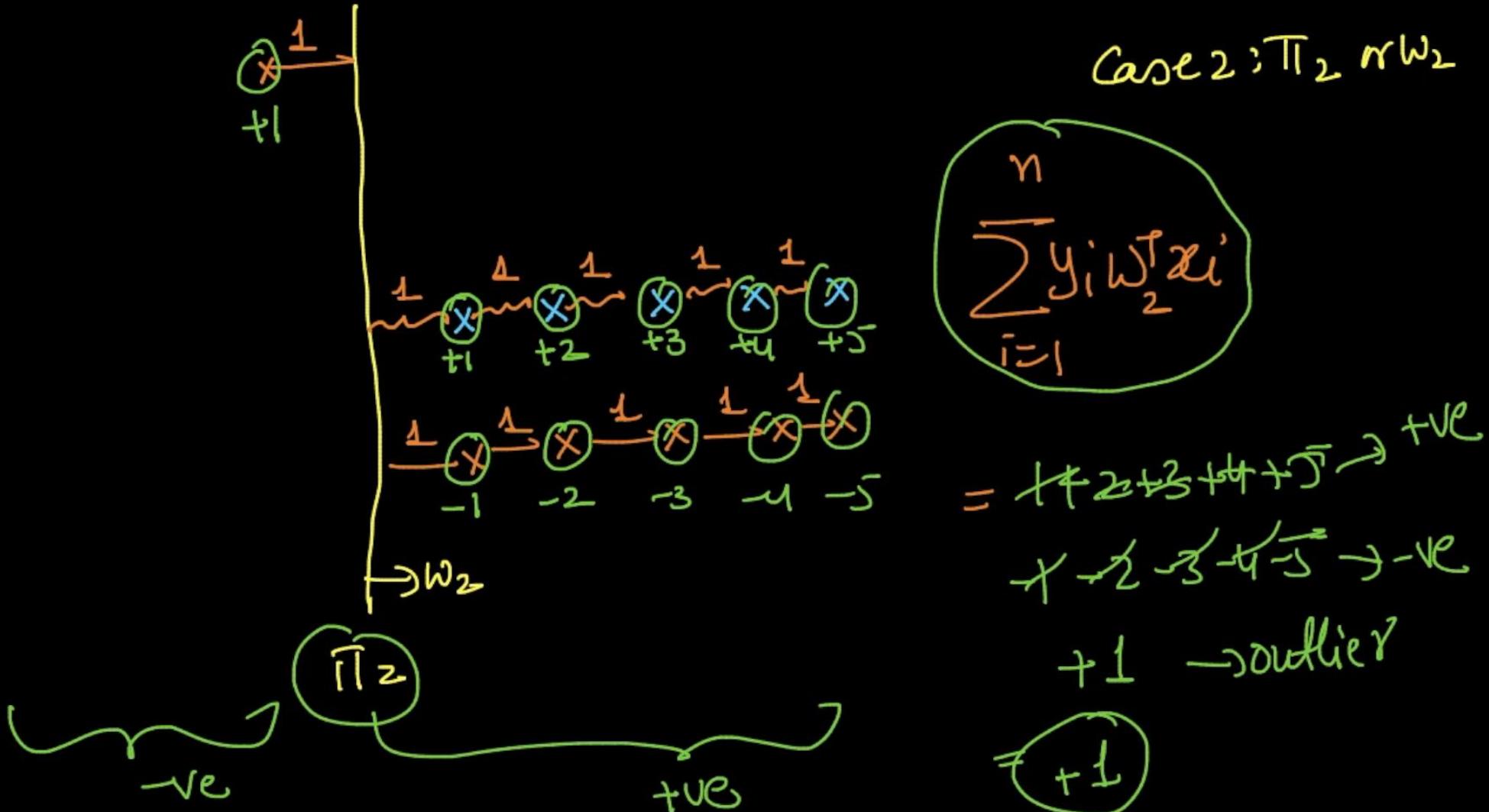
Case 1: Π_1 is my separation

Case 2: $\Pi_2(\omega_2)$

find a best
 $\omega(\pi)$
 Max. sum of
 Signed dist

$$\sum_{i=1}^n y_i \omega_i^T x_i = (+1 + 1 + 1 + 1 + 1) \rightarrow +ve \\ + (-1 - 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1) \rightarrow -ve \\ = -100 \\ = -90$$





obj:- find ω that Maximizes sum of Signed dist
 $\iff \max_{\omega} \sum_i y_i \omega^T x_i$

+1 > -90
✓ $\underline{\pi_2}$ as my classifier

intuitively:- π_1 is better than π_2



one single extreme/outlier pt is changing my
model → hyperplane)

↓
~~Very bad~~

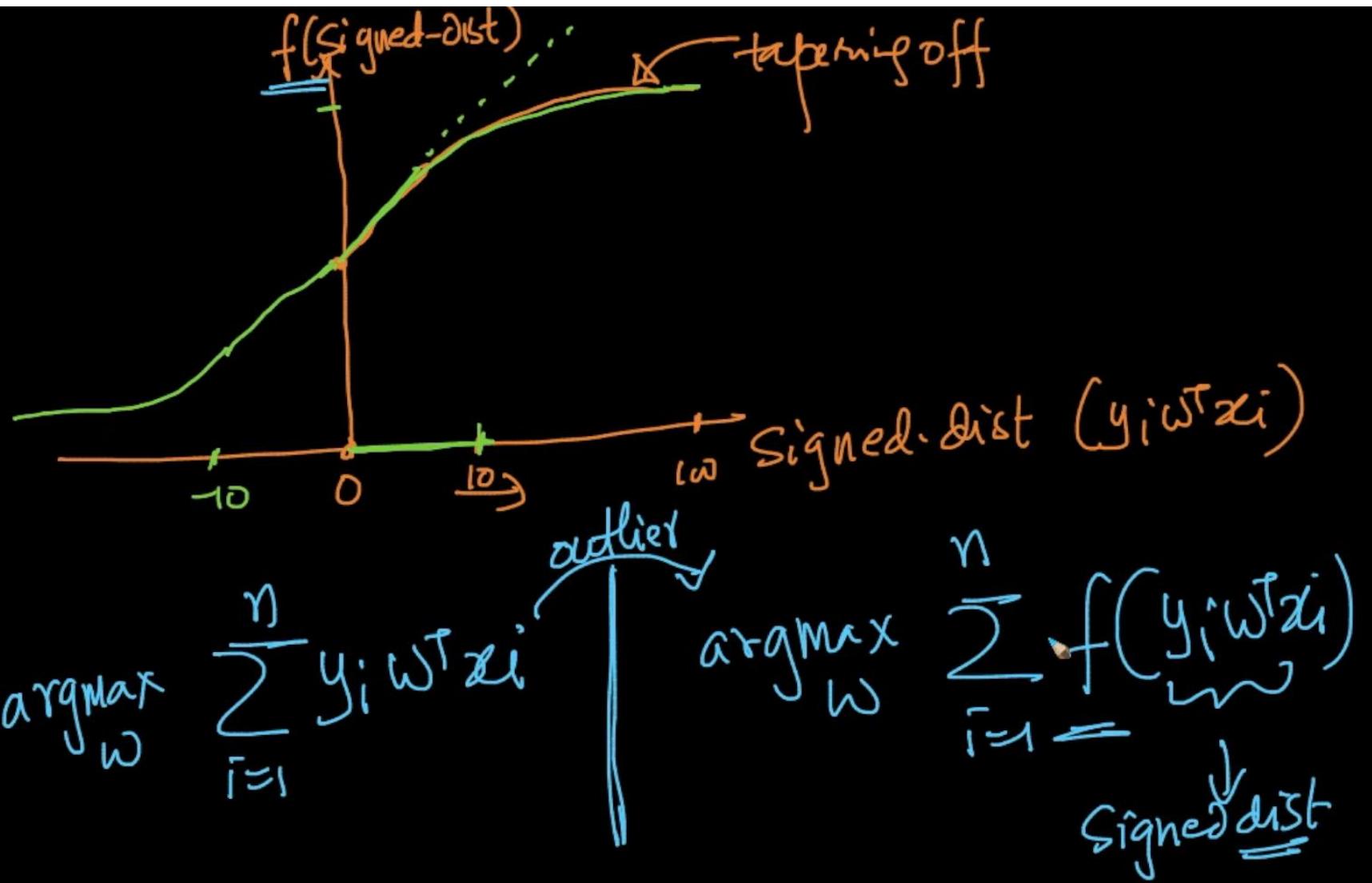
max-sum-of signed-dist

not outlier prone

Squashing:

- Idea:- instead of using signed-dist
- if signed-dist is small :- use it as is
- if signed-dist is large: make it a smaller value





Sigmoid fn:- $\sigma(x)$

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

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

probabilistic interpretation
 $p(y_i = 1)$

30:42

Contents - Google Docs NBayesLogReg.pdf plot(1/(1+e^-x)) - Google Search

Secure | https://www.google.co.in/search?ei=p7QoWtf6PMTqvASm5JiwCg&q=plot%281%2F%281%2Be%5Ex%29%29&oq=plot%281%2F%281%2Be%5Ex%29%29

Google plot(1/(1+e^-x))

All Maps Images Videos News More Settings Tools

About 16,10,00,000 results (0.35 seconds)

Graph for $1/(1+e^{-x})$

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

✓ Max : 1 }
✓ Min : 0 }
 $\sigma(0) = 0.5$

$x = \text{signed dist}$

plot x e^-x, x^2 e^-x, x=0 to 8 - Wolfram|Alpha

Contents - Google Docs NBayesLogReg.pdf plot(1/(1+e^-x)) - Google Search

Secure | https://www.google.co.in/search?ei=p7QoWtf6PMTqvASm5JiwCg&q=plot%281%2F%281%2Be%5Ex%29%29&oq=plot%281%2F%281%2B%2Be%5Ex%29%29

Google plot(1/(1+e^-x))

All Maps Images Videos News More Settings Tools

About 16,10,00,000 results (0.35 seconds)

Graph for $1/(1+e^{-x})$

$f(0) = 0.5$

$w^T z_i = 0$

$p(y_i = 1) = 0.5$

$w^T z_i$ is v. large

$p(y_i = 1) = 1.0$

plot x e^-x, x^2 e^-x, x=0 to 8 - Wolfram|Alpha

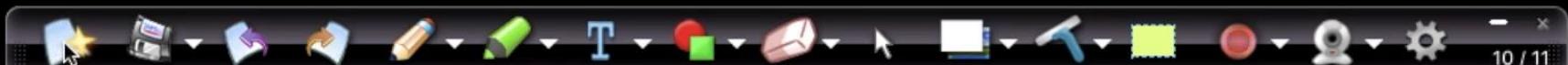
2 / 2

max. sum. of signed dist \rightarrow outlier



$\sigma(\alpha)$ $\xrightarrow{\text{sigmoid}}$
 \hookrightarrow tapering linear
 \hookrightarrow prob.

max. sum. of $\underbrace{\text{transformed signed dist}}_{\sigma}$



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

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

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

→ less impacted
by outliers

✓ max. sum. of signed dist \rightarrow outlier

{ easy to differentiate
prob. interpretation

$$\sigma(\alpha)$$

sigmoid

tapering linear

prob.

max. sum. of transformed signed dist

Deriving

log. regn.

using

geometry

✓ detail

probability \rightarrow brief, pointers

loss-fn. based

✓ detail



Desiring

log. regn.

using

geometry

✓ detail

probability \rightarrow brief, pointers

loss-fn. based

✓ detail



monotonic fn:- $g(x)$

$x \uparrow; g(x) \uparrow \rightarrow$ monotonically incr-fn

if $x_1 > x_2$ then $g(x_1) > g(x_2)$
then $g(x)$ is said to be mon. incr. fn



Contents - Google Docs x NBayesLogReg.pdf x plot(x^2) - Google Search x plot(log(x)) - Google Search x plot(log(x^2)) - Google Search x Chekuri Srikan...

Secure | https://www.google.co.in/search?q=plot(log(x))&oq=plot(log(x)&aqs=chrome.0.69i59j69i57j0l4.3944j0j7&sourceid=chrome&ie=UTF... :

Google plot(log(x))

All Images Maps Videos News More Settings Tools

About 44,10,00,000 results (0.44 seconds)

Graph for $\log(x)$

x: -0.459316465 y: UNDEFINED

More info

Semilogarithmic plot - MATLAB semilogx - MathWorks

opt. prob:

$$x^* = \arg \min_x x^2$$

best(x)

Mimima & maxima

$$f(x) = x^2$$

$$x^* = \arg \min_x f(x)$$



mon. incr

$$g(x) = \underline{\log(x)}$$

optimization

$$\textcircled{1} \rightarrow \boxed{0} = \textcircled{x^*} = \arg \min_{\underline{x}} f(x) \quad ; f(x) = \underline{x^L}$$

$$\textcircled{2} \rightarrow x' = \arg \min_{\underline{x}} g(f(x))$$

$$\log(x^2)$$

claim: $\textcircled{x=x'}$; $g(x)$ is a monotonic fn.



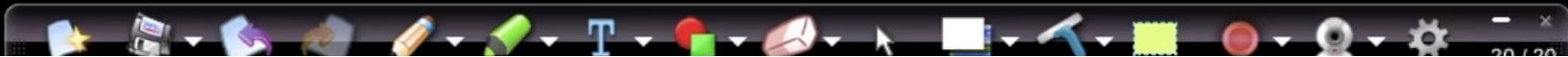
if $g(x)$ is a

monotonic fn.

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

$$\arg \min_x f(\underline{x}) = \arg \min_x g(f(\underline{x}))$$

$$\arg \max_x f(\underline{x}) = \arg \max_x g(f(\underline{x}))$$



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

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

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



$$\arg \max_x f(x) = \arg \min_x -f(x)$$

$$\log(e^x) = x$$

$$\arg \max_x -f(x) = \arg \min_x f(x)$$

opt L2

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

Signed-dist



$$\checkmark \arg \min_w \sum_{i=1}^n \log \left(1 + \exp(-y_i w^\top x_i) \right) \xrightarrow{\text{opt-LR}}$$

$$\arg \min_w \sum -y_i w^\top x_i$$

$$\times \arg \max_w \sum y_i w^\top x_i$$

Sum of signed hist
huge outlier problem

geometry

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

Prob methods

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

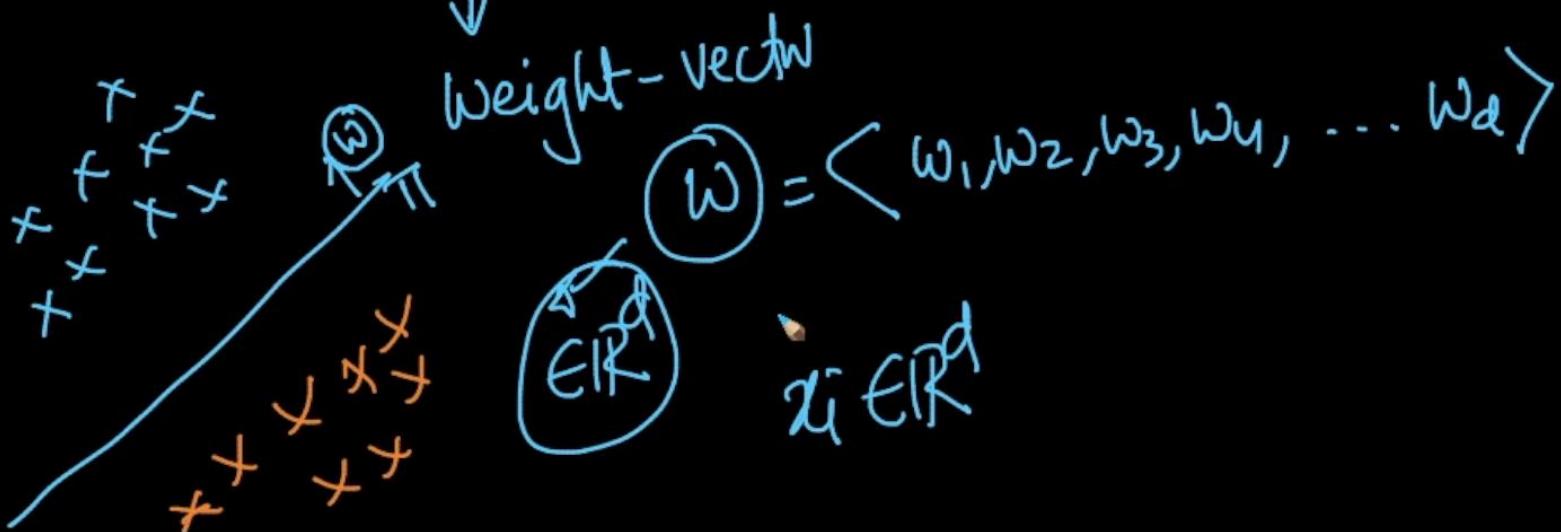
$$p_i = \sigma(w^\top x_i)$$

loss fn



Weight - Vector

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



$$\omega = \langle \omega_1, \omega_2, \omega_3, \dots, \omega_d \rangle$$

decision:

$$x_w \rightarrow y_w \quad \left\{ \begin{array}{l} \text{if } \omega^T x_w > 0 \\ \text{if } \omega^T x_w < 0 \end{array} \right. \quad \begin{array}{l} \text{then } y_w = +1 \\ \text{then } y_w = -1 \end{array}$$

prob.interp:

$$\overline{\sigma}(\omega^T x_w) = P(y_w = +1)$$



interpretation of w :

Case (1)

If $w_i = +ve$, $x_{qj} \uparrow \Rightarrow (w_i x_{qj}) \uparrow$

$$\Rightarrow \left(\sum_{i=1}^k w_i x_{qi} \right) \uparrow$$
$$\Rightarrow \sigma(w^T x_q) \uparrow$$
$$\Rightarrow P(y_q = +1) \uparrow$$

case 2: if $w_i = -\sqrt{c}$

$$f_i$$

$$w_i = -\sqrt{c}$$

$$x_{wi}$$

$$(w_i x_{wi}) \downarrow$$

$$\left(\sum_{i=1}^a w_i x_{wi} \right) \downarrow$$

$$\Rightarrow \sigma(w^T x_w) \downarrow$$

$$\Rightarrow p(y_w = +1) \downarrow$$

$$\Rightarrow p(y_w = -1) \uparrow$$

$$p(y_w = +1) = 1 - p(y_w = -1)$$



Applied | Applied | (70) You | Guadac | (70) SEE | New Tab | (70) You | Applied | DDR 60 | The Boy | Process | + | - | X

appliedaicourse.com/lecture/11/applied-machine-learning-online-course/3014/weight-vector/3/module-3-foundations-of-natural-language-processing-and... | 1 ABP | 51 SBI Bank PO Exa... | 7 Steps to Masterin... | Deep learning work...

Sathvikchiramana 6 Votes

How can we say that if $\text{sigmoid}(W^T \cdot X_q)$ increases then $p(y=+1)$ increases....because $y^*W^T \cdot X_q$ is high for all the correctly classified points so is the $\text{sigmoid}(y^*W^T \cdot X_q)$ so for all correctly classified points $\text{sigmoid}(y^*W^T \cdot X_q)$ is high i.e for both +ve and -ve points..how can we say that if $\text{sigmoid}(y^*W^T \cdot X_q)$ is high then only $p(y=+1)$ is high but not $p(y=-1)$?

Reply    

Jun 02, 2019 18:37 PM

Applied AI Tech Admin

It is because in the problem formulation we have assumed that If the value of product $y_i^*w^T \cdot x_i$ is positive, then the point is correctly classified. Otherwise it is not correctly classified. But here along with checking whether the data points are classified correctly or not, we also have to ensure our model predicts a label for those data points. Hence we need a function to estimate the labels. For this purpose we are taking the sigmoid function into consideration. The possible range of values for result of sigmoid function are [0,1]. Hence we are using a threshold of 0.5. If $\text{sigmoid}(y_q^*W^T \cdot X_q) > 0.5$, then we predict the label as 1. Otherwise we predict the label as 0. Also we have chosen sigmoid as it has got nice probabilistic interpretation and also it is differentiable easily.

If $\text{sigmoid}(y_q^*W^T \cdot X_q) = 0.6$, then we can say $y_q = 1$ with a low confidence as 0.6 is nearer to the threshold 0.5. Whereas if $\text{sigmoid}(y_q^*W^T \cdot X_q) = 1$, we can say $y_q = 1$ with high confidence as 1 is far away from the threshold 0.5 and on the positive side.

If $\text{sigmoid}(y_q^*W^T \cdot X_q) = 0.4$, then we can say $y_q = 0$ with a low confidence as 0.4 is nearer to the threshold 0.5.

Whereas if $\text{sigmoid}(y_q^*W^T \cdot X_q) = 0$ we can say $y_q = 0$ with high confidence as 0 is far away from the threshold 0.5 and on the negative side.

On the basis of these two reasons, we can say that

- As the value of $\text{sigmoid}(y_q^*W^T \cdot X_q)$ increases, the value of $P(y_q=1|x_q)$ increases and $P(y_q=0|x_q)$ decreases

Solving Optimization Problems

Interview Questions on Logistic Regression and Linear Regression

Type here to search | O | Microsoft Edge | e | Google Chrome | C | 23:00 | ENG | 16-08-2019

Screenshot of a browser window showing multiple tabs open, including one for an applied machine learning course.

Address bar: appliedaicourse.com/lecture/11/applied-machine-learning-online-course/3014/weight-vector/3/module-3-foundations-of-natural-language-processing-and...

Open tabs include: Applied, Applied, (70) You, Guapdat, (70) SEE, New Tab, (70) You, Applied, DDR 60, The Boy, Process, Apps, Empire, Why, Where And H..., Solutions to Machi..., The different uses o..., Stranger Things: S1..., 51 SBI Bank PO Exa..., 7 Steps to Masterin..., Deep learning work...

 Applied AI Tech Admin

It is because in the problem formulation we have assumed that If the value of product $y_i^T w - T^T x_i$ is positive, then the point is correctly classified. Otherwise it is not correctly classified. But here along with checking whether the data points are classified correctly or not, we also have to ensure our model predicts a label for those data points. Hence we need a function to estimate the labels. For this purpose we are taking the sigmoid function into consideration. The possible range of values for result of sigmoid function are [0,1]. Hence we are using a threshold of 0.5. If sigmoid($y_q^T W^T X_q$) > 0.5, then we predict the label as 1. Otherwise we predict the label as 0. Also we have chosen sigmoid as it has got nice probabilistic interpretation and also it is differentiable easily.

If sigmoid($y_q^T W^T X_q$) = 0.6, then we can say $y_q = 1$ with a low confidence as 0.6 is nearer to the threshold 0.5. Whereas if sigmoid($y_q^T W^T X_q$) = 1, we can say $y_q = 1$ with high confidence as 1 is far away from the threshold 0.5 and on the positive side.

If sigmoid($y_q^T W^T X_q$) = 0.4, then we can say $y_q = 0$ with a low confidence as 0.4 is nearer to the threshold 0.5.

Whereas if sigmoid($y_q^T W^T X_q$) = 0 we can say $y_q = 0$ with high confidence as 0 is far away from the threshold 0.5 and on the negative side.

On the basis of these two reasons, we can say that

- As the value of sigmoid($y_q^T W^T X_q$) increases, the value of $P(y_q=1|x_q)$ increases and $P(y_q=0|x_q)$ decreases.
- As the value of sigmoid($y_q^T W^T X_q$) decreases, the value of $P(y_q=0|x_q)$ increases and $P(y_q=1|x_q)$ decreases.

 Reply 

Jun 02, 2019 21:27 PM

 Mr. Chauhan

Nice explanation thanx!!



Type here to search



23:00 16-08-2019 ENG

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 $\boxed{z_i = y_i w^T x_i}$

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

$$\sum_{i=1}^n \log \left(1 + \underbrace{\exp(-z_i)}_{\geq 0} \right)$$

• $\exp(-z_i) > 0$

$$\log \left(1 + \underbrace{\exp(-z_i)}_{\geq 0} \right) > 0$$

$$\log(1) = 0$$

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

$$\left\{ \begin{array}{l} \log(1+\delta) > \log(1) \\ \delta > 0 \end{array} \right\}$$

Contents - Google Drive X NBayesLogReg.pdf X cs229-notes1.dvi X Multicollinearity - Wikipedia X plot(exp(-x)) - Google Search X plot(log(1+x)) - Google Search X Chekuri Srikan...

Secure | https://www.google.co.in/search?q=plot(exp(-x))&oq=plot(exp(-x))&aqs=chrome..69i57j0l5.4686j0j4&sourceid=chrome&ie=UTF-8

Google plot(exp(-x))

All Maps Images Videos News More Settings Tools

About 48,80,00,000 results (0.46 seconds)

Graph for $\exp(-x)$

$x: -2.93308038 \quad y: 18.7854076$

$z_i \rightarrow \infty$

$\exp(-z_i) \rightarrow 0$

More info

Graphing Exponential Functions: More Examples - Purplemath

$$w^* = \arg \min_w \left[\sum_{i=1}^n \log(1 + \exp(-z_i)) \right] \geq 0$$

minimal value of $\sum_{i=1}^n \log(1 + \exp(-z_i))$ is "0"

[it occurs when $z_i \Rightarrow \infty$ for all i]

$$\omega^* = \arg \min_{\omega} \sum_{i=1}^n \log(1 + \exp(-z_i))$$

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

then $\exp(-z) \rightarrow 0$

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

$$z_i = \underbrace{y_i}_{\text{+ve}} \underbrace{(w^\top x_i)}_{\text{-ve}}$$

$$\mathcal{D} = \left\{ \langle x_i, y_i \rangle : i=1 \dots n \right\}$$

Correctly-classified

$$z_i \rightarrow +\infty$$

\hookrightarrow modify my w in such a way that each $z_i \rightarrow +\infty$

①

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

②

$$z_i \rightarrow +\infty;$$

- if I pick my ω s.t $n \rightarrow \infty$

(a) all ^{training} points are correctly classified

outliers

(b) $z_i \rightarrow +\infty$

then best ω

$$\omega = [w_1 \ w_2 \ w_3]$$

① overfitting

perfect job in training data

② $w_i \rightarrow +\infty$
 $n(w_i) \rightarrow -\infty$

$w_i^i \rightarrow \infty$

$\sum w_i^i \rightarrow \infty$

One key aspect

minima $\rightarrow 0$

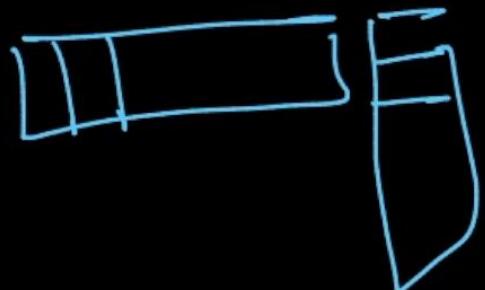
w_i is a normal

$$w^T w = 1$$

regularization

$$\hat{w}^t = \arg \min_w \sum_{i=1}^N \log(1 + \exp(-y_i w^T x_i)) + \lambda \|w\|_2^2$$

$$\begin{cases} w_j \rightarrow \infty \\ w_j \rightarrow -\infty \end{cases}$$



$$\lambda \sum_{j=1}^d w_j^2$$

regularization

$$\hat{w}^t = \arg \min_w \left(\sum_{i=1}^N \log (1 + \exp(-y_i w^T x_i)) + \lambda \|w\|_2^2 \right)$$

loss-term ✓

$w_j \rightarrow \infty$

$w_j \rightarrow -\infty$

regularization term

$$\text{MIN} \leftarrow \sum_{i=1}^n \log(1 + \exp(-z_i)) + \lambda \sum_{j=1}^d w_j^2$$

loss-term

"0"

$z_i \rightarrow +\infty$

reg-term

$w_j \rightarrow 0$

V.V.large

The diagram illustrates a loss function and its regularization term. The loss function is represented by a sum of terms, each involving a logarithmic function and a sigmoid-like term. The regularization term is a sum of squared weights. A 'loss-term' box encloses the first part, and a 'reg-term' box encloses the second. A 'MIN' circle points to the start of the sum. A '0' circle points to the regularization term. A 'V.V.large' circle points to the regularization term. A separate circle shows weights approaching zero.

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

λ : hyperparam in LR

$\begin{cases} \lambda = 0 \Rightarrow \text{overfit to the training side} \\ \lambda = \text{large} \Rightarrow \text{underfit} \end{cases}$

hyper param

K-NN : K

Laplace Smoothing (NB)

α : CV

$$\min \left(\text{loss-fn over training data} + \text{reg. f.} \right)$$

→ cross-validation

$\{ \begin{array}{l} D=0 \Rightarrow \text{overfit} \rightarrow \text{high variance} \\ D=\text{v. large} \Rightarrow \text{underfit} \rightarrow \text{high bias} \end{array} \}$

L_1 regularization and Sparsity

$$\hat{w} = \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i \tilde{w}^T \tilde{x}_i)) + \lambda \|\tilde{w}\|_2$$

$\underbrace{\hspace{10em}}$ logistic-loss $\underbrace{\hspace{2em}}$ $L_2\text{-reg.}$

$\tilde{z} \rightarrow +\infty$
 $\tilde{w}_i \rightarrow +\infty$
 $\tilde{w}_i \rightarrow -\infty$
overfit

Question: alternatives to $L_2\text{-reg.}$

popular:- $L_1\text{-reg}$

$\|w\|_2^2$ for reg

$w_i \rightarrow +\infty$
or $w_i \rightarrow -\infty$

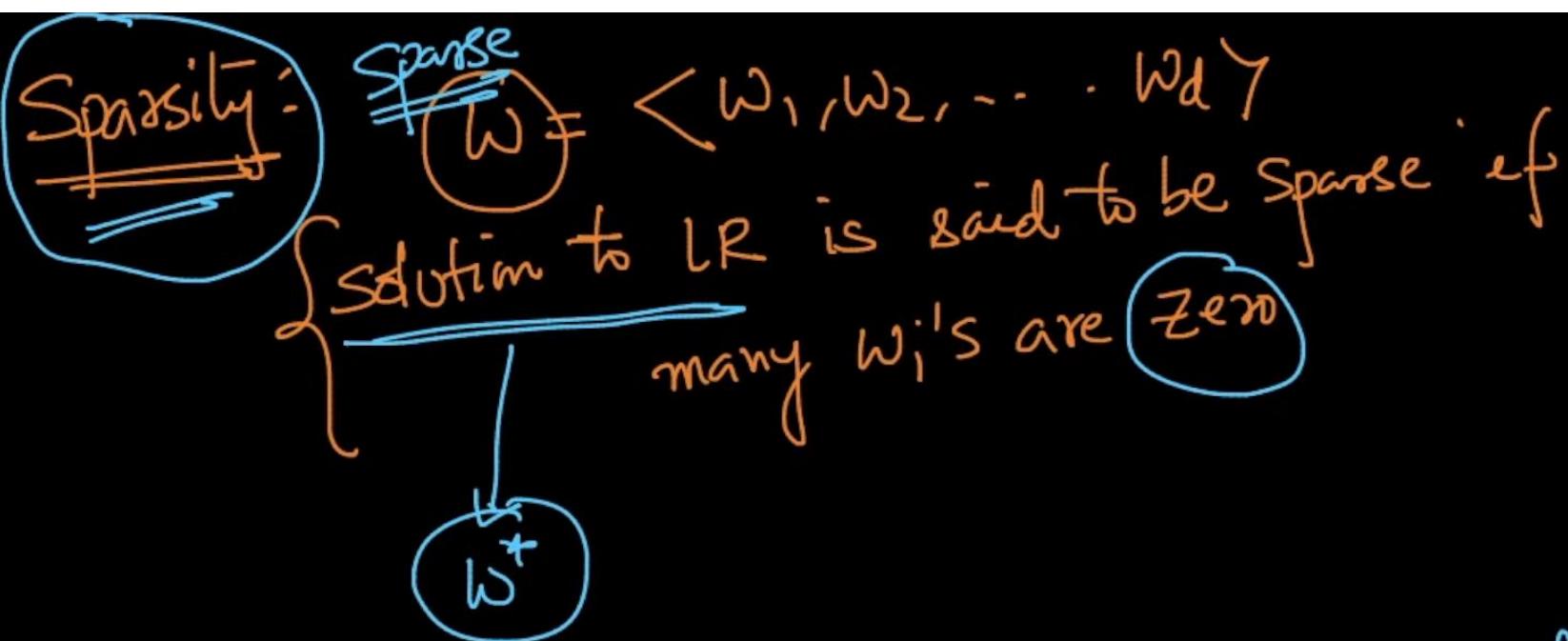
\downarrow
 $\|w\|_1$ for reg.

$$\|w\|_1 = \sum_{i=1}^d |w_i|_{\text{abs}}$$

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

$\checkmark L_1\text{-reg}$

will avoid $w_i \rightarrow +\infty$ or $-\infty$



If we use L_1 reg in LR, all the unimportant or
 less-important features ~~wi's~~ become zero
 =


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

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

\downarrow

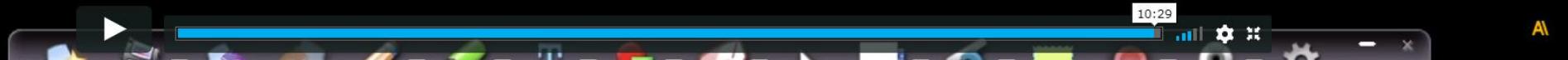
zero - if f_i is used
if L_2 reg is used; w_i becomes small values but not necessarily zero

elastic-net:

either L₁ or L₂

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

Two hyper-param: - λ_1 ~~λ_2~~ λ_2



Probabilistic interpretation of Logistic regression

(Logistic regression) ① Geometry & simple algebra ✓

 ↳ ② Probability

 ↳ ③ loss minimization

(Textbook)

prob. derivation of (LR)

Naive Bayes

$$p(x_i | y_i) \sim N(\mu_i, \sigma_i)$$

→ ① features are real valued:-

↳ Gaussian dist

→ ② $y_i \in \{0, 1\}$
Bernoulli r.v (coin-toss)

~~geom:~~

$$\hat{w}^t = \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \frac{\gamma}{2} \|w\|^2$$

$\Rightarrow \checkmark$

$+1 \text{ or } -1$

~~prob:~~

$$\hat{w}^t = \arg \min_w \sum_{i=1}^n -y_i \log p_i - (1-y_i) \log(1-p_i) + \frac{\gamma}{2} \|w\|^2$$

$\Rightarrow \checkmark$

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

$+1 \text{ or } 0$

Case 1:

y_i : +ve

geom:- $y_i = +1$
prob:- $y_i^- = +1$

geom:- $\log(1 + \exp(-w^T x_i))$
prob:- $\circled{-1} \cdot \log\left(\frac{1}{1 + \exp(-w^T x_i)}\right)$
 $= \log\left(1 + \exp(-w^T x_i)\right)$

$$p_i^- = \sigma(w^T x_i)$$
$$\boxed{-\log(x) = \log(1/x)}$$

Case 2:

$$y_i = -ve_j$$

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

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

geom:

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

prob:

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

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

$$\begin{aligned}
 & -\log(x) \\
 & = \log(1/x)
 \end{aligned}$$



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

divide numerator & den by $\exp(-w^T x_i)$

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

$$\boxed{\frac{1}{e^{-x}} = e^x}$$

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

Loss minimization interpretation of LR

✓ geometric & probabilistic

loss-minimization

$$\left\{ \begin{array}{l} f(x_i) = \underline{\omega^T x_i} \end{array} \right.$$

$$\left\{ \begin{array}{l} \omega^T = \arg \min_{\omega} \sum_{i=1}^n \log \left(1 + \exp(-y_i \omega^T x_i) \right) \end{array} \right.$$

$$\underline{z_i} = y_i \omega^T x_i = \underline{y_i \cdot f(x_i)}$$



ideal - optimzn. model:

classfn: $w^* = \min_w$ loss-fn

0-1 lossfn.



$$\underline{z_i = y_i \cdot w^\top x_i}$$

number of
incorrectly classified
pls

ideal - loss - function

+1 = incorrectly classified

0 = correctly classified

(0-1 loss-fn)

incorrectly
classified

correctly classified

min. loss ✓
Max. profit ✓

$$z_i = \underbrace{y_i \cdot w^\top x_i}_{y_i \cdot f(x_i)}$$



Ideal

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

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



solve optimization problems in ML

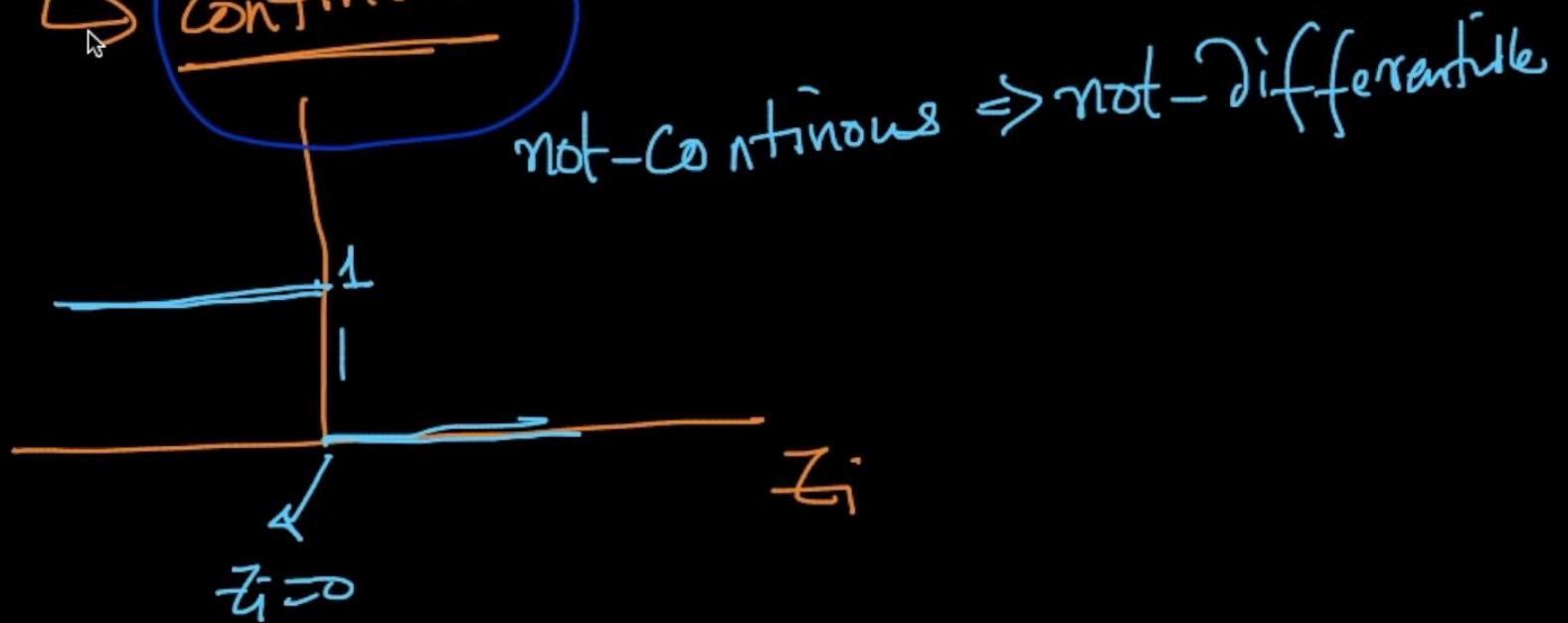
↳ "differentiation" in Calculus.

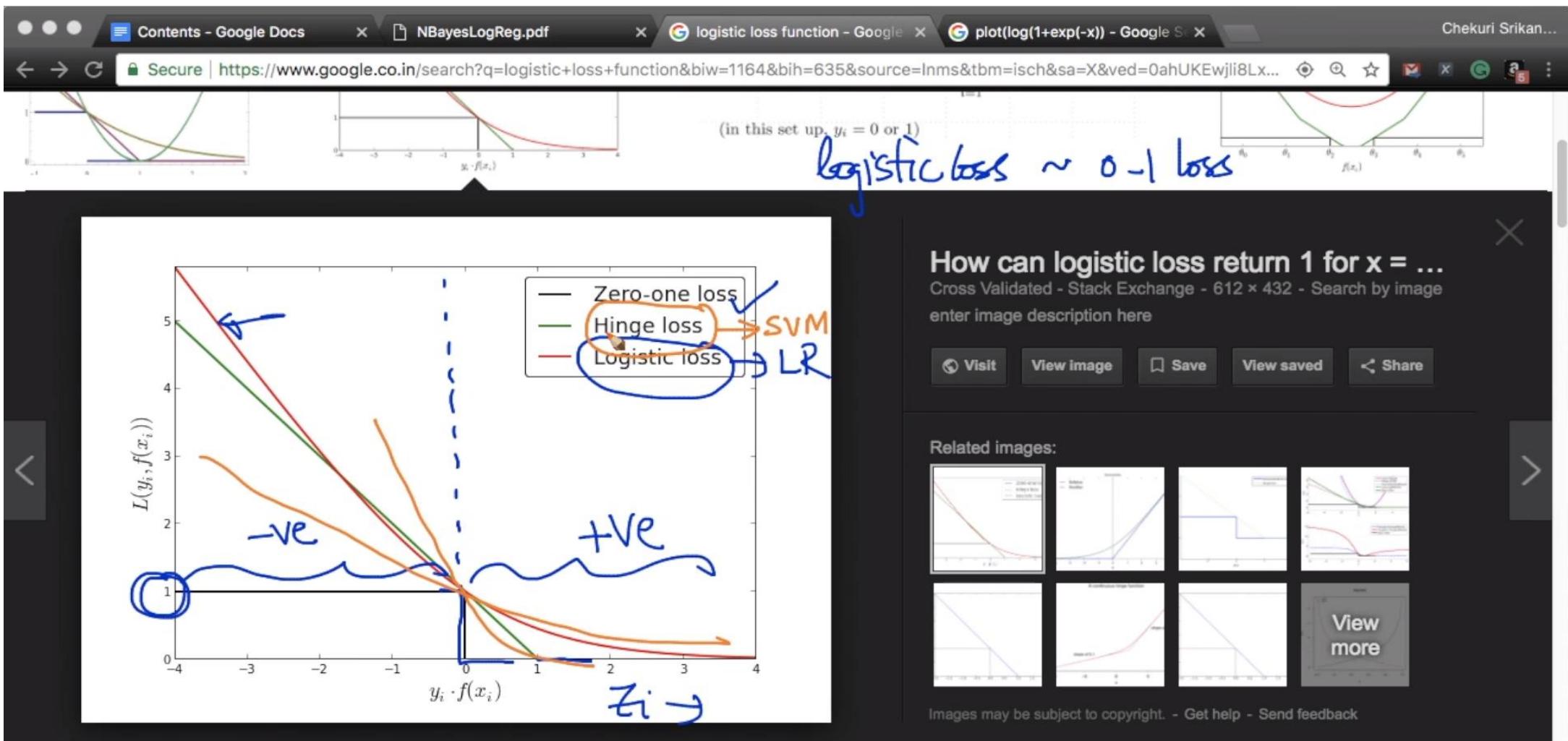
↳ next chapter



$f(x)$ is differentiable

↳ "continuous"





loss-minimization interpretation

↳ (loss-fn) - logistic loss \rightarrow LR

hinge loss \rightarrow SUM

exp-loss \rightarrow AdaBoost

Sq-loss \rightarrow linear Reg



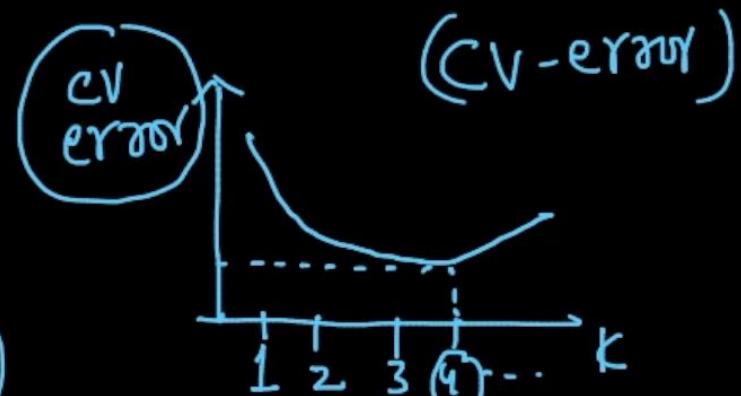
Hyperparameter search optimization

λ : hyperparam

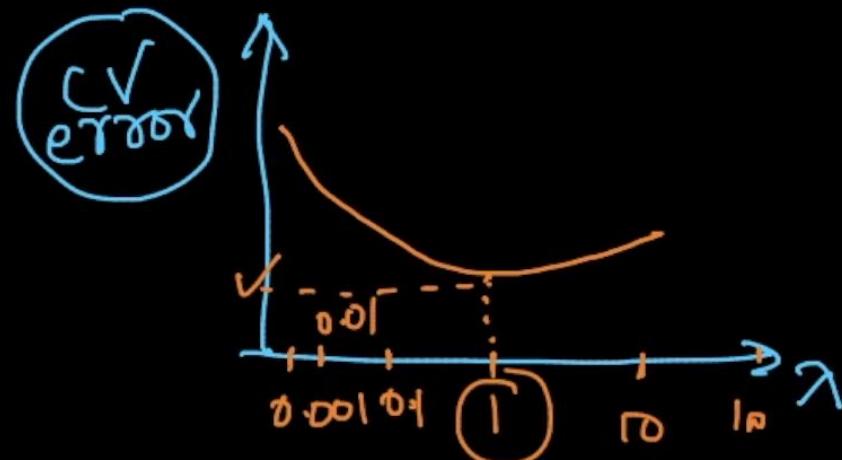
$\lambda = 0 \Rightarrow$ overfitting
 $\lambda = \infty \Rightarrow$ underfitting

(Q) how we determine the best λ

- λ : LR
- K : K-NN
- α : NB (Laplace Smoothing)



Grid Search (Brute force)



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

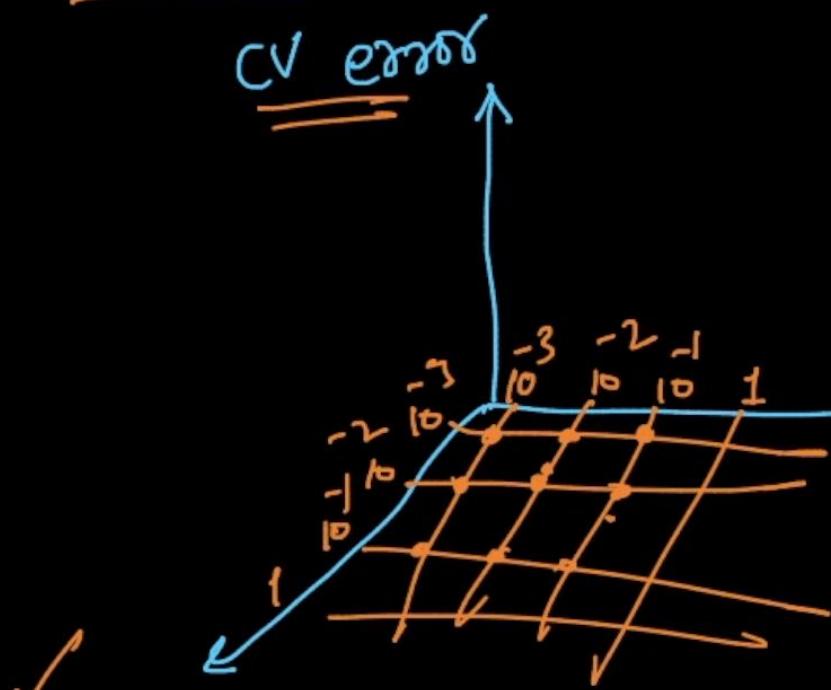
② $\lambda = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]$

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

~~2.2 x 10^2~~

elastic net:-

$$\underline{\lambda_1} \|\underline{w}\|_1 + \underline{\lambda_2} \|\underline{w}\|_2^2$$



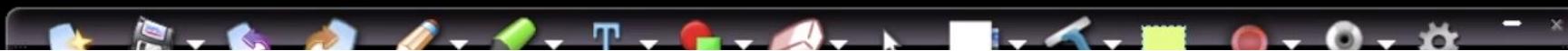
$$\checkmark \lambda_2 = [10^{-3}, 10^{-2}, \dots, 10^3]$$

m_2

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

m_1

$$m_1 \times m_2$$



Grid Search

$$m_1 = m_2 = m_3 = \dots = m$$

deep learning

λ_1 : 1 hyperparam. - m_1 times train

λ_1, λ_2 : 2 hyperparams - $m_1 \times m_2$

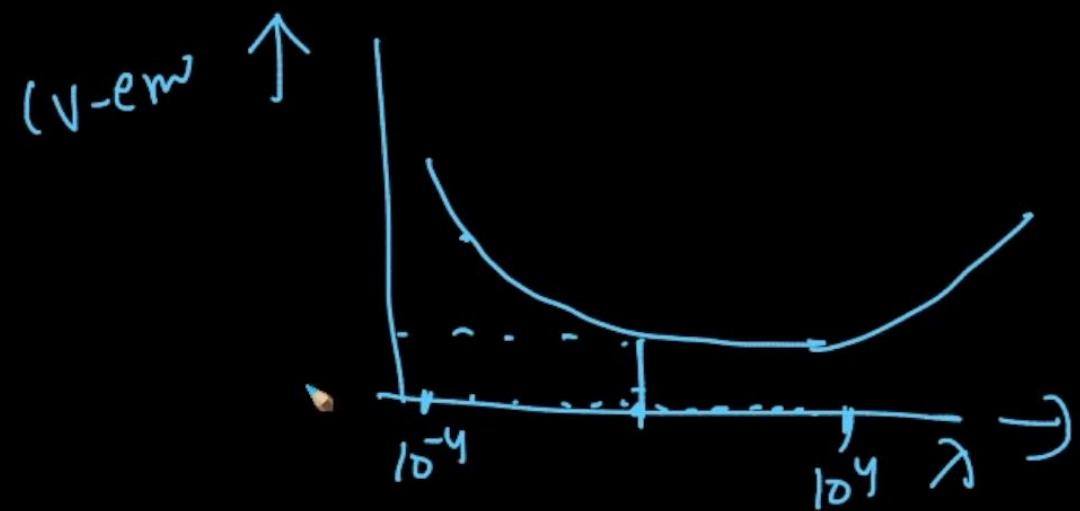
$\lambda_1, \lambda_2, \lambda_3$: -
K-hyperparam: - m^k

$$m_1 \times m_2 \times m_3 \quad m^3$$

as # hyperparam increase, the # times
the model needs to be trained increases
exponentially. $K \uparrow$; $m^K \uparrow$

"Random" Search: \rightarrow is almost as good as grid-search
if hyperparam is large

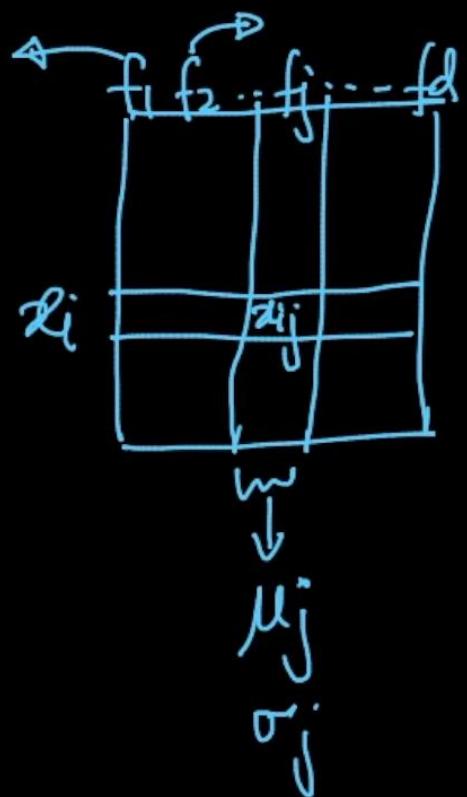
$\lambda \in [10^{-4}, 10^4]$ ← randomly pick values
in the given interval



λ : → hyperparam Search/optimization
↳ Grid Search ($10^{-4}, 10^{-3}, 10^{-2}, \dots$)
↳ Random Search
Sklearn



Column | feature standardization



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

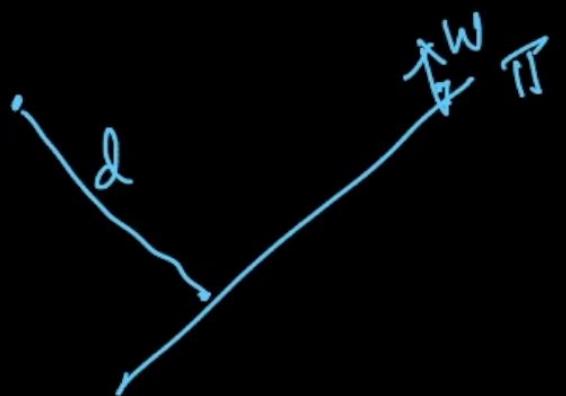
k-NN:- distances
(standardize our feature)

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

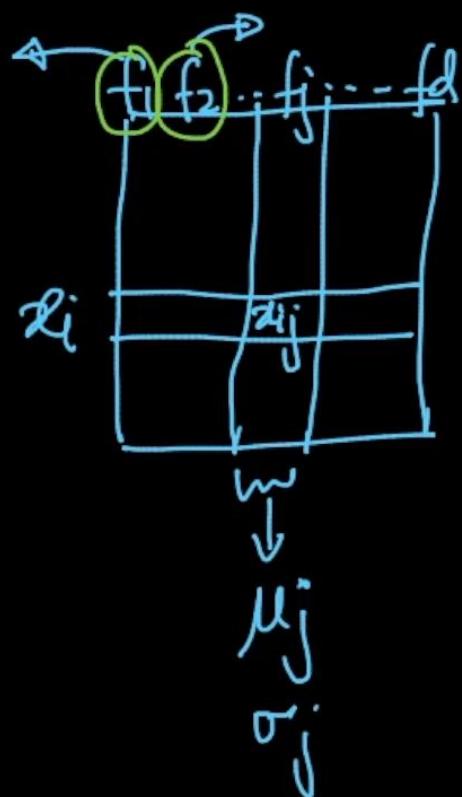


Logistic regression :-

- "mandatory" to perform feature standardization before training on your \mathcal{D}



Column | feature standardization



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

$$\left\{ \begin{array}{l} x_{ij}' = \frac{x_{ij} - \mu_j}{\sigma_j} \\ \text{centering} \\ \text{scaling} \end{array} \right\} \text{: standardization}$$

✓ k-NN; - distances

Standardize our feature

mean-centering
scaling

Logistic regression :-

"mandatory" to perform
feature standardization
before training on
your \mathcal{D}



Feature importance & Model Interpretability

$f_1, f_2, \dots, f_i, \dots, f_d$
 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $w_1, w_2, \dots, w_i, \dots, w_d$

LR : GNB + Bernoulli

assume: if all features are independent (Naive Bayes)

feature-imp:

w_i^*

$\|w_j\|$ = absolute value of weight
corresponding to f_j

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

Case 1: $w_j = +ve$ & large;

$$\sum_{j=1}^d w_j \cdot x_{qj} \uparrow$$

$w^T x_q$

$P(y_q = +)$

Case 2: w_j : -ve & large;

$$\sum_{j=1}^d w_j \cdot x_{qj} \uparrow$$

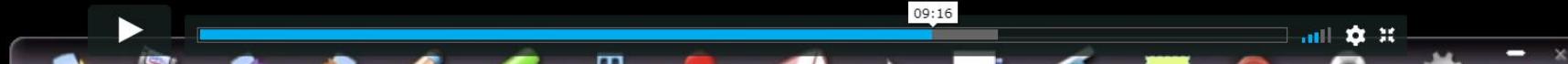
$P(y_q = -)$

determine the important features in LR

e.g: predict gender : male & female
 w_j

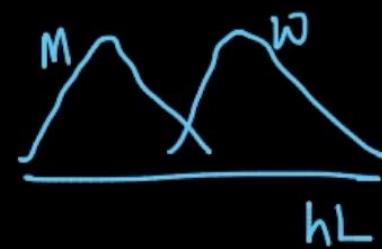
w_{hL} is large
 $w_{hL} \uparrow ; P(y_{qj} = -1) \uparrow$

w_{hL} : -ve
large negative weight value
 w_{hL} hair length

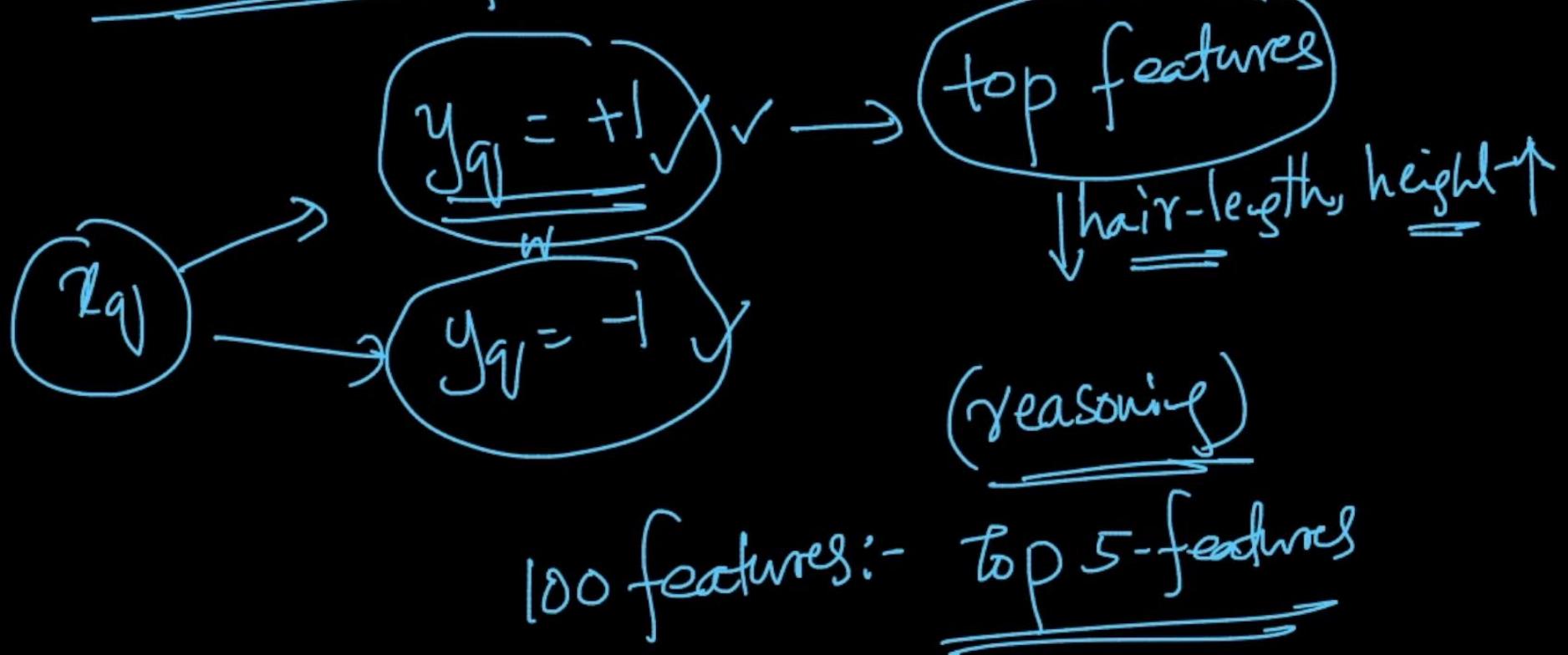


$\{h\}$ height ↑ ; $p(Y_g = +1) \uparrow$
~~male~~

$w_h : +ve$
↳ medium · positive weight



Model interpretability



"Collinearity" (or) Multicollinearity

f.I :- features are indep; $|w_j|$ as f.I values

Collinearity:

\tilde{f}_i, \tilde{f}_j

$$\text{s.t. if, } \tilde{f}_i = \alpha \tilde{f}_j + \beta$$

then f_i & f_j are collinear

Multicollinearity: if f_1, f_2, f_3 & f_4 are

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

then f_1, f_2, f_3 & f_4 are said to be
multicollinear

(Q) Why does (w_j) not be useful as \overline{f} . If
features are collinear \Rightarrow

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

$$\omega^* = \langle \underbrace{1, 2, 3}_{f_1 f_2 f_3} \rangle ; \quad \chi_q = \langle x_{q1}, x_{q2}, x_{q3} \rangle$$

$$\omega^T \chi_q = x_{q1} + 2x_{q2} + 3x_{q3}$$

if $f_2 = 1.5f_1 \Rightarrow f_1 \text{ & } f_2 \text{ are collinear}$

$$\omega^T \chi_q = \cancel{x_{q1}} + 3x_{q1} + 3x_{q3} = 4x_{q1} + 3x_{q3} \\ = \langle 4, 0, 3 \rangle$$

$\checkmark \vec{w}^t = \langle 1, 2, \boxed{3} \rangle$ f_3 is the most imp

$\checkmark \vec{w} \approx \langle 4, 0, 3 \rangle$ f_1 is the most imp
not all imp

Same classifier $\therefore f_2 = 1.5 f_1$

$x_a \rightarrow \vec{w}^T \vec{x}_a$

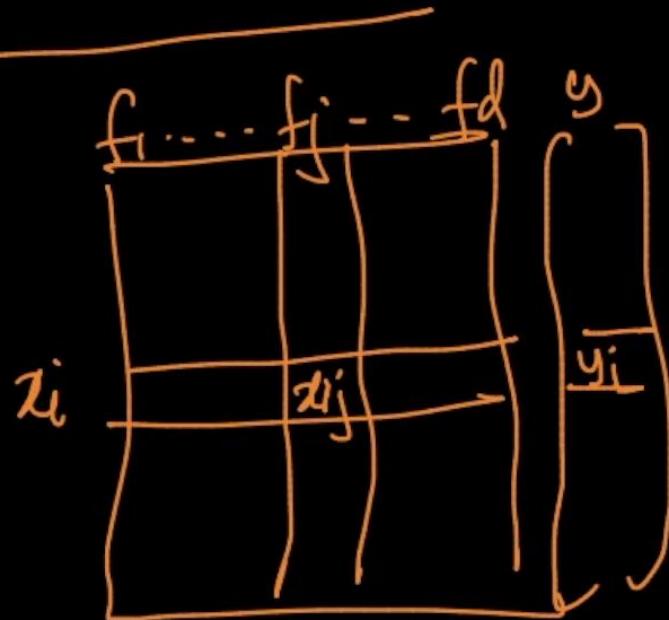
$\{w_j\}$ as $f \cdot I$

determine. If features are multicollinear:

→ perturbation technique:

$$x_{ij} + \epsilon$$

Small noise
 $N(0, \frac{1}{n} I)$



before perturbation: $\omega = \langle \omega_1, \omega_2, \dots, \omega_j, \dots, \omega_d \rangle$

after perturbation: $\tilde{\omega} = \langle \tilde{\omega}_1, \tilde{\omega}_2, \dots, \tilde{\omega}_j, \dots, \tilde{\omega}_d \rangle$

If $\omega_i \neq \tilde{\omega}_i$ significantly
then your features are collinear

↓
Cannot use $|w_j|$'s as f. I

Train & Runtime Space & time complexity

$$n = \# \text{ pts in } D_{\text{train}}$$
$$d = \dim$$
$$\underline{\mathcal{O}(nd)}$$

Train LR:-



$$\underline{w^t = \operatorname{argmin}_w (\text{logistic loss} + \text{reg})}$$

\uparrow next chapter optimization

$$\underline{w^t = \langle w_1, w_2, \dots, w_d \rangle}$$

Run Time (LR)

$$\left\{ \begin{array}{l} \text{Space: } \mathcal{O}(d) \\ \text{Time: } \mathcal{O}(d) \end{array} \right.$$

$$\sum_{j=1}^d w_j x_{qj}$$

$\left\{ \begin{array}{ll} w^T x_{qj} > 0 & \rightarrow \text{+ve} \\ w^T x_{qj} < 0 & \rightarrow \text{-ve} \end{array} \right.$



- if d is small

$$d=30$$

$\rightarrow \underline{LR}$ is v.v-good for low-latency applications

$$x_{ij} \rightarrow 1\text{ms} \quad y_{ij} \checkmark$$

favorite algo
internet
Computers

(30 multiplications
+ 29 additions)

\rightarrow memory efficient

$$w^p$$

$$x_{ij} \rightarrow \boxed{\quad} \rightarrow y_{ij}$$



if d is large

$$d \approx 1000$$

$\nwarrow \uparrow$; underfit

Bias van
✓

$w^T x_q$: - 1000 mult & addition

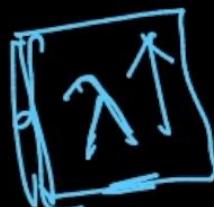
$\rightarrow L_1$ reg: - Sparsity (w_j 's corresponding to less important features = 0)

w^T :

0	0	0	1	0	...	0
1	2	3	...	1000		

\downarrow : reasonably

~~bias, var, latency~~



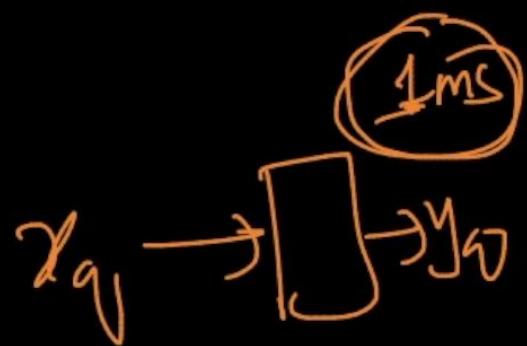
Sparsity \uparrow
more of $w_j = 0$

1 ms

mult & 50 adds

latency reg

Bias vs Latency



Cases:

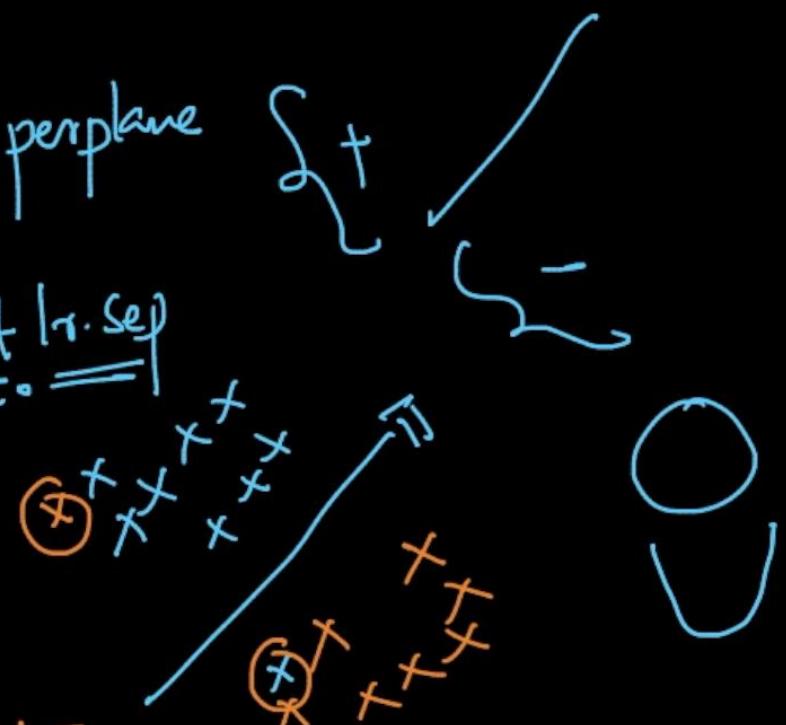
decision surface:- Linear / hyperplane

assumption :- data is linearly
separable (or) almost lin. sep

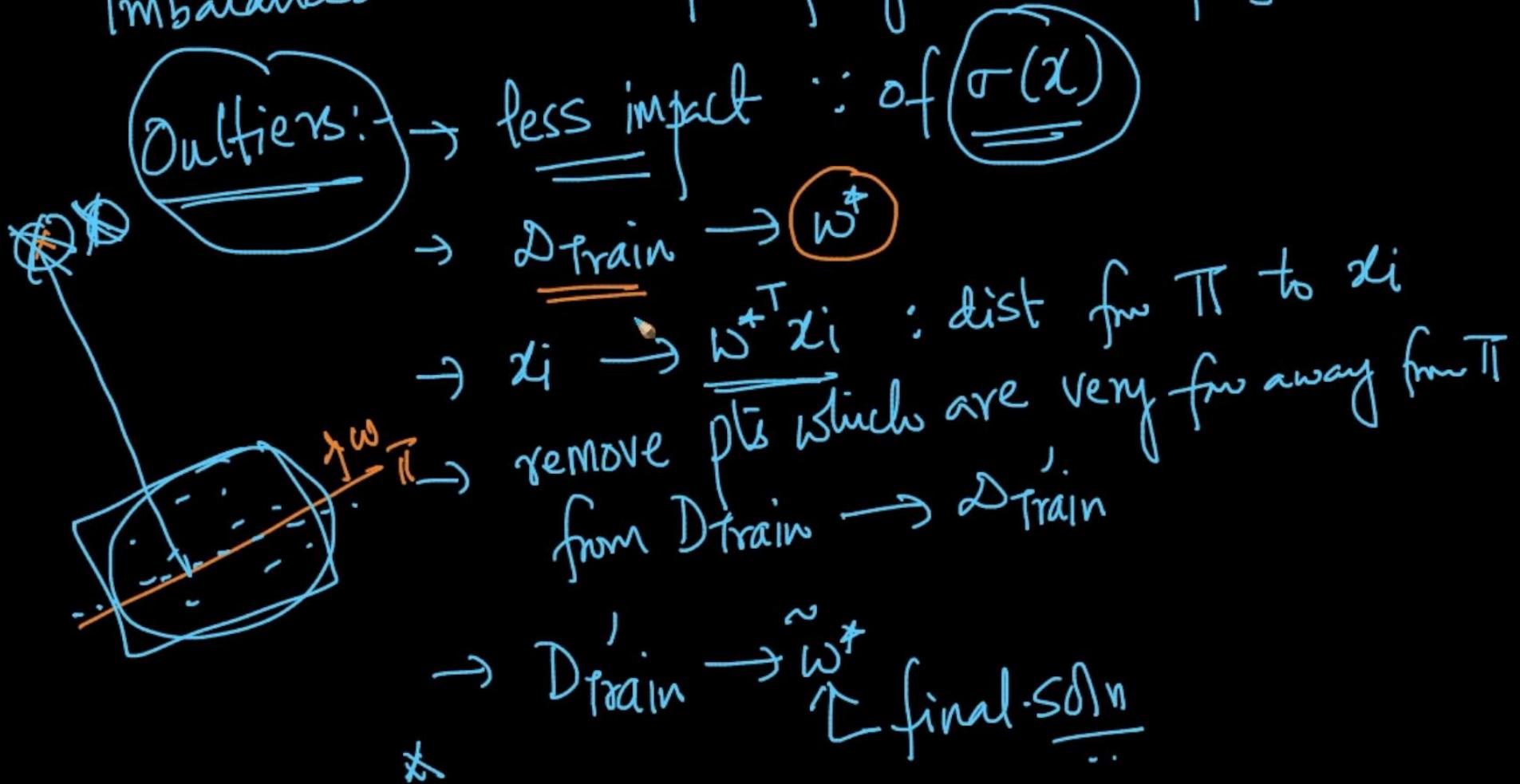
feature imp & Interpretability

$\sum |w_j|$ if features
are not multicollinear

$x_{qj} \rightarrow y_{qj} = +1$



Imbalanced data:- upsampling & down sampling



Missing values:- imputation → mean, median
→ model

Multi-class:

One vs Rest ← Typically
MaxEnt models ← extensions to LR
Softmax classifier → deep learning
Multinomial LR

Similarity Matrix:

extension to LR → Kernel LR
SVM

Best & Worst cases:- \rightarrow almost
separable
 \rightarrow low-latency requirement (L-reg)
 \rightarrow very fast to train

Large 'd':

d is large; $\left\{ \begin{array}{l} \text{chance data is lr-separable} \\ \text{is high} \end{array} \right.$
 \downarrow
low-latency \rightarrow L-reg

Logistic Regression with imbalanced data: geometric intuition

(Q) How does imbalanced data impact Logistic regression? Explain geometrically.

