(1) [ Regression ] or [ Classification ] or [ Generative model ]

given $\boxed{\bar{x}}$ predict $\hat{y}$
features
columns

given $\boxed{\bar{x}}$ predict $\hat{y}$
features
label +1 /-1

given weight & height → age

given movie genre
movie actors $\}$ → yes/no like

(2) · Choose parameters → $\theta$ · the dimension of $\theta$ = the dimension of $x$ (features).

determines the wgt. of each features (weightage).

$$\boxed{\hat{y} = \vec{\theta} \cdot \vec{x} + \theta_0 = \vec{\theta}^T \vec{x} + \theta_0}$$   bias  general eqn formula
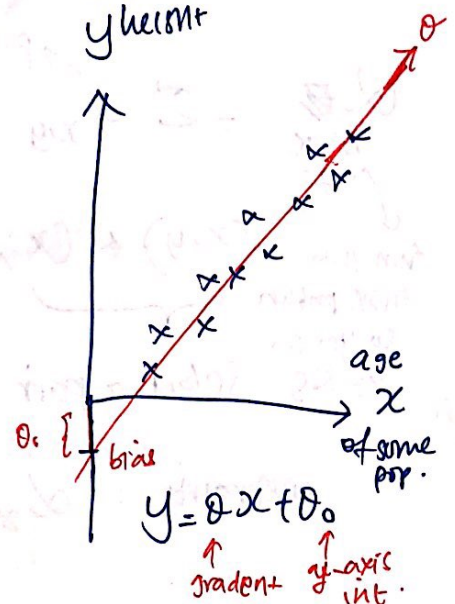
$\downarrow h(x,\theta)$   matrix mult

$\bar{x}$ :
genre . $\theta_1$
actors · $\theta_2$
time · $\theta_3$
country · $\theta_4$
theme · $\theta_5$
scriptwriter· $\theta_6$
director · $\theta_7$
$+$
$\theta_0$

$\hat{y} = $ yes /no like . $\}$ C .

$\hat{y} = $ Rating $\}$ R

(R)



y height
age
x
of some pop.
$\theta_0 \{$ bias
$y = \theta x + \theta_0$
↑ gradent   ↑ y-axis int.

(3) want to find out the $\boxed{best}$ $\theta$

$\hookrightarrow$ how to grade what is good $\theta$? or bad $\theta$?

↓

$\boxed{loss function}$ (1) $\min\limits E_1 = \sum\limits_{i=1}^{N}(y_i - \hat{y})^2 = \frac{1}{2N} \sum\limits_{i=1}^{N}(y_i - (\vec{\theta} \cdot \vec{x}_i + \theta_0))^2$

true value of y from train set for train loss

min least sq.

For training loss, use training set.
↳ test loss   use test set.

(unbounded) linear regression .

$(x_i, y_i)$ → data given

(2) $\min\limits E_2 = \sum\limits_{i=1}^{N} \max \{(1 - (\hat{y}(\theta \cdot x))), 0\}$
$(\theta = -1/+1)$
predicted y

perceptron

(C)



y
x
a good $\theta$ will give $\theta \cdot x$ as the same sign as y

they

(3) $\min E_3 = \frac{1}{n}\sum_{i=1}^{N} \log\left[\frac{1}{1+e^{-y(\theta \cdot x)}}\right]$ } logistic regression.
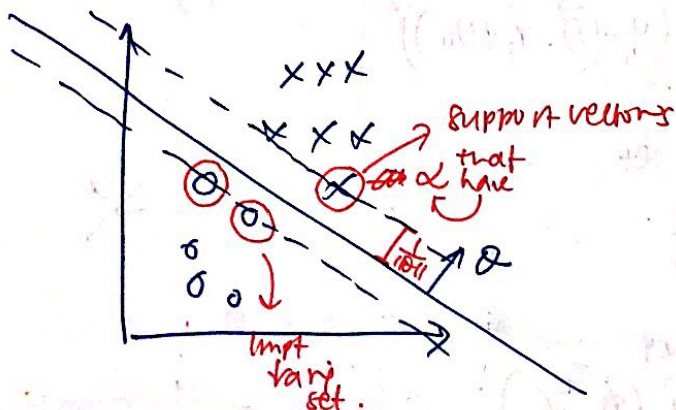
bounded b/w 0 to 1

$\frac{dE_3}{d\theta} \ldots$ ?

Rating.

(better).

Classifier

(4) $\min E_4 \Rightarrow$ (SVM), parameter is called $\alpha$.

dimension of $x < dOS$.

max margin because margin $= \frac{1}{\|\theta\|}$

error fn $\Rightarrow \min \boxed{\frac{1}{2}\|\theta\|^2} = E_4$.

$\min$ / constraint $y(\theta^T x) \geq 1$.

transform to max problem. cannot SD.

→ params     params

$\underset{max}{\alpha} = \sum \alpha_{x,y} - \frac{1}{2}\sum_{x,y}\sum_{x',y'} \alpha_{x,y}\alpha_{x',y'}\, y\, y'(x^T x')$.

turn into max problem so you can use $\theta$s

$(x,y)$ & $(x',y') \Rightarrow$ training set

select a pair

a number

$[x_1 \ldots x_d]\begin{bmatrix} x_1 \\ \vdots \\ x_d\end{bmatrix}$.

constraint: $\alpha_{x,y} \geq 0$.

If I have N items in training set, I have N $\alpha$s.

$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$

test / train features     features

$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & & \cdots & x_{2d} \\ & & & \\ x_{N1} & & \cdots & x_{Nd} \end{bmatrix}$

# of train/test data



support vectors
that have
and $\alpha$   have

$\frac{1}{\|\theta\|}$   $\alpha$

imp't var set.

$Y = \begin{bmatrix} -1 \\ \vdots \\ 1 \end{bmatrix}$ } N    $\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$

④ • After getting loss fn : <u>want to min loss.</u>
↓

differentiate the loss fn.

$\quad\quad$ ↳ ① Exact soln.

$$\frac{dE^{E_5}}{d\theta_1} \propto \left(y_i - \underbrace{\theta_1 x_1 + \theta_2 x_2 + \theta_0}_{}\right)^2$$

$$\frac{dB}{d\theta_1} = \downarrow \; 2x_1$$

Step ① $\quad \dfrac{dE}{d\theta}$ and $\dfrac{dE}{d\theta_0}$ } find.

the error of all train set. ←
$$\boxed{E_1} = \sum_{i=1}^{N} \left(y_i - (\theta \cdot x_i + \theta_0)\right)^2$$

the error of $\theta$ one ←
$$\frac{dE_1}{d\theta} = \sum_{i=1}^{N} 2\left(y_i - [\theta \cdot x_i + \theta_0]\right) \cdot -x_i$$

do the same for <u>ALL</u> parameters.

method ① <u>Exact soln</u> : equate $\dfrac{dE}{d\theta}$ s $= 0$ and find $\theta_s$.

method ② gradient descent :

$\quad\quad$ initialize all $\theta_s$ params to <u>some value</u>. (any value)
repeat for k iterations or until converge
$\quad\quad$ Check current Error → sub $\theta$ to $E(\theta, x, y)$ → sum all errors from TS.
$\quad\quad$ ↱ training sets (ALL)
$\quad\quad$ if $E > e$ → some small value → 0.0001
$\quad\quad\quad\quad$ (if still have error)
$\quad\quad\quad$ ↳ update $\theta$

$$\theta^{new} = \theta^{old} \;\textcolor{red}{\bigcirc{-}}\; \eta \boxed{\frac{dE}{d\theta}}$$

$\quad\quad$ minus of
$\quad\quad$ the error
$\quad\quad$ of this $\theta$ from
$\quad\quad$ the "old" $\theta$ value to
$\quad\quad$ make up the new $\theta$ value.

sub in $\theta$ old values and sub in $x, y$ from training set.

constant , Given, e.s: 0.01 or 0.05 ...

$\quad\quad$ else, (if no more error).
$\quad\quad\quad$ end.

⑤ Backprop (NN).

③ <u>Stochastic GD</u> → sub $\theta$ to $E(\theta, x_i, y_i)$ → choose 1 point, no <u>sum</u> (Random)

Fast ④ Quadratic solve if have constraint.

Conv. feature trick → to code nicer.

$$\vec{\theta} \cdot \vec{x} + \theta_0 = y$$

$\downarrow$

↗ $i^{th}$ data.

$$\begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \cdot \left.\begin{bmatrix} ^i x_1 \\ \vdots \\ ^i x_d \end{bmatrix}\right\} \text{dimension} + \theta_0 = \hat{y}$$

↙ fold

$$\underset{\vec{\theta}}{\begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \\ \theta_0 \end{bmatrix}} \cdot \underset{\vec{x}}{\left.\begin{bmatrix} x_1 \\ \vdots \\ x_d \\ 1 \end{bmatrix}\right\}} = \hat{y}$$

↖ append.

$\underbrace{\qquad\qquad}$

nicer to code.

$$np.dot(\theta, x)$$

$$X = \begin{bmatrix} x_{11} & \cdots & x_{1d} & 1 \\ \vdots & & & 1 \\ \vdots & & & 1 \\ x_{N1} & \cdots & x_{Nd} & 1 \end{bmatrix}$$

$$\vec{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \\ \theta_0 \to bias \end{bmatrix}$$

Scanned with CamScanner

$$\frac{dE_2}{d\theta} \begin{cases} 0 & \text{if} \quad y(\theta \cdot x) > 1 \\ -yx & \text{if} \quad y(\theta \cdot x) < 1. \end{cases} \Bigg\} \text{ correct prediction.}$$

$$\underbrace{\qquad\qquad}_{\text{wrong prediction}}$$

In GD, update $\theta$:

$$\theta^{new} \leftarrow \theta^{old} - \eta \boxed{\left[\frac{dE}{d\theta}\right]} \rightarrow \text{will be } 0 \text{ if } y(\theta^{old} \cdot x) > 1$$

---

**kernel** $\rightarrow$ transform from ~~one data~~ one vector space to another vector space? where In that new vector space, the data is linearly separable.

↳ gets the $x^T \cdot x$ in the new VS straight away :(cheat).



a circle drawn by in Cartesian VS ~~board~~.
  ↳ transform data..

$$(x_1, y_1)$$
$$(x_2, y_2)$$
$$\vdots$$
$$(x_N, y_N)$$

$$x^2 + y^2 = R^2$$

$$x_1^2 + y_1^2 = \boxed{R_1^2} \quad \text{my new data.}$$
$$\vdots$$
$$x_N^2 + y_N^2 = \boxed{R_N^2}$$

label is +ve -ve
same var.