# Knn:-

$L_p$

$$L_2 = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

$$L_1 = |x_1 - x_2| + |y_1 - y_2|$$

$$L_\infty = \max(|x_1 - x_2|, |y_1 - y_2|)$$

## Scaling

Standard :- $\dfrac{x_i - \mu}{\sigma}$

maximum $= \dfrac{x_i}{max}$

max-min $= \dfrac{x_i - min}{max - min}$

## is it metric?

$d_{p,p} = 0$              $d_{p,q} \geq 0$

$d_{p,q} = d_{q,p}$        $d_{p,q} + d_{q,r} \geq d_{p,r}$

# Linear Regression:-

### predict:

$$\hat{y} = w_1 x_1 + w_2 x_2 + b$$

### RSS:

$$Loss_i = (\hat{y}_i - y_i)^2$$

$$J = \frac{1}{2m} \sum (\hat{y} - y)^2$$

### Closed-form:

$$J = (Xw - y)^T (Xw - y)$$

$$w = (X^T X)^{-1} X^T y$$

### Gradiant descent:-

Step 1: initialize $w, b$

Step 2: calculate $\hat{y}$ and cost $\longrightarrow$ RSS

Step 3: calculate $dw, db$

step 4: update $w, b$

Repeat $2 \to 4$ until $|J^K - J^{K+1}| \leq \epsilon$

$$dw_i = \frac{1}{m} (\hat{y} - y) x_i$$

$$db = \frac{1}{m} (\hat{y} - y)$$

$$new \ w_i = w_i - l \ dw_i$$

$$new \ b = b - l \ db$$

underfitting: high Bais (simple)

- more complex $\left\{ \begin{array}{l} \text{add more features} \\ \text{polynomial} \end{array} \right.$

overfitting :- high Variance (complex)

- add more data

- Perform Regularization

$\hat{y} = 6_0 + 4_0 x_1 - 20 x_2$

$\|w\|_2 = \sqrt{4_0^2 + 2_0^2}$

$\|w\|_2^2 = 4_0^2 + 2_0^2$

high $\lambda$ $\Rightarrow$ under fitting
Low $\lambda$ $\Rightarrow$ overfitting

## Regularization :-

① Ridge :- add $\lambda \|w\|_2^2$

closed-form :~

$$J = (Xw - y)^T (Xw - y) + \lambda w^T w$$

$$w = (X^T X + \lambda I)^{-1} X^T y$$

Gradiant :-

$$J = \frac{1}{2m} \sum (\hat{y} - y)^2 + \lambda \|w\|_2^2$$

Steps :

Same as above

but: $dw_i = \frac{1}{m} \sum (\hat{y} - y) \cdot X_i + \lambda \|w\|_2^2$

feature Selection

② Lasso :- add $\lambda \|w\|_1$

$$Loss = \sum (\hat{y} - y)^2 + \lambda \|w\|_1$$

③ Elastic net :- $loss = \sum (\hat{y} - y)^2 + \lambda (\alpha \|w\|_2^2 + (1 - \alpha) \|w\|_1)$

# Linear classifier (Perceptron):-

$y = -1$ or $1$

$\hat{y} = -1$ or $1$

## decision boundary:-

$$w_1 x_1 + w_2 x_2 + b = 0$$

## goal:

$$y^i (\Sigma w x) > 0$$

## update:

when: $y^i (\Sigma w x) \leq 0$

new w: $w \leftarrow w + y^i x^i$

$b \leftarrow b + y^i$

## Loss:

* $y(wx + b) > 0$   no loss

* $y(wx + b) \leq 0$    $loss = -y(wx + b)$

## # of iterations:

$$\left(\frac{R}{r}\right)^2$$

## classification Rule:-

$$y^i = sign(\Sigma w x_i)$$

- ❑ Perform training and check the error:
- ❑ If error is high: (underfitting)
  - ➤ Add more features
  - ➤ More complex model by polynomial
- ❑ If overfitting:
  - ➤ Add more data
  - ➤ Perform regularization

# Logistics Regression

$$y = 0 \text{ or } 1$$
$$\hat{y} = \sigma(\Sigma xw + b)$$

testing:-

$$\sigma(\Sigma wx + b) = \frac{1}{1 + e^{-(\Sigma wx + b)}}$$

$$= 0.5 \quad \Rightarrow \text{ undetermind}$$
$$\geqslant 0.5 \quad \Rightarrow \text{ +ve}$$
$$< 0.5 \quad \Rightarrow \text{ -ve}$$

## Gradiant decent:-

$$\text{Cost} = \frac{1}{m} \Sigma - y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$$y = 0 \text{ or } 1$$
$$\hat{y} = \sigma(\Sigma wx + b)$$

steps:

Same as above

but: $dw_i = \frac{1}{m} (\hat{y} - y) \cdot x$

$db = \frac{1}{m} (\hat{y} - y)$

## General form:-

$$P(y|x) = \frac{1}{1 + e^{-y(wx + b)}}$$

Softmax:- $\dfrac{e^{z_i}}{\Sigma e^{z}}$

# R and P

$$\text{accuracy} = \frac{TP}{all}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

Actual

|  | +ve | -ve |
|---|---|---|
| Predict +ve | TP | FP |
| -ve | FN | TN |

Percision (TP, FP)

recall (TP, FN)

# $F_B$-Score

$$F_B = (1 + B^2)\, \frac{P \times R}{B^2 P + R}$$

$$F_1 = 2\, \frac{P \times R}{P + R}$$

# hard SVM

decision boundry:-

$$w_1 x_1 + w_2 x_2 + b = 0$$

$$w_1 x_1 + w_2 x_2 + b = \pm 1 \longrightarrow SV \text{ boundry}$$

goal:

Primal      min $\frac{1}{2}\|w\|^2$    s.t.      $y^i(\sum wx) \geq 1$

Dual   max $\sum \alpha_j - \frac{1}{2}\sum \alpha_i \alpha_j y^i y^j (x^i . x^j)$   s.t.    $\sum \alpha_i y_i = 0$

$$y^i(wx+b) = 1$$

$$w = \sum \alpha_i y^i x^i$$

$$b = \sum \alpha_i y^i$$

$\boxed{\begin{array}{l} \text{high } c \Rightarrow \text{overfitting} \\ \text{low } c \Rightarrow \text{underfitting} \end{array}}$

# Soft SVM

decision boundry:-

$$w_1 x_1 + w_2 x_2 + b = 0$$

$$w_1 x_1 + w_2 x_2 + b = \pm 1 \longrightarrow SV \text{ boundry}$$

goal:

Primal      min $\frac{1}{2}\|w\|^2 + C\sum \xi_i$ s.t.      $y^i(\sum wx) \geq 1 - \xi_i$

Dual   max $\sum \alpha_j - \frac{1}{2}\sum \alpha_i \alpha_j y^i y^j (x^i . x^j)$   s.t.    $\sum \alpha_i y_i = 0$

that voilate   $C > \alpha > 0 \longleftarrow y^i(wx+b) = 1$

on the boundry   $C = \alpha \longleftarrow y^i(wx+b) = 1 - \xi_i$
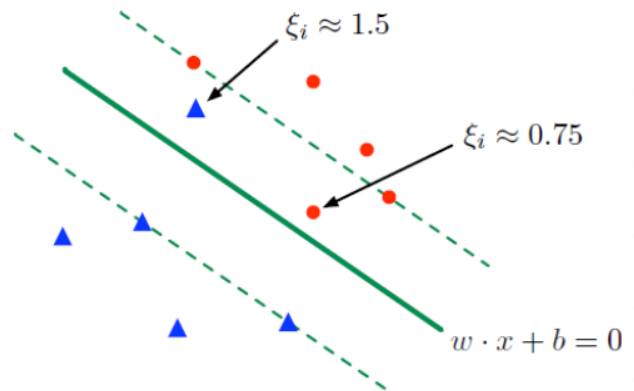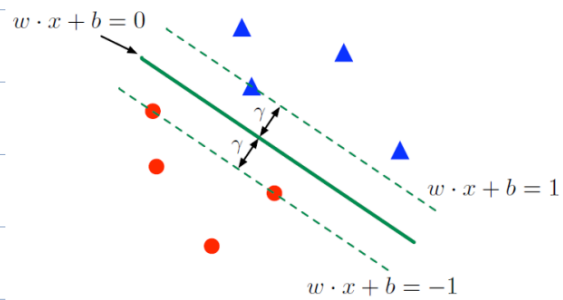
$$w = \sum \alpha_i y^i x^i$$

$$b = \sum \alpha_i y^i$$

test:-

$$\hat{y} = \text{sign}\left( \sum \alpha \, y^i \, x^i . z + b \right)$$



$w \cdot x + b = 0$

$w \cdot x + b = 1$

$w \cdot x + b = -1$



$\xi_i \approx 1.5$

$\xi_i \approx 0.75$

$w \cdot x + b = 0$

measure the
similarity

# Kernal:-

$$\left( 1 + X.z \right)^p$$

# Jaccard Index:-

$$1 - \frac{X \cap z}{X \cup z}$$

high $\sigma$ ⟹ underfitting

small $\sigma$ ⟹ overfitting

# Radial Basis Function:-

$$e^{-\frac{|X-z|^2}{\sigma^2}}$$

more data ⟹ $\sigma$ decreasing

# NN

each neuron $= \sigma(\sum wx + b)$

cross-entropy $= -y \log \hat{y} - (1-y) \log(1-\hat{y})$

$w \leftarrow w - \eta \left(\dfrac{\partial J}{\partial w}\right)$ ①

$b \leftarrow b - \eta \dfrac{\partial J}{\partial b}$ ②

$\dfrac{\partial J}{\partial w_j} = \dfrac{1}{m} \sum (\hat{y}-y) \, x_j$

$\dfrac{\partial J}{\partial b} = \dfrac{1}{m} \sum (\hat{y}-y)$

**Loss function is not convex!**

# Three Problems with activation functions:-

① non Zero centered   (slower in converge)

② Saturated          (no learning)

③ computational cost

$\max(0, x)$          $\dfrac{1}{1+e^{-z}}$

|  | ReLU | Tanh | Sigmoid |
|---|---|---|---|
| zero-centered | Not zero-centered | Zero-centred | Not zero-centered |
| Saturation | Dose not saturated | saturated | saturated |
| Computational Cost | efficient | slow | slow |
| Drivative | $\begin{cases} 0, & \text{if } x \le 0 \\ 1, & \text{other} \end{cases}$ | $1 - \tanh^2$ | $\sigma(x)(1-\sigma(x))$ |

for multiple classes :-

What we change

$\big\{$
\# of output Layer

activation of output Layer   SoftMax

Loss

$$\left( \sum_{i=1}^{K} -y \log(\hat{y}) \right)$$

↳ \# of classes

↳ one-hot-encoding

for Regression:-

What we change

$\big\{$
activation in the output Layer (no activation)

Loss

$$MSE = \frac{1}{2m} \sum (\hat{y} - y)^2$$

# Decision Tree



error = $1 - \max(P_i)$

Gini = $1 - \sum (P_i)^2$

Entropy = $-\sum_i P_i \log_2(P_i)$

$$IG(D_p) = I(D_p) - I(D)_L \frac{N_L}{N_p} - I(D)_R \frac{N_R}{N_p}$$

Cost of tree = error rate + $\lambda$ (# of leafs)

# Ensemle

## AdaBoost

Weighted error rate = $\dfrac{\text{total weight of mistakes}}{\text{total weight of all data}}$

x $S_1$   0.1
  $S_2$   0.6
x $S_3$   0.9
  $S_u$   0.5

$\epsilon = \dfrac{1.0}{2.1}$

weight of classifier = $\dfrac{1}{2} \ln\left(\dfrac{1-\epsilon}{\epsilon}\right)$

④ re w.t the samples   $\propto e^{-wt}$ (origin w.t)   , $\propto e^{+wt}$ (if it's correct)

$$\alpha_i^{j+1} = \begin{cases} \alpha_i^{)} \cdot e^{-w.t} & , \text{ if correct} \\ \alpha_i^{)} \cdot e^{+w.t} & , \text{ if incorrect} \end{cases}$$

## Bagging

## RF

# Clustering

Loss func = $\Sigma$ of squared distance from center

## ① K-means

assign to nearest cluster
update the center

K-means ++ = Smart initialization

LBG = Start with 2 then split until K
MacQueen = update as you assign

## ② EM

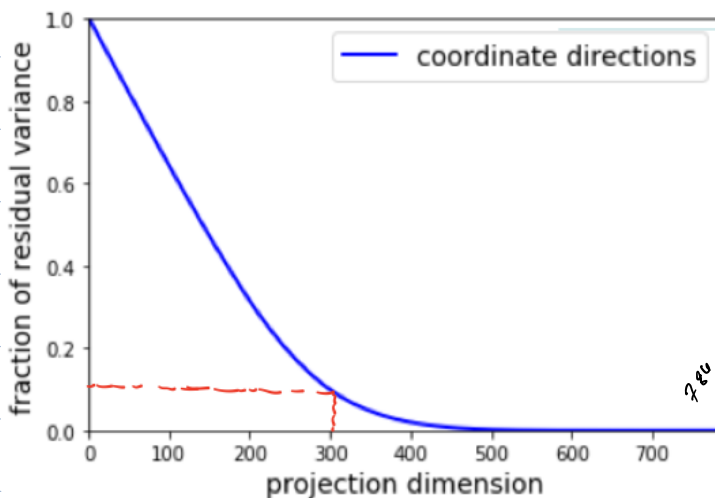## ③ Agglomerative hierarchical clustering

# Dimensionally Reduction

① Remove lowest variance

② PCA

if you don't have
exact # of axis
you fix % of
$\sigma$ then get $\lambda_1 \dots \lambda_n$
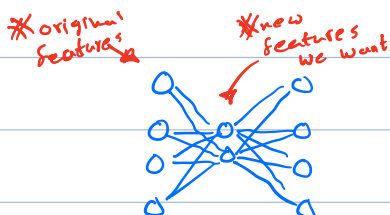that give you % you want

How $\sigma$ I will get if I select 3

$$\sigma = \frac{\lambda_1 + \lambda_2 + \lambda_3}{\Sigma \lambda}$$



if we keep 300 features
we will lose just 10% of
the variance

③ T-sne (neural network)



*original features

*new features we want

loss = squared loss