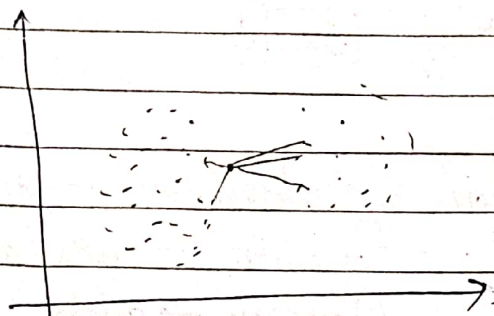


Non-parametric Learning Algo.
(Training data required for inference)
KNN (K-Nearest Neighbours) (learning faster, inference slower) → opposite for Parametric

Usually k is odd.



Calculate all dist.

↓ arrange in ascending order
consider the first k values.

↓ majority of their labels → our prediction

→ euclidean distance = $\sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
(elliptic in 2D, spherical shape in 3D)

for 2 pts 2D → $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

→ manhattan distance = $\sum_{i=1}^n |x_i - y_i|$

for 2 pts 2D.

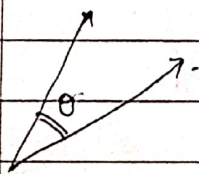
minkowski distance = $\left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$ $D = |x_2 - x_1| + |y_2 - y_1|$

$p=1$: Manhattan Distance

$p=2$: Euclidean Distance

$p \rightarrow \text{infinity}$: Chebychev Distance

cosine distance/similarity.



$$\cos \theta = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|} = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n b_i^2}}$$

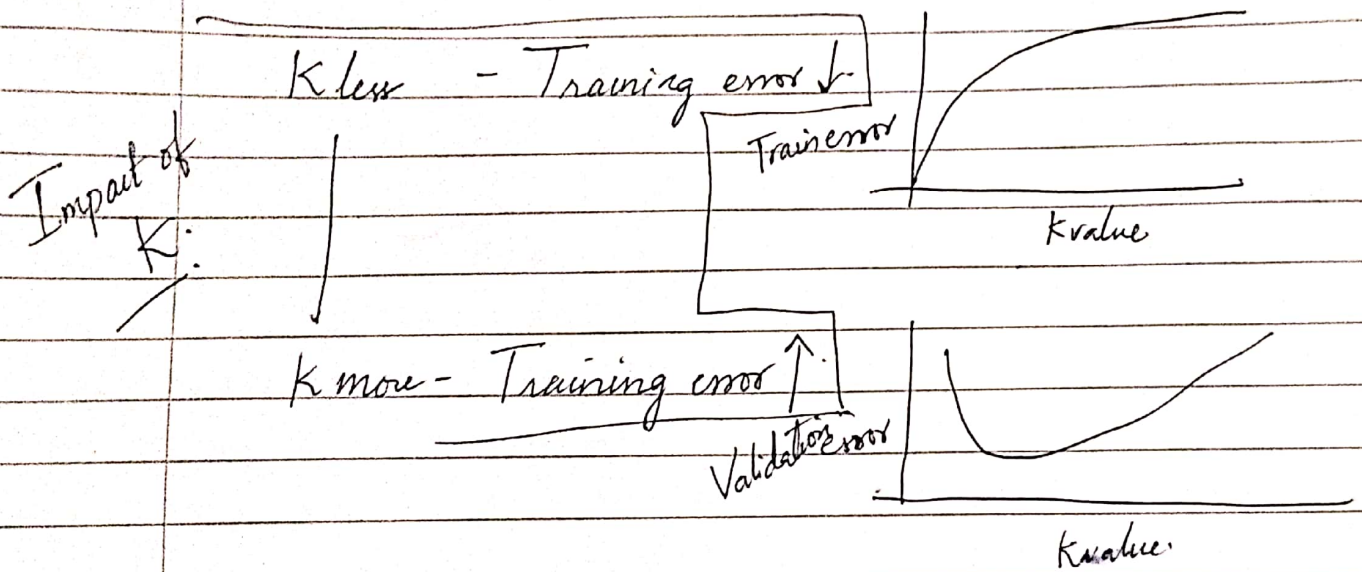
Hamming Distance = 01101010

⊕ 11011011

10110001 → 4 1s, (differ at 4 places).
 $d(01101010, 11011011) = 4$

Assumptions in traditional KNN.

- Equal weight to all attributes (to compute distance)
- Scale of attribute range same
- Equal variance
- Classes - spherical.
- Attributes - very less noise
- Equal importance



⇒ Distance-weighted KNN

• ~~This large value of K~~

if K is small - also sensitive to outliers

$$Pred_{test} = \frac{\sum_{i=0}^K w_i \cdot class_i}{\sum_{i=0}^K w_i}$$

$$w_k = \frac{1}{Dist(x_k, x_{test})}$$

K-large - may include unnecessary values from other classes.

→

Weighted distance function.

$$D(x_i, x_j) = \sqrt{\sum_{i=1}^n w_i (x_{im} - x_{jm})^2}$$

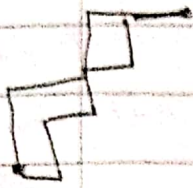
every feature has
difference
weight

$w_i \rightarrow \text{high}$

$w_i \rightarrow \text{low}$

$w_i = 0$

→ Mean Absolute Error



$$MAE(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} |y_i - \hat{y}_i|$$

does a poor job when scale of the data is high

→ Mean Squared Error

$$MSE(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} (y_i - \hat{y}_i)^2$$

Squaring the value, sum won't be zero
" emphasizes larger differences

→ Median Absolute Error

$$MedAE(y, \hat{y}) = \text{median}(|y_1 - \hat{y}_1|, \dots, |y_n - \hat{y}_n|)$$

Scale of prediction too high
logs scales it down

→ does not support multiclass

Median AE prevents outliers contributing more error to model evaluation

→ Mean Squared Logarithmic Error

$$MSLE(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} (\log_e(1 + y_i) - \log_e(1 + \hat{y}_i))^2$$

when targets have exponential growth → population counts, average sales of commodity over years

Method penalizes an under predicted estimate greater than over predicted estimate

regression +
classification
metric

→ Max Error.

$$\text{Max Error}(y, \hat{y}) = \max(|y_i - \hat{y}_i|)$$

aims to find max error made by model for single sample →
finding best model that covers all samples.

Similarity Measures.

→ Explained Variance Score (y, \hat{y}) Best possible score → 1.0
lower values are worse

$$= 1 - \frac{\text{Var}\{y - \hat{y}\}}{\text{Var}\{y\}}$$

Estimates deviation between prediction and actual values

→ R^2 score, the coeff of determination

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

calculate unadjusted
 R^2 without correcting
for bias in sample variance of y

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \epsilon_i^2$$

Classification Metrics

	actual	predicted	
True +ve	1	1	11 → 1
False +ve	0	1	01 → 0
False -ve	1	0	10 → 0
True -ve	0	0	00 → 1

Error/Dissimilar
Measure
(less the better)

Accuracy/Similar
Measure
(more the better)

$$A\bar{B} + \bar{A}B = \text{EXOR}$$

$$AB + \bar{A}\bar{B} = \text{EXNOR}$$

→ Accuracy = $\frac{TP + TN}{TP + TN + FP + FN}$ } all true values / total

→ Precision = $\frac{TP}{TP + FP}$ } fraction of predicted +ve that are actually +ve

Precision, Recall, Confusion matrix

Recall (Sensitivity) = $\frac{TP}{TP + FN}$
 TP ← correct
 FN → 0 then recall max
 all correct

Total days - 30 (one month)

actual rain - 10

Model - 9 days (5✓, 4x)

Precision = 5/9 ; Recall = 5/10

Accuracy :- $\frac{(TP + TN)}{\text{total}}$ → $\frac{TP + FP + TN + FN}{\text{total}}$

True +ve Rate :- $\frac{TP}{\text{actual yes}}$ → $\frac{FN + TP}{\text{actual yes}}$
 (Sensitivity / Recall)

False +ve Rate :- $\frac{FP}{\text{actual no}}$ → $\frac{TN + FP}{\text{actual no}}$

True -ve Rate :- $\frac{TN}{\text{actual no}}$ → $\frac{FP + TP}{\text{actual no}}$
 (Specificity)

Precision :- $\frac{TP}{\text{predicted yes}}$ → $\frac{FP + TP}{\text{predicted yes}}$

Misclassified Rate :- $\frac{FP + FN}{\text{Total}}$ → $\frac{TP + FP + TN + FN}{\text{Total}}$
 False +ve Rate = $\frac{FN}{\text{actual yes}}$ → $\frac{FN + TP}{\text{actual yes}}$
 (miss rate)

→ $F_1 = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}} = \frac{2 * (\text{precision} * \text{recall})}{\text{precision} + \text{recall}}$
 50% 50% focus

→ $F_\beta \text{ score} = \frac{(1 + \beta^2) * (\text{precision} * \text{recall})}{(\beta^2 * \text{precision} + \text{recall})}$

$\beta = 0$ P ↑

$\beta = \infty$ R ↑

- $F_{0.5}$ (beta = 0.5) :- Precision ↑ Recall ↓
- F_1 (beta = 1.0) :- Precision 50% Recall 50%
- F_2 (beta = 2.0) :- Precision ↓ Recall ↑

ROC Curve and ROC AUC Score

ROC AUC curve
overfitting
underfitting

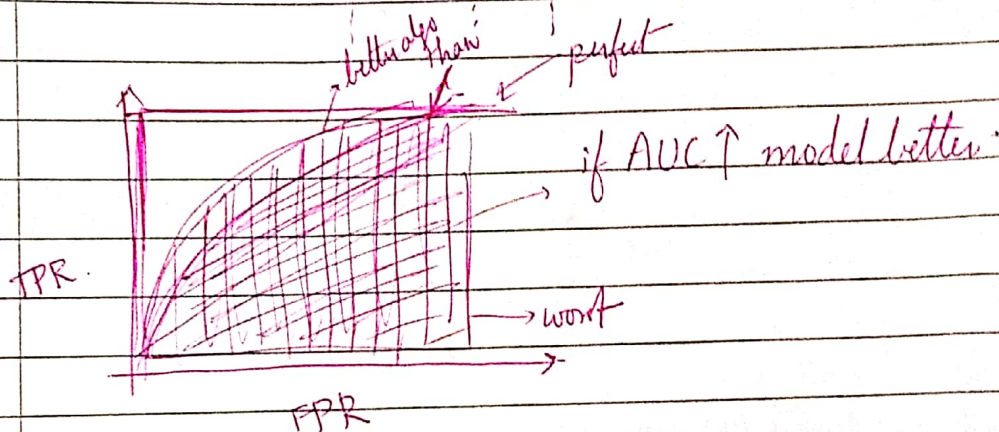
helps in understanding balance b/w TPR & FPR

- thresholds → all unique prediction probabilities in 'descending order'
 - FPR = $\frac{FP}{(FP+TN)}$
 - TPR = $\frac{TP}{(TP+FN)}$
- for each threshold

Higher the AUC, better the model is at predicting 0s as 0s and 1s as 1s.

eg:-
th₁ 0.31
else 0
th₂ 0.41
else 0

fps tps
th₁ x₁ y₁
th₂ x₂ y₂
th₃ x₃ y₃



Overfitting and underfitting

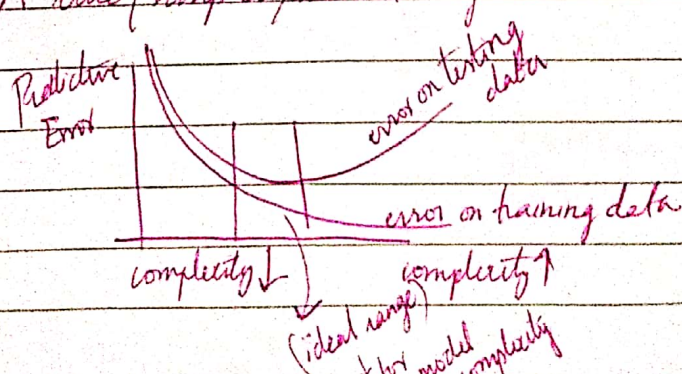
$x_1 \propto y$ $x_2 \propto \frac{1}{y}$ $x_3 \nrightarrow$ no impact y
 $w_1 \uparrow$ $w_2 \downarrow$ $w_3 = 0$
always tries to get +ve impact on 'y'

Training Error
Training Loss
Underfit
Bias

Testing Error
Testing Loss
Overfit
Variance

→ A complex model $y = \text{high order polynomial in } x$

→ A true (simpler) model $y = aX + b + \text{noise}$



training acc test acc

60

58

60

40

98

70

98

96

HB
LV
HB
HB
LB
HB
LB
LV