

Supervised Learning

Practical Machine Learning

Kannan Singaravelu

May, 2022

In this lecture...

- Understanding classification setting
- Logit functions & algorithms for classification
- Classification loss functions
- Model assessment ~~X~~
- Support Vector Networks ~~Top~~
- Lagrangian Dual Formulation
- Three steps to find nearest neighbours ~~- KNN~~
- Overview of hyperparameters tuning

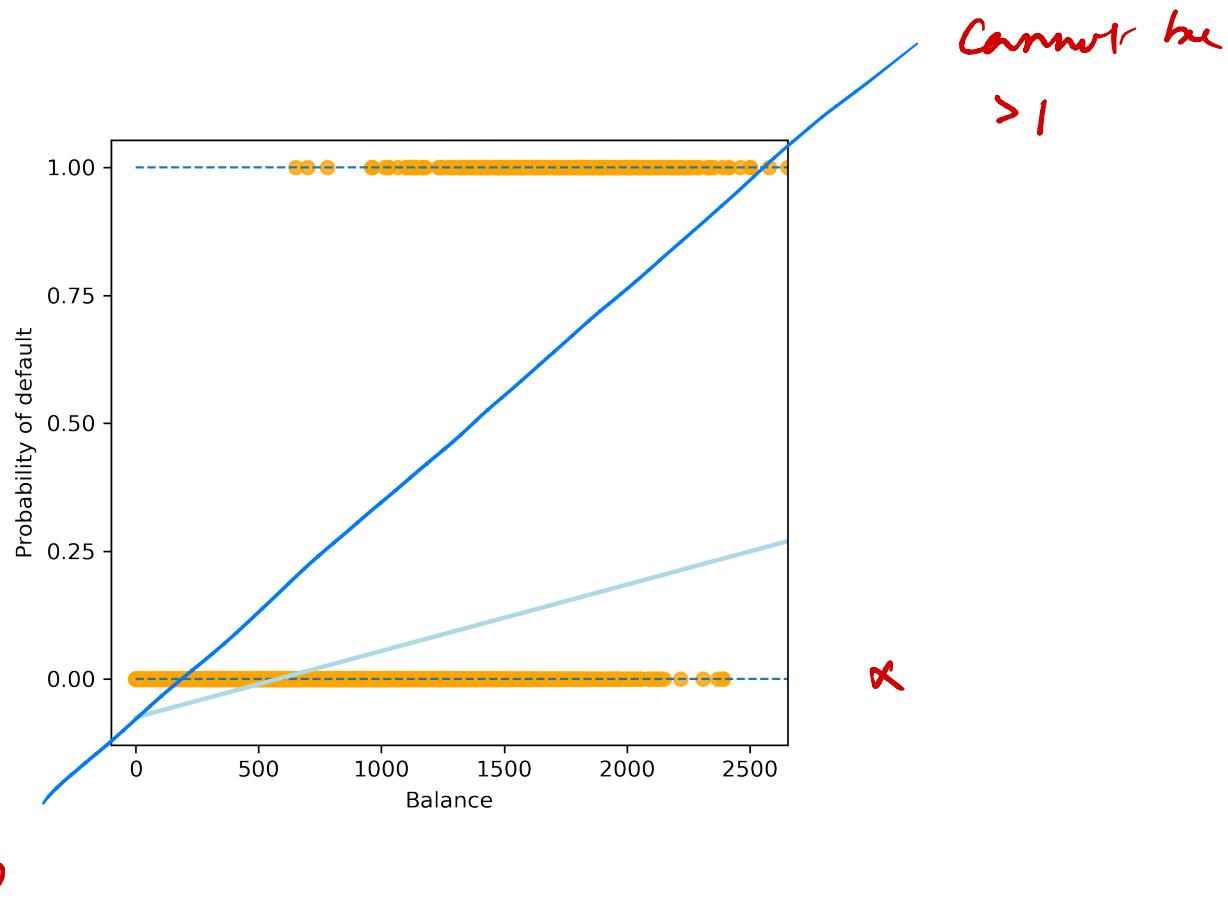
9/10 use Conv

- a) Hard margin
- b) SVC
- c) SVM

Different

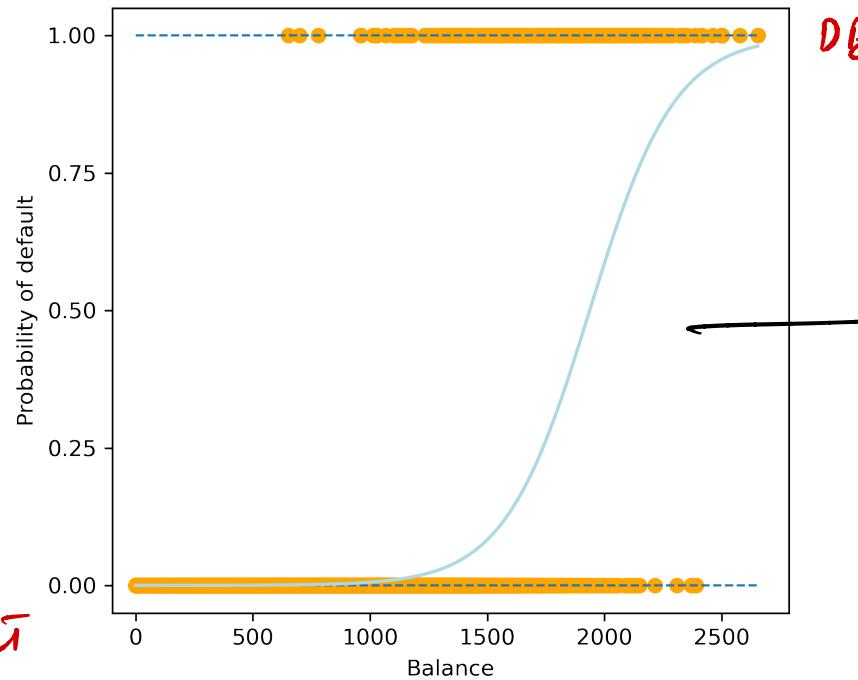
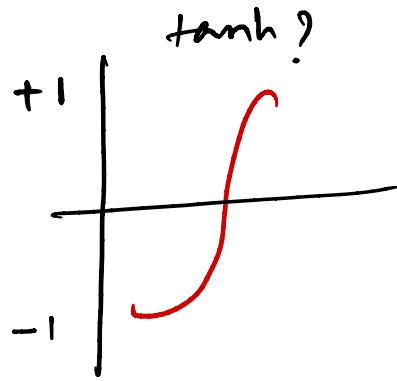


Why Not Linear Regression?





Why Not Linear Regression?



$\Sigma \log(\text{odds})$
logit fn
=

The Classification Problem

- Classification models are models which predict a qualitative response.
- Estimate f on training observations $(x_1, y_1), \dots, (x_n, y_n)$ where y_1, \dots, y_n
- The training error is given as

$$\frac{1}{n} \sum_{j=1}^n I(y_i \neq \hat{y}_i)$$

where, *Indicator loss*

$$y_i = \hat{y}_i \Rightarrow L = 0$$

$$y_i \neq \hat{y}_i \Rightarrow L = 1$$

"0-1 loss"

y_i is the predicated class label,

$I(y_i \neq \hat{y}_i)$ is an indicator variable that equals 1 if $y_i \neq \hat{y}_i$ and equals 0 if $y_i = \hat{y}_i$

- The test error rate is then given as $\text{Ave } I(y_i \neq \hat{y}_i)$

► Bayes Classifier

	Class 1	Class 2	
1	40%	60%	✓
2	80% ✓	20%	

- The test error rate can be minimised if we classify the classes correctly given its predictors.
- Assign test observation with predictor x_0 to the class j for which
 $Pr(Y = j | X = x_0)$ is the largest.
- Classifier based on this conditional probability is popularly known as the Bayes Classifier.
- For binary classification of class 1 and class 2, predicting class 1 corresponds to
 $Pr(Y = 1 | X = x_0) > 0.5$ class 2 otherwise.
? threshold
- Bayes error rate produces lowest possible error rate and is given as

$$\underline{1 - E \left(\max_j Pr(Y = j | X) \right)}$$

1 60% prob. $\Rightarrow 1 - 60\% = 0.40$
2 80% prob. $\Rightarrow 1 - 80\% = 0.20$

Error
↓ rate

► Logistic Regression

- One of the most widely used algorithms for classification.
- Map quantitative data onto categorial variables.
- Probabilistic model used to assign class labels.
- Two main type of classification problems
 - ✓ Binary or Binomial
 - ✓ Multi-class or Multinomial
- Also referred as the shallow neural network

► Logistic Regression

- Find a function of the predictor variables that relates to 0 and 1 as outcome.
- Instead of y as outcome variable (like in linear regression), we use a function of y called the logit.
- Logit can be modelled as a linear function of the predictor and can be mapped back to a probability which in turn can be mapped to a class.
- Typically use probability with a cutoff value.

► Logit

- $\text{Logit} = \log(\text{odds}) = \underline{w_0 + w_1x_1 + w_2x_2 + \dots + w_qx_q}$

- The odds of an event is defined as

$$Odds = \frac{p}{1 - p}, \text{ where } p \text{ is the probability of event}$$

- Probability of an outcome or event can be computed as

$$p = \frac{Odds}{1 + Odds}$$

- If $p = 0.5$, $Odds = 1$ and $\log(Odds) = 0$

- Log odds can range between $-\infty$ to $+\infty$ and the midpoint is when $p = 0.5$

Probability Calculation

$$\text{logit} \left(\frac{P}{1-P} \right) = \log(\text{odds}) = w_0 + w_1 x_1 + w_2 x_2 + \dots = z$$

$$\text{odds} = \frac{P}{1-P}, \quad P \Rightarrow \text{prob.}$$

$$\log \left(\frac{P}{1-P} \right) = \text{linear combination} = z$$

SIGMOID

$$P = \frac{1}{1+e^{-z}}$$

$$\frac{P}{1-P} = e^z$$

taking
Isolate P , reciprocal

$$\frac{1-P}{P} = \frac{1}{e^z}$$

$$\frac{1}{P} = 1 + \frac{1}{e^z} \Rightarrow \frac{1+e^z}{e^z} \therefore P = \frac{e^z}{1+e^z} \times e^{-z}$$

Probability Calculation

When the response has 3 levels

$$\ln\left(\frac{P_1}{P_2}\right) = z_1, \quad \ln\left(\frac{P_3}{P_2}\right) = z_2$$

$\left(\begin{array}{l} = \text{prob of losing} \\ \text{w.r.t. draw } (P_2) \end{array} \right)$

$\left(\begin{array}{l} \text{prob of gaining} \\ \text{w.r.t. draw } (P_2) \end{array} \right)$

P_1 = Down or Loss

P_2 = Draw or Neutral

P_3 = Gain or VP

$$P_1 + P_2 + P_3 = 1$$

$$\underline{P_1 + P_3} = 1 - P_2$$

Falling exp,

$$\textcircled{1} \quad \frac{P_1}{P_2} = e^{z_1}$$

$$\textcircled{2} \quad \frac{P_3}{P_2} = e^{z_2}$$

$$\Rightarrow \frac{P_1 + P_3}{P_2} = e^{z_1} + e^{z_2}$$

(A)

$$\frac{1 - P_2}{P_2} = e^{z_1} + e^{z_2}$$

$$\frac{1}{P_2} = 1 + e^{z_1} + e^{z_2} \Rightarrow$$

$$P_2 = \frac{1}{1 + e^{z_1} + e^{z_2}}$$

► Probability Calculation

from ①

$$P_1 = P_2 \cdot e^{z_1}$$

$$P_1 = \frac{e^{z_1}}{1 + e^{z_1} + e^{z_2}},$$

from ②

$$P_3 = P_2 \cdot e^{z_2}$$

$$P_3 = \frac{e^{z_2}}{1 + e^{z_1} + e^{z_2}}$$



Probability Calculation

▶ Probability Calculation

- Function that converts log-odds to probability is the logistic function

$$\log\left(\frac{p}{1-p}\right) = w_0 + w_1x_1 + w_2x_2 + \dots + w_px_p = z$$

$$\frac{p}{1-p} = e^z$$

$$\frac{1}{p} - 1 = \frac{1}{e^z}$$

$$\frac{1}{p} = \frac{1 + e^z}{e^z}$$

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

▶ Probability Calculation

- If response has 3 levels,

$$\log\left(\frac{p_1}{p_2}\right) = z_1 \quad \log\left(\frac{p_3}{p_2}\right) = z_2$$

$$\frac{p_1}{p_2} = e^{z_1} \quad (1)$$

$$\frac{p_3}{p_2} = e^{z_2} \quad (2)$$

$$\frac{p_1 + p_3}{p_2} = e^{z_1} + e^{z_2} \quad (3)$$

$$p_1 + p_2 + p_3 = 1 \Rightarrow p_1 + p_3 = 1 - p_2$$

$$\text{rewriting (3), } \frac{1 - p_2}{p_2} = e^{z_1} + e^{z_2} \Rightarrow \frac{1}{p_2} - 1 = e^{z_1} + e^{z_2} \Rightarrow \frac{1}{p_2} = 1 + e^{z_1} + e^{z_2}$$

$$p_2 = \frac{1}{1 + e^{z_1} + e^{z_2}}$$

$$p_1 = \frac{e^{z_1}}{1 + e^{z_1} + e^{z_2}} \quad \text{from (1)}$$

$$p_3 = \frac{e^{z_2}}{1 + e^{z_1} + e^{z_2}} \quad \text{from (2)}$$

► Logistic Regression

- Model the probabilities of the output classes given a function that is linear in x .
- Mathematically,

$$P_j = g \left(w_0 + \sum_{j=1}^p x_{ij} w_j \right) = g \left(w^T x \right)$$

~~where $g(z) = \frac{1}{1 + e^{-w^T x}} = \frac{1}{1 + e^{-z}}$~~

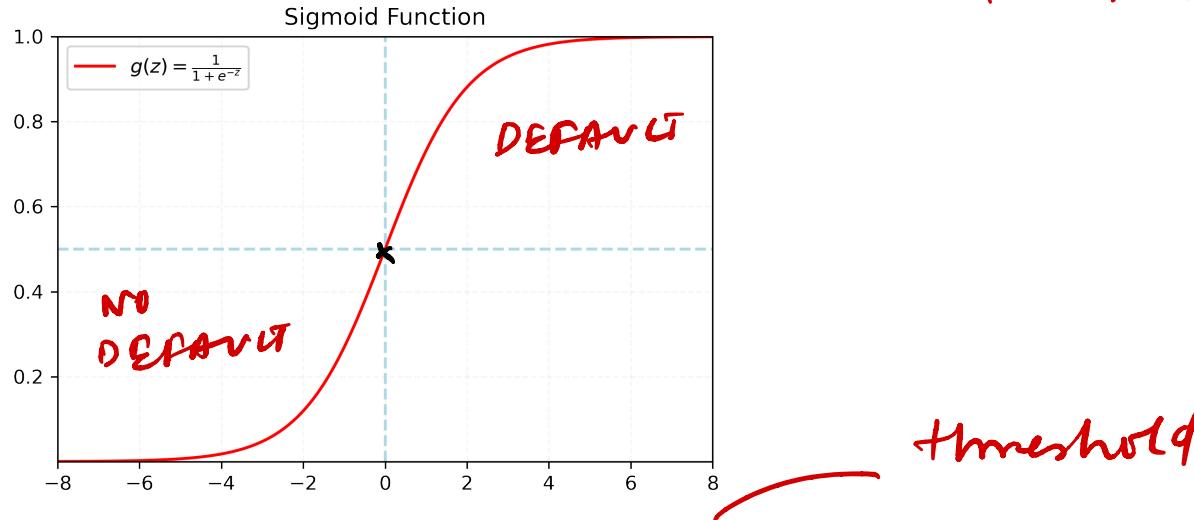
sigmoid = α

- Sigmoid function maps values that range from positive to negative infinity to values that only range from 0 to 1.
- Use probability with a threshold value for predictive classification.

► Logistic Regression

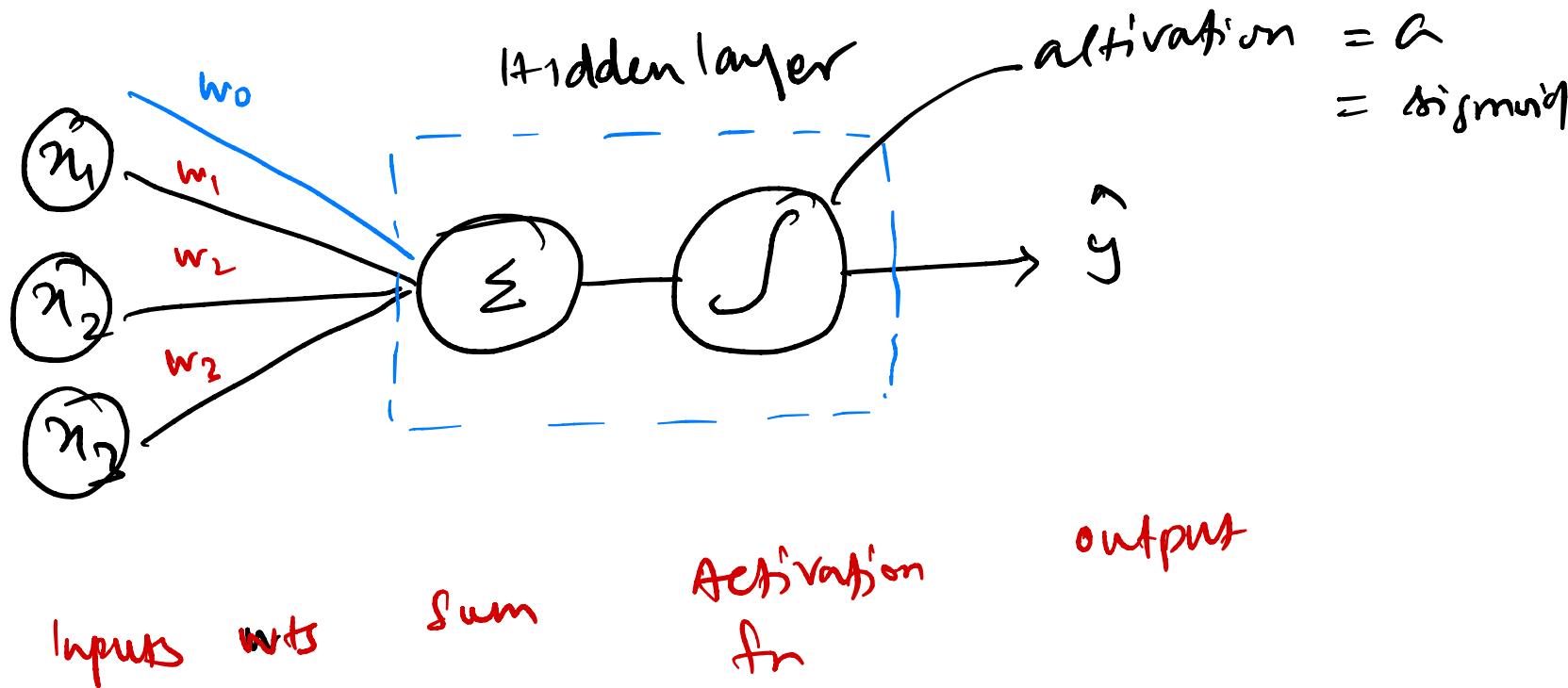
- Prediction of the probability is then given as

$$\hat{y} = \begin{cases} 0 & \text{if } \hat{p} < 0.5 \\ 1 & \text{if } \hat{p} \geq 0.5 \end{cases} \quad \text{where, } \hat{p} = g(z) = \text{SIGMOID - Activation function.}$$



- $g(z) < 0.5$ when $z < 0$ and $g(z) \geq 0.5$ when $z \geq 0$
- Model predicts 1 if $g(z)$ is ≥ 0.5 , and 0 if it is < 0.5

► Shallow Neural Networks



feed forward NN \Rightarrow Single Neuron
Perceptron

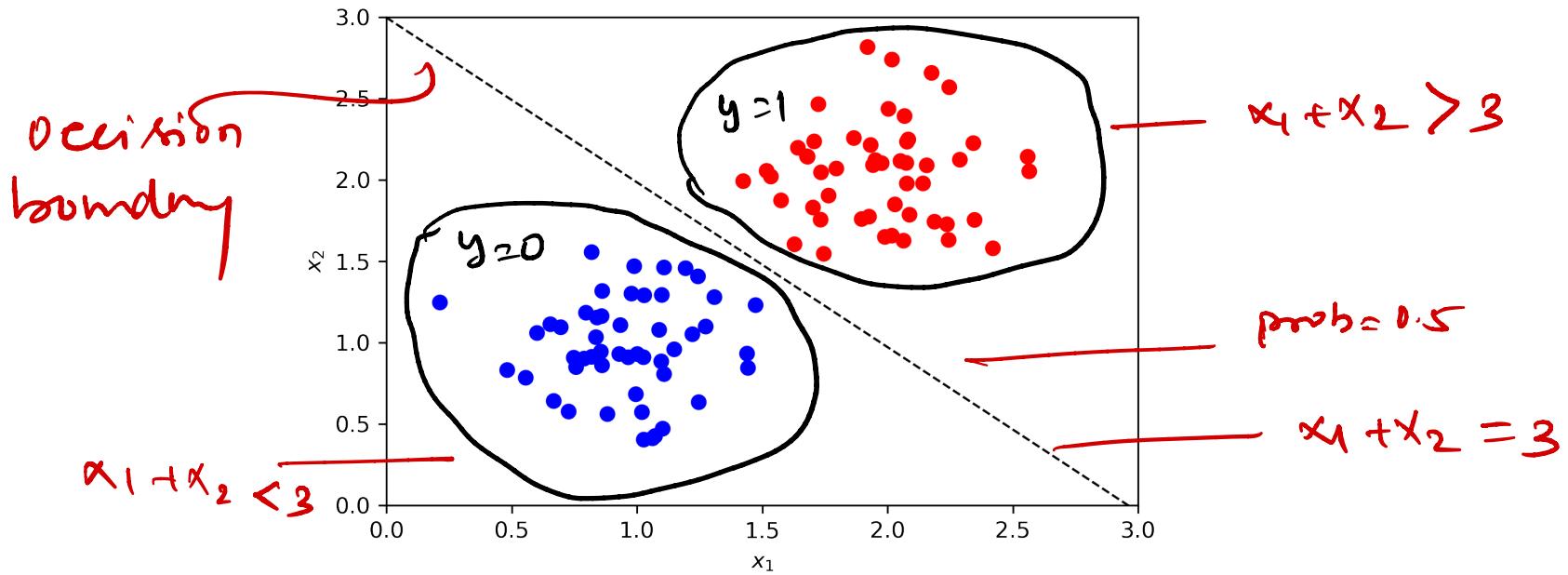
Shallow NN.

Linear Decision Boundary

$$f(x) = g(w_0 + w_1 x_1 + w_2 x_2)$$

-3 1 1

$$-3 + x_1 + x_2 \geq 0$$





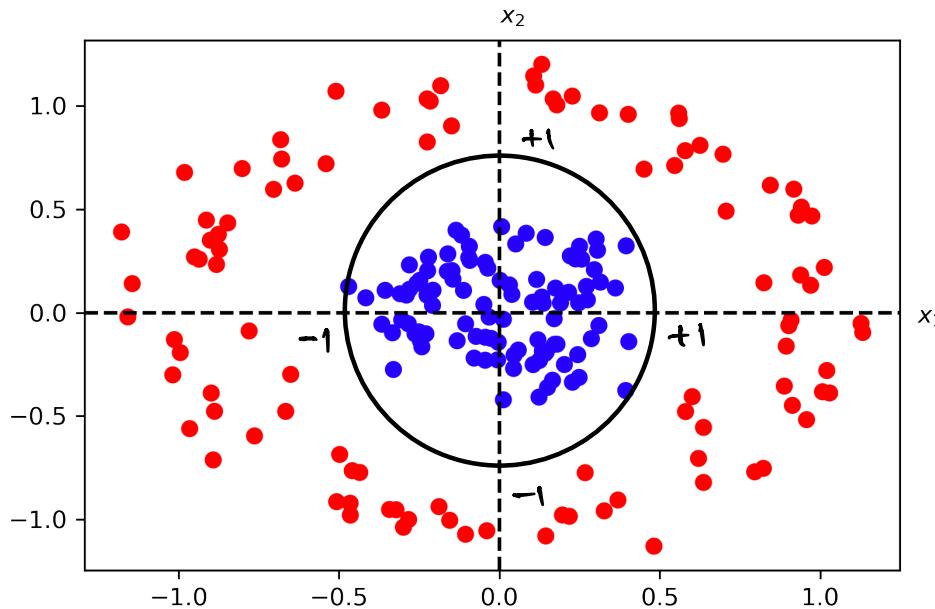
Non-Linear Decision Boundary

$$f(n) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2$$

-1 0 0 1 1

$$-1 + x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 1$$



► Decision Boundary

- There is a linear relationship between the features and the input of the sigmoid function.
- Decision boundary is a property of the hypothesis including parameters w_0, w_1, w_2 and not the property of the dataset.
- Adding polynomial features will result in more complex decision boundaries.

► Loss Functions for Classification

- In general, the **0-1 loss function** for classification take the form $l(\hat{y}, y) = 1(\hat{y} \neq y)$, where 1 is the indicator function.
- **Score Function**

$$\text{Predict} = \begin{cases} +1, & \text{if } f(x) > 0 \\ -1, & \text{if } f(x) < 0 \end{cases}$$

$$\hat{y} = \begin{cases} +1 & f(x) > 0 \\ -1 & f(x) \leq 0 \end{cases}$$

where $f(x)$ is the score for the input x and f may be called a score function.

- Magnitude of the score represents the confidence of our prediction.

► Loss Functions for Classification

- Margin

- Functional margin is $\hat{y}y$
- It is a measure of how correct our predictions are
 - If y and \hat{y} are the same sign, prediction is correct and margin is positive.
 - If y and \hat{y} have different sign, prediction is incorrect and margin is negative.
- Objective is to maximise the margin.

- Margin-based

- Most classification losses depend only on the margin.
- Geometric margin - a related concept used in SVM.

$$\hat{y}y$$

$$y=1, \hat{y}=1, \hat{y} \cdot y = +1$$

$$y=1, \hat{y}=-1, \hat{y} \cdot y = -1 \times$$

► Loss Functions for Classification

■ Zero-One Loss

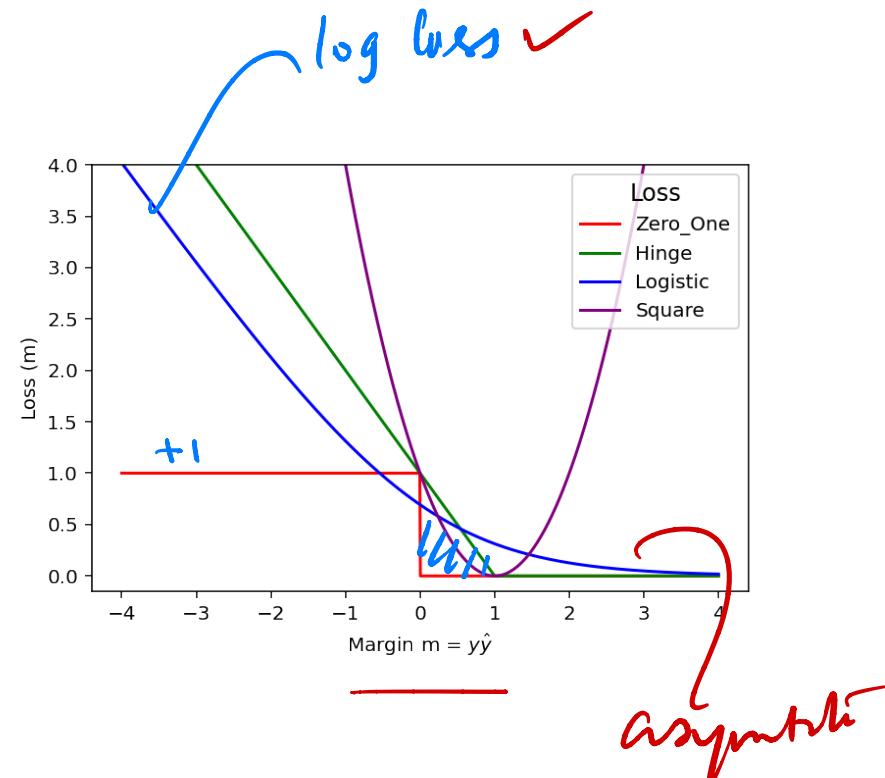
- $l_{0-1} = 1(m \leq 0)$
- Positive margin is correct classification

■ Hinge Loss

- $l_{Hinge} = \max(1 - m, 0) = (1 - m)_+$
- Margin error when $m < 1$
- Not differentiable at $m = 1$

■ Log Loss

- $l_{Logistic} = \log(1 + e^{-m})$
- Differentiable, wants more margin (loss never zero).



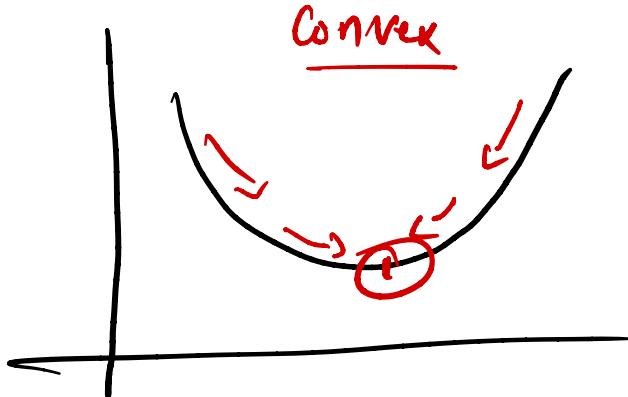
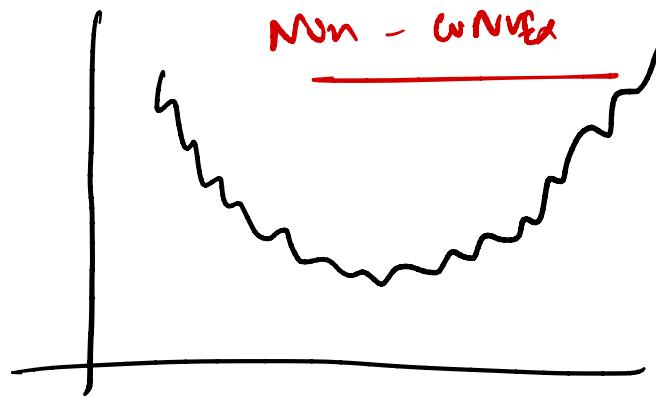


Why Not Square Loss?

$$(\hat{y} - y)^2 = (\underline{g(x)} - y)^2$$

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

Since $\frac{1}{1 + e^{-x}}$ (σ) is non linear in the above eqn., $J(w)$ (cost fn) ends up being a non-convex fn. if we define it as a square cost fn.



► Why Not Square Loss? Logit cost fn.

$$P(y|x) = \begin{cases} \hat{y} & \text{if } y=1 \\ 1-\hat{y} & \text{if } y=0 \end{cases}$$

where \hat{y} is the chance that $y=1$

$$P(y|x) = \hat{y}^y \cdot (1-\hat{y})^{(1-y)} \quad \text{--- (1)}$$

$$\begin{aligned}\hat{y} &\Rightarrow y=1 \\ 1-\hat{y} &\Rightarrow y=0\end{aligned}$$

Taking log,

$$\log(P(y|x)) = y \log \hat{y} + (1-y) \log(1-\hat{y})$$

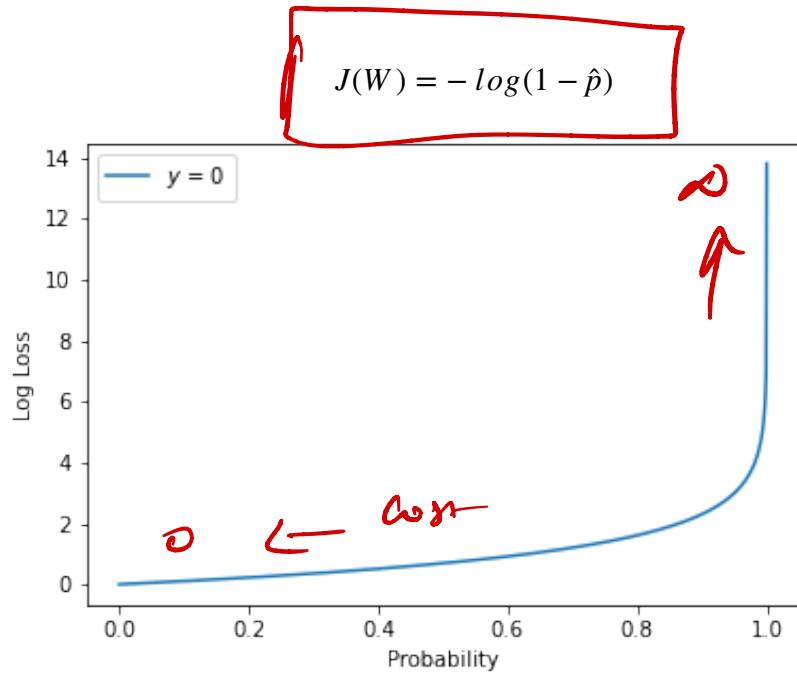
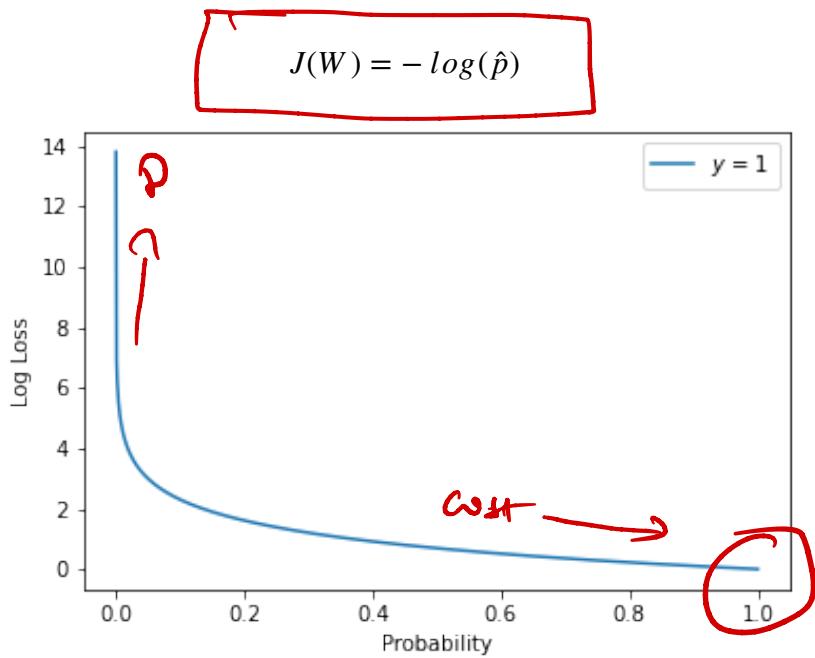
(2)

► Why Not Square Loss?

maximizing the log of prob is minimizing
the loss fn.

$$\therefore \log P(y|x) = -f(\hat{y}, y)$$

► Evaluation Metrics : Classification



- If $y = 1$, and $\hat{p} = 1$, then cost = 0
If $y = 1$, and $\hat{p} = 1$, then cost = 0
- $\hat{p} \rightarrow 0$, cost $\rightarrow \infty$
- If $y = 0$, and $\hat{p} = 0$, then cost = 0
If $y = 0$, and $\hat{p} = 0$, then cost = 0
- $\hat{p} \rightarrow 1$, cost $\rightarrow \infty$

Evaluation Metrics : Cost Function

$P(\text{Labels in training set}) = \prod_{i=1}^n P(y_i | x_i)$

MLE

, taking log.

$$\log P(y_1, y_2, \dots | x_1, x_2, \dots) = \sum_{i=1}^n \log P(y_i | x_i)$$

From ②,

$$= - \sum (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i))$$

$$J(w) = \frac{1}{n} \sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)$$

B

BINARY cross entropy

▶ Evaluation Metrics : Cost Function

- For a single training pair (x, y)

$$\hat{y} = P(Y = 1 | X = x) = \sigma(\mathbf{w} \cdot \mathbf{x}), \quad \text{where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$P(y|x) = \begin{cases} \hat{y} & \text{if } y = 1 \\ 1 - \hat{y} & \text{if } y = 0 \end{cases}$$

If the labels are 0 and 1, then Y is a Bernoulli random variable $Y \sim Ber(p)$,

where, $p = \sigma(\mathbf{w} \cdot \mathbf{x})$,

$$\therefore P(Y = y | X = x) = \sigma(\mathbf{w} \cdot \mathbf{x})^y [1 - \sigma(\mathbf{w} \cdot \mathbf{x})]^{(1-y)}$$

$$P(y|x) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

$$\log P(y|x) = y \log \hat{y} + (1 - y) \log(1 - \hat{y}) = -L(\hat{y}, y)$$

- Maximising the log of probabilities is minimising the loss

▶ Evaluation Metrics : Cost Function

- **Binary Cross Entropy** - used with models that output a probability between 0 and 1, the cost function is given as below,

$$\begin{aligned} J(W) = \log \prod_{n=1}^n P(y^{[i]} | x^{[i]}) &= -\frac{1}{n} \sum_{n=1}^n y^{[i]} \log(\hat{y}^{[i]}) + (1 - y^{[i]}) \log(1 - \hat{y}^{[i]}) \\ &= -\frac{1}{n} \sum_{n=1}^n y^{[i]} \log(f(x^{[i]}; W)) + (1 - y^{[i]}) \log(1 - f(x^{[i]}; W)) \\ &= -\frac{1}{n} \sum_{n=1}^n y^{[i]} \log(\hat{p}^{[i]}) + (1 - y^{[i]}) \log(1 - \hat{p}^{[i]}) \end{aligned}$$

- Cost function is convex
- Use Gradient Descent to find the global minimum

Confusion Matrix

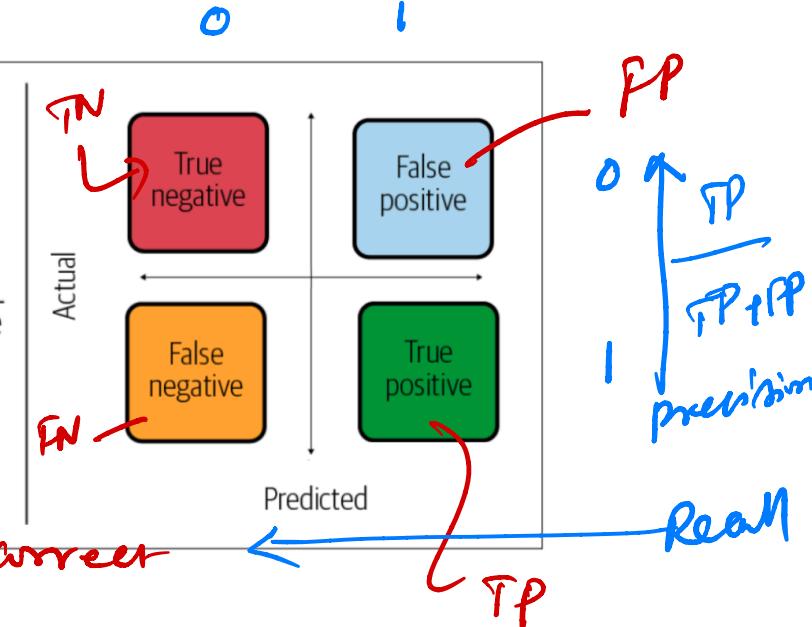
- Confusion Matrix

$$\text{Precision} = \frac{\text{True positive}}{\text{Actual results}} \quad \text{or} \quad \frac{\text{True positive}}{\text{True positive} + \text{False positive}}$$

$$\text{Recall} = \frac{\text{True positive}}{\text{Predictive results}} \quad \text{or} \quad \frac{\text{True positive}}{\text{True positive} + \text{False negative}}$$

$$\text{Accuracy} = \frac{\text{True positive} + \text{True negative}}{\text{Total}}$$

$$\text{F1 Score} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{\text{TP}}{\text{TP} + 0.5 * (\text{FP} + \text{FN})}$$



\Rightarrow Harmonic mean.

how many true relevant items
are selected

Accuracy Paradox

what Model ?

student not default

student default

Predicted

Model - 1 ✓		
Actual	0	1
0	5000	70
1	30	20

TN

FN

FP

TP

Predicted

Model - 2 ✓		
Actual	0	1
0	20	30
1	70	5000

$$\text{Accuracy} = \frac{5000+20}{5000+20+70+30} = 98.05\%$$

$$\text{Misclassification} = 1 - \text{Accuracy} = 1.95\%$$

$$M_1 = \frac{5020}{5720} = 98$$

$$M_2 = \frac{5020}{5120} = 98$$

Balanced Accuracy

- Machine Learning algorithms work best when the number of samples in each class are about equal.

$$\text{Balanced Accuracy} = \frac{\text{sensitivity} + \text{specificity}}{2}$$

$$\frac{\text{tpr} + \text{tnr}}{2}$$

- Sensitivity is also known as true positive rate or recall that measures the proportion of real positives that are correctly predicted out of total positive prediction made by the model.
- Sensitivity = $\frac{\text{TP}}{(\text{TP} + \text{FN})}$
- Specificity is also known as true negative rate, it measures the proportion of correctly identified negatives over the total negative prediction made by the model.
- Specificity = $\frac{\text{TN}}{(\text{TN} + \text{FP})}$

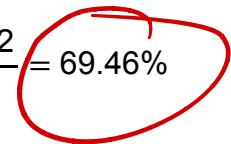
Balanced Accuracy

		Predicted
		Model - 1
		0
Actual		1
0	5000	70
1	30	20

$$\text{Sensitivity} = \frac{20}{20+30} = 40\%$$

$$\text{Specificity} = \frac{5000}{5000+70} = 98.62\%$$

✓ $\text{Balanced Accuracy} = \frac{40 + 98.62}{2} = 69.46\%$



- A model can have high accuracy with bad performance, or low accuracy with better performance, which can be related to the accuracy paradox.
- Accuracy can be a useful measure only if we have a balanced dataset.

▶ Class Imbalances

- Machine Learning algorithms work best when the number of samples in each class are about equal.
- Observation in one class is higher than the observation in other classes leads to class imbalances - which is a common problem in Machine Learning.
- Most algorithms are designed to maximize accuracy and reduce errors.
- Imbalance data can hamper model accuracy big time.
 - high accuracy by predicting the **majority class** (not that model is any good)
 - fail to capture the **minority class**
- Algorithm selection, under/over sampling and penalising algorithms are some of the most common approaches to address imbalance in data.

arising class nts

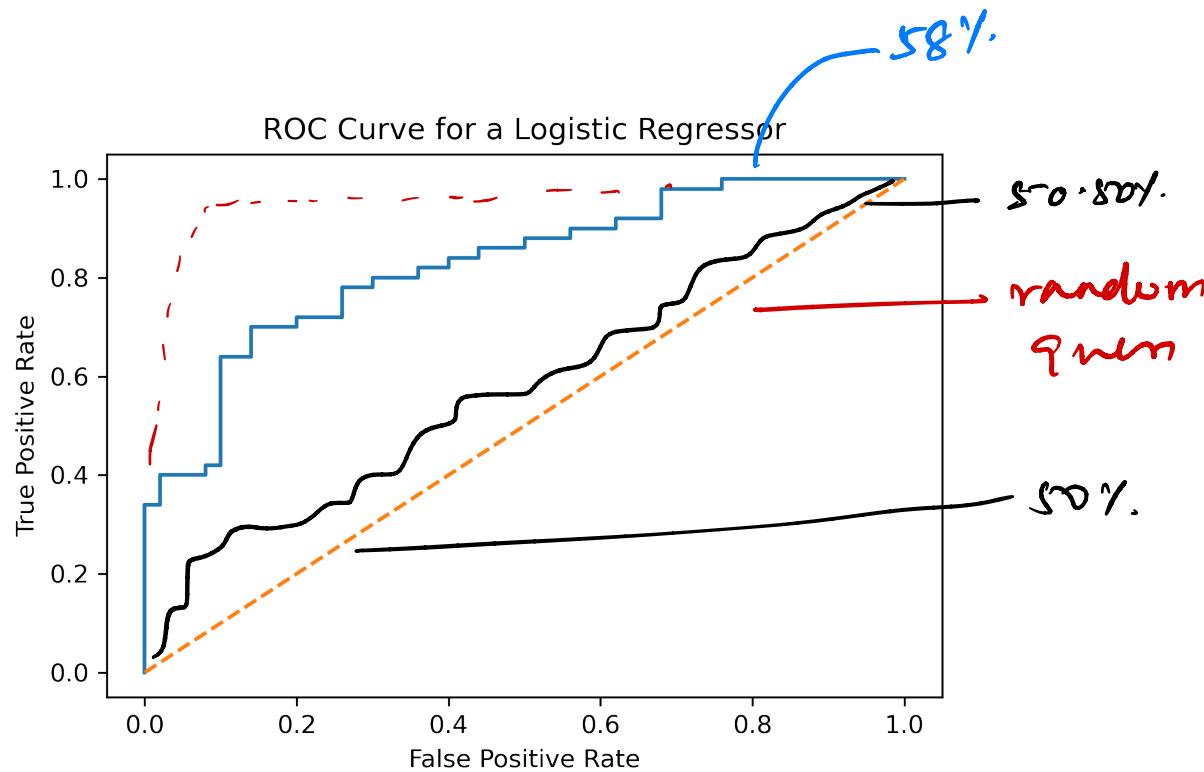
class nts (in mas)
1 min more mis

► ROC Curves

- Receive Operating Characteristic curve is a plot of the false positive rate versus the true positive rate.
- True positive rate is the recall and describes the goodness of the model in predicting the positive class when the actual outcome is positive. It is also referred to as sensitivity.
- False positive rate is defined as the inverted specificity (ie., 1-Specificity).
- ROC curves are used where there is equal number of observations of each class.

► ROC Curves

- The area under the ROC curve (AUC) is a measure of how well a model can distinguish between two classes.



Multiclass Classification

- Logistic Regression model can be extended to multiple classes.
- Given an instance x , the model computes a score $s_k(x)$ for each k ,
- Probability of each class is then estimated by applying softmax function (normalised exponential) to the scores
- Mathematically, this is represented as

$\hat{p}_k = \frac{f_k}{\sum_k f_k}$, where f_k is the output of classifier k and \hat{p}_k is the probability of the

observation being in class k and sum of \hat{p}_k for all the values of k is equal to 1.

$$\hat{p}_k = \sigma(s(x))_k = \frac{\exp^{s_k(x)}}{\sum_{j=1}^K \exp^{s_j(x)}} \text{ and for } k \text{ classes, we train } \frac{k(k-1)}{2} \text{ classifiers.}$$

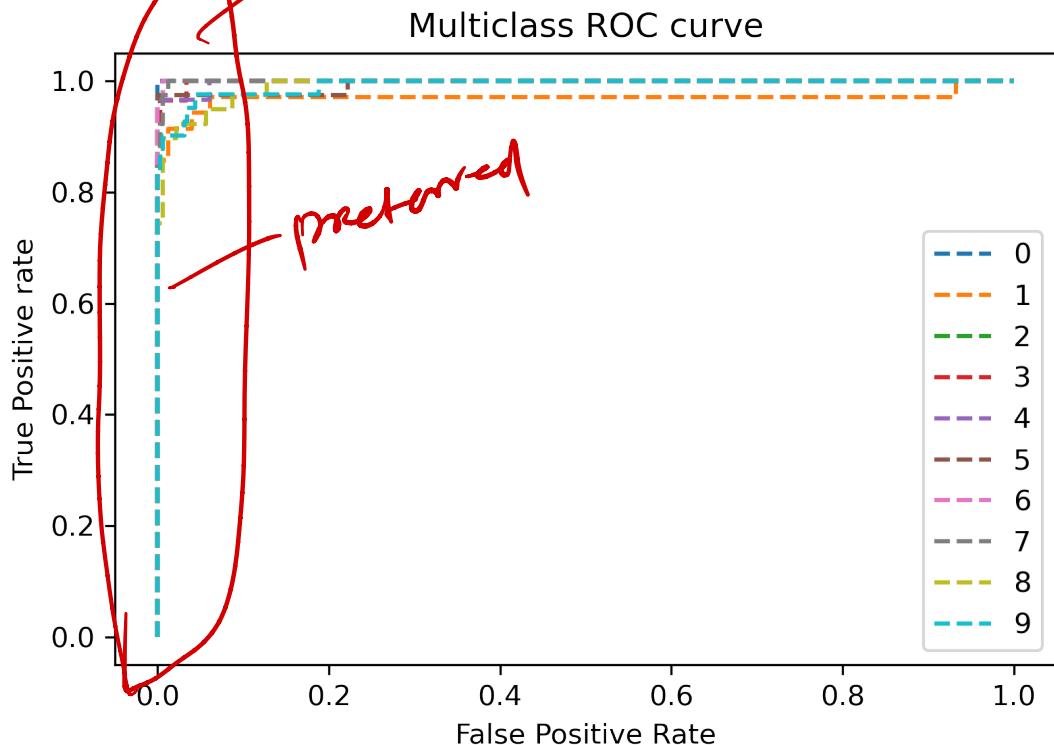
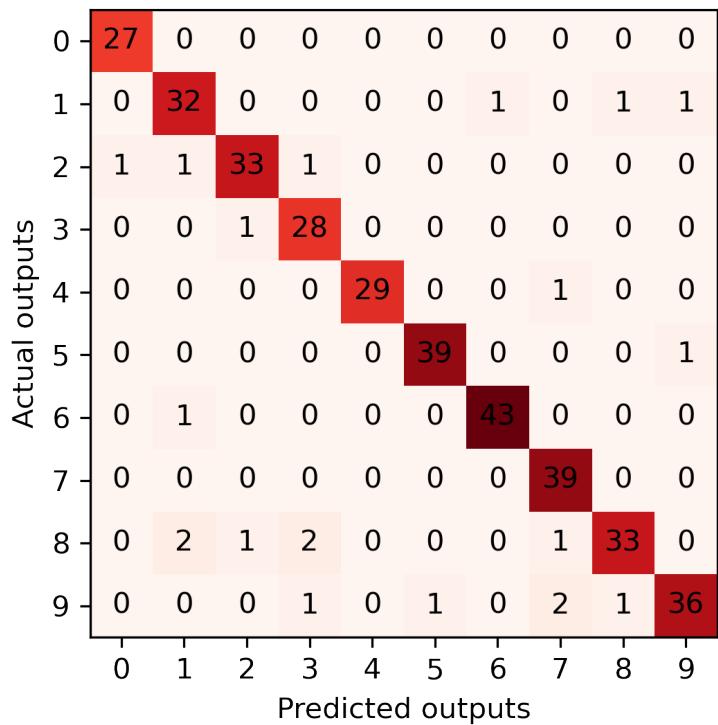
Digit Classification

MNIST

60k train
10k test

Stephan
is prepared

- Image recognition of handwritten digits.



► Support Vector Network

VAPNIL

NIPS
kernel ≈ 2
trick

- One of the popular Machine Learning algorithms.
- Supervised Machine Learning for
 - Linear and non-linear classification
 - Regression and classification problem
- Work well with small or medium sized datasets.

(a)

- Maximal Margin Classifier **HARD MARGIN**

b

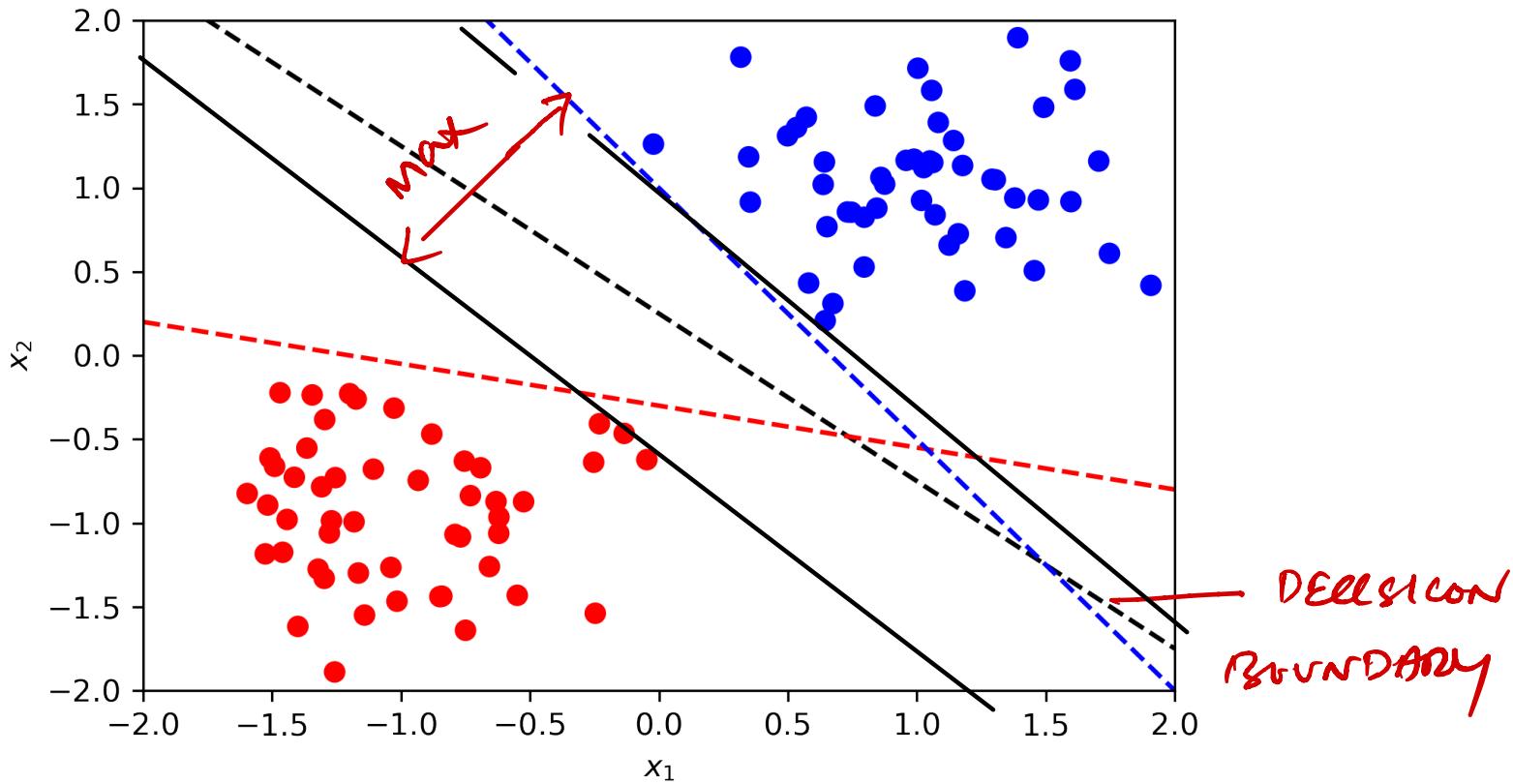
- Support Vector Classifier **SOFT MARGIN**

c

- Support Vector Machine **KERNELS TRICKS**

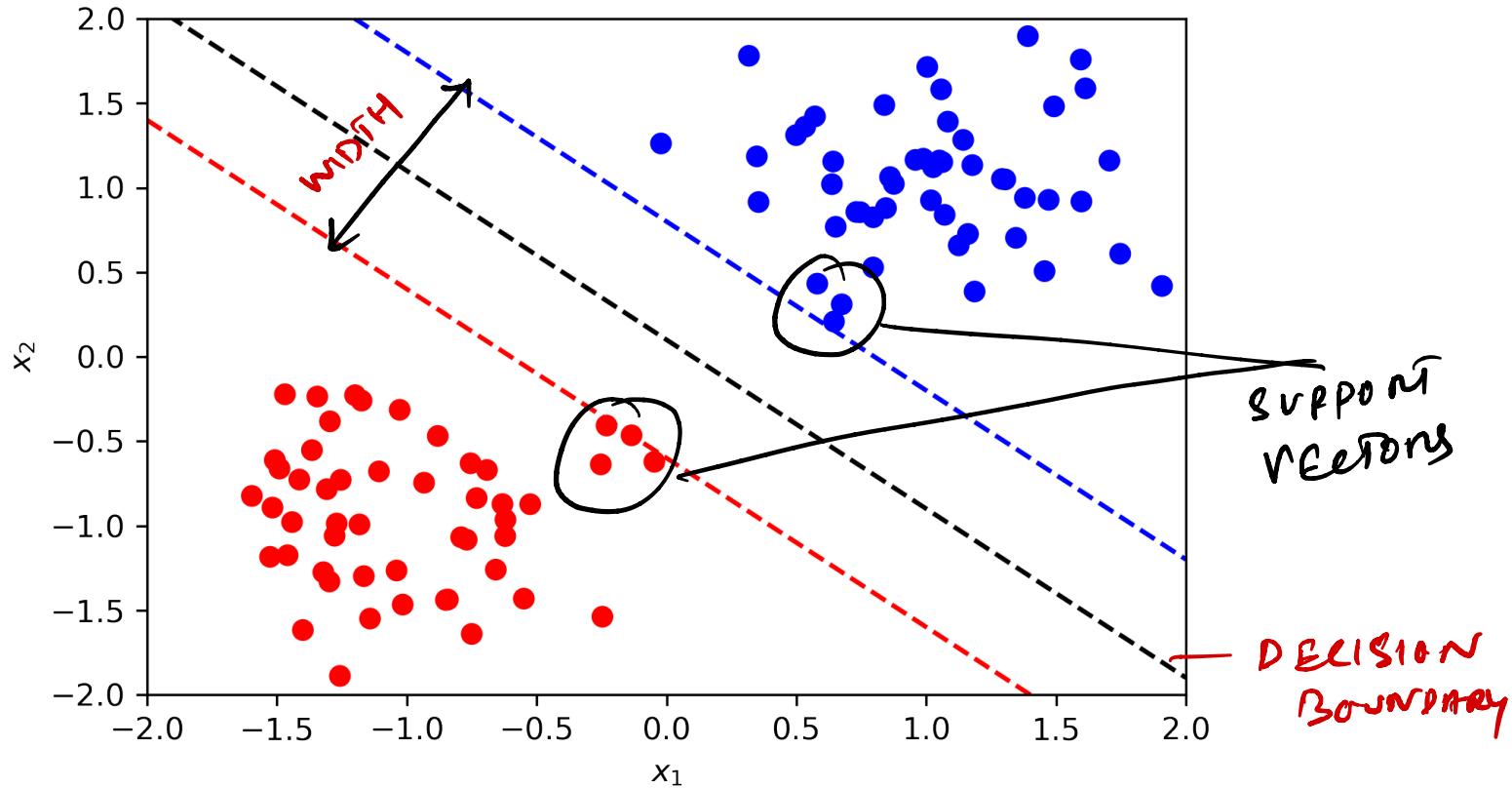


The Hyperplane



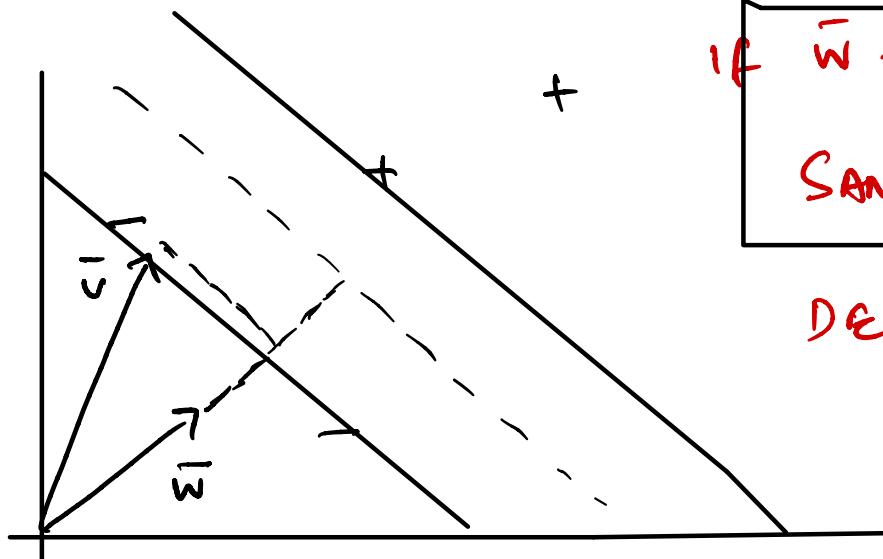


Maximal Margin





Maximal Margin



$$\bar{w} \cdot \bar{v} \geq c \quad c = -b$$

if $\bar{w} \cdot \bar{v} + b > 0$, Then true
SAMPLE

DECISION RULE - ①

Ⓐ

$$\bar{w} \cdot \bar{x}_+ + b \geq 1$$

$$\bar{w} \cdot \bar{x}_- + b \leq -1$$

ⓐ MATHEMATICAL CONVENTION
we use the same fn.

y_i s.t $y_i = +1$ for + samples
 $= -1$ for - samples.

ⓑ MARGIN - BASED APPROXIMATION
from Ⓩ

► Maximal Margin

$$y_i (\bar{w} \cdot \bar{x}_{+i} + b) \geq 1 \quad y_i = 1 \quad 1 \geq 1$$

$$y_i (\bar{w} \cdot \bar{x}_{-i} + b) \geq -1 \quad y_i = -1 \quad -1 \geq -1$$

$$y_i (\bar{w} \cdot \bar{x}_i + b) - 1 \geq 0$$

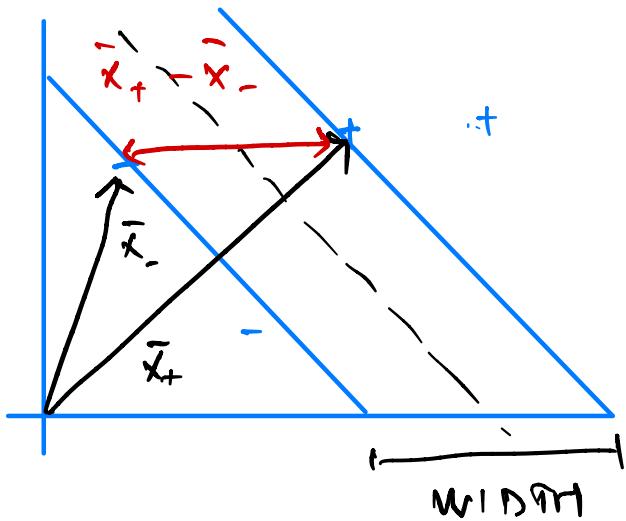
For support vectors, we have

$$y_i (\bar{w} \cdot \bar{x}_i + b) - 1 = 0$$

- ②



Maximal Margin



$$\begin{aligned}
 \text{WIDTH} &= \frac{(\bar{x}_+ - \bar{x}_-). \bar{w}}{\|\bar{w}\|} \\
 &= \frac{\bar{x}_+ \cdot \bar{w} - \bar{x}_- \cdot \bar{w}}{\|\bar{w}\|} \\
 &= \frac{1 - b - (-1 - b)}{\|\bar{w}\|} \\
 \text{WIDTH} &= \frac{2}{\|\bar{w}\|}
 \end{aligned}$$

$$\begin{aligned}
 \text{MAX } \frac{2}{\|\bar{w}\|} &\approx \text{MIN } \frac{\|\bar{w}\|}{2} \stackrel{\text{MIN }}{\approx} \frac{1}{2} \|\bar{w}\|^2 \quad (3)
 \end{aligned}$$



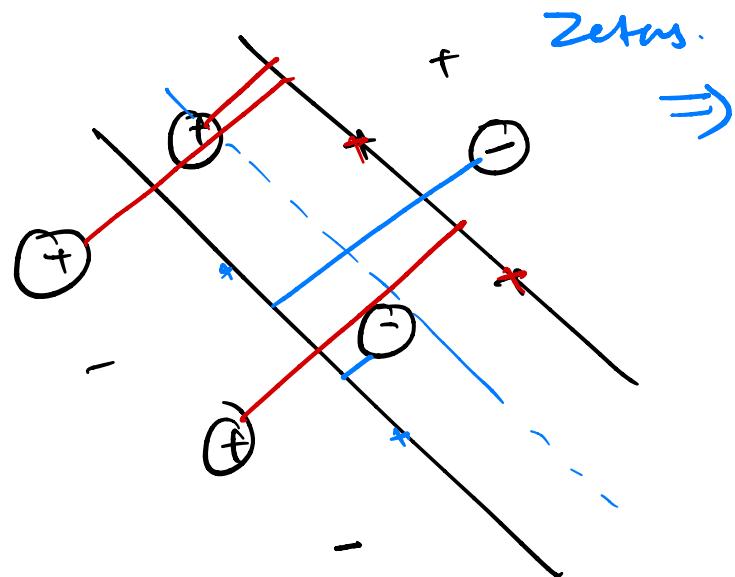
Hard Margin Classifier

- Maximise the width of the margin is equal to

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

s.t

$$y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0$$



Support Vector Classifier

$$\text{SVC} = \frac{\text{Margin}}{\text{Error}} + \frac{\text{Classification Error}}{\text{Margin}}$$

$$\text{Err} = \frac{\min}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \zeta_i$$

- Also called as soft margin classifier

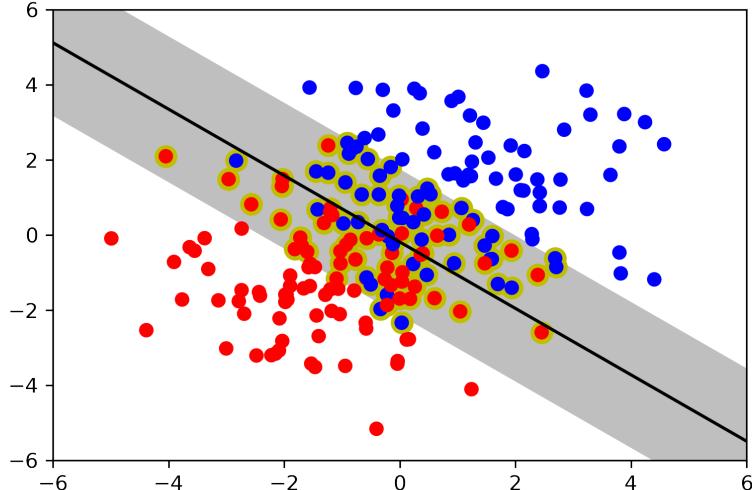
$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \zeta_i$$

hp zeta

subject to,

$$y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq (1 - \zeta_i)$$

where $\zeta_i \geq 0$ and C is a non-negative tuning parameter



- $\zeta_i = 0$ when vector does not reside inside the margin and is on the correct side of the hyperplane.
- $0 > \zeta_i < 1$ when vector resides inside the margin.
- $\zeta_i = 1$ when vector resides on the hyperplane.
- $\zeta_i > 1$ when vector is on the wrong side of the hyperplane.

► Support Vector Classifier

$$\max(0, 1 - m)$$

- Choose slack variables as

$$\zeta_i = \underbrace{\max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))}_{}$$

- We can rewrite SVC formulation as a simplified optimization problem

$$\min \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{red line}} + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$

- SVC (soft margin) can be formulated with hinge loss

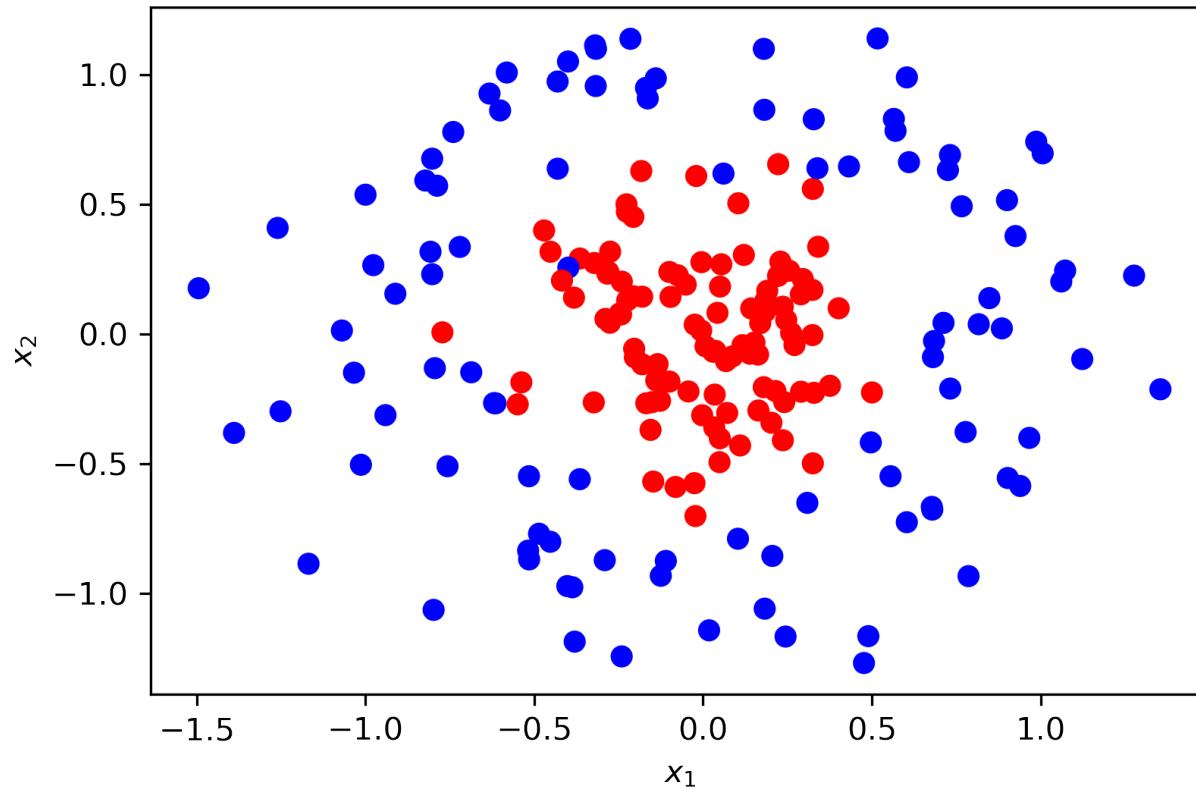
► Support Vector Machines

- Extension of support vector classifier.
- Enlarge features space using kernels.
- Accommodate non-linear boundaries between the classes.
- Implicitly create features to project into higher dimensional space.

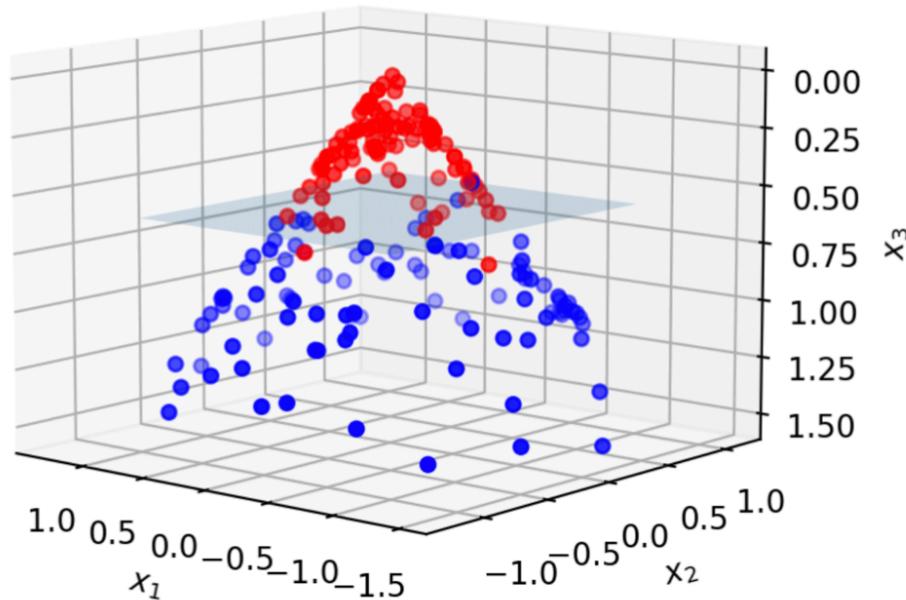
- Hyperplane divide the classes in the higher dimension space.
- Decision boundary on the original dimensional space is non-linear.



Support Vector Machines



► Support Vector Machines



Lagrange Dual Formulation

- In finding extremum of a function with constraints, we can use Lagrange multiplier.
- Lagrange multiplier will give us new expression which we can maximise or minimise without thinking about the constraints.
- Basic idea is to convert a constrained problem into a form such that derivative (can still be applied) of an unconstrained problems.
- Finding the extremum means, find ∇ and set them to 0.

Lagrangian Dual Formulation

$$L = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i (\bar{w} \cdot \bar{x}_i + b) - 1]$$

$$\begin{aligned}\frac{\partial L}{\partial \bar{w}} &= \bar{w} - \sum \alpha_i y_i \bar{x}_i = 0 \\ \Rightarrow \bar{w} &= \sum \alpha_i y_i \bar{x}_i \quad \textcircled{1}\end{aligned}$$

$$\frac{\partial L}{\partial b} = - \sum \alpha_i y_i = 0 \Rightarrow \sum \alpha_i y_i = 0 \quad \textcircled{2}$$

Plugging the values in 2, we have

$$L = \frac{1}{2} \left(\sum \alpha_i y_i \bar{x}_i \right) \left(\sum \alpha_j y_j \bar{x}_j \right) - \sum \alpha_i y_i \bar{x}_i \cdot \left(\sum \alpha_j y_j \bar{x}_j \right)$$

$$- \sum_{\substack{\alpha_i y_i \\ = 0}} b + \sum \alpha_i$$

Lagrangian Dual Formulation

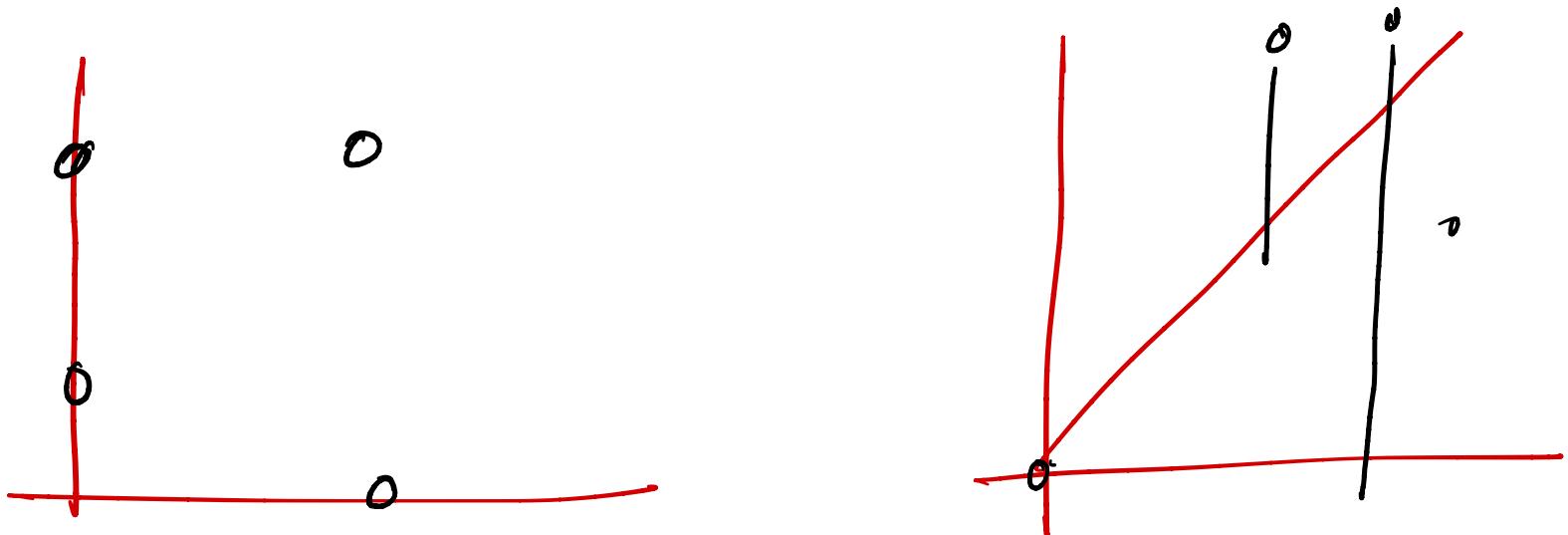
$$L = \sum \lambda_i - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \bar{x}_i \cdot \bar{x}_j \quad (4)$$

Decision rule

$$\sum \lambda_i y_i \bar{x}_i \bar{v}_i + b \geq 0.$$

DEPENDS ON THE
Dot product of the
samples.

Lagrangian Dual Formulation



transformation $\phi(\bar{x})$

Kernel $\phi(\bar{x}_i), \phi(\bar{x}_j)$

if we have f_m , $K(x_i, x_j) = \phi(\bar{x}_i) \cdot \phi(\bar{x}_j)$

Lagrangian Dual Formulation

- Step 1 : Apply method of Lagrange multipliers

$$\min \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to constraints} \quad y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0$$

- Step 2 : Differentiate vectors

$$\frac{\partial \|\mathbf{w}\|^2}{\partial \mathbf{w}} = 2\mathbf{w} \quad \text{and} \quad \frac{\partial \mathbf{x} \cdot \mathbf{w}}{\partial \mathbf{w}} = \mathbf{x}$$

- Step 3 : Find derivatives of the Lagrangian L

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_i^l \alpha_i y_i \mathbf{x}_i = 0$$

$$\frac{\partial L}{\partial b} = \sum_i^l \alpha_i y_i = 0$$

Lagrangian Dual Formulation

- Step 4 : Do the algebra

$$L = \sum_i^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

- Minimisation depends only on dot products of sample vectors

Kernel Trick

- Kernel is a function that creates the implicit mapping to a higher dimensional space.
- Transform non-linear on \mathbf{x} into linear on $\phi(\mathbf{x})$.
- No explicit feature generation (hence no increase in the size of dataset).
- Solution involves only the inner products of the observations of the general form $K(x_i, x_j)$ which is equal to $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$
- Choice of kernels for best fit hyperplane.

choices of HP

- Linear $\frac{(\mathbf{x}_i \cdot \mathbf{x}_j)}{(r + \mathbf{x}_i \cdot \mathbf{x}_j)^d}$
- Polynomial $e^{(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)}$
- Gaussian RBF $e^{(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)}$
- Sigmoid $\tanh(\gamma \mathbf{x}_i \cdot \mathbf{x}_j + r)$

SVM for Regression

- Train model that includes as many vectors as possible inside the margin.
- No penalty for vectors that reside inside the margin.
- Penalise for vectors that lie outside the margin.
- Cost function would be the well-loss.
- Optimisation of width takes the form of

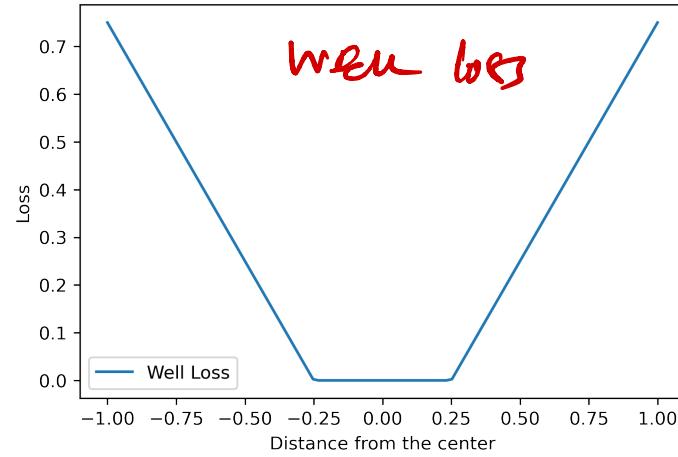
$$\min \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n (\zeta_i + \zeta_i^*)$$

subject to,

$$y_i - \mathbf{w} \cdot \mathbf{x}_i - b \leq \epsilon + \zeta_i$$

$$\mathbf{w} \cdot \mathbf{x}_i + b - y_i \leq \epsilon + \zeta_i^*$$

$$\zeta_i, \zeta_i^* \geq 0$$



- One of the simplest ML algorithms used for both regression and classification.
- Lazy algorithm as it doesn't technically train a model to make a predictions.
- Observed is predicted to be in the class of that of the largest proportion of k - nearest observations.
- Direct attempt at approximating the conditional expectation using actual data.

- Non-parametric method and it does not assume any function as there are no parameters to be estimated.
 - k close to 1 gives the most flexibility (low bias), but highest variability (high variance)
 - Large k gives the least flexibility (high bias), but lowest variability (low variance)
 - $k = N$ will assign all new test observations to a single class.
- Odd number k is preferred for binary classification.
- Best way to select k is through
 - Cross validation
 - Elbow method

► KNN Regression

- For regression, the predicted value is the mean of K neighbours and the estimator is given as

$$\hat{f}(x) = \text{Average} [y_i | x_i \in \mathcal{N}_k(x)]$$

where, $\mathcal{N}_k(x)$ is a neighbourhood of x containing the k closest observations.

KNN Classification

- For classification, the predicted value is the class with the plurality, i.e., the class most represented among the neighbours.
- This is equivalent to taking a majority vote among the k nearest neighbours.
- For each class $j = 1, \dots, K$, we then compute the empirical (conditional) probability

$$Pr(G = j | X = x_0) = \frac{1}{k} \sum_{i \in \mathcal{N}_k(x)} I(y_i = j)$$

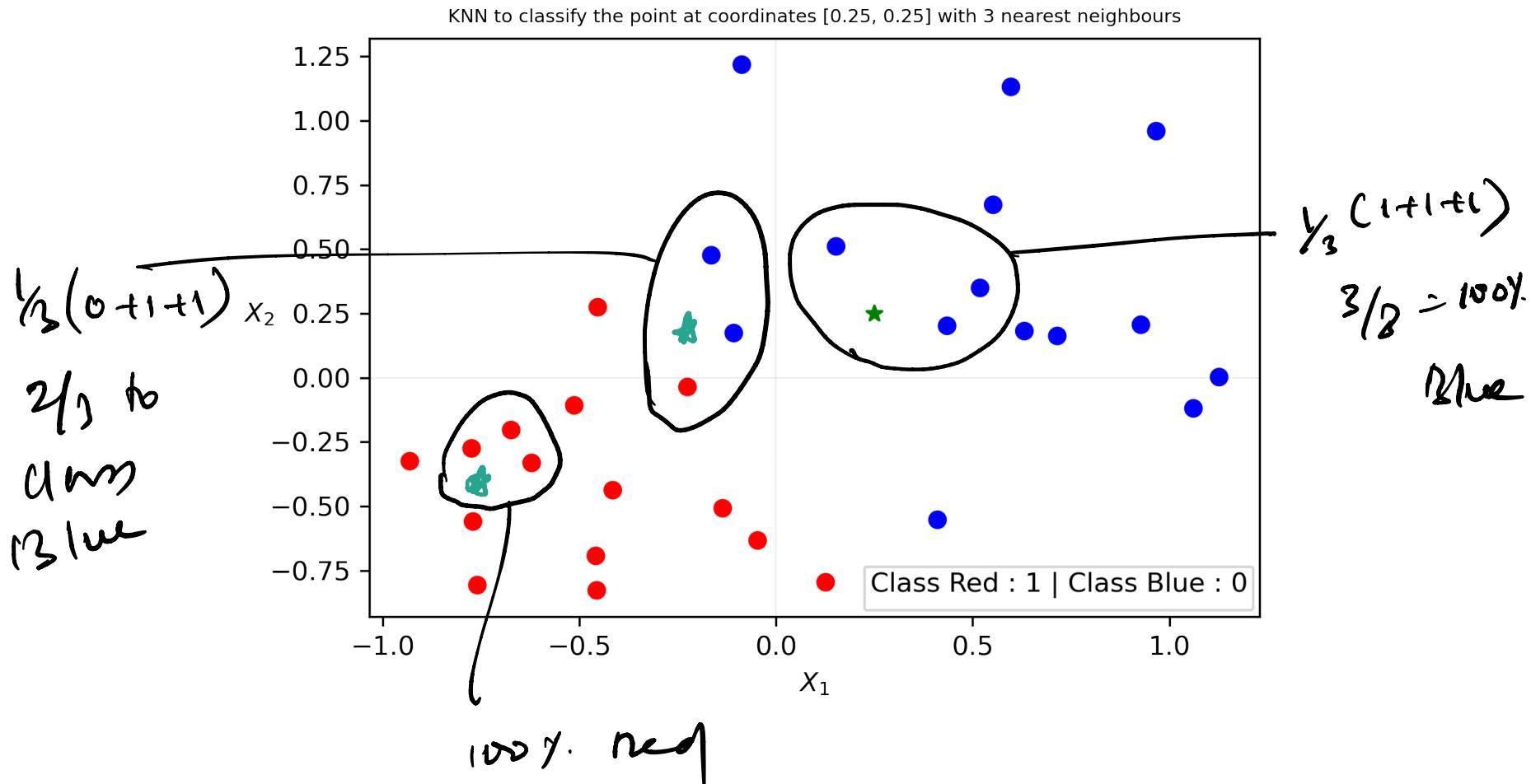
$I(y_i = j)$

and assign the observation to the class with the highest empirical probability. Here, $I(y_i = j)$ is the indicator function returning 1 if $y_i = j$ and 0 otherwise.

- KNN output would be biased if the dataset is imbalanced or on presence of outliers.

KNN Example

$K=3$



The Three Steps

- Once, k is specified, finding nearest neighbours is a 3-step process.

- Compute the distance, usually Euclidian

L_2

- Sort by ascending distance to the find the k NN

- Compute the average or probability of k NN observations

L_{reg}

L_{clf}

► Nearest Neighbours

- KNN uses distance as the metric in determining its neighbours.
- Distance can be thought of as a measure of similarity.
- Generalised distance metric is called the Minkowski distance

$$d = \left(\sum_{n=i}^n |x_i - y_i|^p \right)^{1/p}$$

P=2 default

where x_i and y_i are the two observations for which distance d is being calculated with a hyper parameter, integer p .

- $p = 1$ is the Manhattan distance
- $p = 2$ is the Euclidean distance

KNN Example

	X1	X2	Y	Class	Distance
0	-0.621282	-0.331191	1	Red	1.047337
1	0.410545	-0.552438	0	Blue	0.818341
2	-0.046838	-0.632821	1	Red	0.931389
3	-0.455832	-0.826784	1	Red	1.287502
4	-0.087054	1.217604	0	Blue	1.024628
5	-0.135597	-0.507716	1	Red	0.850188
6	0.518747	0.349231	0	Blue	0.286482
7	1.061250	-0.119117	0	Blue	0.891276
8	-0.225160	-0.035936	1	Red	0.554560
9	0.152672	0.510864	0	Blue	0.278429
10	-0.453268	0.274211	1	Red	0.703685
11	0.631439	0.181521	0	Blue	0.387537
12	-0.513032	-0.107380	1	Red	0.842579
13	0.597201	1.131221	0	Blue	0.947153
14	-0.164562	0.476409	0	Blue	0.472359
15	0.434864	0.202033	0	Blue	0.190986
16	-0.415179	-0.437181	1	Red	0.956390
17	0.966692	0.959089	0	Blue	1.008194
18	-0.770452	-0.559163	1	Red	1.302331
19	-0.931161	-0.324453	1	Red	1.313445
20	-0.107496	0.173775	0	Blue	0.365532
21	-0.458115	-0.692723	1	Red	1.179048
22	0.927720	0.206114	0	Blue	0.679139
23	-0.673497	-0.203217	1	Red	1.028714
24	-0.759423	-0.806535	1	Red	1.461233
25	0.126097	-0.795368	1	Red	1.052686
26	0.552044	0.672569	0	Blue	0.519418
27	-0.774020	-0.274782	1	Red	1.150658
28	1.126512	0.002594	0	Blue	0.910760
29	0.715034	0.162489	0	Blue	0.473197

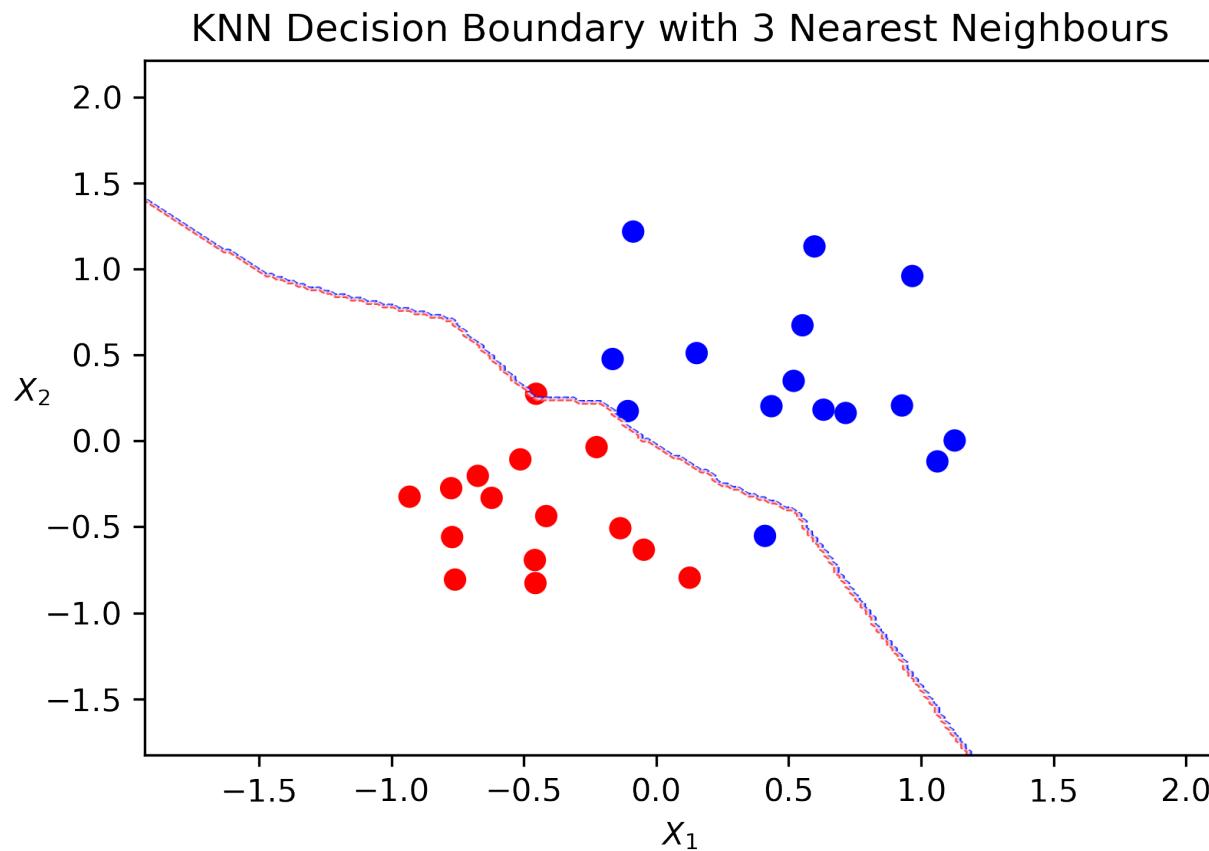
	X1	X2	Y	Class	Distance
15	0.434864	0.202033	0	Blue	0.190986
9	0.152672	0.510864	0	Blue	0.278429
6	0.518747	0.349231	0	Blue	0.286482
20	-0.107496	0.173775	0	Blue	0.365532
11	0.631439	0.181521	0	Blue	0.387537
14	-0.164562	0.476409	0	Blue	0.472359
29	0.715034	0.162489	0	Blue	0.473197
26	0.552044	0.672569	0	Blue	0.519418
8	-0.225160	-0.035936	1	Red	0.554560
22	0.927720	0.206114	0	Blue	0.679139
10	-0.453268	0.274211	1	Red	0.703685
1	0.410545	-0.552438	0	Blue	0.818341
12	-0.513032	-0.107380	1	Red	0.842579
5	-0.135597	-0.507716	1	Red	0.850188
7	1.061250	-0.119117	0	Blue	0.891276
28	1.126512	0.002594	0	Blue	0.910760
2	-0.468383	-0.632821	1	Red	0.931389
13	0.597201	1.131221	0	Blue	0.947153
16	-0.415179	-0.437181	1	Red	0.956390
17	0.966692	0.959089	0	Blue	1.008194
4	-0.087054	1.217604	0	Blue	1.024628
23	-0.673497	-0.203217	1	Red	1.028714
0	-0.621282	-0.331191	1	Red	1.047337
25	0.126097	-0.795368	1	Red	1.052686
27	-0.774020	-0.274782	1	Red	1.150658
21	-0.458115	-0.692723	1	Red	1.179048
3	-0.455832	-0.826784	1	Red	1.287502
18	-0.770452	-0.559163	1	Red	1.302331
19	-0.931161	-0.324453	1	Red	1.313445
24	-0.759423	-0.806535	1	Red	1.461233

Classification

[0.0, 100.0]

Blue

KNN Example



▶ Hyper-parameter Tuning

- Parameters that are not directly learnt within estimators.
- Passed as arguments to the constructor of the estimator classes (in scikit-learn).
- Recommended to search the hyper-parameter space for the best cross validation score.
- Typical examples include tuning
 - Alpha for Lasso, Ridge
 - C, kernel and gamma for SVC
 - n_neighbors for KNN

► Hyper-parameter Tuning

- Some of the popular approaches to hyper parameter tuning are
 - Manual
 - Grid Search
 - Random Search



GridSearch

- The conventional way of performing hyper-parameter optimisation.
- Exhaustive search through a manually specified subset of the hyper-parameter space of a learning algorithm.
- Algorithm must be guided by some performance metric
 - cross-validation on the training set
 - evaluation on a validation set
- Implements a “fit” and a “score” method among other methods.
- The parameters of the estimator used to apply these methods are optimised by cross-validated grid-search over a parameter grid.



RandomSearch

- Randomised search narrow down the results.
- Random combinations of the hyper-parameters are used to find the best parameters.
- Yields high variance during computing
- Works best under the assumption that not all hyperparameters are equally important.
- Given the random search pattern, changes of finding the optimal parameter are comparatively higher.



References

- Scott Fortmann-Roe (2012) , Understanding the Bias-Variance Tradeoff
- Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani (2013), An Introduction to Statistical Learning
- Sebastian Raschka (2015), Python Machine Learning
- Francois Chollet (2017), Deep Learning with Python
- Aurelien Geron (2019), Hands-on Machine Learning with Scikit-Learn, Keras & TensorFlow
- Hariom Tatsat, Sahil Puri, and Brad Lookabaugh (2020), Machine Learning & Data Science Blueprints for Finance
- Motunrayo Olugbenga (2022), Balanced Accuracy: When Should You Use It?
- Scikit-learn User Guide
- Google Developers, Introduction to Machine Learning

Note: Some of the materials from the above resources are adapted for these notes under [CC BY-SA 4.0](#). For detailed interpretation on the subject, refer to the above resources.