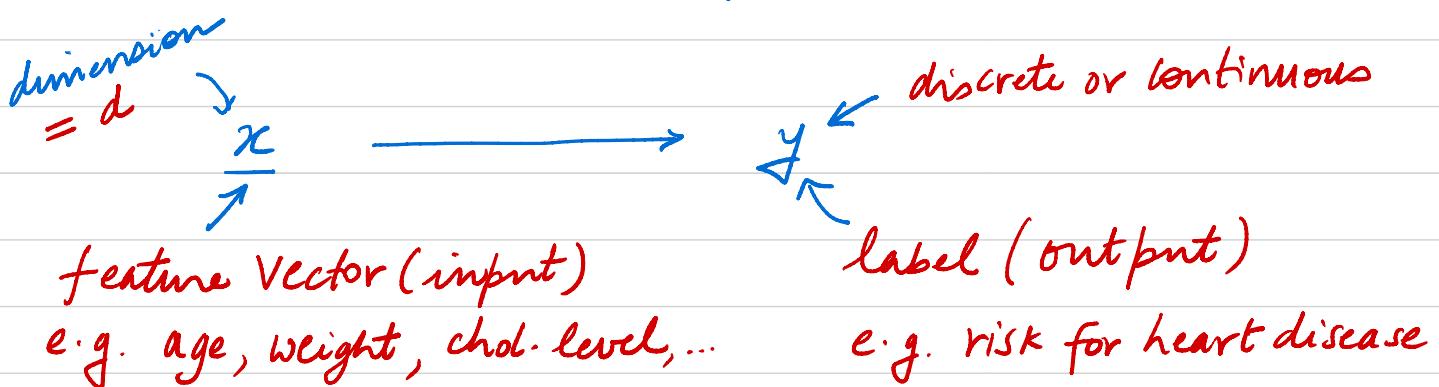


Machine Learning

Goal To learn mapping from \underline{x} to \hat{y}

Supervised Learning Use labelled training data

$$\mathcal{T} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$$

Classification Label y takes on finite # values

If only two values $\{-1, 1\}$ - binary classification

If M values ($M > 2$) - M -ary classification

Classifier uses training data \mathcal{T} to form function f

f is used to predict labels for \underline{x} outside \mathcal{T}

$$\hat{y} = f(\underbrace{\underline{x}}_{\text{feature vector } \notin \mathcal{T}})$$

Measuring classifier Performance

Loss function $L(\hat{y}, y) = L(f(\underline{x}), y)$

$$0-1 \text{ Loss : } L(\hat{y}, y) = \begin{cases} 1 & \text{if } \hat{y} \neq y \\ 0 & \text{if } \hat{y} = y \end{cases}$$

$$\mathcal{T} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$$

Training Error

$$Err_{\text{train}}(\mathcal{T}) = \frac{\sum_{i=1}^N L(f(\underline{x}_i), y_i)}{N}$$

average loss for classifier f over training set \mathcal{T}

For 0-1 loss,

$$Err_{\text{train}}(\mathcal{T}) = \frac{\# \text{ data in } \mathcal{T} \text{ labelled incorrectly by } f}{N}$$

It is easy to design f to have $Err_{\text{train}}(\mathcal{T}) = 0$

Example $y \in \{-1, +1\}$, 0-1 loss

$$f(\tilde{\underline{x}}) = \begin{cases} \tilde{y} & \text{if } (\tilde{\underline{x}}, \tilde{y}) \in \mathcal{T} \text{ OVERFITTING} \\ -1 & \text{otherwise} \end{cases}$$

But f may be terrible on data outside \mathcal{T} .

Prediction Error and Validation Error

$$Err_{\text{pred}} = E[L(f(\underline{x}), y) \mid f \text{ is trained on } \mathcal{T}]$$

\uparrow "average" over all possible (\underline{x}, y)

Estimate using another labelled data set independent
of \mathcal{T} , Validation set $\mathcal{V} = \{(\tilde{\underline{x}}_1, \tilde{y}_1), \dots, (\tilde{\underline{x}}_V, \tilde{y}_V)\}$

$$\hat{Err}_{\text{pred}}(\mathcal{V}) = \frac{\sum_{i=1}^V L(f(\tilde{\underline{x}}_i), \tilde{y}_i)}{V} = Err_{\text{val}}(\mathcal{V})$$

For 0-1 loss,

$$Err_{\text{val}}(\mathcal{V}) = \frac{\# \text{ data in } \mathcal{V} \text{ labelled incorrectly by } f}{V}$$

Displaying Error Performance

- Focus on 0-1 loss
- Errtrain and Errval measure overall error rate on training and validation datasets.
- More detailed measures of error may be needed for some applications (e.g. cancer detection)

Confusion Matrix

		True Class	
		+1	-1
Classifier Output	+1	TP True Positives	FP False Positives
	-1	FN False Negatives	TN True Negatives
total		P	N

Focus on special case of binary classification

$$P = \# \text{ data with } +1 \text{ label}$$

$$N = \# \text{ data with } -1 \text{ label}$$

$$N_{\text{tot}} = N + P$$

$$\text{True Positive Rate } TPR = \frac{TP}{P}$$

$$\text{False Positive Rate } FPR = \frac{FP}{N}$$

$$\text{True Negative Rate} = \text{specificity} = 1 - FPR$$

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

$$\text{Accuracy} = \frac{TP + TN}{N_{\text{tot}}} = 1 - Err$$

K-Nearest Neighbor classifier

Distance / Dissimilarity Measure

$\Delta(\underline{x}_1, \underline{x}_2)$ measures distance between feature vectors \underline{x}_1 and \underline{x}_2

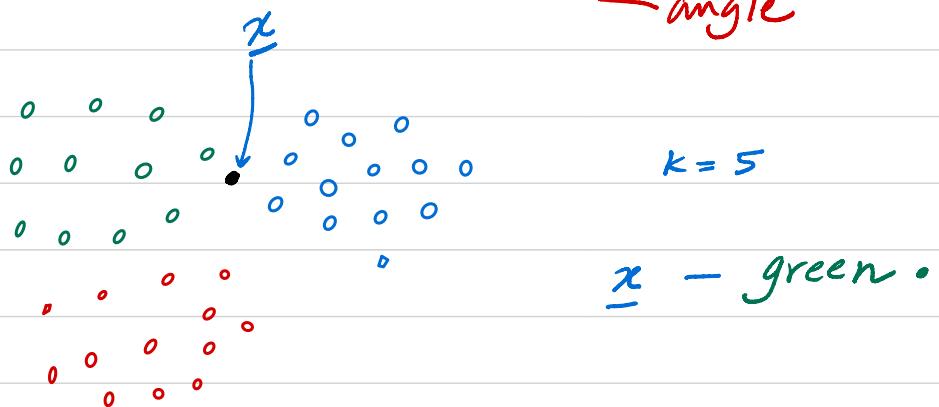
Examples : 1) $\Delta(\underline{x}_1, \underline{x}_2) = \|\underline{x}_1 - \underline{x}_2\|$

$$= \sqrt{(\underline{x}_1 - \underline{x}_2)^T (\underline{x}_1 - \underline{x}_2)}$$

2) $\Delta(\underline{x}_1, \underline{x}_2) = 1 - \cos(L(\underline{x}_1, \underline{x}_2))$

↑ angle

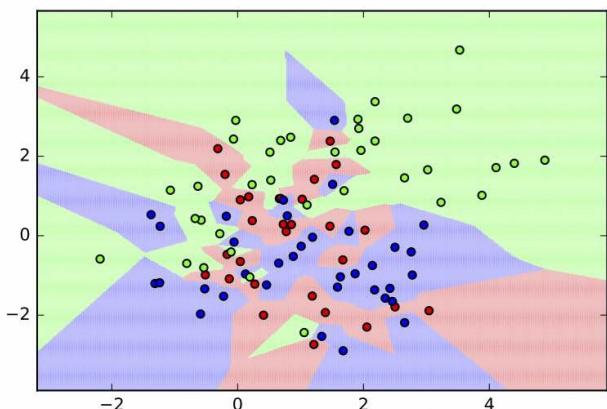
K-NN



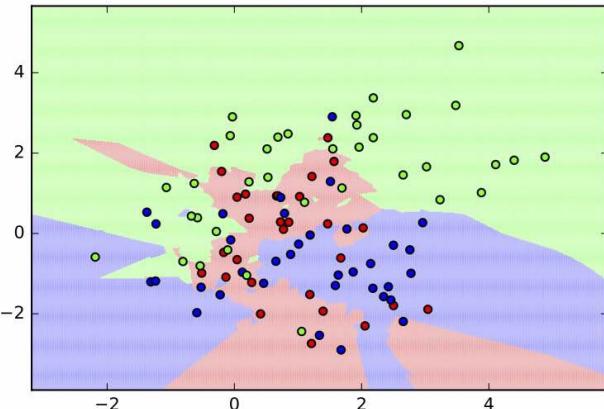
For each \underline{x} pick k nearest neighbors using $\Delta(\cdot, \cdot)$ and choose label as most popular label.

For $k=1$, algorithm is simply called NN classifier

NN



5-NN



Binary Bayes Classifier $\underline{x} \longrightarrow y \in \{+1, -1\}$

Suppose we know:

$$\pi_y \triangleq p(y) = P\{y=y\}, y = +1, -1$$

and $p(\underline{x}|y), y = +1, -1$.

Then, binary classification problem is equivalent to binary hypothesis testing:

$$H_0 : \underline{x} \sim p(\underline{x}|-1) \quad \text{priors}$$

$$H_1 : \underline{x} \sim p(\underline{x}|1) \quad \pi_{-1} \text{ and } \pi_1$$

$$\text{Joint distribution } p(\underline{x}, y) = p(\underline{x}|y) \pi_y$$

For 0-1 loss,

$$\begin{aligned} \text{Err}_{\text{pred}} &= P_e \\ &= \pi_{-1} \underbrace{P\{\text{say } 1 \mid -1 \text{ true}\}}_{P_{FA}} + \pi_1 \underbrace{P\{\text{say } -1 \mid 1 \text{ true}\}}_{P_M} \end{aligned}$$

MAP Rule

$$\frac{p(\underline{x}|1)}{p(\underline{x}|-1)} \stackrel{\text{say } 1}{\geq} \frac{\pi_{-1}}{\pi_1} \equiv \pi_1 p(\underline{x}|1) \stackrel{\text{say } 1}{\geq} \pi_{-1} p(\underline{x}|-1)$$

$$\hat{y}_{\text{MAP}} = \arg \max_y \pi_y p(\underline{x}|y)$$

MAP rule is also called Bayes classifier

$$\hat{y}_{\text{Bayes}} = f_{\text{Bayes}}(\underline{x}) = \hat{y}_{\text{MAP}}$$

f_{Bayes} minimizes Err_{pred} - optimal, but...