

Machine Learning

Goal To learn mapping from x to y

Supervised Learning Use training data

$$\mathcal{T} =$$

Classification Label y takes on finite # values

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

If M values ($M > 2$) -

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

Measuring classifier Performance

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

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

Training Error

$$Err_{\text{train}}(\mathcal{T}) =$$

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

For 0-1 loss,

$$Err_{\text{train}}(\mathcal{T}) =$$

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

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

$$f(\tilde{\underline{x}}) =$$

Prediction Error and Validation Error

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

Estimate using another labelled data set independent
of \mathcal{T} ,

$$\hat{Err}_{\text{pred}}(\mathcal{V}) =$$

For 0-1 loss,

$$Err_{\text{val}}(\mathcal{V}) =$$

Displaying Error Performance

- Focus on
- Error rate and Error measure
on training and validation datasets.
- More detailed measures of error may be needed for some applications

Confusion Matrix

		True Class	
		+1	-1
Classifier Output	+1	True Positives	False Positives
	-1	False Negatives	True Negatives
total		P	N

Focus on special case
of binary classification

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

$$N_{\text{tot}} =$$

$$\text{True Positive Rate } TPR =$$

$$\text{False Positive Rate } FPR =$$

$$\text{True Negative Rate} = \text{specificity} =$$

$$\text{Precision} =$$

$$\text{Accuracy} =$$

K-Nearest Neighbor classifier

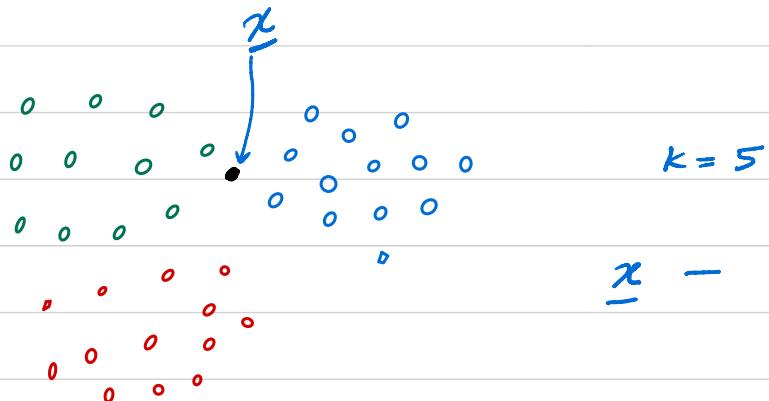
Distance / Dissimilarity Measure

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

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

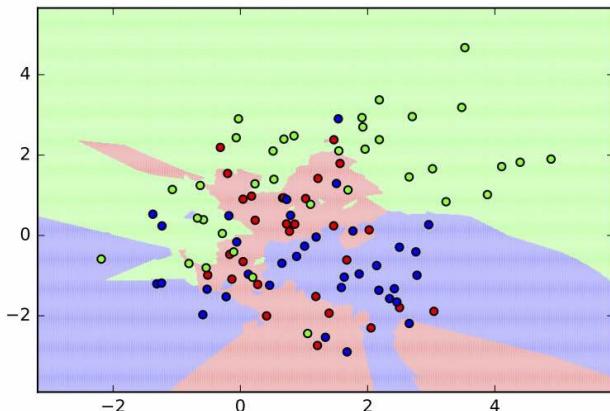
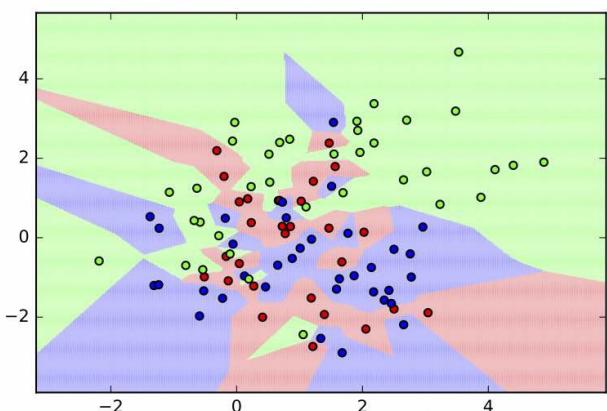
2) $\Delta(\underline{x}_1, \underline{x}_2) =$

k -NN



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

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



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

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:

Joint distribution

For 0-1 loss,

$$\text{Err}_{\text{pred}} =$$

MAP Rule

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

$$\hat{y}_{\text{MAP}} =$$

MAP rule is also called

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

$$f_{\text{Bayes}}$$