

If attribute A is used to split D into v subsets, $\{D_1, D_2, \dots, D_v\}$, the resulting information is

$$\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

Information gain is defined as the difference between the original information (before splitting) and the remaining information (after splitting D by A):

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

PID	Fever	Cough	Sore Throat	Tiredness	Flu
1	no	yes	no	yes	+
2	no	yes	no	no	-
3	mild	yes	no	yes	+
4	yes	mild	no	yes	+
5	yes	no	yes	yes	+
6	yes	no	yes	no	-
7	mild	no	yes	no	+
8	no	mild	no	yes	-
9	no	no	yes	yes	+
10	yes	mild	yes	yes	+
11	no	mild	yes	no	+
12	mild	mild	no	no	+
13	mild	yes	yes	yes	+
14	yes	mild	no	no	-

Disclaimer: Synthetic data

$$\text{Gain}(\text{Fever}) = \text{Info}(D) - \text{Info}_{\text{Fever}}(D) = 0.940 - 0.694 = 0.246 \text{ bits}$$

$$\text{SplitInfo}_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \log_2 \left(\frac{|D_j|}{|D|} \right)$$

$$\text{GainRatio}_A(D) = \frac{\text{Gain}(A)}{\text{SplitInfo}_A(D)}$$

Gini Index:

$$\text{Gini}(D) = 1 - \sum_{i=1}^2 p_i^2 \quad \text{Gini}_A(D) = \frac{|D_1|}{|D|} \text{Gini}(D_1) + \frac{|D_2|}{|D|} \text{Gini}(D_2)$$

$$\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D)$$

Evaluating Classifier Performance:

		predicted class		
actual class		Yes	No	Total
	Yes	TP	FN	P
	No	FP	TN	N
	Total	P'	N'	P+N

$$\text{Accuracy} = \frac{TP+TN}{P+N} \quad \text{Error rate} = \frac{FP+FN}{P+N} = 1 - \text{Accuracy}$$

$$\text{Sensitivity} = \frac{TP}{P} \quad \text{Specificity} = \frac{TN}{N}$$

$$\text{Accuracy} = \text{Sensitivity} \times \left(\frac{P}{P+N} \right) + \text{Specificity} \times \left(\frac{N}{P+N} \right)$$

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

$$\text{Recall} = \frac{TP}{TP+FN} = \frac{TP}{P} = \text{TPR} \quad \frac{FP}{N} = \text{FPR}$$

$$F = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Dissimilarity and Similarity Measures:

1. Minkowski distance: $d(\mathbf{x}, \mathbf{y}) = (\sum_{k=1}^n |x_k - y_k|^r)^{1/r}$
2. Manhattan distance ($r = 1$): $d(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^n |x_k - y_k|$
3. Euclidean distance ($r = 2$): $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$
4. Cosine similarity: $\cos(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$
5. Infinity (Sup) Distance: $d(\mathbf{x}, \mathbf{y}) = \max_{1 \leq j \leq d} |x_j - y_j|$

SVM:

$$y_i = \text{sign}(\mathbf{w} \cdot \mathbf{x}_i + b). \quad d = \frac{2}{\|\mathbf{w}\|}$$

$$\min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} \quad \text{subject to} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N.$$

– Dual optimization problem

$$\mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i, \quad \sum_{i=1}^N \lambda_i y_i = 0.$$

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

Neural Network:

– Activation Functions:

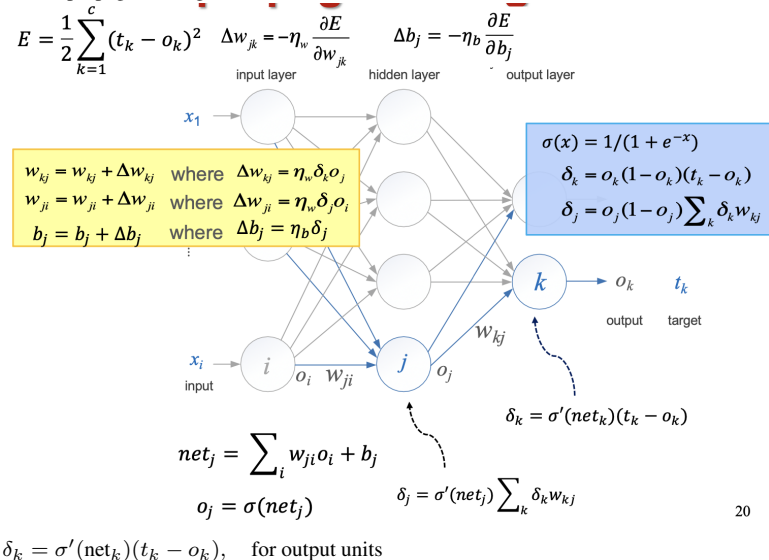
$$\text{Linear: } \sigma(x) = x \quad \text{Sigmoid: } \sigma(x) = \frac{1}{1+e^{-ax}} \quad \text{Tanh: } \sigma(x) = \tanh(\gamma x) = \frac{e^{2\gamma x} - 1}{e^{2\gamma x} + 1}$$

$$\text{Sign: } \sigma(x) = \text{sign}(x) = \begin{cases} +1, & x \geq 0 \\ -1, & x < 0 \end{cases} \quad \text{ReLU: } \sigma(x) = \begin{cases} x, & x \geq 0 \\ 0, & x < 0 \end{cases} = \max(0, x)$$

$$\text{Leaky ReLU: } \sigma(x) = \begin{cases} x, & x \geq 0 \\ ax, & x < 0 \end{cases} = \max(ax, x), \quad \text{where } a \ll 1$$

– Gradient Descent: $\mathbf{w}' = \mathbf{w} - \eta \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$

– Back propagation Algorithm:



$$\delta_j = \sigma'(\text{net}_j) \sum_k \delta_k w_{kj}, \quad \text{for hidden units}$$

$$\frac{\partial E}{\partial w_{kj}} = \frac{\partial E}{\partial o_k} \cdot \frac{\partial o_k}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial w_{kj}} = \frac{\partial E}{\partial o_k} \cdot \sigma'(\text{net}_k) \cdot \frac{\partial (\sum_j w_{kj} o_j + b_k)}{\partial w_{kj}} =$$

$$-(t_k - o_k) \cdot \sigma'(\text{net}_k) \cdot o_j = -\delta_k \cdot o_j$$

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{ji}} = \frac{\partial E}{\partial o_j} \cdot \sigma'(\text{net}_j) \cdot \frac{\partial (\sum_i w_{ji} o_i + b_j)}{\partial w_{ji}} =$$

$$-(\sum_k \delta_k w_{kj}) \cdot \sigma'(\text{net}_j) \cdot o_i = -\delta_j \cdot o_i$$

CNN:

– Stride: steps per moving. – Zero padding : pads the input with zeros around the border. –

Pooling: Max: max one within filter size; Average: average within filter size.

– Regularization: $J'(\mathbf{w}) = J(\mathbf{w}) + \alpha R(\mathbf{w})$

L1 Regularization (LASSO): $R(\mathbf{w}) = \|\mathbf{w}\|_1 = \sum_k |w_k|$

L2 Regularization (Ridge): $R(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_k (w_k)^2$

Elastic Net Regularization: $R(\mathbf{w}) = \|\mathbf{w}\|_1 + \beta \|\mathbf{w}\|_2^2$

Also can be done by early stopping.

Clustering:

– Use Euclidean Distance: $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$

– K-Means: 1. Initialize K random points as K clusters' centers. 2. Assign every point to its cluster by which center it nearest. 3. Calculate each clusters' average again to set as new center. 4. Repeat 2-3, until no points assignments. Initialization influence results.

– HAC: –Single Linkage: the minimum distance between any pair of two data samples from each cluster. –Complete Linkage: the maximum distance between any pair of two data samples from each cluster. –Average Linkage: the average distance between all pairs of two data samples from each cluster. –Centroid Distance: the distance between the means of data samples (i.e., centroids) from each cluster.

• What is the limitations of K-Means algorithm? Need to choose K. Can stuck at poor local minimum. Need good metric. • What are the limitations of HAC algorithm? Memory- and computationally-intensive.

Regression:

$$f_{w,b}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

$$\text{Minimize the } l_2 \text{ loss: } \min_{\mathbf{w}, b} \hat{L}(f_{w,b}) = \min_{\mathbf{w}, b} \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i + b - y_i)^2$$

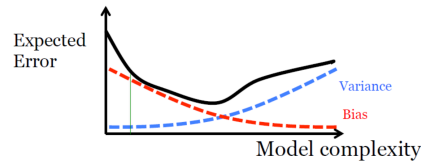
Loss function: mean squared error between $\mathbf{w}^T \mathbf{x}_i + b$ and y_i .

Bias and Variance:

– Bias: Error caused by the wrong assumptions made in the learning algorithms or models.

– Variance: Error due to the learning sensitivity to small fluctuations in the training set.

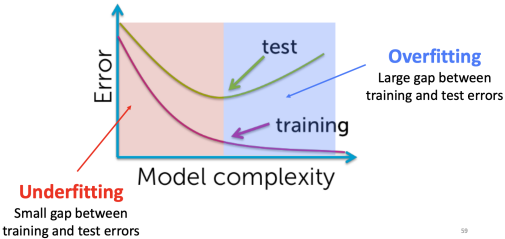
Training a classifier $f_\theta(x)$



Expected error of a classifier \approx bias + variance (+noise)

– Underfitting: High bias and low variance.

– Overfitting: Low bias and high variance.



PCA:

• \mathbf{x}_i (black): the original data.

• \mathbf{v} (red): PCA subspace.

• $(\mathbf{v}^T \mathbf{x}_i) \mathbf{v}$ (blue): projected data.

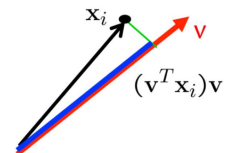
• Green: $\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v}$: projection error (MSE).

• Minimizing MSE \Leftrightarrow Maximizing Projected Variance

• Blue² + green² = black²

• Black is fixed (given data)

• Maximizing blue (variance) is equivalent to Minimizing green (MSE).



Bayes' Theorem:

$$P(A | B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(B) = \sum_i P(B | A_i) P(A_i)$$

Naïve Bayes:

$$P(a_1, \dots, a_d | v_j) = P(a_1 | v_j) \cdots P(a_d | v_j) = \prod_{i=1}^d P(a_i | v_j)$$