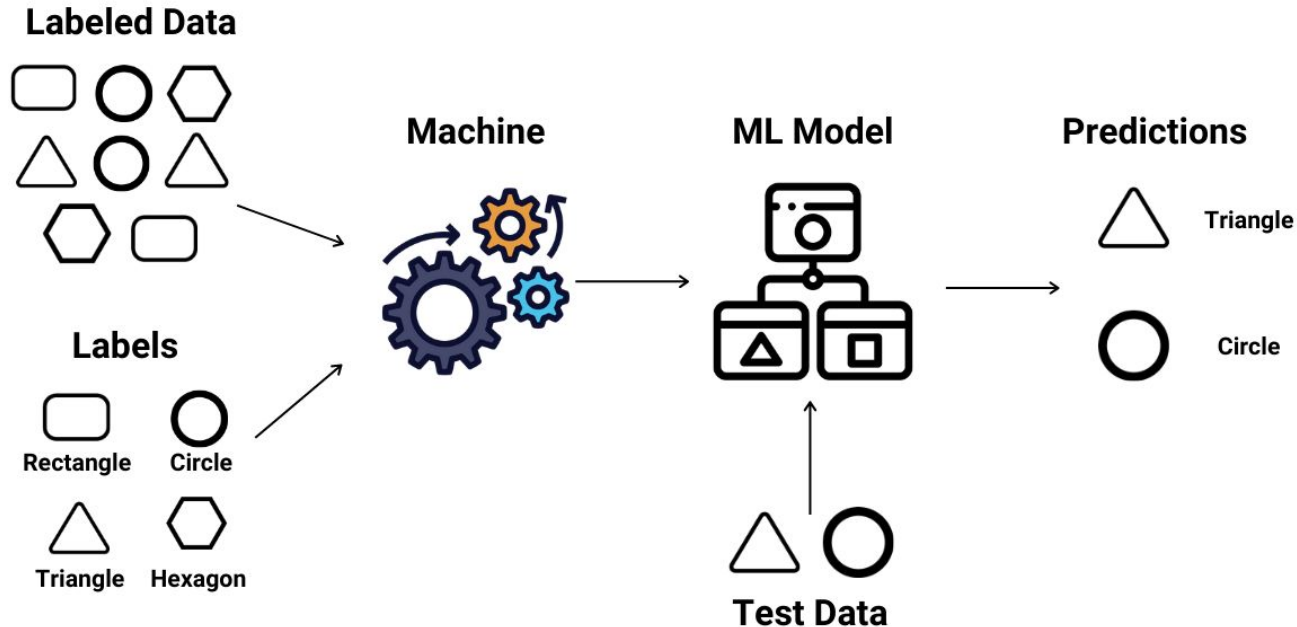


Metrics&Scaling for Classification/Clustering

Supervised Learning Algorithm

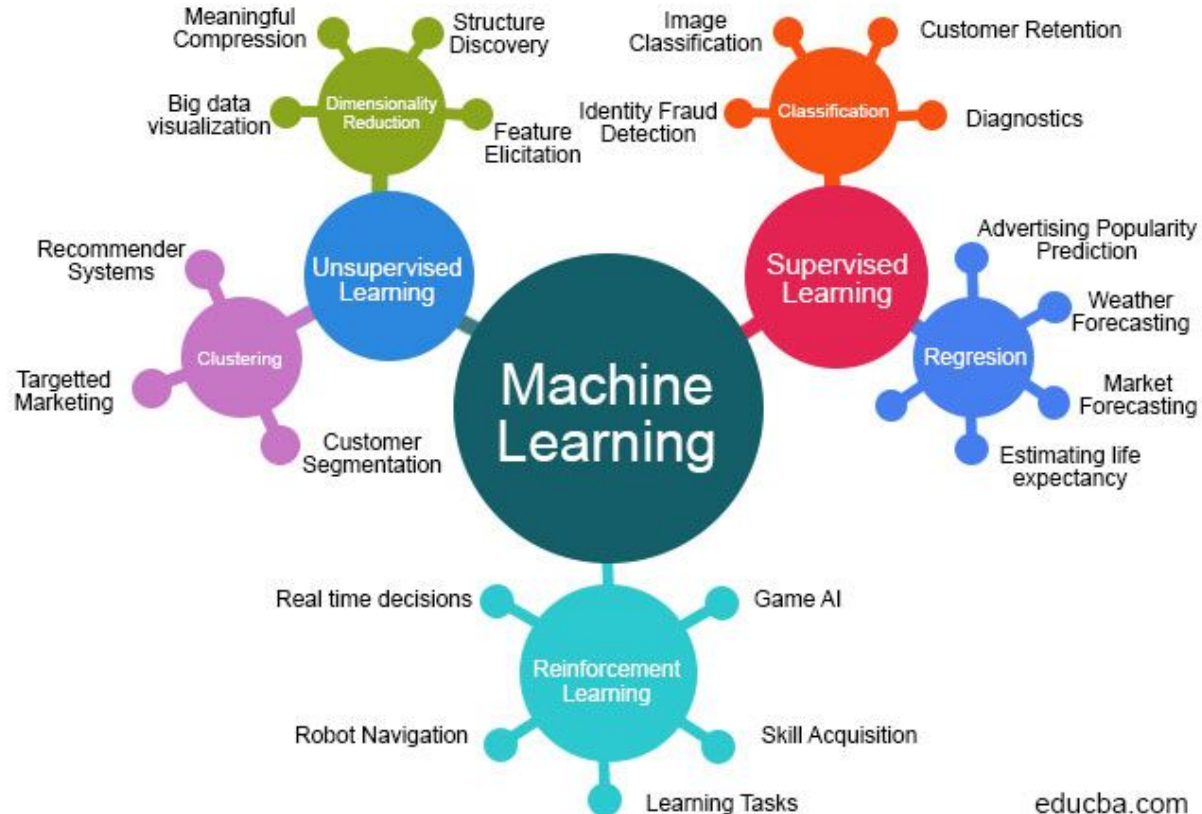


Supervised Learning

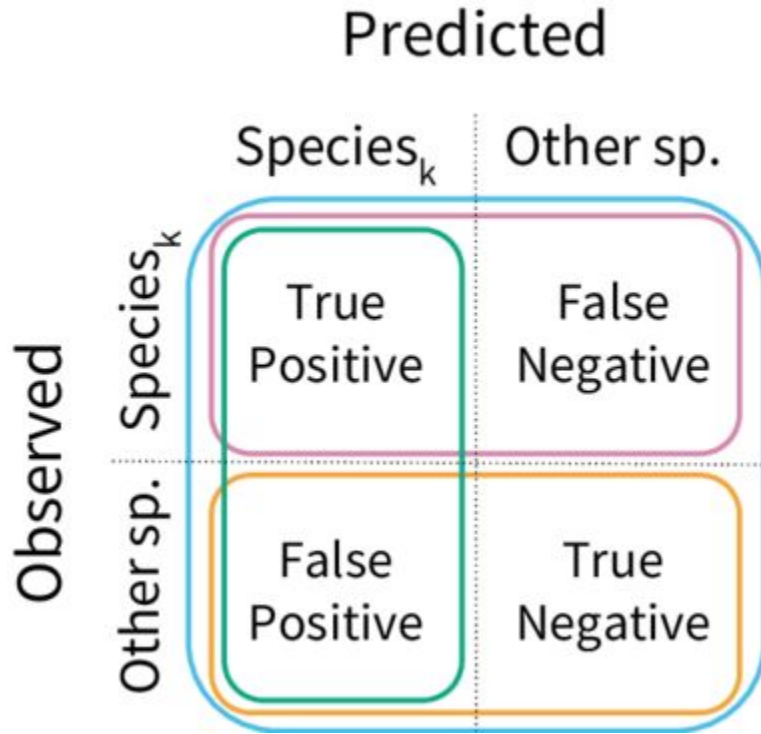



Big Picture


Machine Learning Algorithms





Metrics



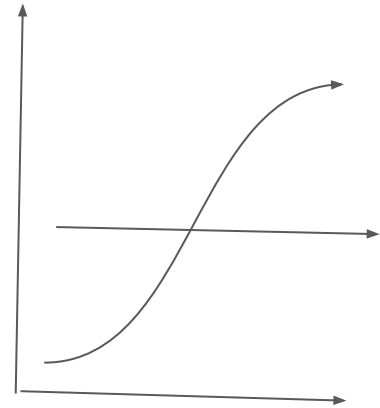
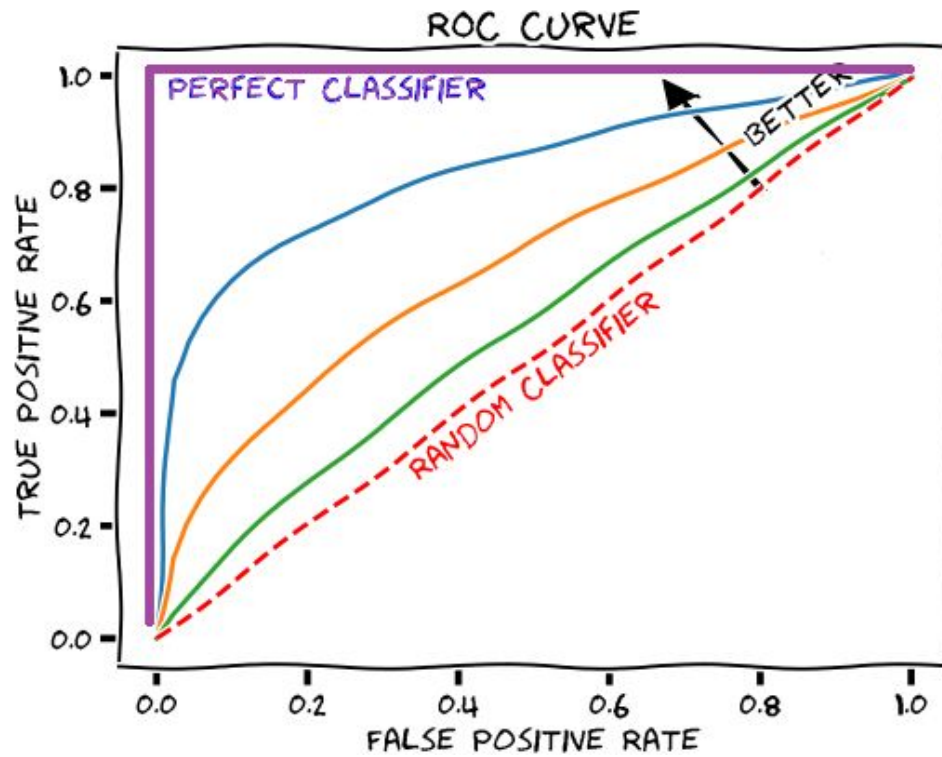
 Accuracy = $\frac{TP + TN}{TP + TN + FP + FN}$

 Specificity = $\frac{TN}{TN + FP}$

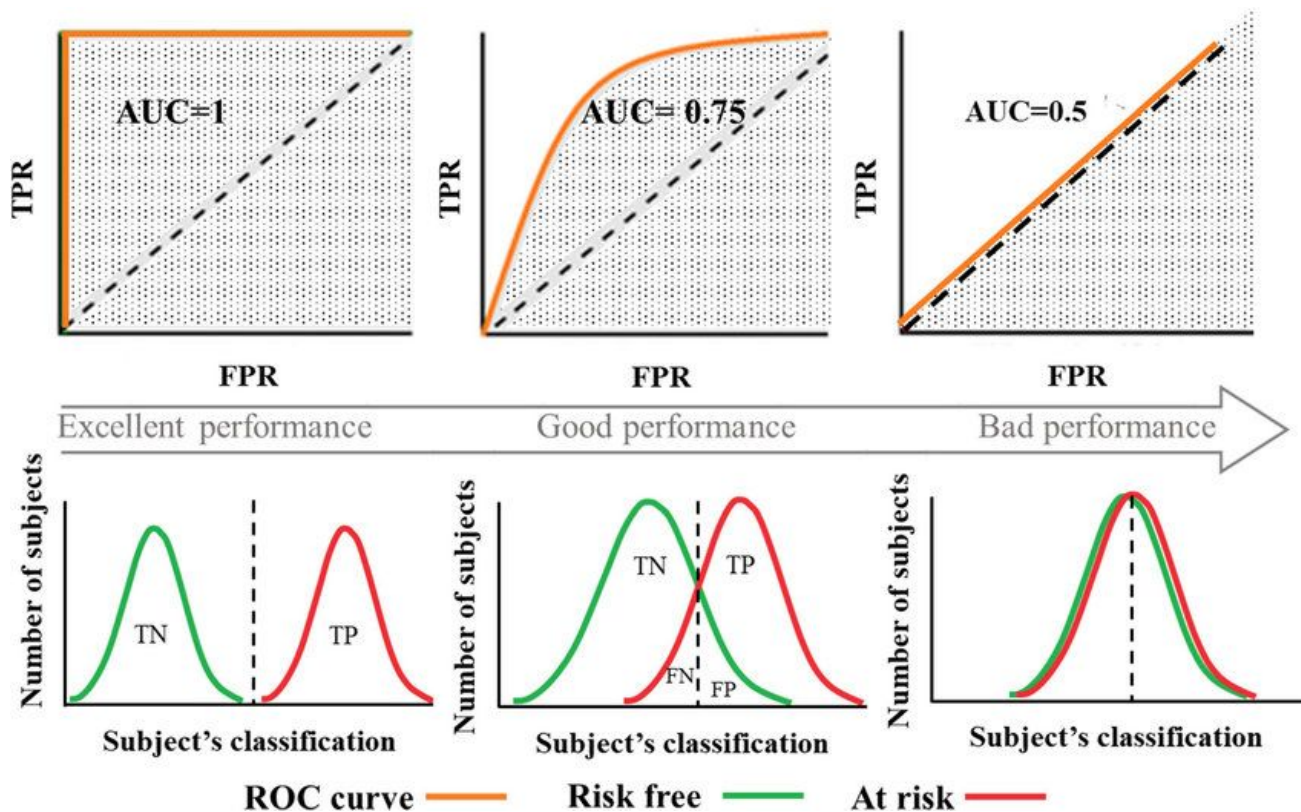
 Precision = $\frac{TP}{TP + FP}$

 Recall = $\frac{TP}{TP + FN}$

Metrics



Metrics



Feature Scaling

Normalization or Standardization

- **Feature Scaling** means scaling features to the same scale.
- **Normalization** scales features between 0 and 1, retaining their proportional range to each other.

Normalization

$$X' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

Diagram annotations: A red arrow labeled "new value" points to X' . A red arrow labeled "original value" points to x .

- **Standardization** scales features to have a mean (μ) of 0 and standard deviation (σ) of 1.

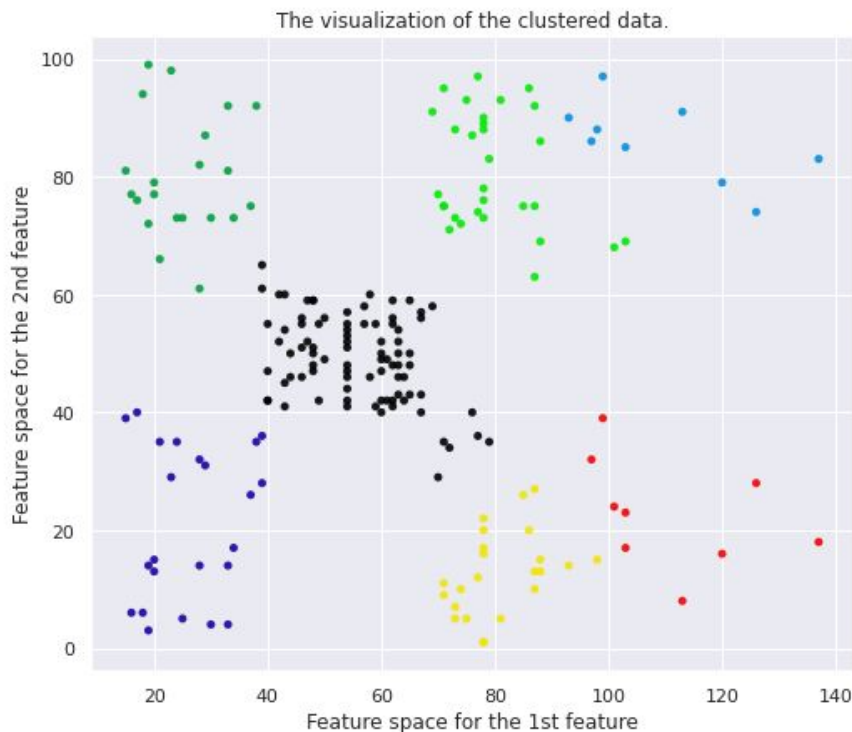
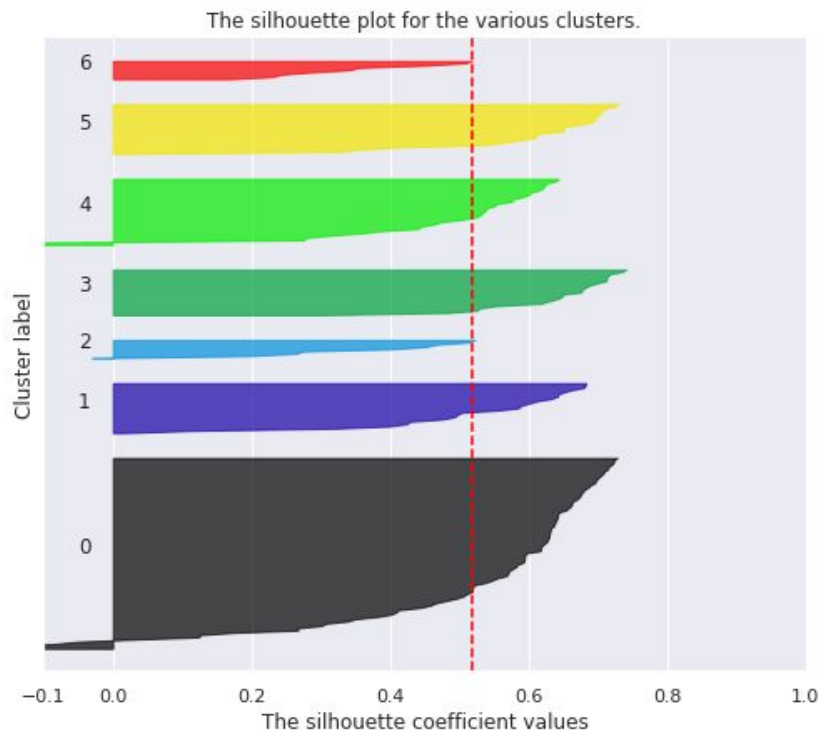
Standardization

$$X' = \frac{x - \mu}{\sigma}$$

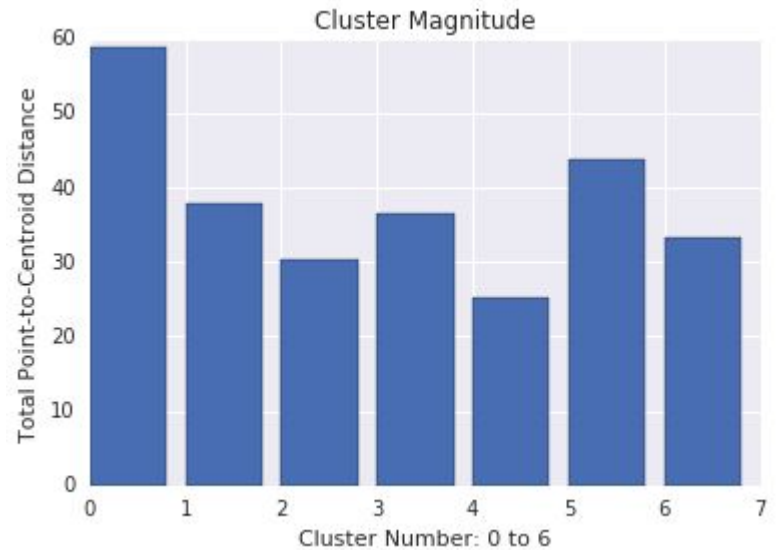
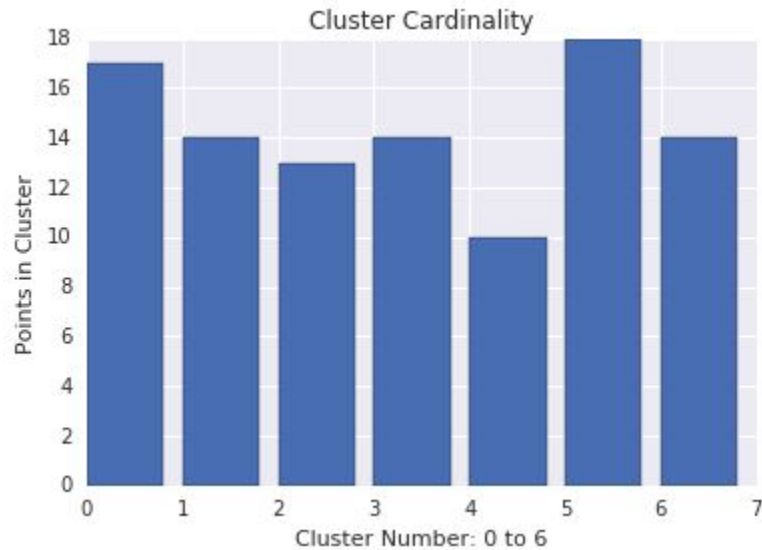
Diagram annotations: A red arrow labeled "new value" points to X' . A red arrow labeled "original value" points to x . A red arrow labeled "mean" points to μ . A red arrow labeled "standard deviation" points to σ .

Metrics for Clustering - Silhouette

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 7$



Metrics for Clustering - K-means



Metrics for Clustering - K-means

In cluster analysis, the elbow method is a heuristic used in determining the number of clusters in a data set. The method consists of plotting the explained variation as a function of the number of clusters, and picking the elbow of the curve as the number of clusters to use

$$\text{minimize} \left(\sum_{k=1}^k W(C_k) \right)$$

where C_k is the k^{th} cluster and $W(C_k)$ is the within-cluster variation

