# INT104 ARTIFICIAL INTELLIGENCE

## Review II
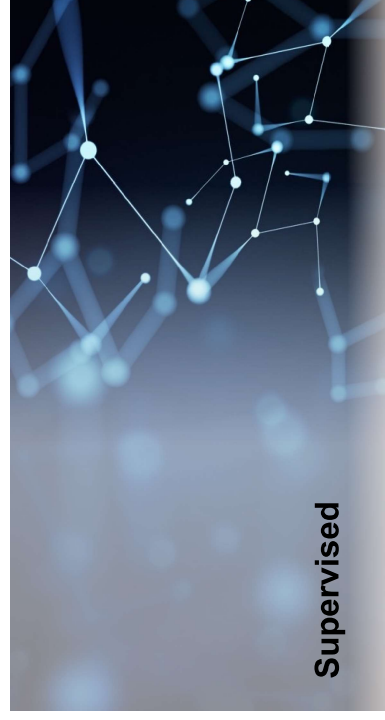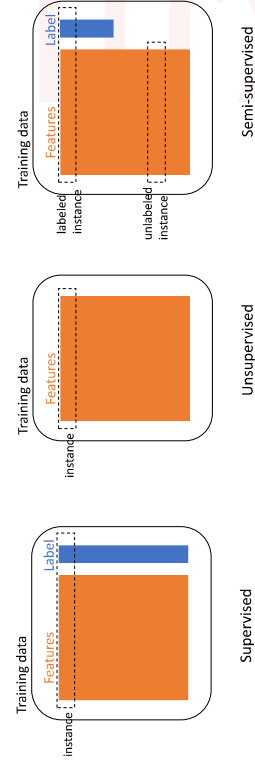
Fang Kang

Fang.kang@xjtlu.edu.cn

Xi'an Jiaotong-Liverpool University
西交利物浦大学

---

## CONTENT

➢ Supervised methods
- ◆ Classification and Regression
- ◆ SVM
- ◆ Decision Tree
- ◆ Random Forest

➢ Unsupervised methods
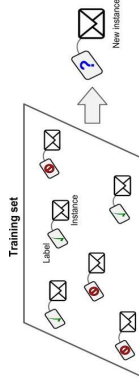- ◆ K-means
- ◆ Hierarchical clustering
- ◆ GMM

---

## Supervised



---

# Supervised vs. unsupervised

Training data — Features — instance

**Supervised** (Training data, Features, Label)

**Unsupervised** (Training data, Features)

**Semi-supervised** (Training data, Features, Label, labeled instance, unlabeled instance)

---

# Regression

Regression attempts to determine the strength and character of the relationship between one dependent variable (usually denoted by Y) and a series of other variables (known as independent variables).

$$x \xrightarrow{\ f(x)\ } \hat{y} = f(x)$$

Value

Value?

New instance

Feature 1

---

# Classification

**Classification:** Classification algorithms find a function that determines which category the input data belongs to.

**Binary Classification** is a supervised learning algorithm that classifies new observations into one of two classes.

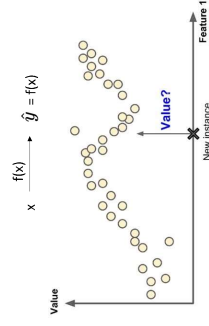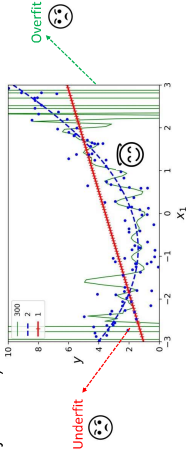Training set — Label — Instance — New instance

**Multiclass/Multilabel Classification**
- **Multiclass classification** refers to classification tasks that can distinguish between more than two classes.
- **Multilabel classification** refers to classification system that outputs multiple binary tags.

# Learning Curves

$$\hat{y} = ax_2 + bx_1 + c$$
$$x_2 = x_1^2$$

If you perform high-degree Polynomial Regression you will likely fit the training data much better than with plain Linear Regression. (Is high-degree polynomial always better?)



Underfit
Overfit

**Bias**: refers to the error from erroneous assumptions in the learning algorithm. (inability to capture the underlying patterns in the data).
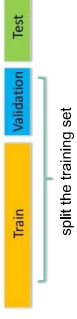
**Variance**: refers an error from sensitivity to small fluctuations in the training data. (difference in fits between data sets)

# Cross Validation

- Train/test/validation split

- To avoid selecting the parameters that perform best on the test data but maybe not the parameters that generalize best, we can further split the training set into training fold and validation fold

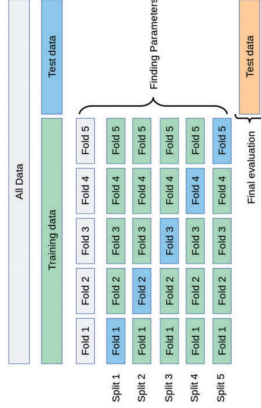- Can maximize the accuracy on the training data



| Train | Validation | Test |

split the training set

- **Training fold**: used to fit the model
- **Validation fold**: used to estimate prediction error for model selection
- **Test set**: used for assessment of the prediction error of the final chosen model
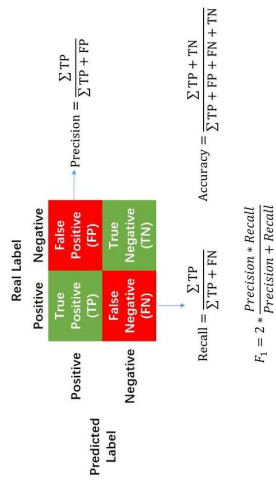
# K-fold Cross-Validation



All Data

Training data | Test data

Split 1 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5
Split 2 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5
Split 3 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5
Split 4 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5
Split 5 | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5

Finding Parameters

Final evaluation | Test data

# Confusion Matrix



Real Label

| | | Positive | Negative |
|---|---|---|---|
| Predicted Label | Positive | True Positive (TP) | False Positive (FP) |
| | Negative | False Negative (FN) | True Negative (TN) |

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

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

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

$$F_1 = 2 * \frac{Precision * Recall}{Precision + Recall}$$

# Support Vector Machine (SVM)

## Linear SVM Classification



Some Linear Classifier

SVM Classifier

- Linear separability
- Fitting widest possible "street" between classes
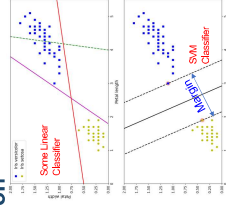  *Performs better with new data*
- **Large Margin Classification**
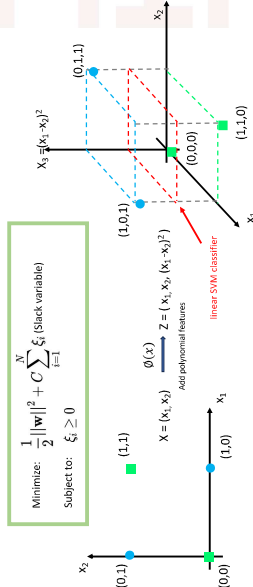- Margin, Support Vectors

**Hard Margin SVM**

All instances being off the street and on the right side

**Soft Margin SVM**

Allow margin violations

**Support Vectors**

- Decision boundary is not affected by more training instances
- It is determined by support vectors ( instances located on the edge of street )

# Nonlinear SVM Classification



$X = (x_1, x_2)$

$\emptyset(x)$ — Add polynomial features

$Z = (x_1, x_2, (x_1 \cdot x_2)^2)$ — linear SVM classifier

Minimize: $\frac{1}{2}||\mathbf{w}||^2 + C\sum_{i=1}^{N}\xi_i$ (Slack variable)

Subject to: $\xi_i \geq 0$

Nonlinear transformation $\emptyset(x)$ : not only one form

Cover's theorem: High-dimensional space is more likely to be linearly separable than in a low-dimensional space.

https://en.wikipedia.org/wiki/Cover's_theorem

## Decision Tree Definition

- A tree-like model that illustrates series of events leading to certain decisions
- Each node represents a test on an attribute and each branch is an outcome of that test

**Who to loan?**

- Not a student
- 45 years old
- Medium income
- Fair credit record
- **Yes**

- Student
- 27 years old
- Low income
- Excellent credit record
- **No**

**Depth:** the length of the longest path from the root node to a leaf node

---

## Best attribute = lowest Gini impurity

**In practice, we compute $gini(X)$ only once!**

$gini(X_{color=brown}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 \approx 0.444$

$gini(X_{color=white}) = 0.5$

$gini(X, color) = \frac{3}{7} \cdot 0.444 + \frac{4}{7} \cdot 0.5 \approx 0.476$

$gini(X_{fly=yes}) = 0$

$gini(X_{fly=no}) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \approx 0.375$

$gini(X, fly) = \frac{3}{7} \cdot 0 + \frac{4}{7} \cdot 0.375 \approx \boxed{0.214}$

---

## Random Forests



---

## Nonlinear SVM: Kernel Trick

**Input Space:** dimension n

**High-dimensional Feature Space:** dimension N >> n

$N \gg n$

Expensive operation and requires large memory

**Kernel Trick**

Universal approximator,
Corresponding feature space
$\phi(x)$ is infinite dimensional space

non-linearly separable data

infinite-dimensional space

**Common kernels:**

Linear: $K(\mathbf{a}, \mathbf{b}) = \mathbf{a}^T \mathbf{b}$

Polynomial: $K(\mathbf{a}, \mathbf{b}) = (\gamma \mathbf{a}^T \mathbf{b} + r)^d$

Gaussian Radial Basis Function: $K(\mathbf{a}, \mathbf{b}) = \exp(-\gamma \|\mathbf{a} - \mathbf{b}\|^2)$

Sigmoid: $K(\mathbf{a}, \mathbf{b}) = \tanh(\gamma \mathbf{a}^T \mathbf{b} + r)$

---

## Best attribute = highest information gain

**In practice, we compute $entropy(X)$ only once!**

$entropy(X) = -p_{mammal} \log_2 p_{mammal} - p_{bird} \log_2 p_{bird} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \approx 0.985$

$entropy(X_{color=brown}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \approx 0.918$

$entropy(X_{color=white}) = 1$

$gain(X, color) = 0.985 - \frac{3}{7} \cdot 0.918 - \frac{4}{7} \cdot 1 \approx 0.020$

$entropy(X_{fly=yes}) = 0$

$entropy(X_{fly=no}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \approx 0.811$

$gain(X, fly) = 0.985 - \frac{3}{7} \cdot 0 - \frac{4}{7} \cdot 0.811 \approx \boxed{0.521}$

---

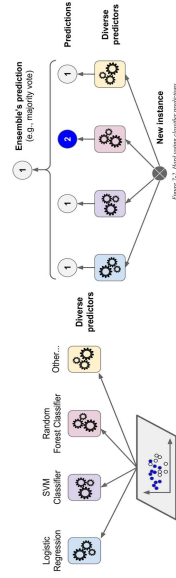## Ensemble Learning

**Ensemble :** A group of predictors

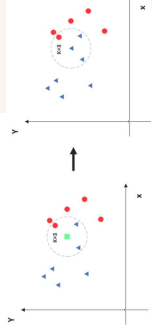**Voting Classifier**

**Hard Voting**

# K Nearest Neighbors (KNN)

As in the general problem of classification, we have a set of data points for which we know the correct class labels

When we get a new data point, we compare it to each of our existing data points and find similarity

Take the most similar k data points (k nearest neighbours)

From these k data points, take the majority vote of their labels. The winning label is the label / class of the new data point

---

## Ensemble method

- **Random Forests** are one of the most common examples of ensemble learning.

- Other commonly-used ensemble methods:
  - ➢ **Bagging**: multiple models on random subsets of data samples.
  - ➢ **Random Subspace Method**: multiple models on random subsets of features.
  - ➢ **Boosting**: train models iteratively, while making the current model focus on the mistakes of the previous ones by increasing the weight of misclassified samples.
  - ➢ **Stacking**: instead of using  hard voting to aggregate the predictions of all predictors in an ensemble, train a model to perform this aggregation.

---

# K-means clustering algorithm

Goal: Assign all data points to k clusters

**Step 1**: Pick k *random* initial cluster centroids

**Step 2**: Paint the data points that are closer to red centroid red, and those closer to blue centroid blue
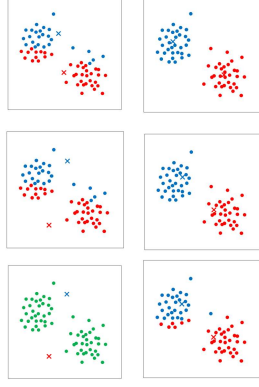
**Step 3**: Update the positions of centroids
Red centroid := average of current red points
Blue centroid := average of current blue points

Repeat

Until no more pointes need to be repainted, i.e., the centroids no longer change

Clustering is done

Euclidean distance $d(p,q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \cdots + (p_n - q_n)^2}$.

---

## Unsupervised

---

# Agglomerative clustering example
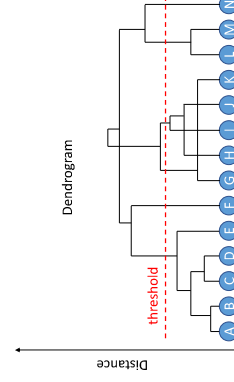
Dendrogram

Distance

threshold

---

## Hierarchical clustering

- Hierarchical Clustering is a set of clustering methods that aim at building a hierarchy of clusters
  - ➢ A cluster is composed of smaller clusters

- There are two strategies for building the hierarchy of clusters:
  - ➢ **Agglomerative** (bottom-up): we start with each point in its own cluster and we merge pairs of clusters until only one cluster is formed.
  - ➢ **Divisive** (top-down): we start with a single cluster containing the entire set of points and we recursively split until each point is in its own cluster.

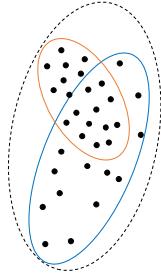- The most popular strategy in practical use is bottom-up (agglomerative)!

# Gaussian mixture model (GMM)

K-means make <u>hard</u> assignments to data points: $x^{(i)}$ must belong to one of the clusters $1, 2, \cdots, K$

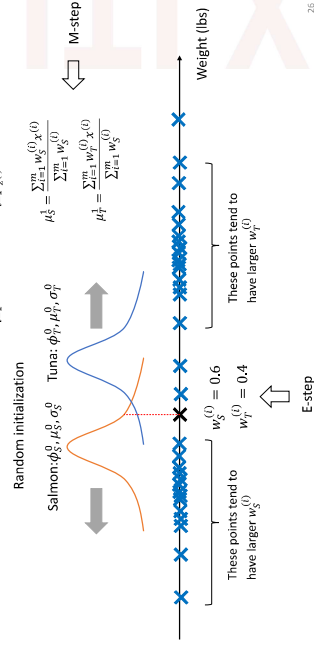Sometimes, one data point can belong to multiple clusters

- Clusters may overlap
- Hard assignment may be simplistic
- Need a <u>soft</u> assignment:
  data points belong to clusters with different ***probabilities***



---

# Demonstration with $k = 2$, 1-D Gaussian

Maximize likelihood of the whole data: $\mathcal{L}(\theta) = p(X|\Theta) = \prod_{i=1}^{m} p(x^{(i)}|\theta) = \prod_{i=1}^{m} \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}|\theta) = \prod_{i=1}^{m} \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}|\theta)}{Q_i(z^{(i)})}$

Random initialization

Salmon: $\phi_S^0, \mu_S^0, \sigma_S^0$    Tuna: $\phi_T^0, \mu_T^0, \sigma_T^0$

$\mu_S^1 = \frac{\sum_{i=1}^m w_S^{(i)} x^{(i)}}{\sum_{i=1}^m w_S^{(i)}}$

$\mu_T^1 = \frac{\sum_{i=1}^m w_T^{(i)} x^{(i)}}{\sum_{i=1}^m w_S^{(i)}}$

$w_S^{(i)} = 0.6$
$w_T^{(i)} = 0.4$

E-step

M-step

These points tend to have larger $w_T^{(i)}$

These points tend to have larger $w_S^{(i)}$

Weight (lbs)