# - Module 9 –
# $k$-Nearest Neighbors

## Outline

- $k$-Nearest Neighbors
  - ➢ Distance Metrics
  - ➢ Effects of $k$
- Classifier Performance Metrics

# *k*-Nearest Neighbors (1)

- *k*-Nearest Neighbors (*k*-NN) is a **non-parametric**, **instance-based** learning algorithm

  - non-parametric:  it does not make explicit assumptions about the form of the mapping function, $h(X)\colon X \to Y$
    (as opposed to parametric techniques, e.g., linear regression, $h(X)\colon Y = X\beta$)

  - instance-based:  it does not explicitly learn a model, i.e., there is no training stage
    (instead, the training samples are only used during the prediction stage)

- *k*-NN can be used for classification and regression
  ➔ for this course, focus on multiclass classification where the input features can have continuous or discrete values
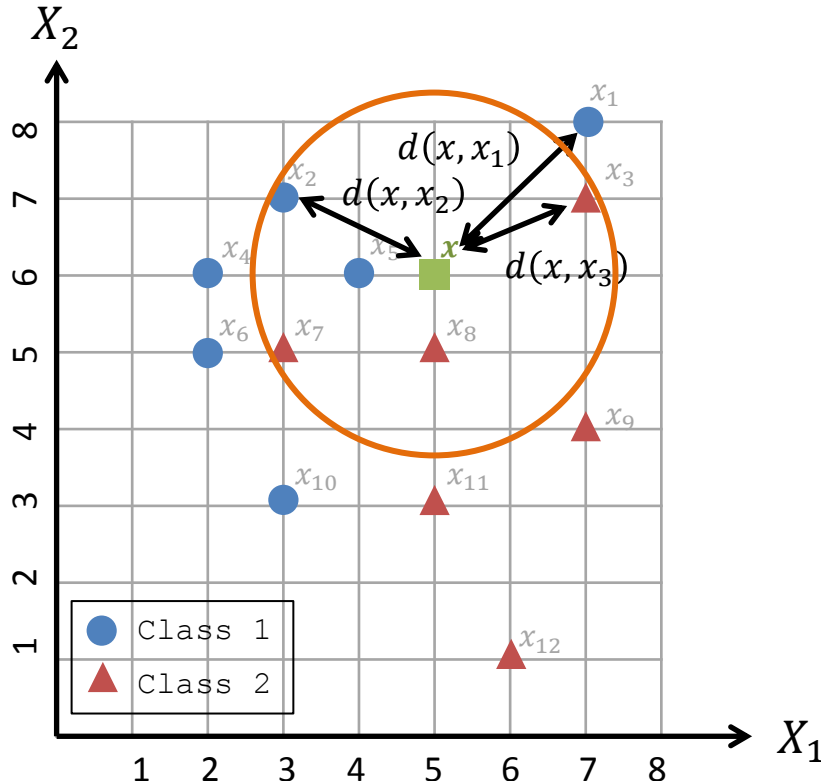
# $k$-Nearest Neighbors (2)

- Steps to use $k$-NN for classification, $kNN(x, \mathcal{T}, k)$, given
  - a test sample, $x$
  - a training set, $\mathcal{T}$      consists of the training samples, $(x_i, y_i)$ for $i = 1, \dots, N$
  - a positive integer, $k$
  - a **distance metric**, $d(x, x_i)$    measures the distance between $x$ and $x_i \, \forall \, i = 1, \dots, N$

  **Step 1:**   For each training sample $x_i \, \forall \, i = 1, \dots, N$, compute the distance, $d(x, x_i)$, between $x$ and $x_i$

  **Step 2:**   Order the training samples by increasing distance

  **Step 3:**   Select the first $k$ training samples that have the smallest $d(x, x_i)$, i.e., $k$ nearest neighbors of $x$

  **Step 4:**   Determine the predicted value $\hat{y}$ of $x$ by majority vote, i.e., most commonly occurring class of the $k$ nearest neighbors

               (Note: $k$ is usually an odd number to avoid ties)

# $k$-Nearest Neighbors: Example (3)

**Example:**

Consider the following <u>training samples</u>.

**Classify test sample $x$ using 5-NN.**



5 nearest neighbors

$d(x, x_1) = 2.83$    $d(x, x_5) = 1.00$    ●

$d(x, x_2) = 2.24$    $d(x, x_8) = 1.00$    ▲

$d(x, x_3) = 2.24$    $d(x, x_2) = 2.24$    ●

$d(x, x_4) = 3.00$    $d(x, x_3) = 2.24$    ▲

$d(x, x_5) = 1.00$    $d(x, x_7) = 2.24$    ▲

$d(x, x_6) = 3.16$    $d(x, x_1) = 2.83$

$d(x, x_7) = 2.24$    $d(x, x_9) = 2.83$

$d(x, x_8) = 1.00$    $d(x, x_4) = 3.00$

$d(x, x_9) = 2.83$    $d(x, x_{11}) = 3.00$

$d(x, x_{10}) = 3.61$    $d(x, x_6) = 3.16$

$d(x, x_{11}) = 3.00$    $d(x, x_{10}) = 3.61$

$d(x, x_{12}) = 5.10$    $d(x, x_{12}) = 5.10$

Majority vote

▲

sort

$x$ ➡ ▲ Class 2

4

# $k$-Nearest Neighbors: Distance Metrics (4)

- The <u>distance metrics</u>, $d(x, x')$, where $x$ and $x'$ are <u>vectors</u> of $p$ <u>features</u> indexed by $j$, for $j = 1, \dots, p$, are computed as follows:
  - for <u>continuous</u> features
    - $L^2$ norm (Euclidean distance):

$$d(x, x') = \sqrt{\sum_{j=1}^{p} (x_j - x_j')^2}$$

    - $L^1$ norm (Manhattan distance):

$$d(x, x') = \sum_{j=1}^{p} |x_j - x_j'|$$

  - for <u>categorical</u> features
    - Hamming distance:

Indicator function:
$I(condition)$
$= \begin{cases} 1, & \text{if } condition = true \\ 0, & \text{otherwise} \end{cases}$

$$d(x, x') = \sum_{j=1}^{p} I(x_j \neq x_j')$$

# $k$-Nearest Neighbors: Distance Metrics (5)
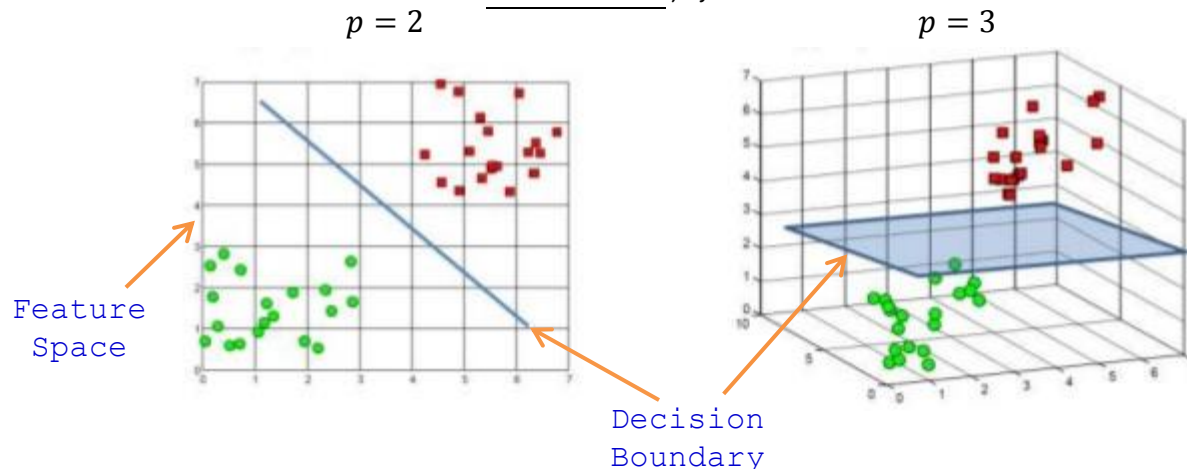
- Note that the <u>distance metric</u> is <u>dependent</u> on the <u>scales</u> of the <u>features</u>

  (e.g., if <u>a feature</u> <u>changes</u> <u>measurements</u> from <u>centimeters</u> to <u>miles</u> while <u>keeping</u> <u>other features</u> the <u>same</u>, a <u>different set</u> of <u>nearest neighbors</u> will be <u>selected</u>)

➔ to <u>avoid</u> this, it is common to <u>standardize</u> each <u>feature</u> $j$, i.e., $x_{i,j}$ is replaced with $(x_{i,j} - \mu_j)/\sigma_j$, <u>before</u> performing $k$-NN

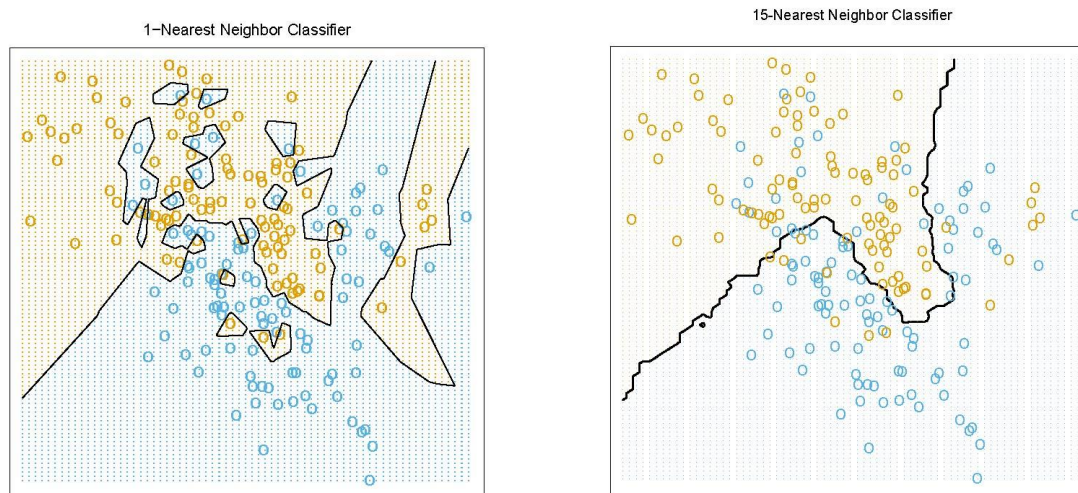# $k$-Nearest Neighbors: Effects of $k$ (6)

- In $k$-NN, the <u>parameter</u>, $k$, is a <u>tuning parameter</u> that controls the <u>smoothness</u> of the **decision boundary** of a $k$-NN <u>classifier</u> and its <u>fit</u> for the given <u>dataset</u>

  {<u>decision boundary</u>:    <u>partition</u> the **feature space** (space spanned by possible values of $X_1, X_2, …, X_p$) into $C$ <u>non-overlapping</u> <u>regions</u>, <u>one</u> for <u>each class</u>

  (e.g., <u>lines</u> for <u>two-dimensional</u> ($p = 2$) <u>input vectors</u> or <u>hyper-surfaces</u> for <u>higher-dimensional</u> ($p > 2$) <u>input vectors</u>)}

$p = 2$                                                   $p = 3$

Feature
Space

Decision
Boundary

# $k$-Nearest Neighbors: Effects of $k$ (7)

- When $k$ is <u>small</u>, $k$-NN has <u>high flexibility</u> but highly <u>sensitive to outliers</u>
  - exhibit <u>overfitting</u> (i.e., <u>low bias</u> and <u>high variance</u>) and the <u>decision boundary</u> is more <u>jagged</u>
- When $k$ is <u>large</u>, $k$-NN is more <u>resilient to outliers</u> but has <u>lower flexibility</u>
  - exhibit <u>underfitting</u> (i.e., <u>high bias</u> and <u>low variance</u>) and the <u>decision boundary</u> is more <u>smooth</u>



1-Nearest Neighbor Classifier        15-Nearest Neighbor Classifier

➔ the <u>best value</u> of $k$ is <u>chosen</u> using <u>cross-validation</u>, i.e., select $k$ that gives the <u>lowest average cross-validation estimate</u> of <u>prediction error</u>

- $K$-fold <u>cross-validation</u> for $k$-NN:

  **Step 1:** <u>Sub-divide</u> the <u>training set</u> into $K$ <u>non-overlapping</u> <u>subsets</u> of <u>equal size</u>

  **Step 2:** For <u>each trial</u> $q = 1, \ldots, K$
  - <u>all subsets</u> except the $q$-th subset is <u>used</u> as the <u>training set</u>, $\mathcal{T}_q$, and the $q$-th subset is <u>used</u> as the <u>validation set</u>, $\mathcal{V}_q$
  - for <u>each sample</u> $x_i$ in $\mathcal{V}_q$
    - <u>perform</u> $kNN(x_i, \mathcal{T}_q, k)$ to obtain the <u>predicted value</u>, $\hat{y}_i$
  - <u>compute</u> the <u>cross-validation estimate</u> of <u>prediction error</u> (or **misclassification rate**) for <u>trial</u> $q$:

$$Err_{CV}^q = \frac{\sum_{i=1}^{|\mathcal{V}_q|} I(y_i \neq \hat{y}_i)}{|\mathcal{V}_q|},$$

  where $|\mathcal{V}_q|$ <u>denotes</u> the <u>number of samples</u> in $\mathcal{V}_q$

# $k$-Nearest Neighbors: Cross-Validation (9)

**Step 3:** Compute the average cross-validation estimate of prediction error, $Err_{CV}$, across $K$ trials

$$Err_{CV} = \frac{1}{K}\sum_{q=1}^{K} Err_{CV}^{q}$$

- To determine the best value of $k$ (in $k$-NN), repeat steps 1 to 3 for all considered $k$ values and select $k$ that gives the lowest average cross-validation estimate of prediction error, $Err_{CV}$

# $k$-Nearest Neighbors: Advantages and Limitations (10)

- **Advantages** of $k$-NN:
  - easy to implement
  - ability to achieve good classification performance when the training set is large
  - works with multiclass datasets

- **Limitations** of $k$-NN:
  - prediction stage is computationally intensive as the algorithm has to iterate through all samples in the training set
    - ➜ often an issue where the training set is very large
  - high memory requirement as it has to store all of the training samples all the time
  - unable to provide information as to whether a feature is more informative than another

# Classifier Performance Metrics (1)

- How to <u>assess</u> the <u>performance</u> of a <u>classifier</u>?
  - **Accuracy:** percentage of <u>correct classifications</u>
  - **Error Rate:** percentage of <u>misclassifications</u> (1 - **Accuracy**)

- However, <u>accuracy</u> or <u>error rate</u> assume <u>equal costs</u> for <u>misclassification</u>, but this <u>assumption</u> is often <u>not valid</u> for <u>real applications</u>
  - **Example:** Spam filter
    - cost of misclassifying a <u>spam e-mail</u> as a <u>legitimate e-mail</u> ➜ **LOW**
    - cost of misclassifying a <u>legitimate e-mail</u> as a <u>spam e-mail</u> ➜ **HIGH**
  - **Example:** Medical screening
    - cost of misclassifying a <u>healthy person</u> as a <u>potential cancer patient</u> ➜ **LOW**
    - cost of misclassifying a <u>potential cancer patient</u> as a <u>healthy person</u> ➜ **HIGH**

- Need to <u>differentiate</u> <u>different kinds</u> of <u>misclassifications</u> ➜ **Confusion Matrix**

# Classifier Performance Metrics (2)

- Confusion matrix
  - **True Positives (TP):** actual class of a sample is positive and the predicted class is also positive

    (e.g., the person has cancer and the classifier also indicates that the person has cancer)
  - **True Negatives (TN):** actual class of a sample is negative and the predicted class is also negative

    (e.g., the person does not have cancer and the classifier also indicates that the person does not have cancer)
  - **False Positives (FP):** actual class of a sample is negative but the predicted class is positive

    (e.g., the person does not have cancer but the classifier indicates that the person has cancer)
  - **False Negatives (FN):** actual class of a sample is positive but the predicted class is negative

    (e.g., the person has cancer but the classifier indicates that the person does not have cancer)

|        |          | Predicted | |
|--------|----------|-----------|-----------|
|        |          | Positive | Negative |
| Actual | Positive | # of **True Positives** (TP) | # of **False Negatives** (FN) |
| Actual | Negative | # of **False Positives** (FP) | # of **True Negatives** (TN) |

(false positive)     (false negative)

You're pregnant     You're not pregnant

# Classifier Performance Metrics (3)

- Different <u>performance metrics</u> can be calculated from the <u>confusion matrix</u>:

  - $Accuracy = \dfrac{TP+TN}{TP+TN+FP+FN} = \dfrac{TP+TN}{total}$

  - $Precision = \dfrac{TP}{TP+FP} = \dfrac{TP}{predicted\ positive}$

  - $Sensitivity/Recall = \dfrac{TP}{TP+FN} = \dfrac{TP}{actual\ positive}$

  - $Specificity = \dfrac{TN}{TN+FP} = \dfrac{TN}{actual\ negative}$

  - $False\ Positive\ Rate = \dfrac{FP}{TN+FP} = \dfrac{FP}{actual\ negative}$
    $$= 1 - specificity$$

  - $False\ Negative\ Rate/MissRate = \dfrac{FN}{FN+TP}$
    $$= \dfrac{FN}{actual\ positive}$$