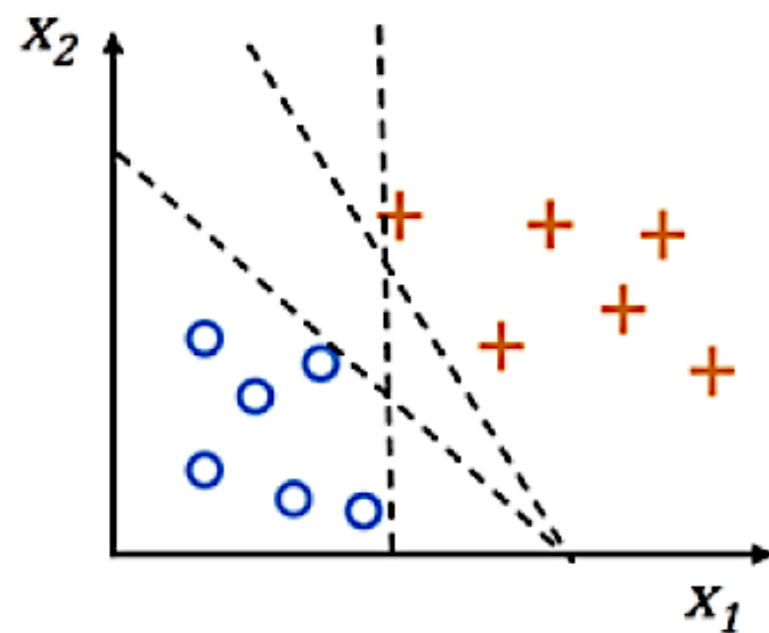# SVMs, Decision Trees, Random Forest, and *k*-NNs

Machine Learning for Engineering Applications
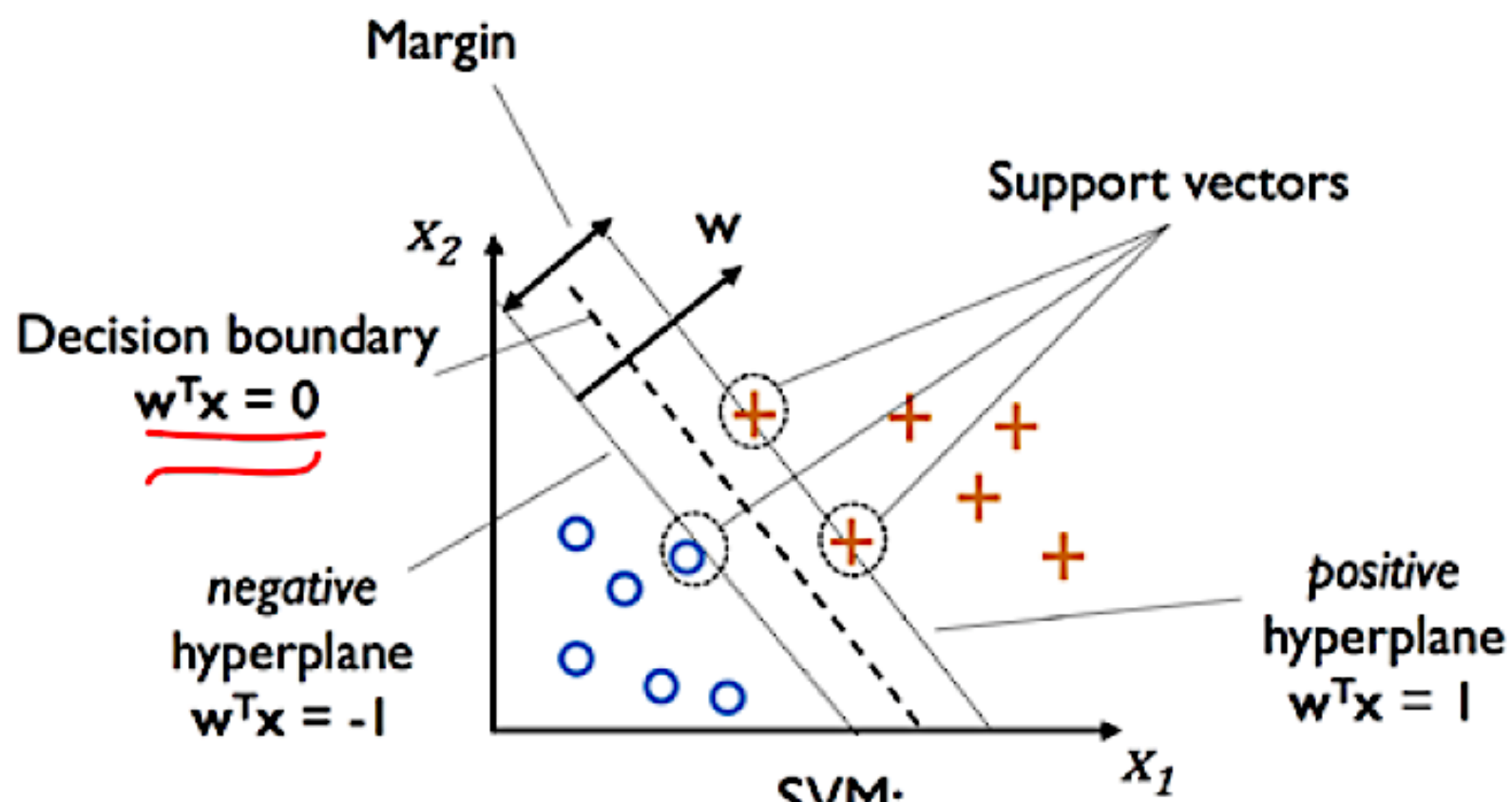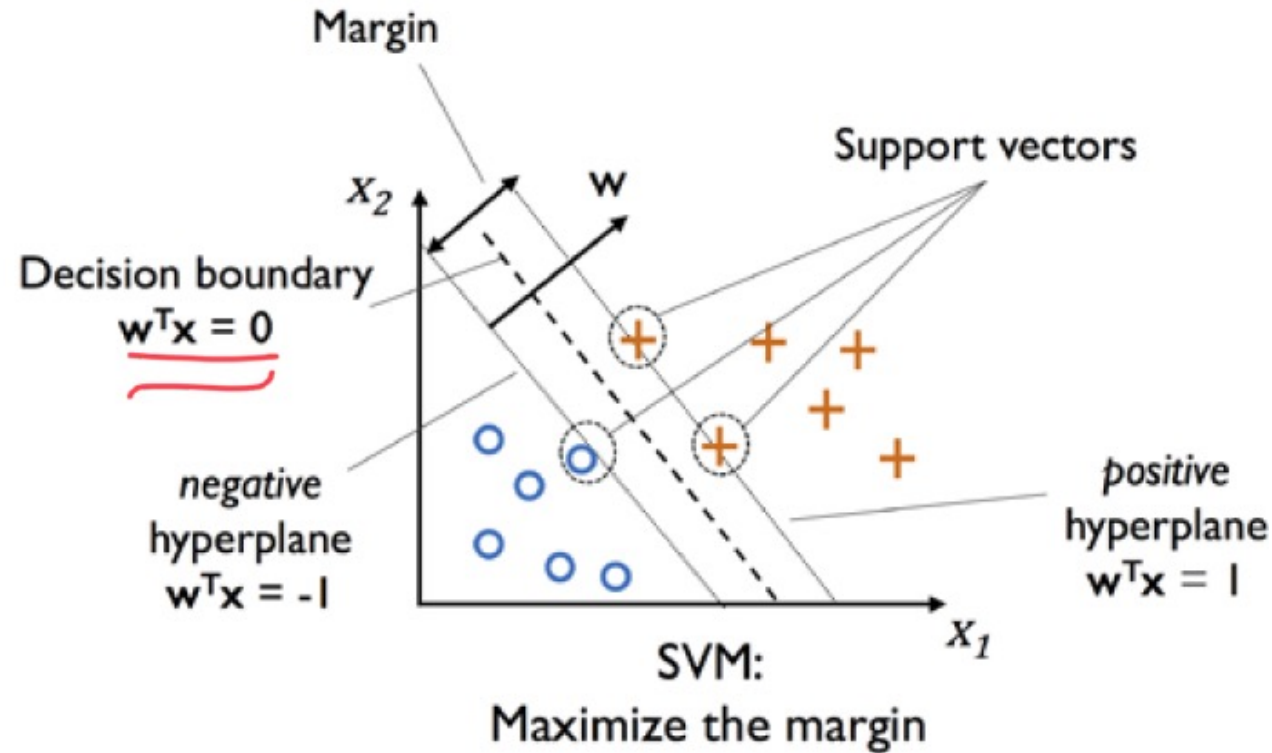
Fall 2023

# Support Vector Machines (SVM)

- **SVM** is one of the most power and most used classifier in the world!

- The Python community has created many SVM libraries that support all types of classification problems

- Non-linear classification is the essential part of this tool

Margin

$x_2$

w

Support vectors

Decision boundary
$w^Tx = 0$

$x_2$

negative
hyperplane
$w^Tx = -1$

positive
hyperplane
$w^Tx = 1$

$x_1$

Which hyperplane?

SVM:
Maximize the margin

$x_1$

- **SVM** determines many parts of a classification

- Identifying the vectors
- Drawing the hyperplanes
- Calculate direction (magnitude)
- Margins
- The distance of the hyperplanes
- So on…

- **Reason**: The best optimized (*larger*) boundary will help the model not to overfit

Margin

Support vectors

$x_2$

w

Decision boundary

$\mathbf{w^Tx} = 0$

negative hyperplane

$\mathbf{w^Tx} = -1$

positive hyperplane

$\mathbf{w^Tx} = 1$

$x_1$

SVM:

Maximize the margin

- **SVM** determines many parts of a classification

$$w_0 + w^T x_{pos} = 1$$

$$\Rightarrow w^T \left( x_{pos} - x_{neg} \right) = 2$$

Margin distance

$$w_0 + w^T x_{neg} = -1$$

- Normalize the weights:

$$\|w\| = \sqrt{\sum_{j=1}^{m} w_j^2}$$

- **You end up with:**

$$\frac{w^T \left( x_{pos} - x_{neg} \right)}{\|w\|} = \boxed{\frac{2}{\|w\|}}$$

Maximize this part!

Margin

- For evaluation, the margin is expressed in its reciprocal:

$$\frac{1}{2}\|w\|^2$$

- **We need a slack variable ($\xi$):**

- $\xi$: This drives the model to be a soft-margin classifier

- Hyperplane effect:

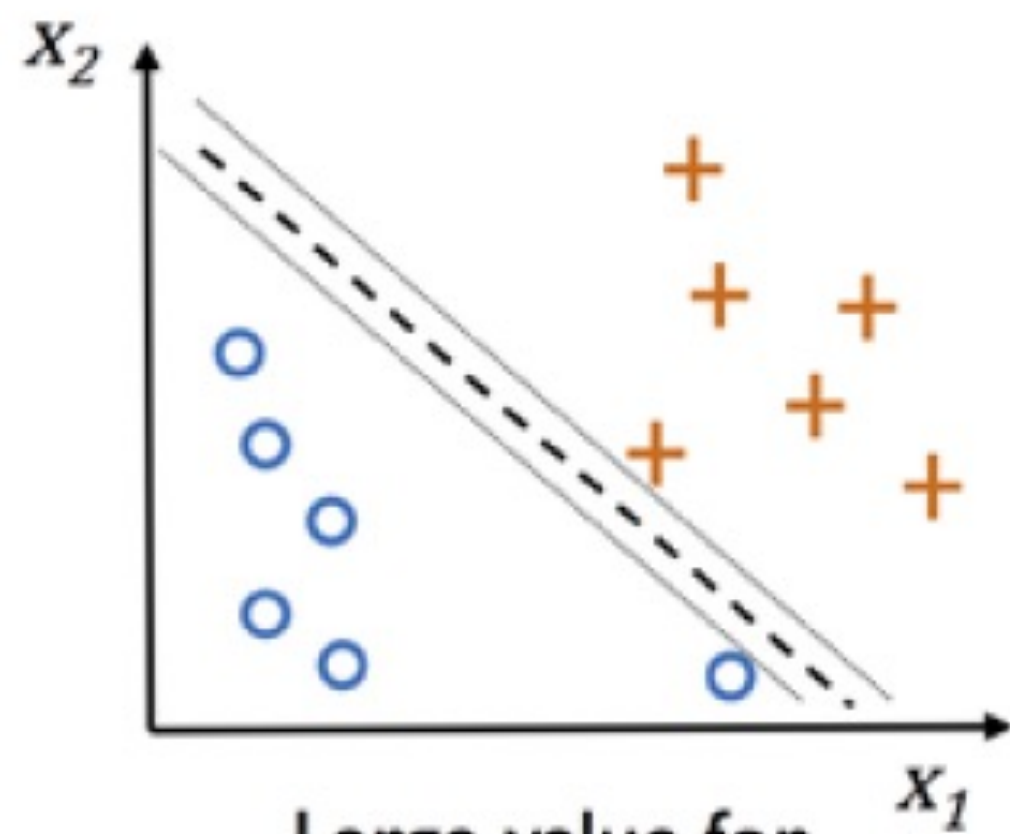$$w_0 + \boldsymbol{w}^T \boldsymbol{x}^{(i)} \geq 1 - \xi^{(i)} \quad if \quad y^{(i)} = 1$$

$$w_0 + \boldsymbol{w}^T \boldsymbol{x}^{(i)} \leq -1 + \xi^{(i)} \quad if \quad y^{(i)} = -1$$

$$\text{for } i = 1 \ldots N$$

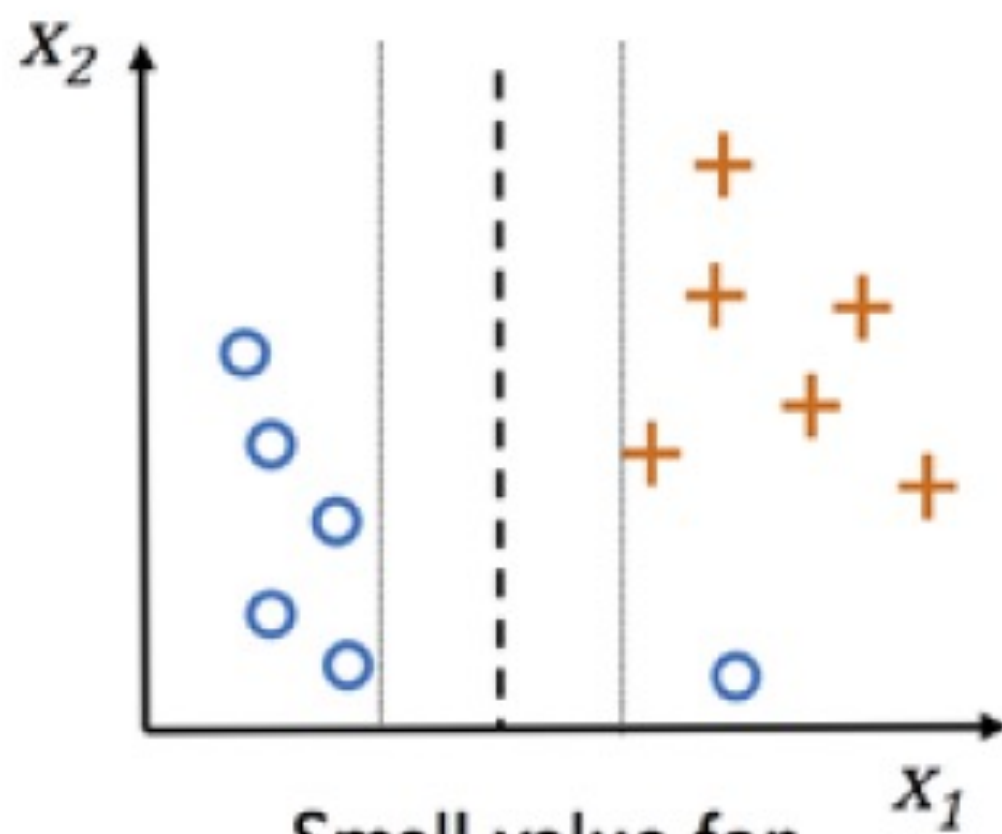- **You now want to minimize the $\xi$:**

$$\frac{1}{2}\|w\|^2 + C\left(\sum_i \xi^{(i)}\right)$$

- The C-variable controls the penalty intensity
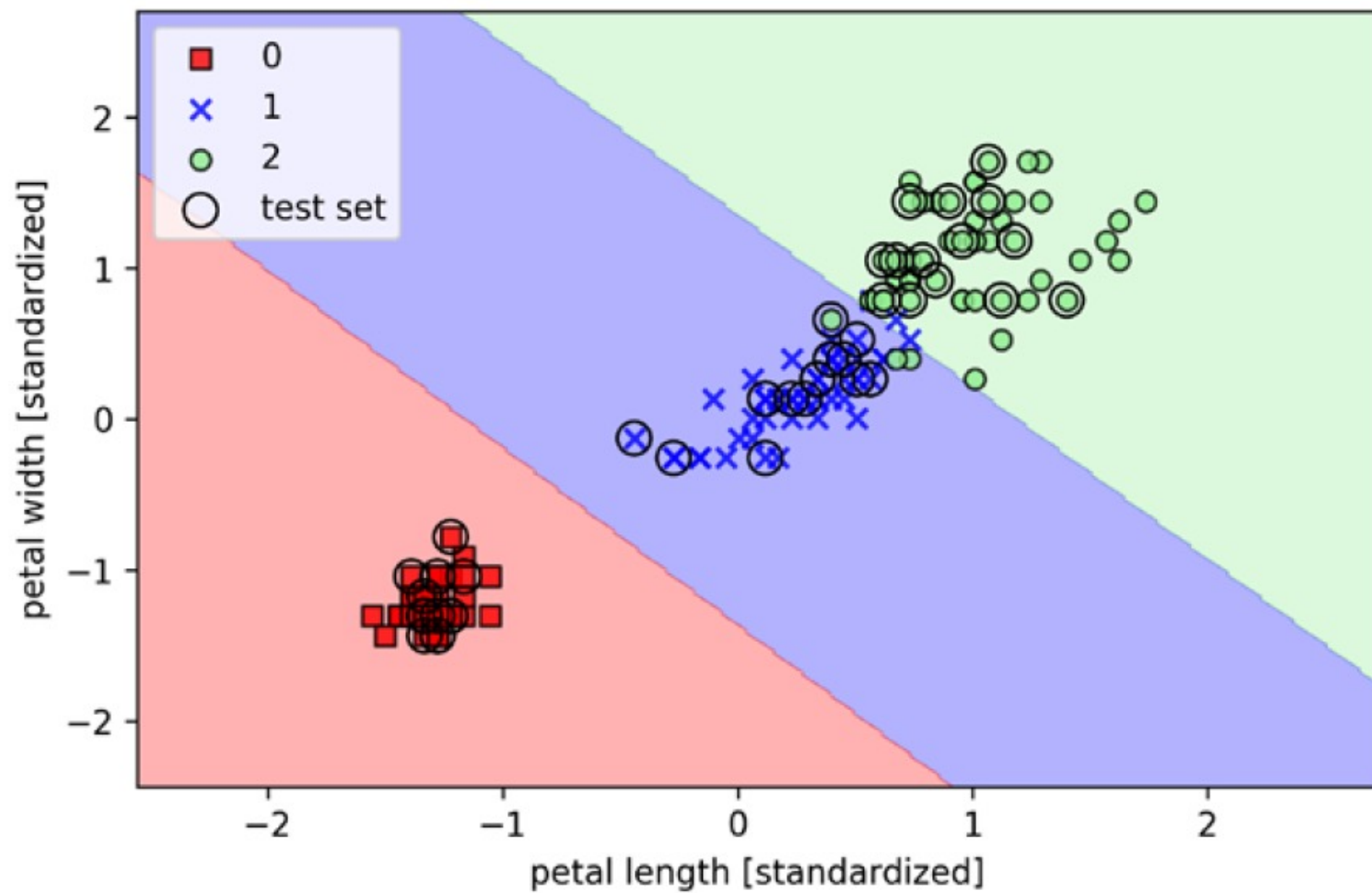
- Large C-values:  Large Error Penalties / Vice-Versa

Large value for
parameter C

Small value for
parameter C

# SVM LEAP
# Demo

# SVM with Non-Linearity
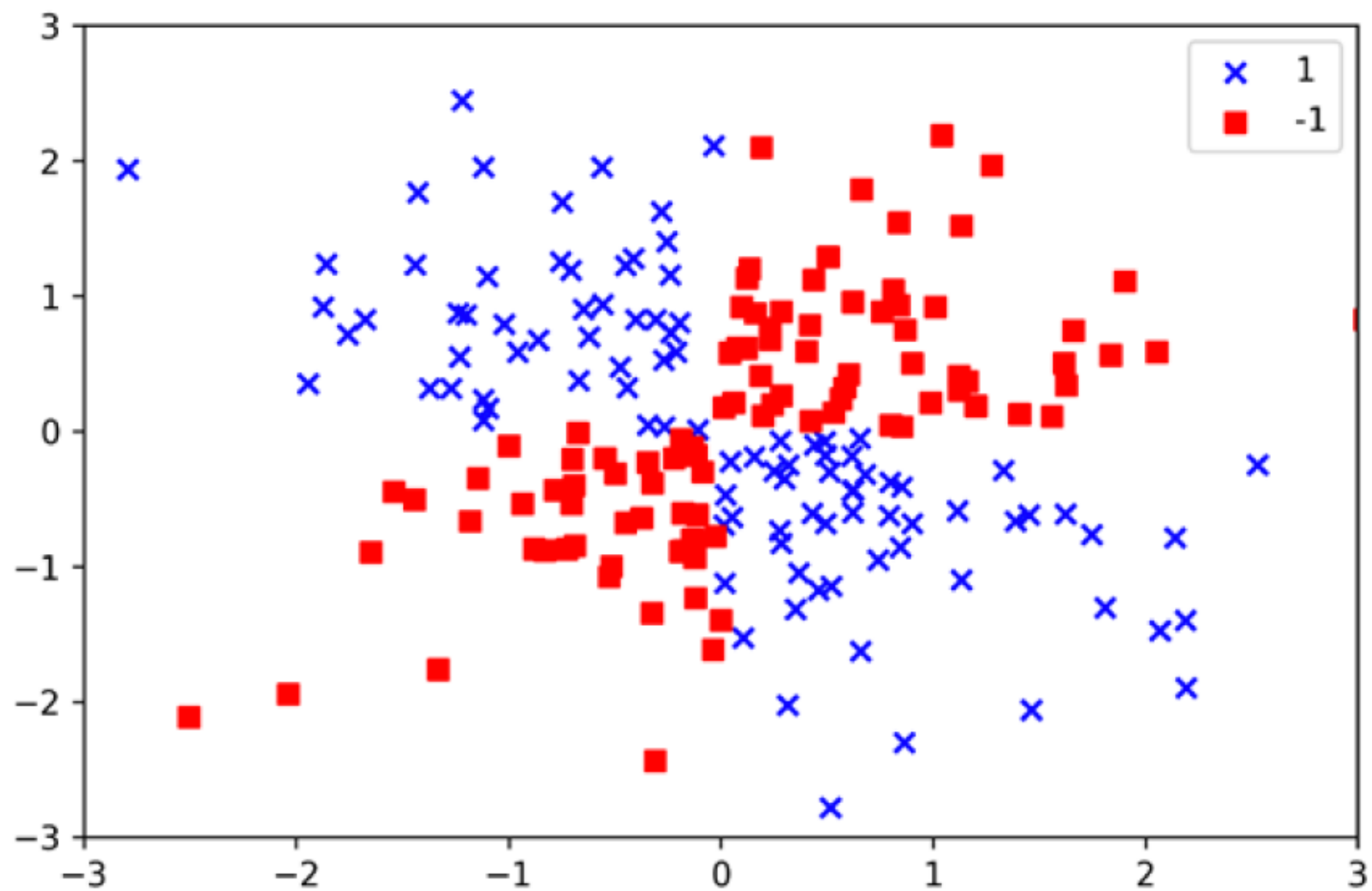
- The XOR non-linear example with random noise

```
import matplotlib.pyplot as plt
import numpy as np

np.random.seed(1)

X_xor = np.random.randn(200, 2)

y_xor = np.logical_xor(X_xor[:, 0] > 0,
    ... X_xor[:, 1] > 0)

y_xor = np.where(y_xor, 1, -1)
```
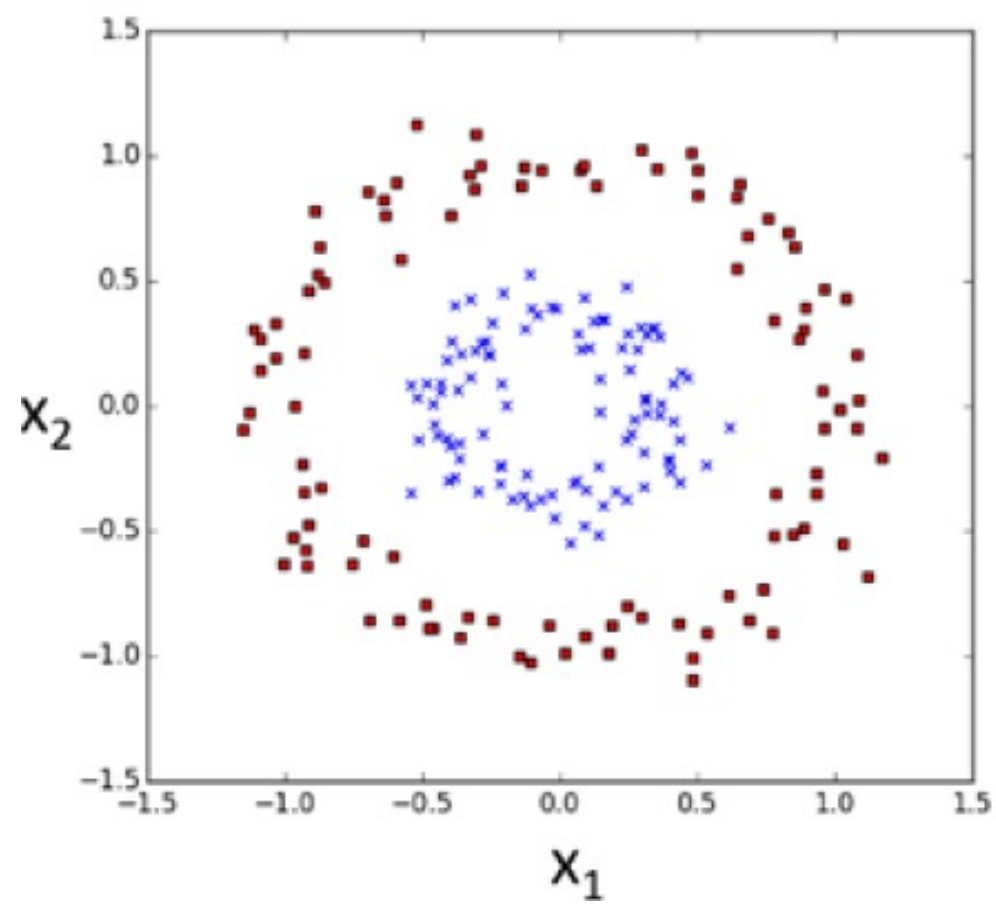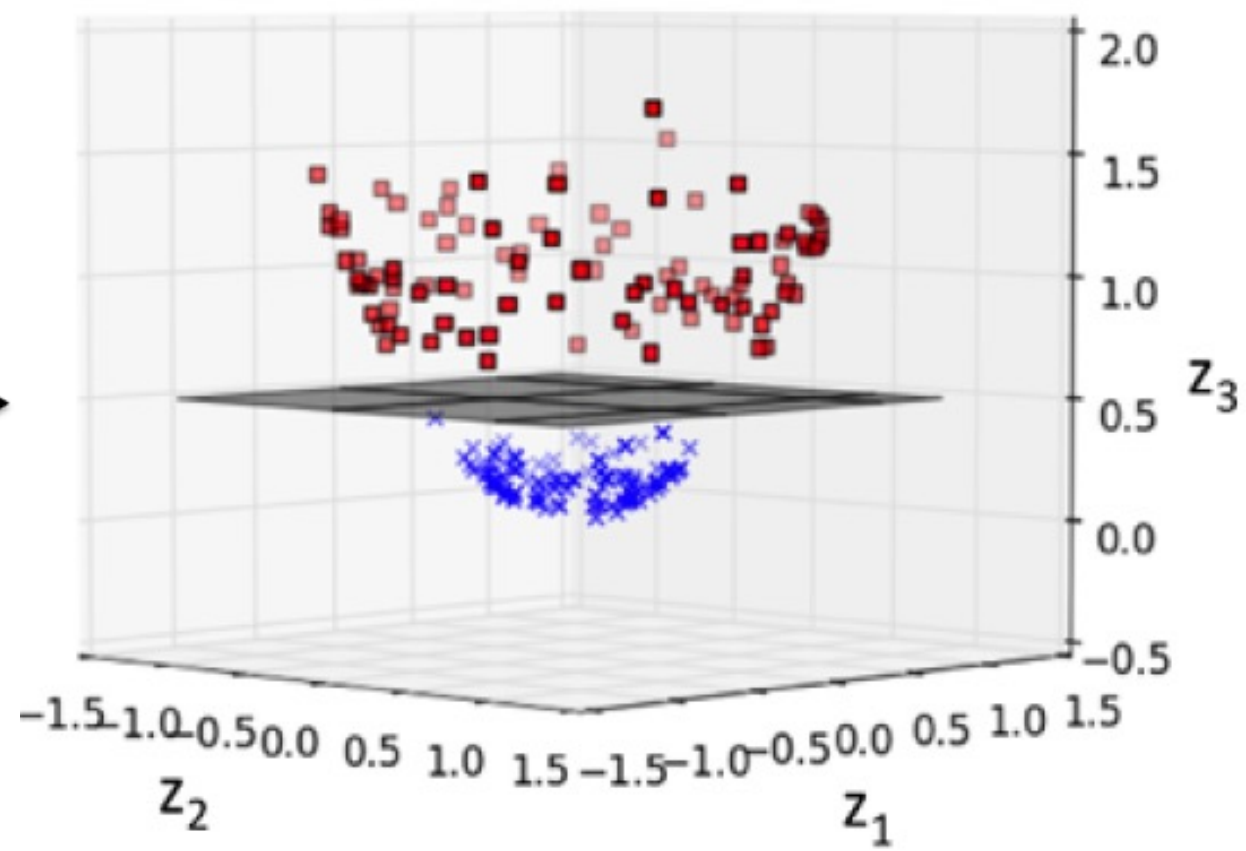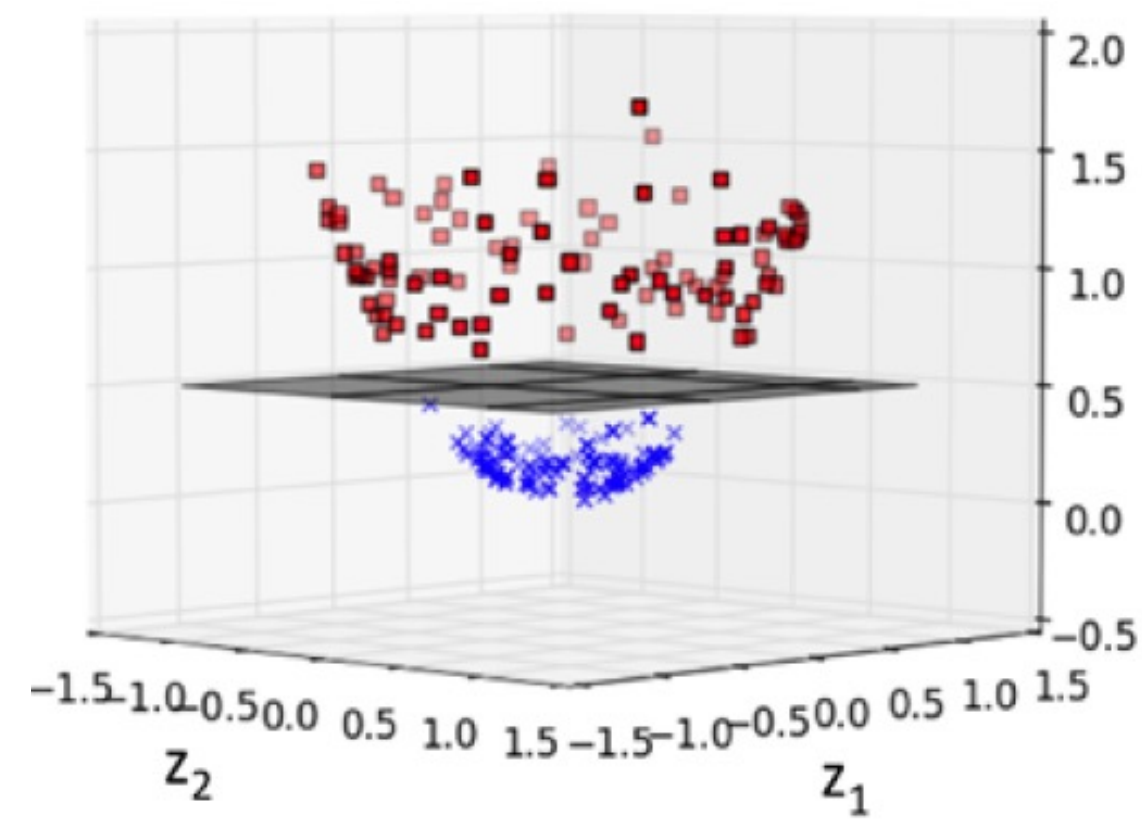
- Kernel methods, or Kernel SVM, helps with this type of non-linear issue

- The kernel method is a transformation of the data: 2D to a 3D space:

$$\phi(x_1, x_2) = (z_1, z_2, z_3) = (x_1, x_2, x_1^2 + x_2^2)$$

- The transformation:

$$x^{(i)T} x^{(j)}$$

$$K\left(x^{(i)}, x^{(j)}\right) = \phi\left(x^{(i)}\right)^T \phi\left(x^{(j)}\right)$$

$$K\left(x^{(i)}, x^{(j)}\right) = \exp\left(-\frac{\left\|x^{(i)} - x^{(j)}\right\|^2}{2\sigma^2}\right)$$

Radial Basis Function (RBF)

$$K\left(x^{(i)}, x^{(j)}\right) = \exp\left(-\gamma\left\|x^{(i)} - x^{(j)}\right\|^2\right)$$

$$\gamma = \frac{1}{2\sigma^2}$$

Low Gamma Values

Large Gamma Values

- The XOR non-linear example with random noise

```
svm = SVC(kernel='rbf', random_state=1, gamma=0.10, C=10.0)

svm.fit(X_xor, y_xor)


plot_decision_regions(X_xor, y_xor, classifier=svm,
                      ...test_idx=range(y_train.size,
                          ..y_train.size + y_test.size))

plt.legend(loc='upper left')
plt.show()
```

# Decision Tree Learning

- Decision Tree classification

- A classifier that **breaks down a decision** of which class a new sample belongs to through Q&As

- Decision Tree is a model developed by deducting the class labels of the samples

- The splits are based on features with the largest **Information Gains** (IG)

- Samples @ each node belong to the same class

- Depending on the data: Trees can have a lot of nodes -> **overfitting** -> pruning the tree helps this issue by setting maximum depths (*hyperparameter*)

- What makes the tree split?

- **Information Gains** – To know the gain, there must be a calculation relationship between the **parent node** and the **child nodes**.

- The overall idea:

$$IG\left(D_p, f\right) = I\left(D_p\right) - \sum_{j=1}^{m} \frac{N_j}{N_p} I\left(D_j\right)$$

- Different between the impurities of the parent node and the aggregated child nodes

- *f* – feature (column) that performs the split

- I – Impurity function

- $D_P$ – dataset of the parent

- $D_j$ – dataset of the $j^{th}$ child node

- $N_P$ - # of samples @ parent node

- $N_j$ - # of samples @ the $j^{th}$ child node

$$IG\left(D_p, f\right) = I\left(D_p\right) - \sum_{j=1}^{m} \frac{N_j}{N_p} I\left(D_j\right)$$

Parent's

Children's

- **Binary Decision Tree**

- Easier to implement by most ML libraries

- Each Parent node <u>can only spawn 2 child nodes</u>:

$$IG\left(D_p, f\right) = I\left(D_p\right) - \frac{N_{left}}{N_p} I\left(D_{left}\right) - \frac{N_{right}}{N_p} I\left(D_{right}\right)$$

Parent's

Child 1

Child 2

- **Types of Binary Decision Trees**:
  1. Gini
  2. Entropy
  3. Classification Error

```python
from sklearn.tree import DecisionTreeClassifier
tree = DecisionTreeClassifier(criterion='gini', max_depth=4,
                                    ... random_state=1)
tree.fit(X_train, y_train)
X_combined = np.vstack((X_train, X_test))
y_combined = np.hstack((y_train, y_test))

plot_decision_regions(X_combined,
                      ... y_combined,
                      ... classifier=tree,
                      ... test_idx=range(y_train.size,
                       .. y_train.size + y_test.size))
plt.xlabel('petal length [cm]')
plt.ylabel('petal width [cm]')
plt.legend(loc='upper left')
plt.show()
```

See textbook code example to visualize a decision tree (Chapter 3)

# Random Forests

- **Random Forest** is random by nature

- "*Easy*" to implement, hard to explain.

- Basically:  A collection of *Decision Trees*

- Each Decision Tree is randomly trained

- **The goal**:  Aggregate the accuracy of all the trees

- **Implementation of the Random Forest:**

1.  Draw a random **bootstrap** sample of size $n$ (randomly choose $n$ samples from the training set with replacement).

2.  Grow a decision tree from the bootstrap sample. At each node:

    a.  Randomly select $d$ features without replacement.

    b.  Split the node using the feature that provides the best split according to the objective function, for instance, maximizing the information gain.

3.  Repeat the steps 1-2 $k$ times.

4.  Aggregate the prediction by each tree to assign the class label by **majority vote**.

- **Random Forest Classification:**

```
from sklearn.ensemble import RandomForestClassifier

forest = RandomForestClassifier(criterion='gini',
    ... n_estimators=25,   #25 random decision trees!
    ... random_state=1)
forest.fit(X_train, y_train)
plot_decision_regions(X_combined, y_combined,
                ... classifier=forest, test_idx= range
                    ..(y_train.size, y_train.size +
                    .. y_test.size))
plt.xlabel('petal length')
plt.ylabel('petal width')
plt.legend(loc='upper left')
plt.show()
```

# *k*-Nearest Neighbors

- *k*-Nearest Neigbor (*kNN*) is a supervised instance-based learning algorithm

- *k* – is a variable

- **How does it work:**
  1. Choose the # of k and distance metric
  2. Find the *k*-NN of the sample that needs to be classified
  3. Assign the category label based on majority vote

- **The Pro:**
  - New data is easily determined by the votes (adaptable)

- **The Con:**
  - The large number of new samples will take over your physical memory (a linear growth)
  - This will depend on the # of features in the dataset

- **k-NN Rules:**

  - If the votes are tied:
    - Distance to the closest
  - If two or more classes have the same distance:
    - Which ever class is listed first in the dataset (*user-determined*)

- The **k-factor** determines great balance, overfitting or underfitting.

- Distance is controlled by the **p-parameter** (*in the code*)

- # Popular Distances:

**Metrics intended for real-valued vector spaces:**

| identifier | class name | args | distance function |
|---|---|---|---|
| "euclidean" | EuclideanDistance | | `sqrt(sum((x - y)^2))` |
| "manhattan" | ManhattanDistance | | `sum(|x - y|)` |
| "chebyshev" | ChebyshevDistance | | `max(|x - y|)` |
| "minkowski" | MinkowskiDistance | p | `sum(|x - y|^p)^(1/p)` |
| "wminkowski" | WMinkowskiDistance | p, w | `sum(|w * (x - y)|^p)^(1/p)` |
| "seuclidean" | SEuclideanDistance | V | `sqrt(sum((x - y)^2 / V))` |
| "mahalanobis" | MahalanobisDistance | V or VI | `sqrt((x - y)' V^-1 (x - y))` |

- p=1  → Manhattan
- p=2  → Euclidean

```python
from sklearn.neighbors import KNeighborsClassifier

# all preprocessing of data goes here…

knn = KNeighborsClassifier(n_neighbors=5, p=2,
                                ... metric='minkowski')


# create the classes through training
knn.fit(X_train, y_train)


# vote on unseen data into the trained classes
plot_decision_regions(X_combined, y_combined,
                      ... classifier=knn,
                      ... test_idx=range(y_train.size,
                       … y_train.size + y_test.size))
```

# For next class…

- **Start reading Chapter 4**

- **HW# 2 will be posted after class**