

# Classification

(Slides largely based on Prof. Sandra Avila's Machine Learning Course)

Prof. Rosa Paccotacya Yanque

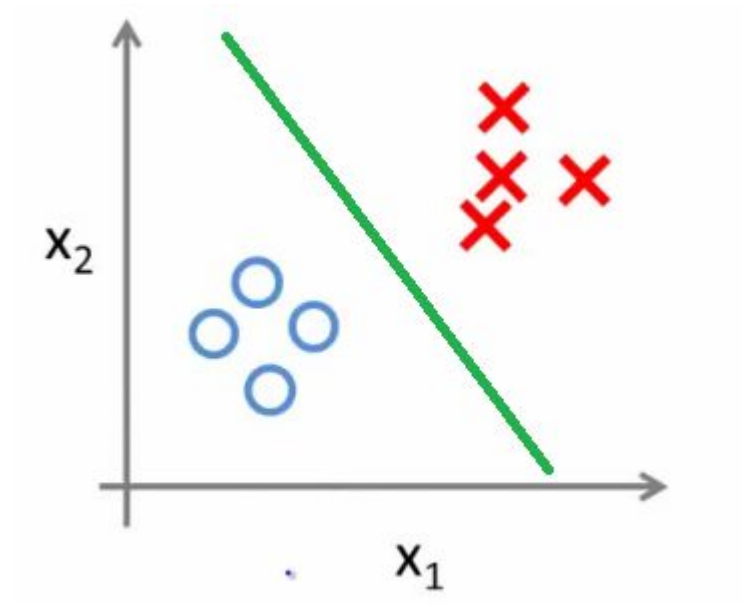
# Agenda

- ▷ Logistic Regression
- ▷ Ensemble Learning
- ▷ Support Vector Machines (SVMs)
- ▷ Decision Tree
- ▷ Random Forest

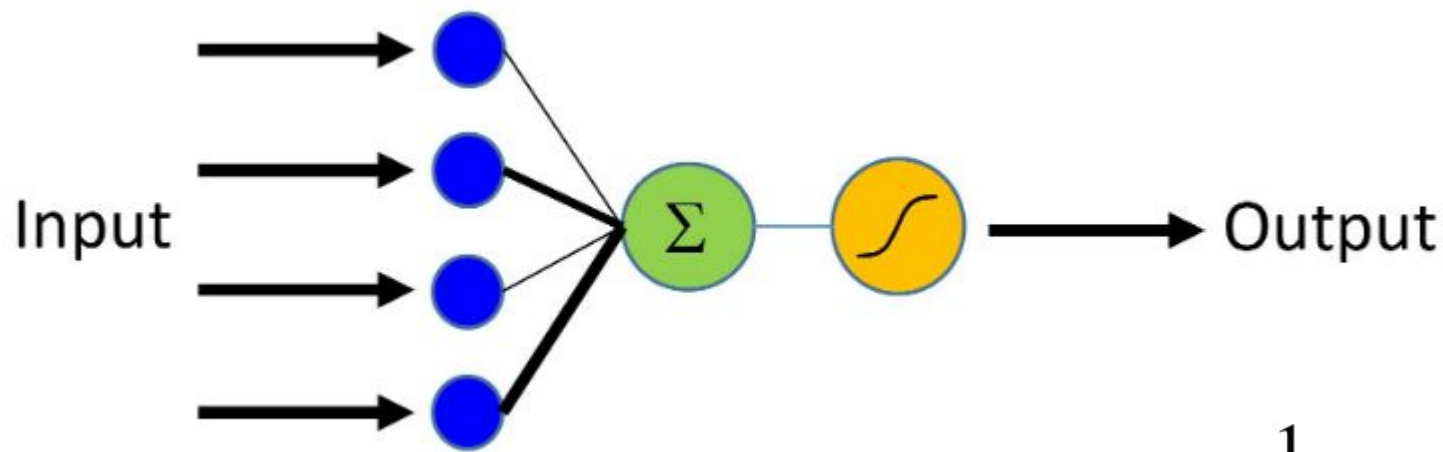
# Logistic Regression

# Binary Classification

- Transactions (fraudulent or no)
- People with disease or no
- Tumors (benign or not)
- E-mail (spam or not)



## Model



$$h_{\theta}(x) = g(\theta^T x) \quad g(z) = \frac{1}{1 + e^{-z}}$$

$$z = \theta^T x$$

# Sigmoid or Logistic Function

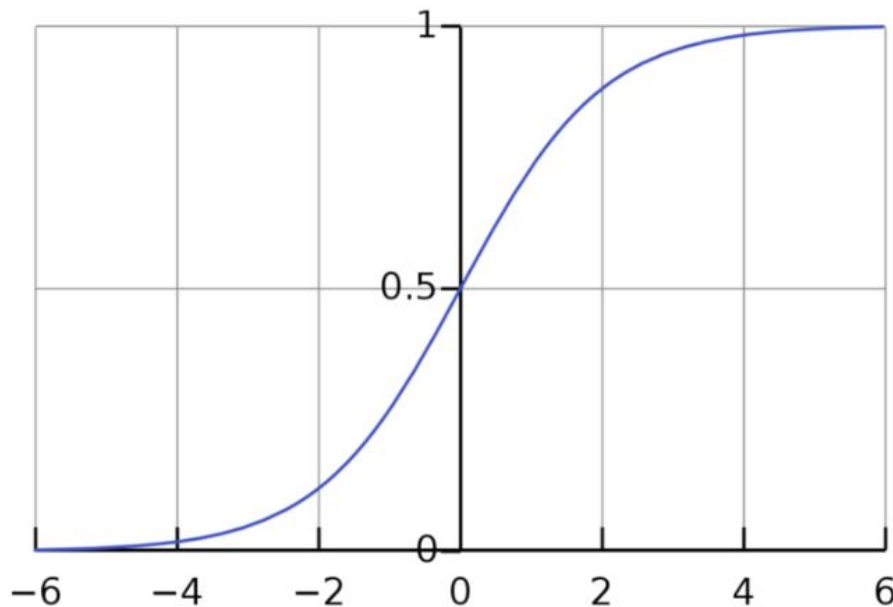
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$z = \theta^T x$$

$$z = 0, e^0 = 1 \Rightarrow g(z) = 1/2$$

$$z \rightarrow \infty, e^{-\infty} \rightarrow 0 \Rightarrow g(z) = 1$$

$$z \rightarrow -\infty, e^{\infty} \rightarrow \infty \Rightarrow g(z) = 0$$



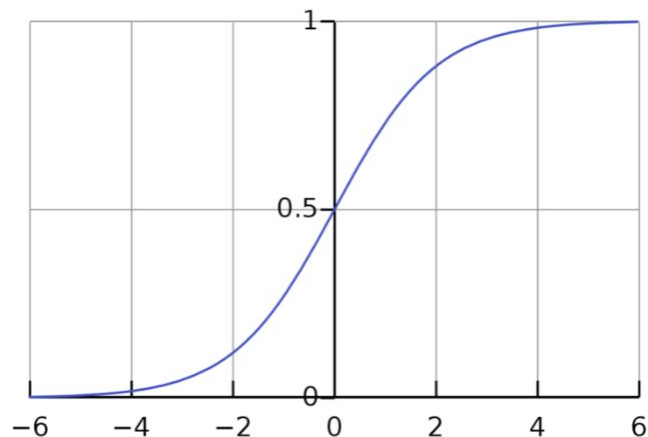
# Interpretation

$$h_{\theta}(x) \geq 0.5 \rightarrow y = 1$$

$$h_{\theta}(x) < 0.5 \rightarrow y = 0$$

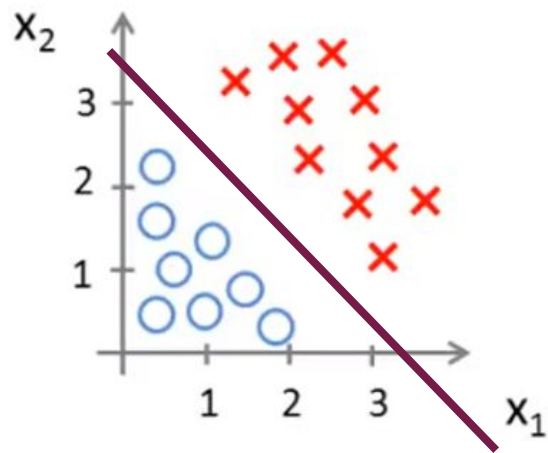
$$h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

$$\text{when } \theta^T x \geq 0$$



$$\theta^T x \geq 0 \Rightarrow y = 1$$

$$\theta^T x < 0 \Rightarrow y = 0$$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

Predict " $y = 1$ " if  $-3 + x_1 + x_2 \geq 0$



Training set:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

m examples  $x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \quad x_0 = 1, y \in \{0, 1\}$

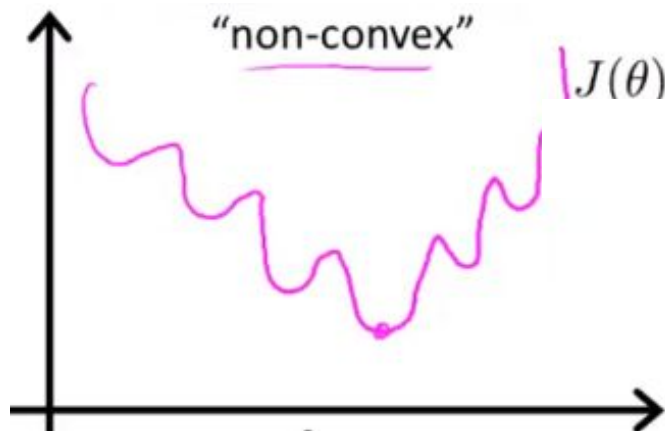
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose parameters  $\theta$  ?

$$\text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Linear regression:  $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Cost Function

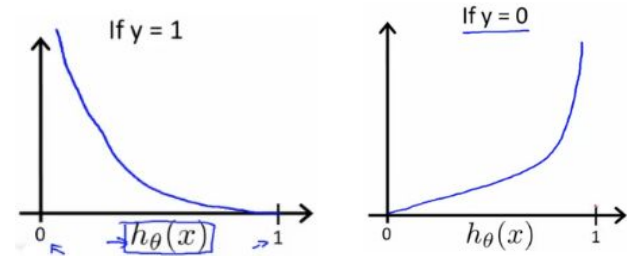
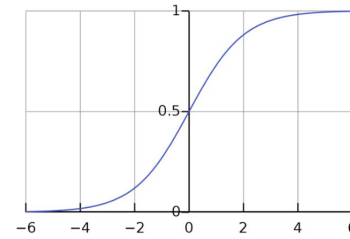


# LR COST FUNCTION

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

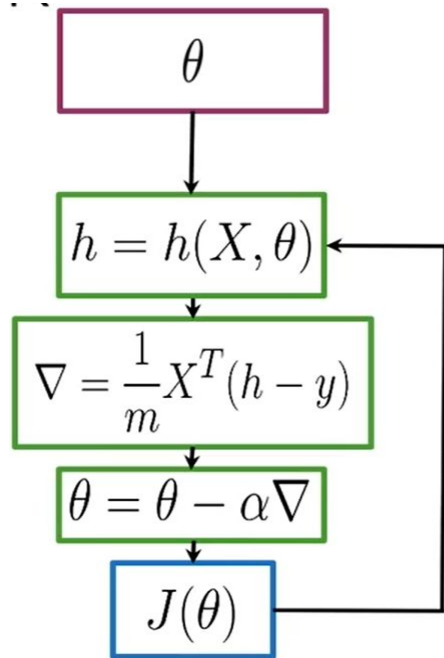
$$\text{Cost}(h_{\theta}(x), y) = -\log(h_{\theta}(x)) \quad \text{if } y = 1$$

$$\text{Cost}(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x)) \quad \text{if } y = 0$$

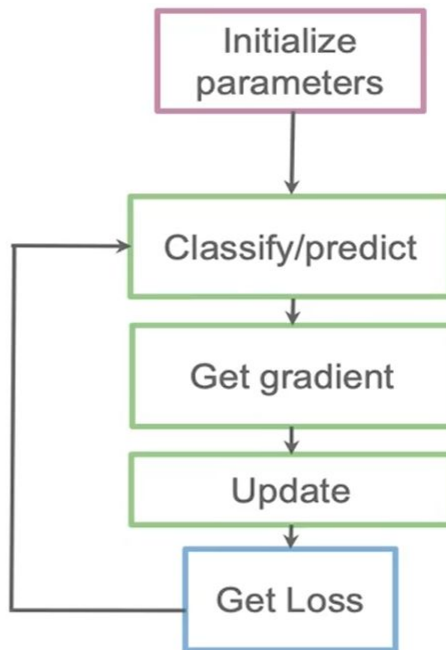


$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$$

# Training: Gradient Descent



Until good enough



$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \end{array} \right\}$$

# Multiclass Classification

- Email tagging: Work, Friends, Family
- Skin Lesion: Melanoma, Carcinoma, Nevus, Keratosis
- Video: Pornography, Violence, Gore scenes, Child abuse

# Classification

Email tagging: Work, Friends, Family

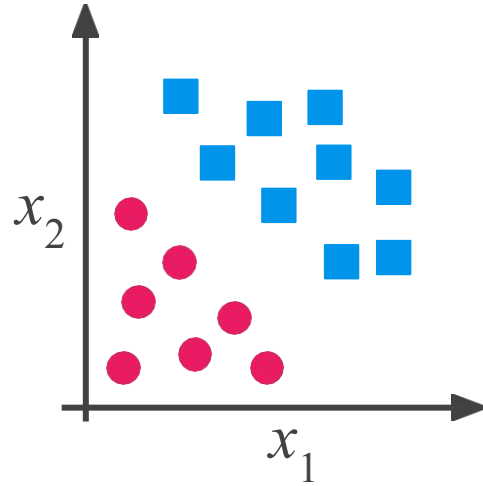
$$y = 1 \quad y = 2 \quad y = 3$$

Skin Lesion: Melanoma, Carcinoma, Nevus, Keratosis

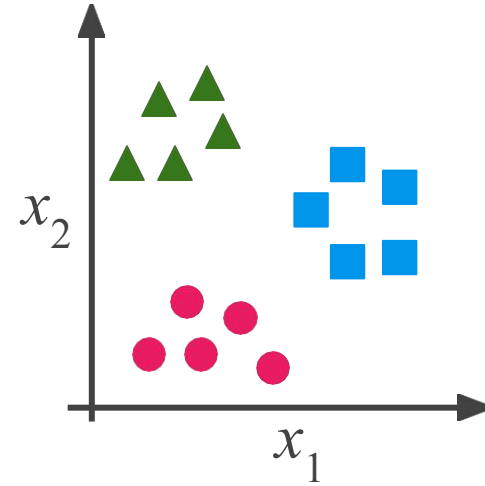
$$y = 1 \quad y = 2 \quad y = 3 \quad y = 4$$

Video: Pornography, Violence, Gore scenes, Child abuse

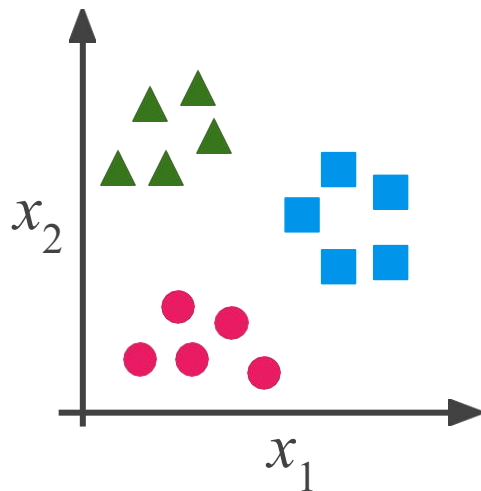
## Binary Classification



## Multi-class Classification



# One-vs-All (One-vs-Rest)



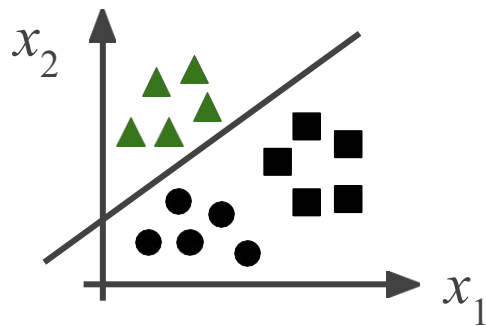
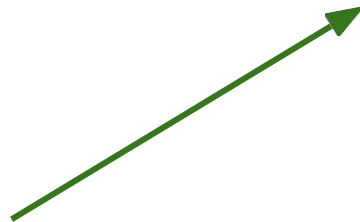
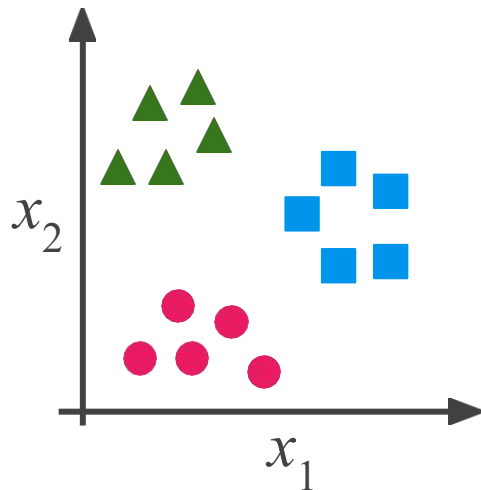
Class 1: ▲

Class 2: ■

Class 3: ●



# One-vs-All (One-vs-Rest)

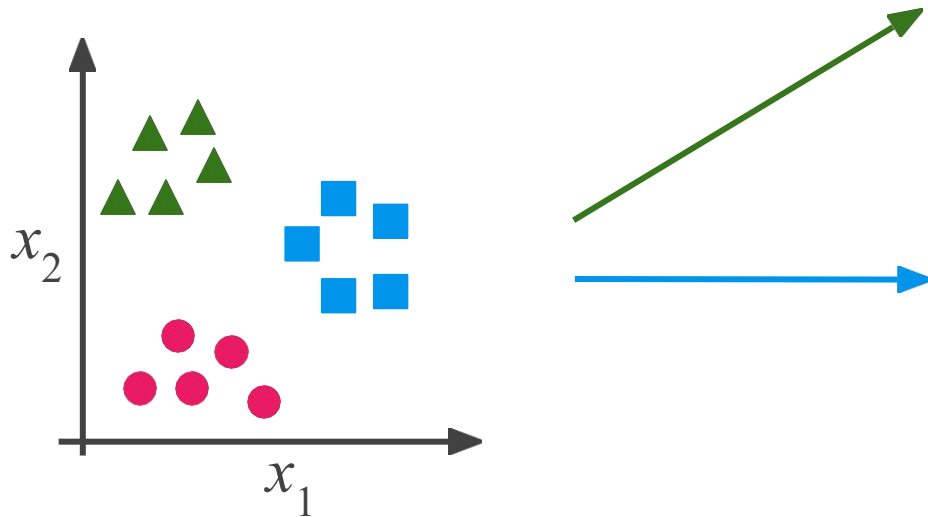


Class 1: ▲

Class 2: ■

Class 3: ●

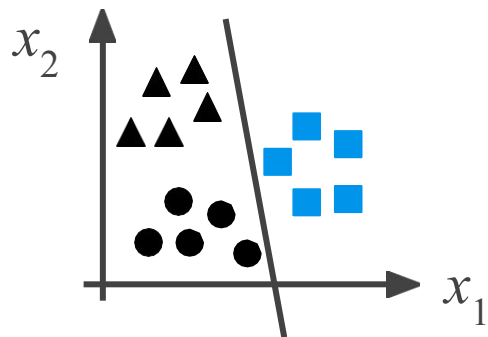
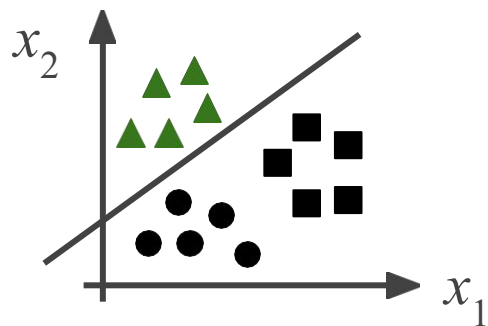
# One-vs-All (One-vs-Rest)



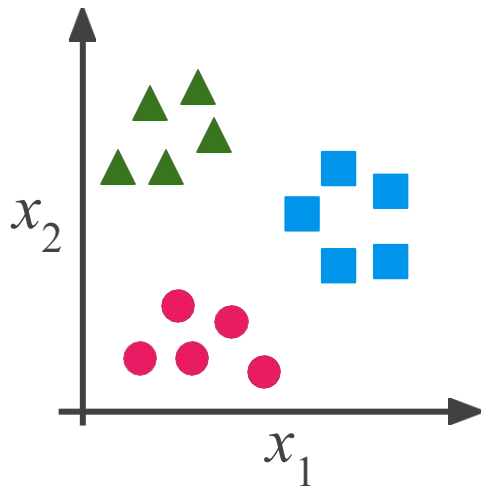
Class 1: ▲




Class 2: ■

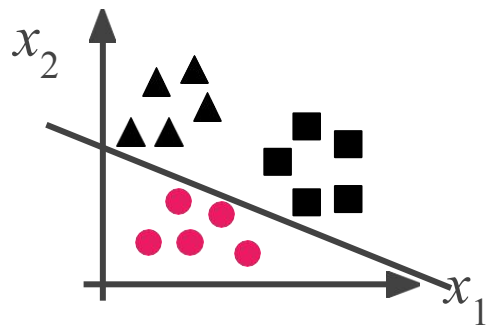
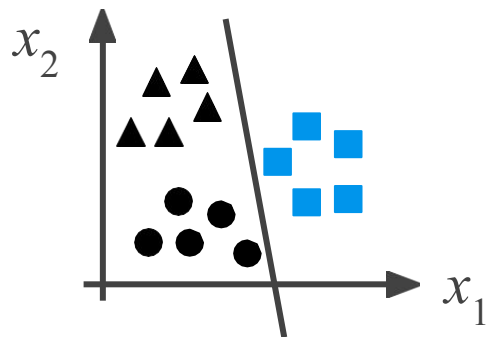
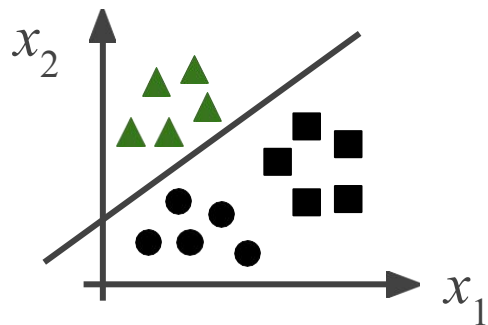
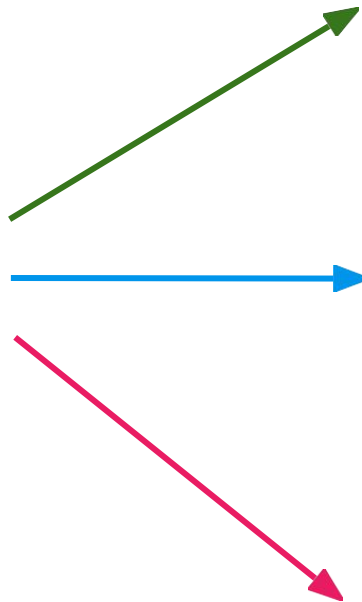
Class 3: ●



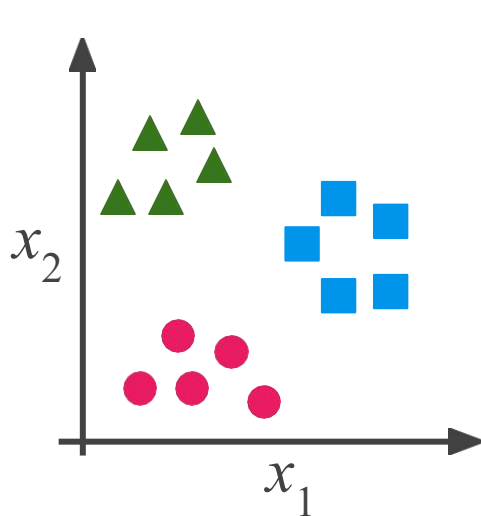
# One-vs-All (One-vs-Rest)



Class 1:   
Class 2:   
Class 3: 



# One-vs-All (One-vs-Rest)

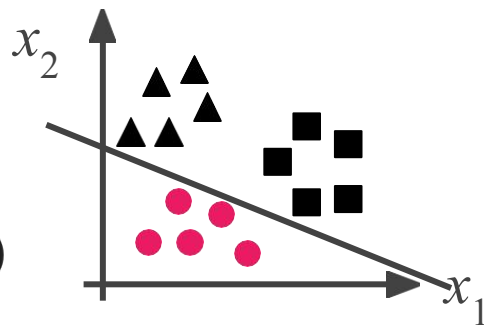
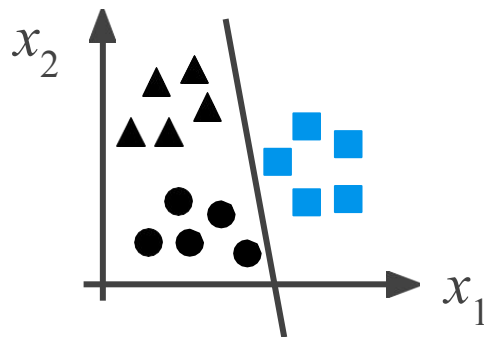
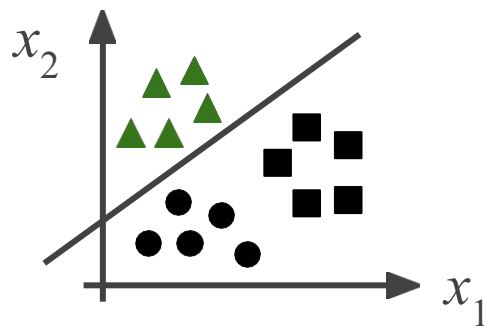
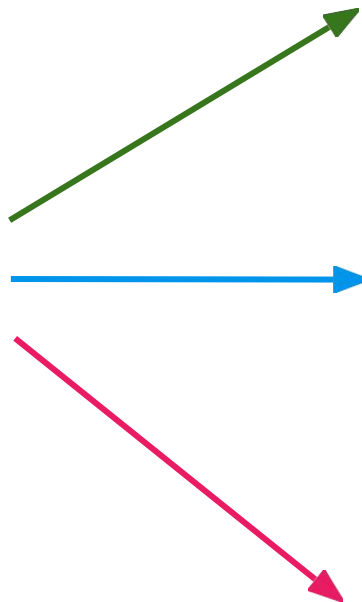


Class 1: ▲

Class 2: ■

Class 3: ●

$$h_{\theta}^{(i)}(x) = P(y = i \mid x; \theta) \quad (i=1,2,3)$$



## One-vs-All (One-vs-Rest)

Train a logistic regression classifier  $h_{\theta}^{(i)}(x)$  for each class  $i$  to predict the probability that  $y = i$ .

On a new input  $x$ , to make a prediction, pick the class  $i$  that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

# References

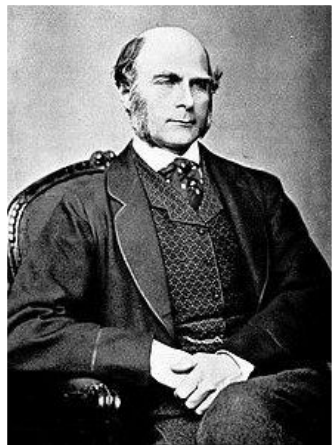
## Machine Learning Books

- Hands-On Machine Learning with Scikit-Learn and TensorFlow, Chap. 4
- Pattern Recognition and Machine Learning, Chap. 4

## Machine Learning Courses

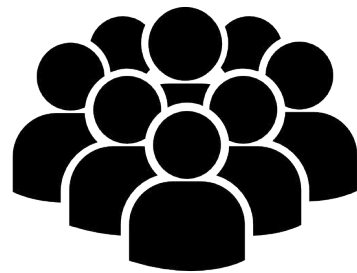
- <https://www.coursera.org/learn/machine-learning>, Week 3
- Logistic Regression — The Math of Intelligence (Week 2):  
<https://youtu.be/D8alok2P468>
- <http://cs229.stanford.edu/notes/cs229-notes1.pdf>

# Ensemble Learning



Francis Galton  
(1822-1909)

Animal's weight?



~800 people

1,197 kg

**1,207 kg**





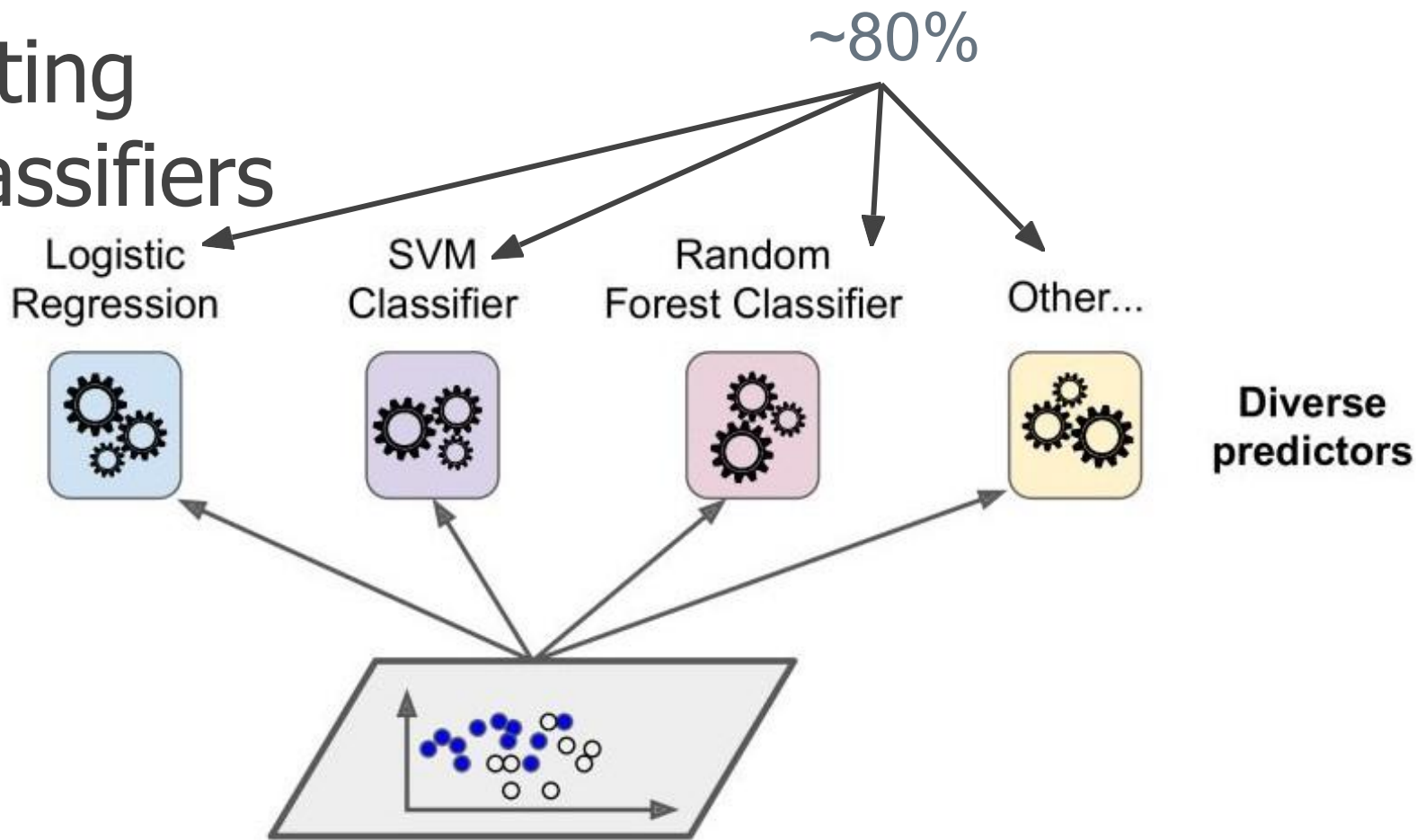
## Wisdom of the Crowd



# Ensemble Learning

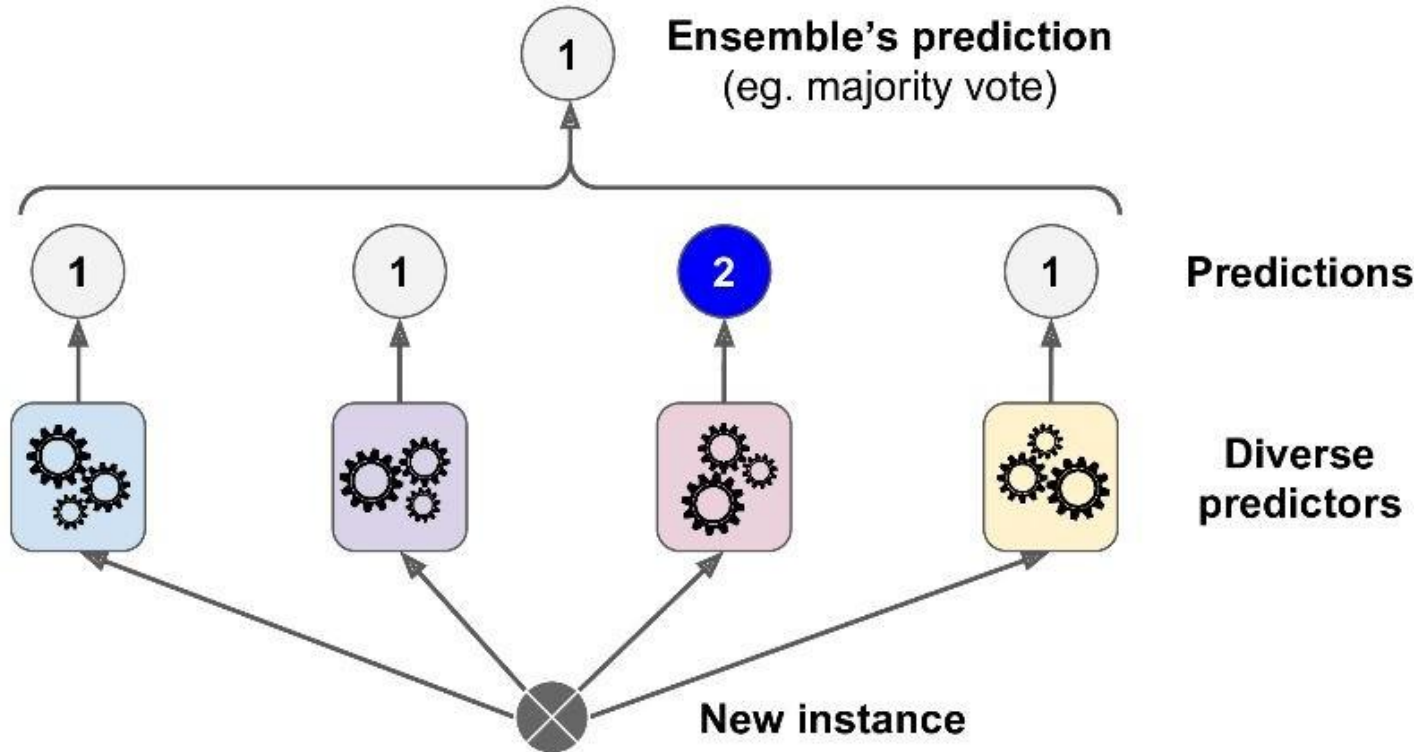
- Multiple learning algorithms **to obtain better predictive performance** than could be obtained from any learning algorithms individually.

# Voting Classifiers



# Voting Classifiers

Hard/Soft voting classifier



# Voting Classifiers

- Voting classifier **often achieves a higher accuracy than the best classifier** in the ensemble.

# Voting Classifiers

- Voting classifier **often achieves a higher accuracy than the best classifier** in the ensemble.
- Even if each classifier is a **weak learner**, the ensemble can still be a **strong learner**, provided there are a sufficient number of weak learners and they are sufficiently diverse.

# Voting Classifiers

- Ensemble methods work best when the predictors are as **independent** from one another as possible.



# Voting Classifiers

- Ensemble methods work best when the predictors are as **independent** from one another as possible.
- One way to get diverse classifiers is to train them using **very different algorithms**: this increases the chance that they will make very different types of errors, improving the ensemble's accuracy.



# Voting Classifiers

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC

log_clf = LogisticRegression()
rnd_clf = RandomForestClassifier()
svm_clf = SVC()

voting_clf = VotingClassifier(
    estimators=[('lr', log_clf), ('rf', rnd_clf), ('svc', svm_clf)],
    voting='hard'
)

voting_clf.fit(X_train, y_train)
```

# Voting Classifiers

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC

log_clf = LogisticRegression()
rnd_clf = RandomForestClassifier()
svm_clf = SVC()

voting_clf = VotingClassifier(
    estimators=[('lr', log_clf), ('rf', rnd_clf), ('svc', svm_clf)],
    voting='hard'
)

voting_clf.fit(X_train, y_train)
```

```
LogisticRegression 0.864
RandomForestClassifier 0.896
SVC 0.888
VotingClassifier 0.904
```

# Ensemble Methods

- **Bagging (and Pasting)**
- Boosting
- Stacking

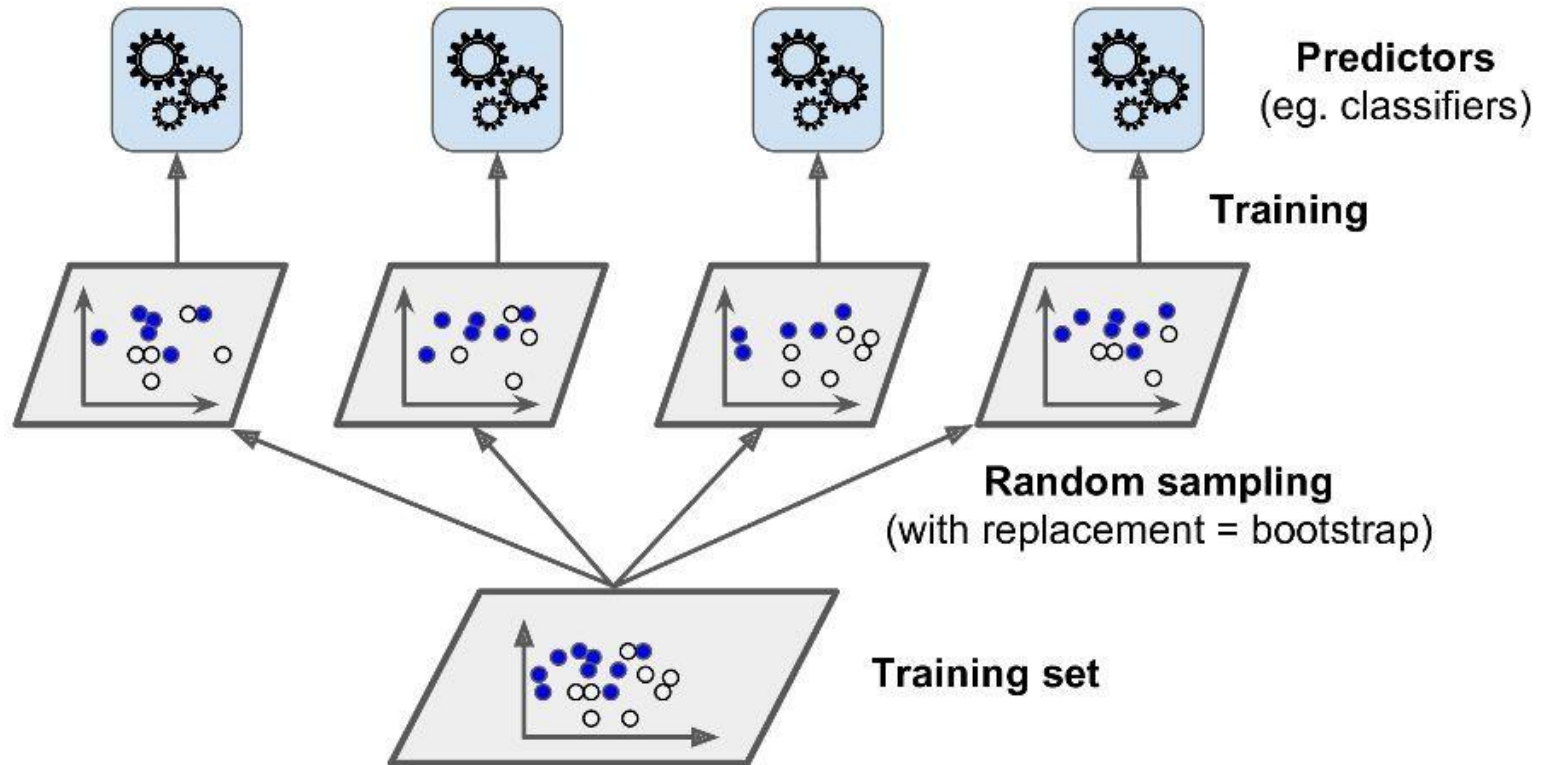
# Bagging and Pasting

- Use **the same training algorithm** for every predictor, but to train them on **different random subsets** of the training set.

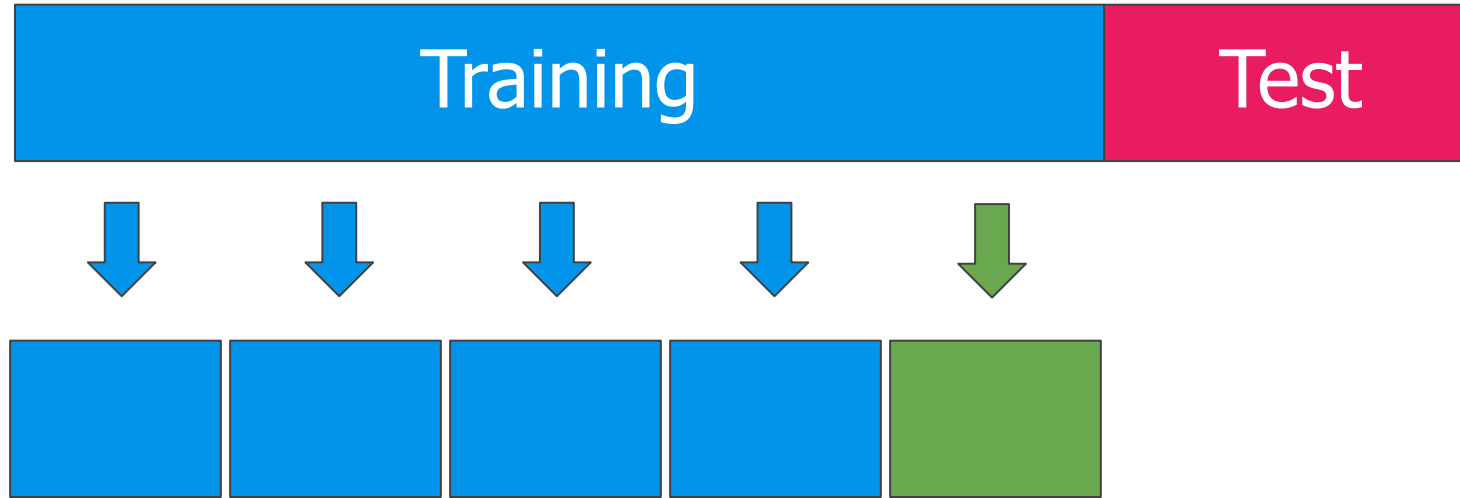
# Bagging and Pasting

- Use **the same training algorithm** for every predictor, but to train them on **different random subsets** of the training set.
- **Bagging** (short for Bootstrap Aggregating): sampling is performed **with** replacement.
- **Pasting**: sampling is performed **without** replacement.

# Bagging and Pasting

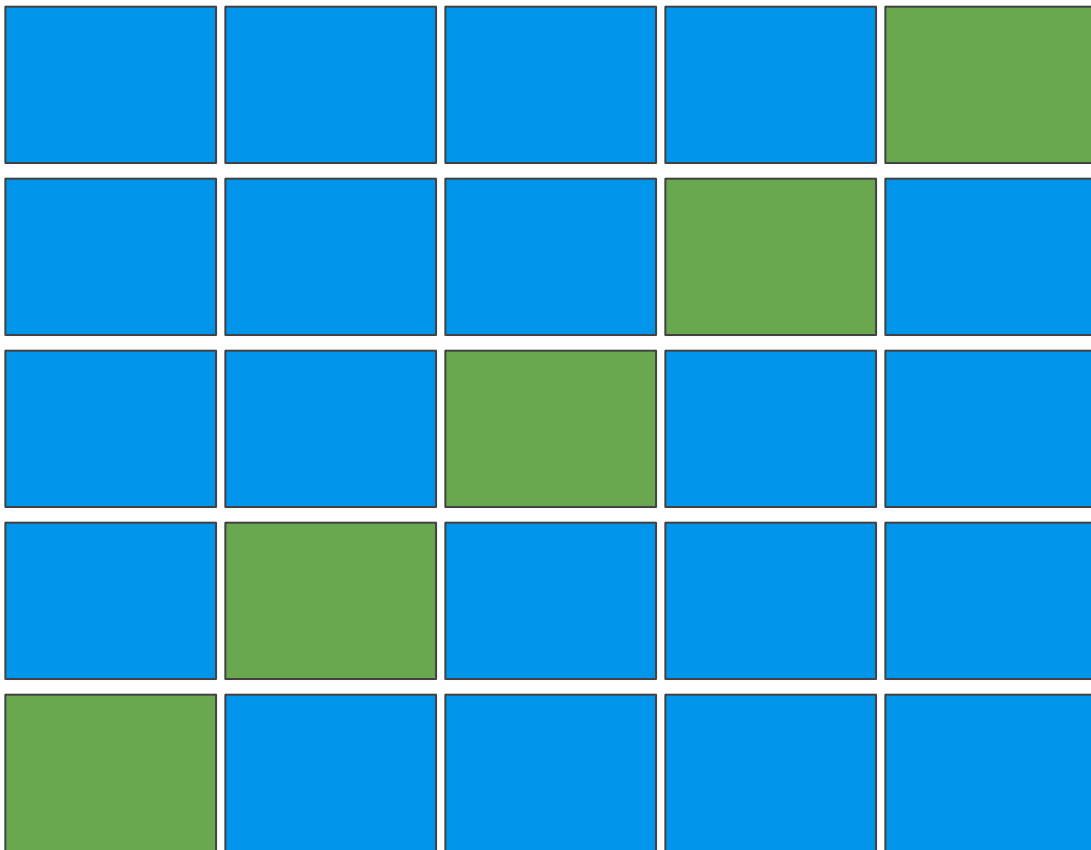


# Cross Validation



Training

Test

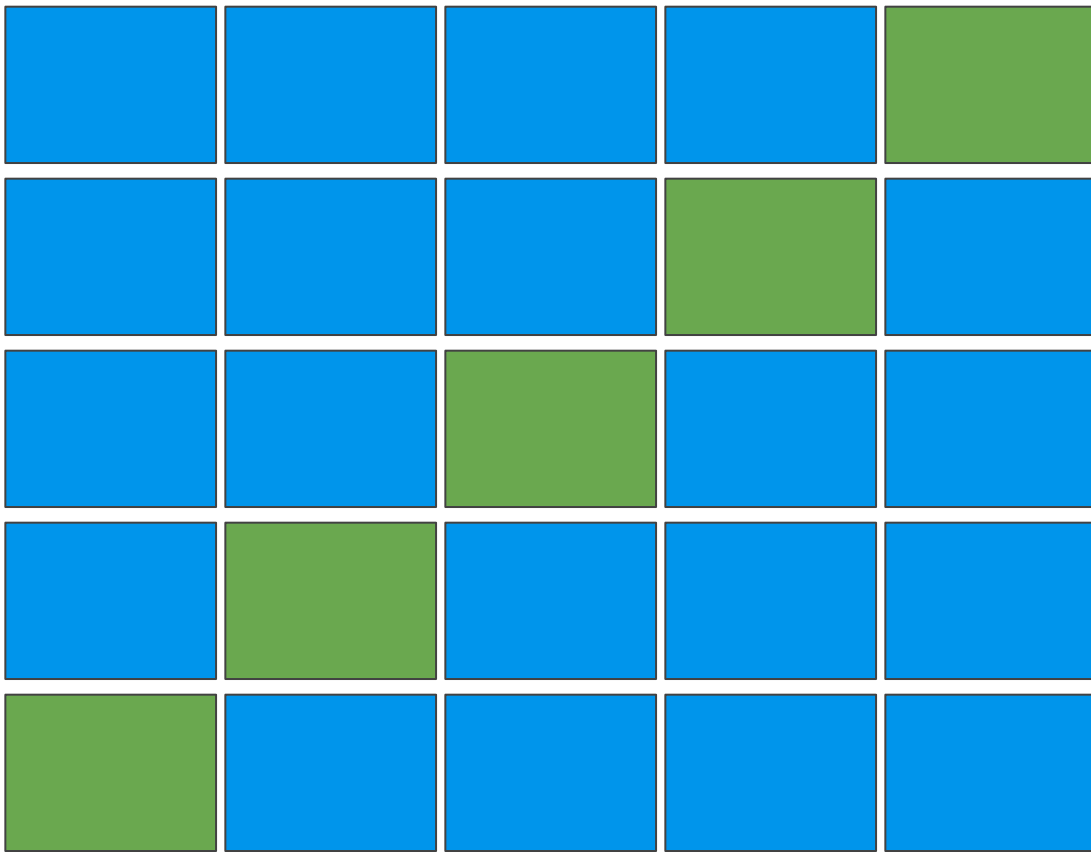


Cross  
Validation



Training

Test



Cross  
Validation  
(one  
model)

Training

Test

Randomsubset

Random subset

Random subset

Random subset

Random subset

Bagging  
(many  
models)

Training

Test



Bagging

# Bagging and Pasting

- Once all predictors are trained, the ensemble can make a prediction for a new instance by simply **aggregating the predictions of all predictors**.
- **Bagging and Pasting scale very well.**

# Bagging and Pasting

```
from sklearn.ensemble import BaggingClassifier
from sklearn.tree import DecisionTreeClassifier
bag_clf = BaggingClassifier(
    DecisionTreeClassifier(),
    n_estimators=500, max_samples=100,
    bootstrap=True, n_jobs=-1
)
bag_clf.fit(X_train, y_train)
y_pred = bag_clf.predict(X_test)
```

# Bagging and Pasting

- **Random Patches Ensemble** method: sampling **both** training instances and **features**.
- This is particularly useful when dealing with high-dimensional inputs.

# Ensemble Methods

- Bagging (and Pasting)
- **Boosting**
- Stacking

# Boosting

- The general idea of most boosting methods is **to train predictors sequentially**, each trying to correct its predecessor.



# Boosting

- The general idea of most boosting methods is **to train predictors sequentially**, each trying to correct its predecessor.
- Most popular: AdaBoost and Gradient Boost.

# AdaBoost [Freund and Schapire, 1997]

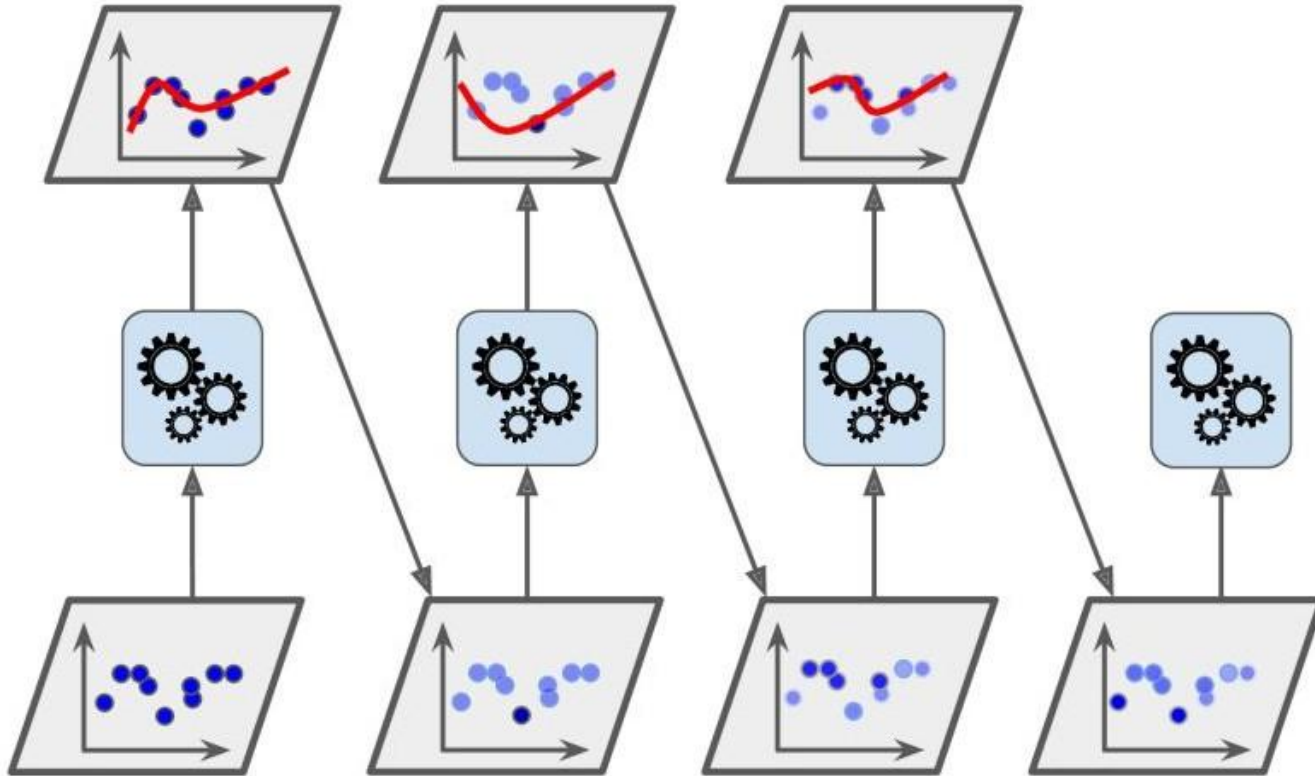
- One way for a new predictor to correct its predecessor is to pay a bit **more attention** to the training instances that **the predecessor underfitted**.

# AdaBoost [Freund and Schapire, 1997]

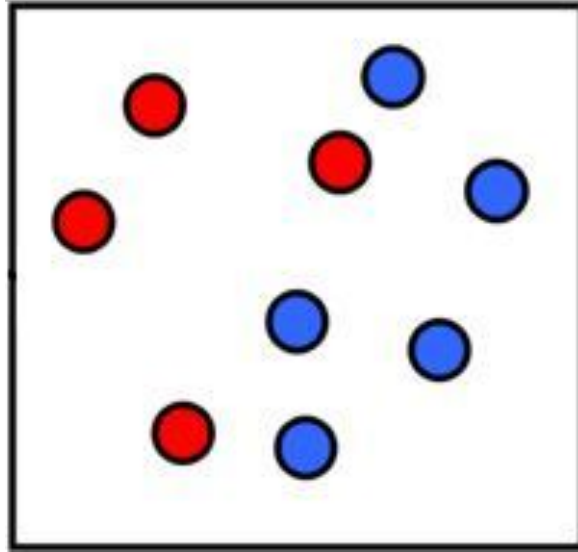
- One way for a new predictor to correct its predecessor is to pay a bit **more attention** to the training instances that **the predecessor underfitted**.
- This results in new predictors focusing more and more on **the hard cases**.

“A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting”, Y. Freund and R. Schapire (1997):  
<http://goo.gl/OlduRW>

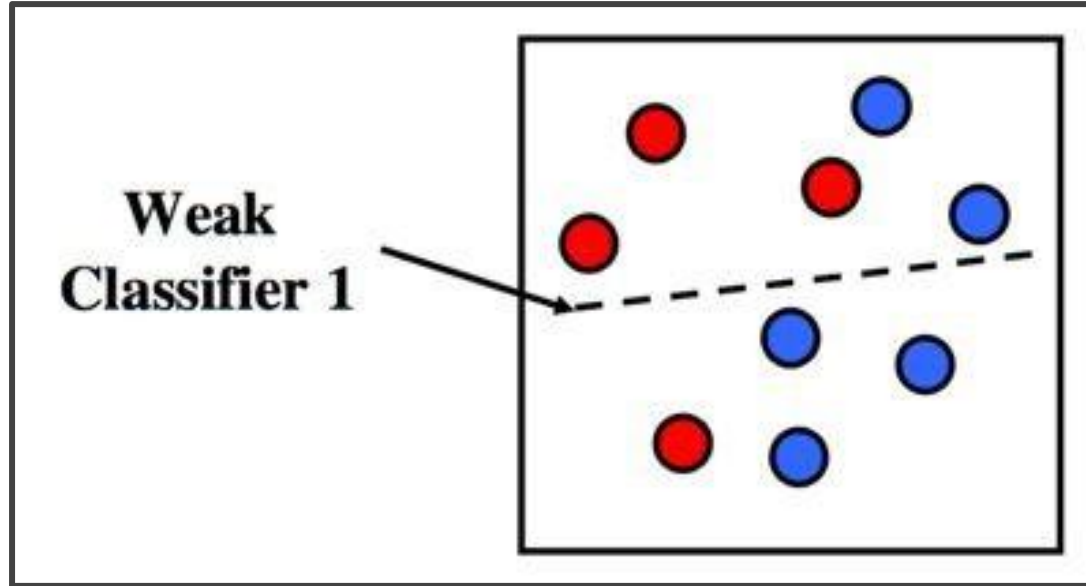
# AdaBoost



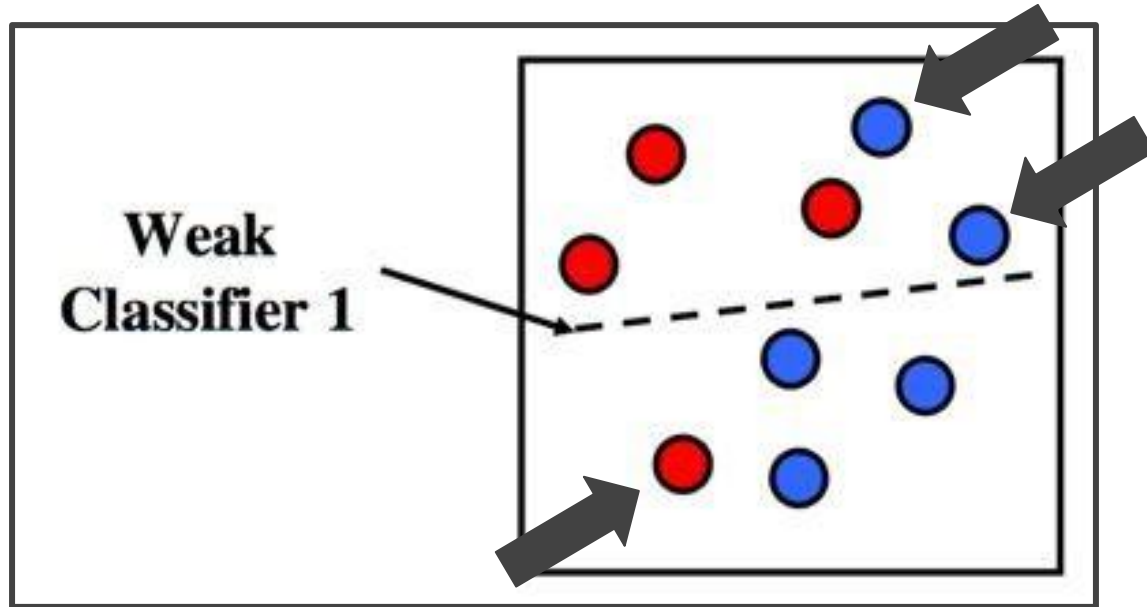
# AdaBoost



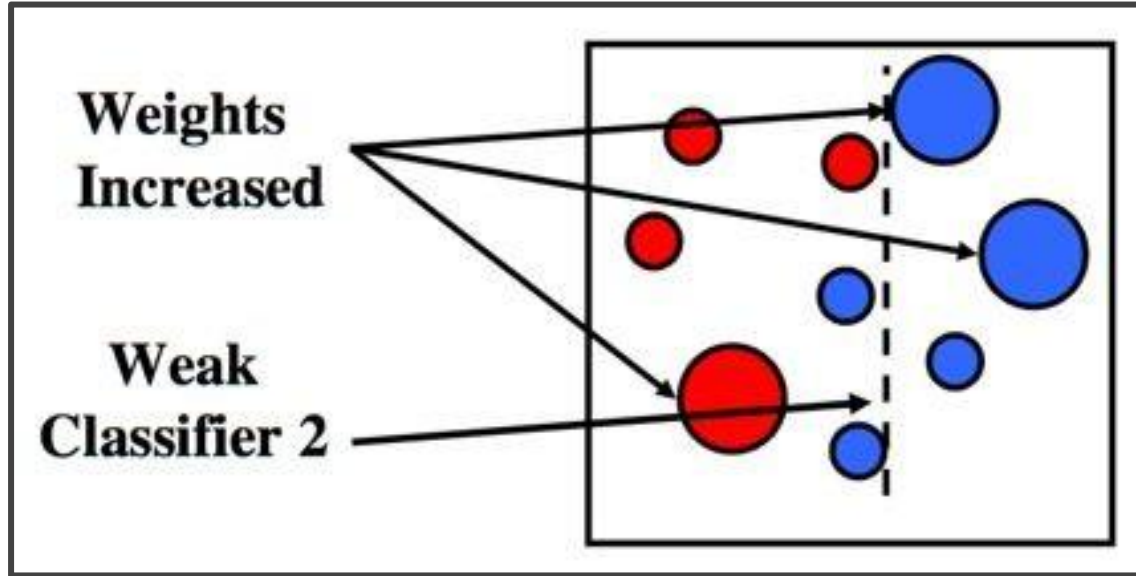
# AdaBoost



# AdaBoost

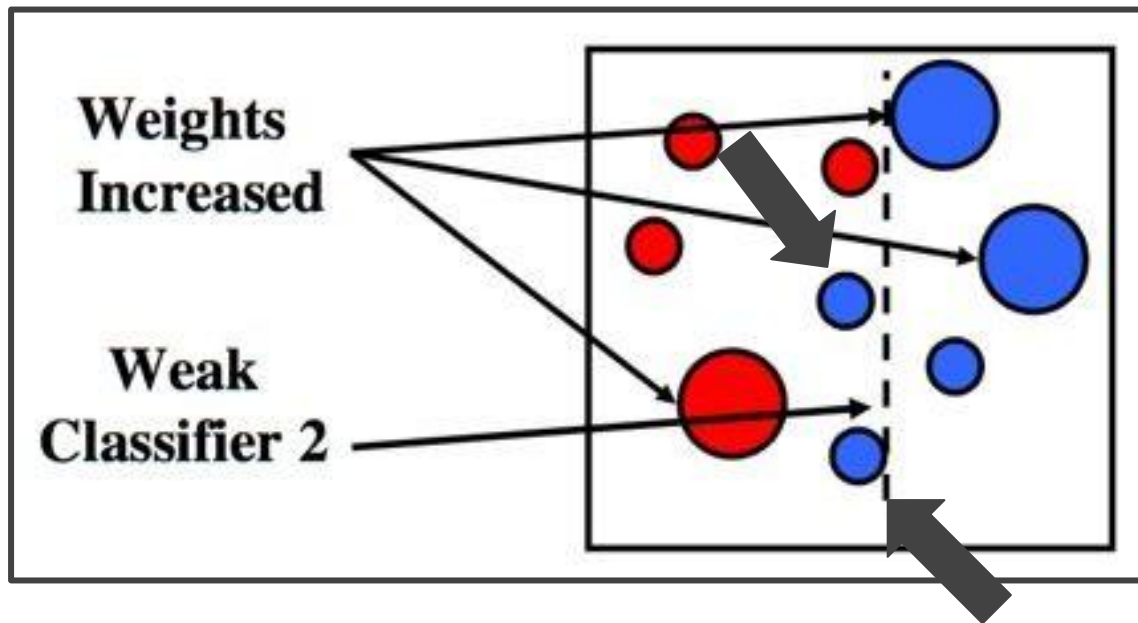


# AdaBoost

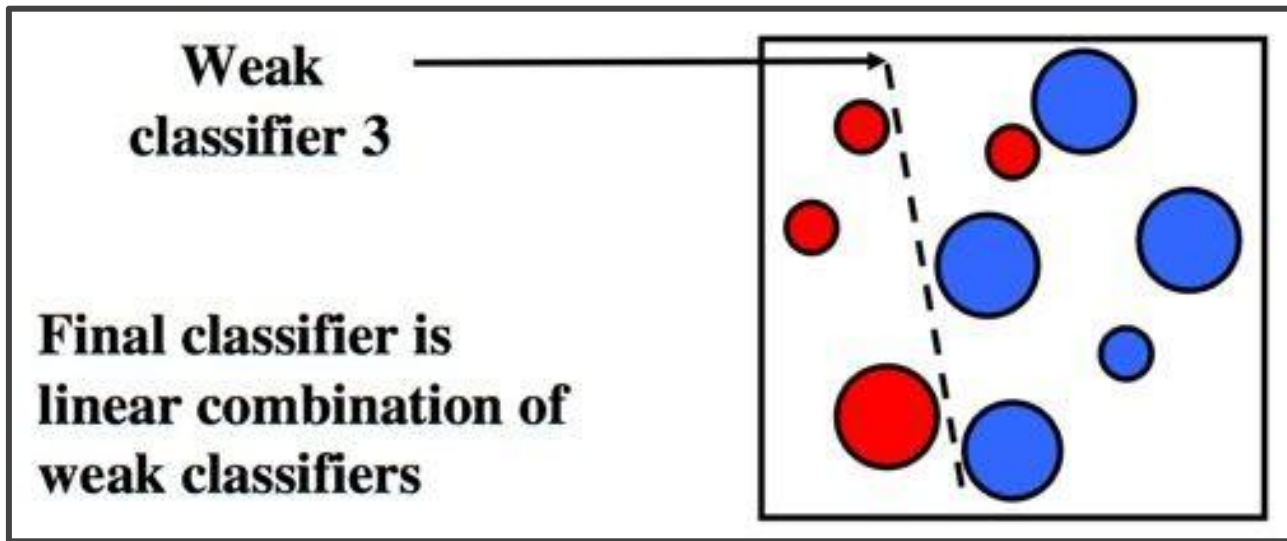




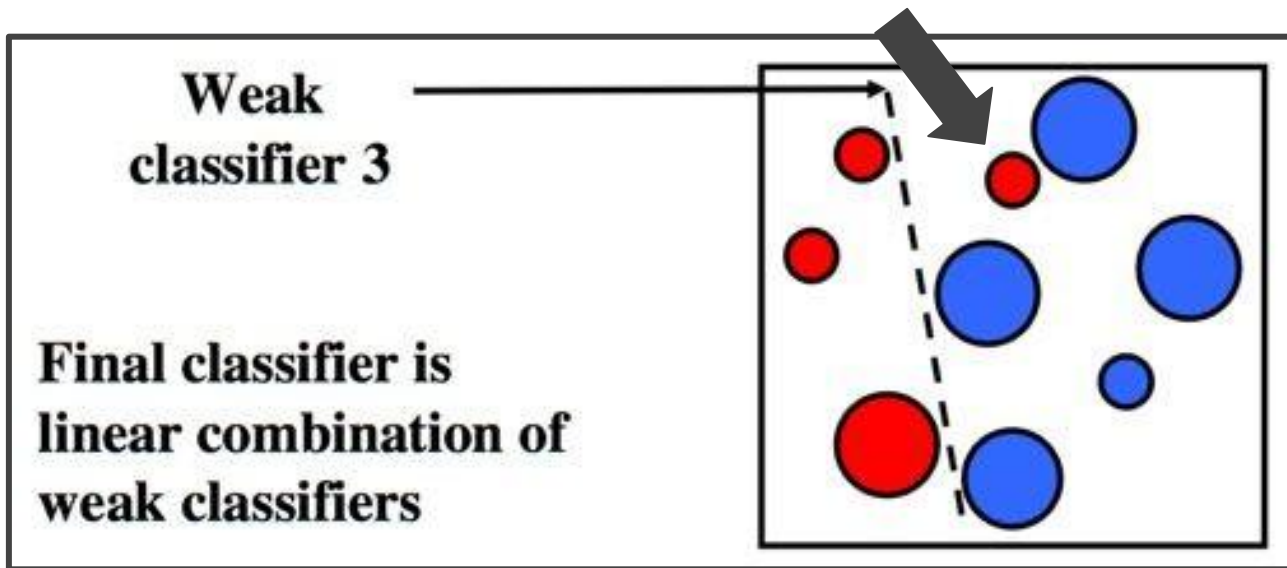
# AdaBoost



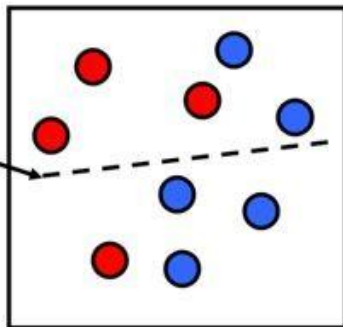
# AdaBoost



# AdaBoost

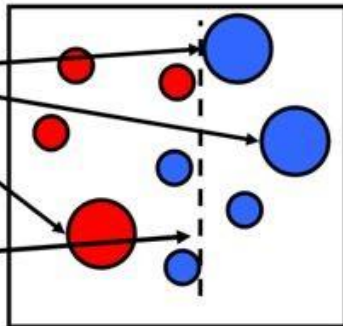


**Weak  
Classifier 1**

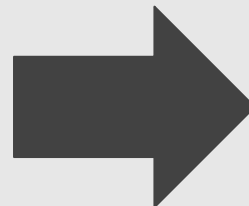
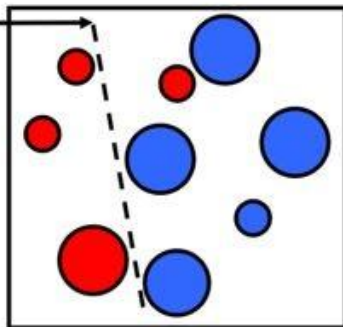


**Weights  
Increased**

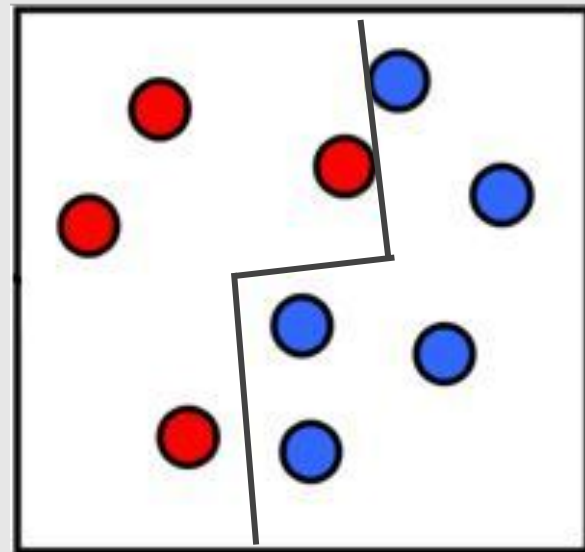
**Weak  
Classifier 2**



**Weak  
classifier 3**



**Final classifier is  
linear combination of  
weak classifiers**



# AdaBoost

1. Assign every observation,  $x_i$ , an initial weight value,  $w_i = \frac{1}{n}$ , where  $n$  is the total number of observations.
2. Train a "weak" model. (most often a decision tree)
3. For each observation:
  - 3.1. If predicted incorrectly,  $w_i$  is increased
  - 3.2. If predicted correctly,  $w_i$  is decreased
4. Train a new weak model where observations with greater weights are given more priority.
5. Repeat steps 3 and 4 until observations perfectly predicted or a preset number of trees are trained.

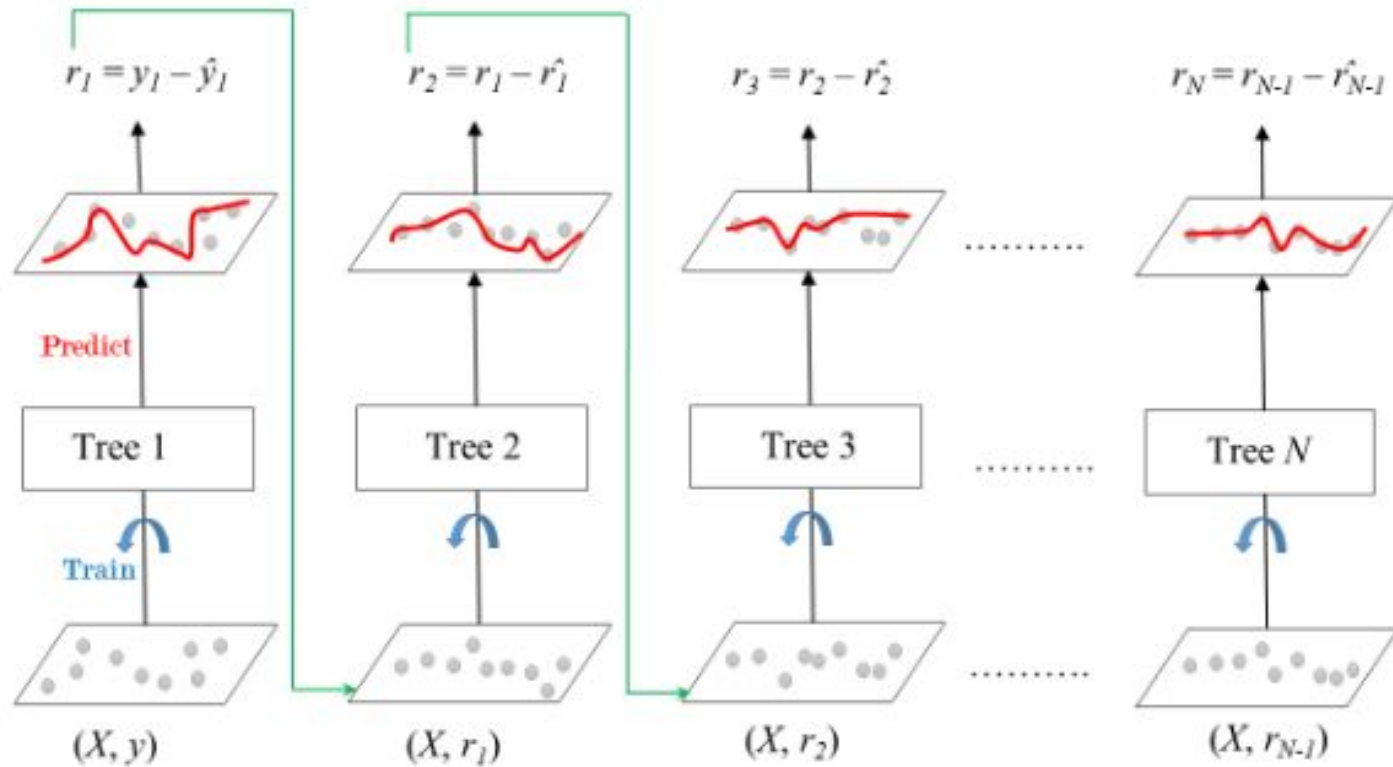
Chris Albon

# Gradient Boosting [Breiman, 1997]

- Instead of tweaking the instance weights at every iteration like AdaBoost does, this method fit the new predictor to the **residual errors** made by the previous predictor.

# Gradient Boosting [Breiman, 1997]

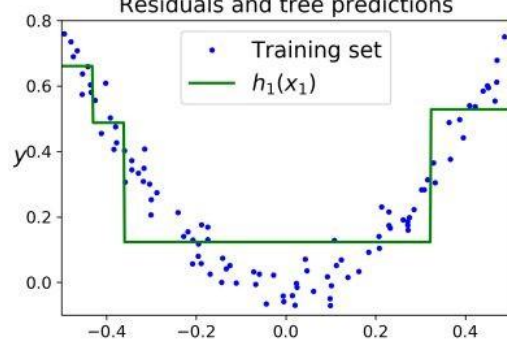
- Instead of tweaking the instance weights at every iteration like AdaBoost does, this method fit the new predictor to the **residual errors** made by the previous predictor.
- Instead of training on a newly sample distribution, the weak learner **trains on the remaining errors**.



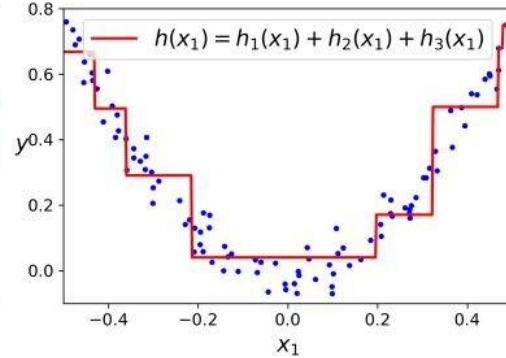
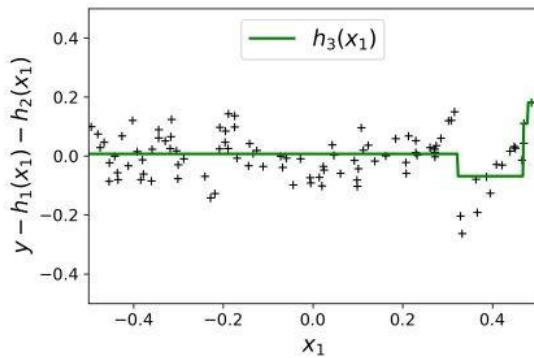
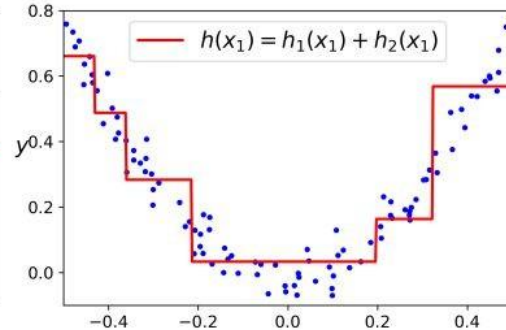
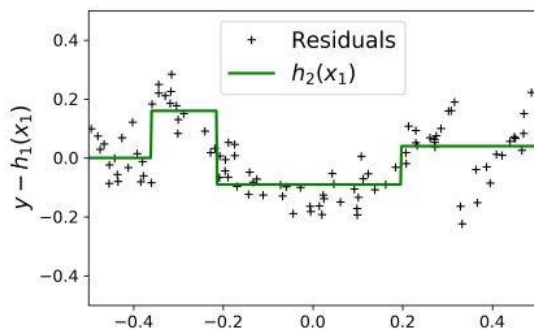
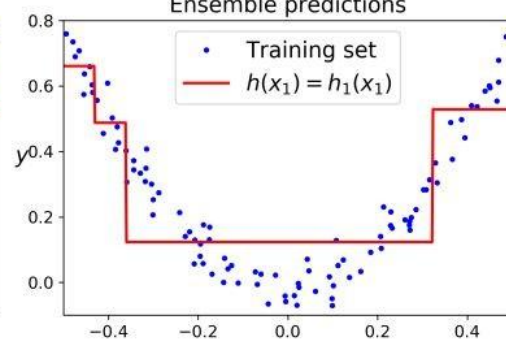
$$y(\text{pred}) = y_1 + (\text{eta} * r_1) + (\text{eta} * r_2) + \dots + (\text{eta} * r_N)$$



Residuals and tree predictions



Ensemble predictions



# Gradient Boosting [Breiman, 1997]

1. Fit a simple linear regressor or decision tree on data  
**[call x as input and y as output]**
2. Calculate error residuals. Actual target value, minus predicted target value  
**[ $e1 = y - y_{\text{predicted1}}$  ]**
3. Fit a new model on error residuals as target variable with same input variables  
**[call it  $e1_{\text{predicted}}$ ]**
4. Add the predicted residuals to the previous predictions  
**[ $y_{\text{predicted2}} = y_{\text{predicted1}} + e1_{\text{predicted}}$ ]**
5. Fit another model on residuals that is still left, i.e. **[ $e2 = y - y_{\text{predicted2}}$ ]** and repeat steps 2 to 5 until it starts overfitting or the sum of residuals become constant.

# Gradient Boosting [Breiman, 1997]

- XGboost [Chen and Guestrin, 2016]:

Extreme Gradient Boosting

<https://github.com/tqchen/xgboost>

It aims at being extremely fast, scalable and portable.

# Ensemble Methods

- Bagging (and Pasting)
- Boosting
- **Stacking**

# Stacking [Wolpert, 1992]

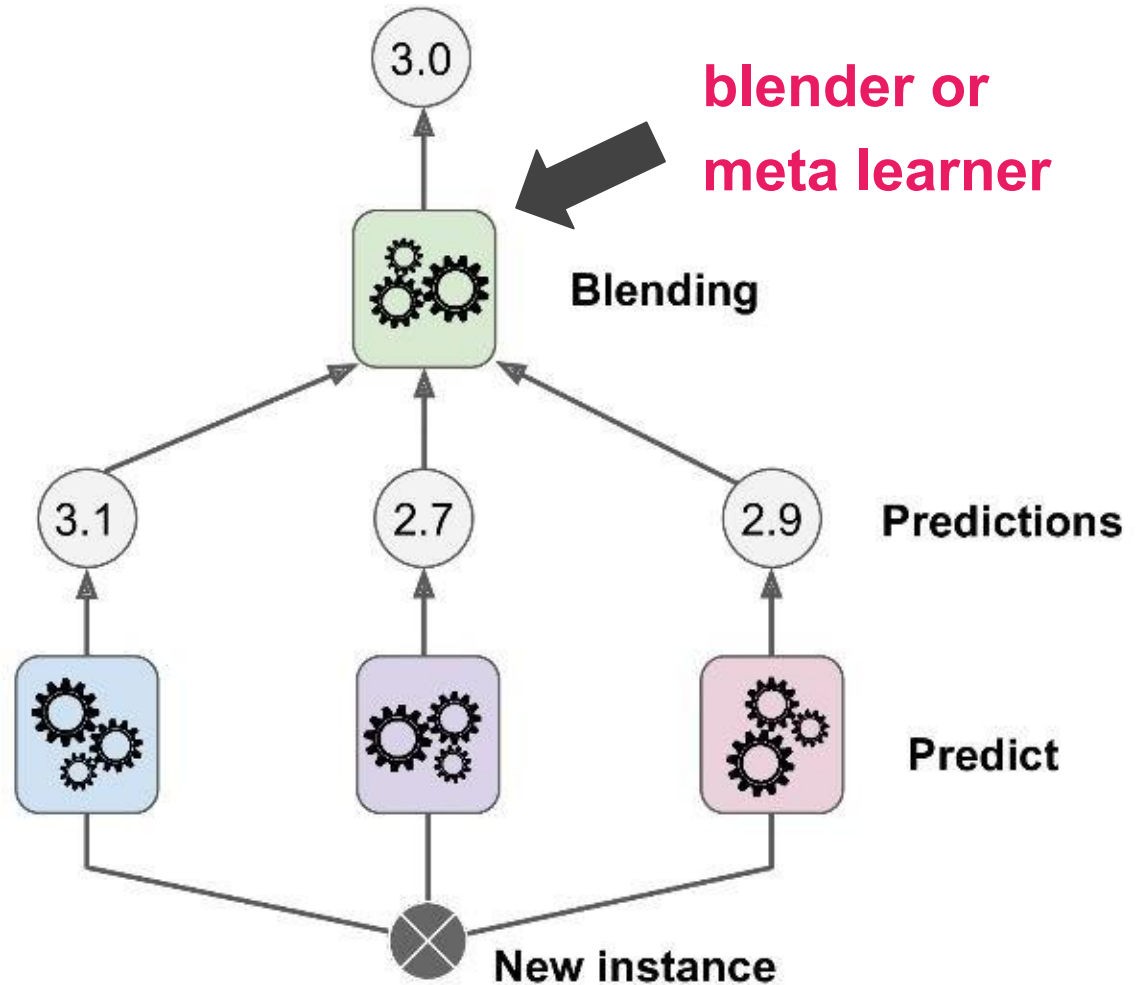
- Stacking (short for Stacked Generalization)

“Stacked Generalization”, D. Wolpert (1992): <http://goo.gl/9l2NBw>

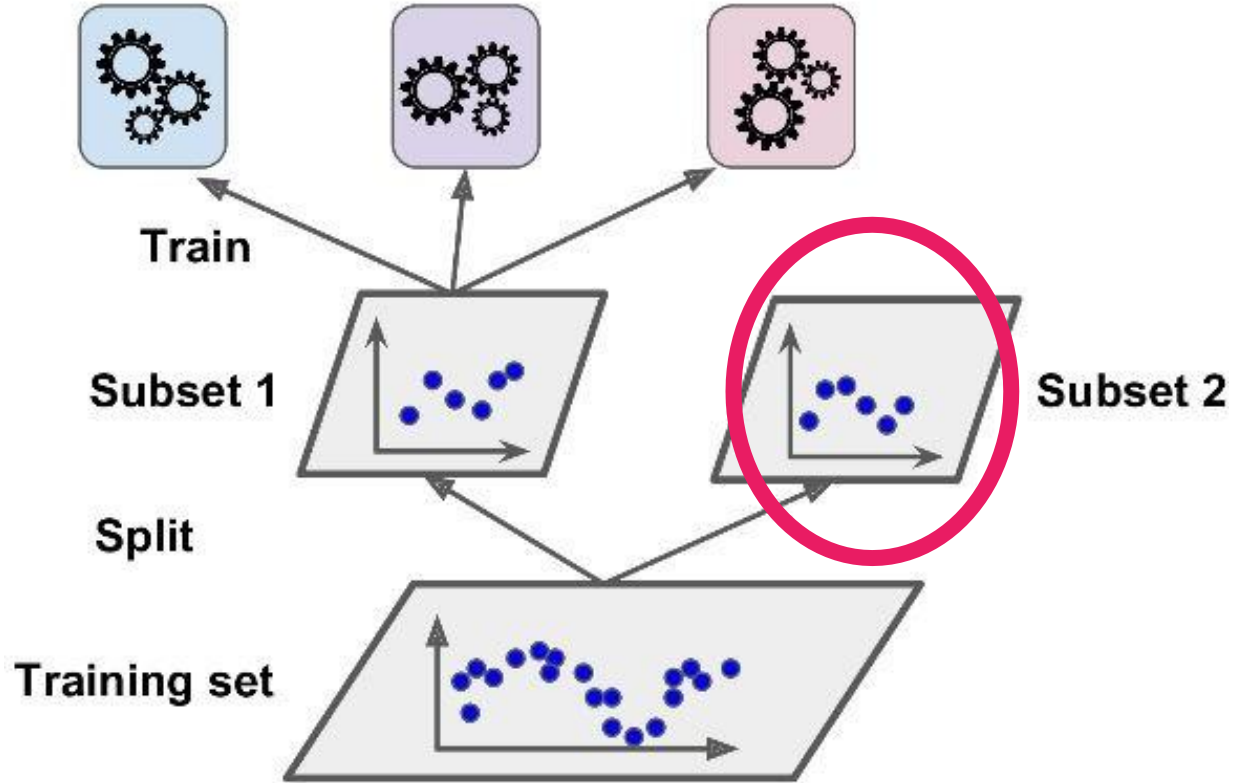
# Stacking [Wolpert, 1992]

- Stacking (short for Stacked Generalization)
- Instead of using trivial functions (such as hard voting) to aggregate the predictions of all predictors in an ensemble, we **train a model to perform this aggregation**.

# Stacking

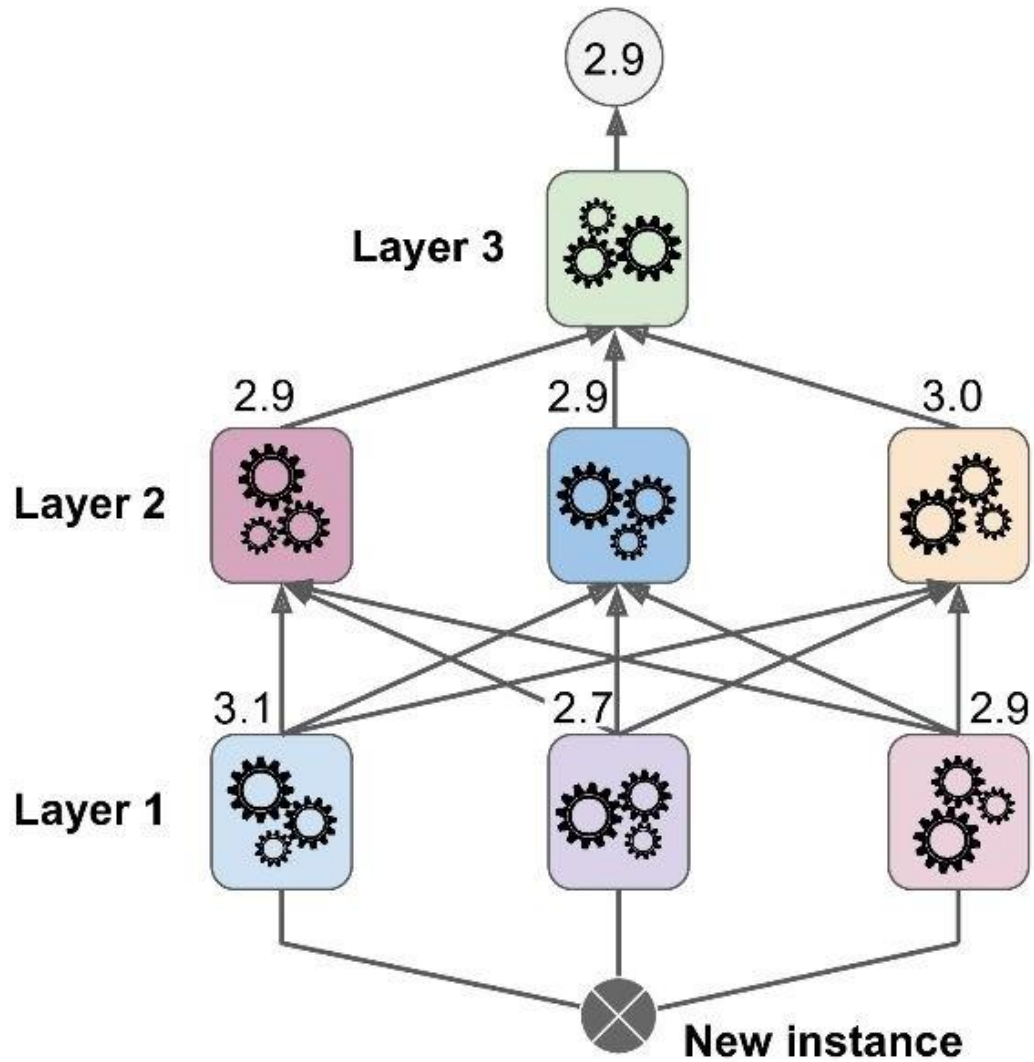


To train the  
blender, a  
common  
approach  
is to use a  
**hold-out set.**





# Multi-layer Stacking Ensemble



# Stacking [Wolpert, 1992]

- Scikit-Learn does not support stacking directly. =(

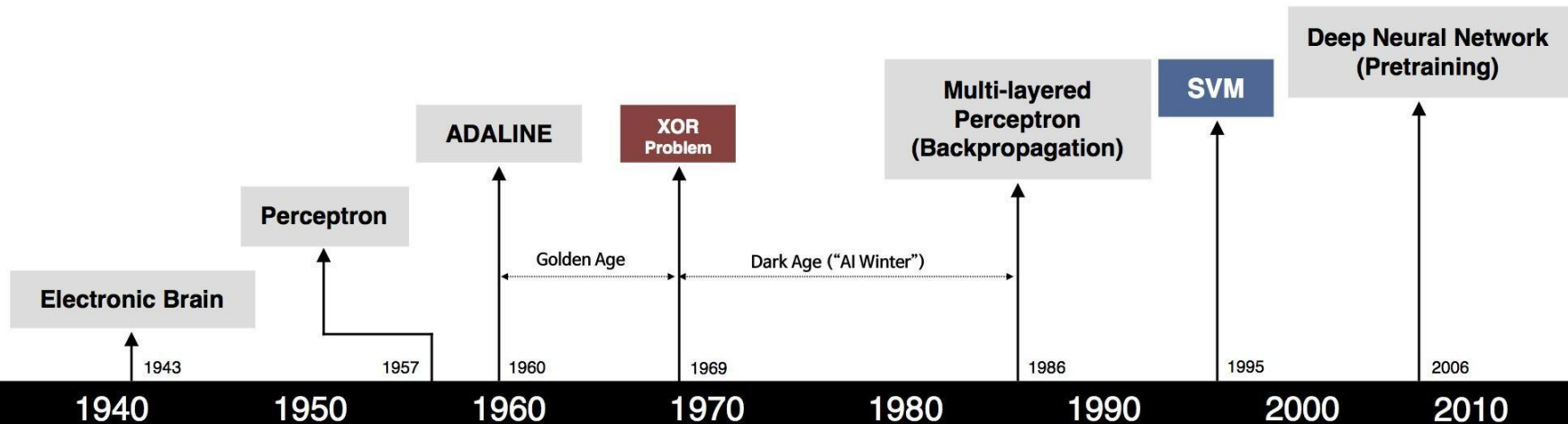
“Stacked Generalization”, D. Wolpert (1992): <http://goo.gl/9l2NBw>

# References

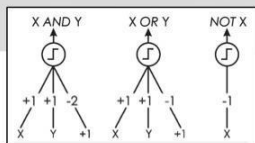
## Machine Learning Books

- Hands-On Machine Learning with Scikit-Learn and TensorFlow, Chap. 6 & 7
- Pattern Recognition and Machine Learning, Chap. 14
- Pattern Classification, Chap 8 & 9 (Sec. 9.5)
- “Scikit Learn Ensemble Learning, Bootstrap Aggregating (Bagging) and Boosting” <https://youtu.be/X3Wbfb4M33w>

# Support Vector Machines (SVMs)



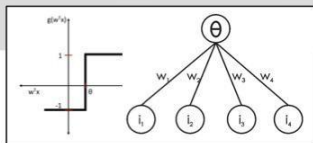
S. McCulloch – W. Pitts



- Adjustable Weights
- Weights are not Learned



F. Rosenblatt



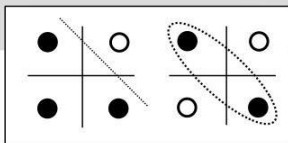
- Learnable Weights and Threshold



B. Widrow – M. Hoff



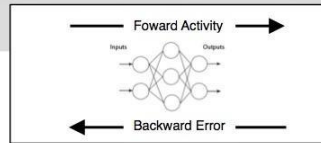
M. Minsky – S. Papert



- XOR Problem



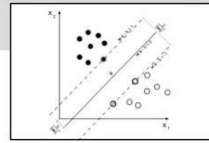
D. Rumelhart – G. Hinton – R. Williams



- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting



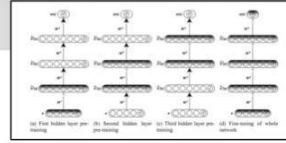
V. Vapnik – C. Cortes



- Limitations of learning prior knowledge
- Kernel function: Human Intervention



G. Hinton – S. Ruslan



- Hierarchical feature Learning

# Traditional Recognition



Classifier



“cat”



Edges



Classifier



“cat”



Edges



Histogram



Classifier



“cat”



Edges



Histogram



K-means  
Sparse  
code



Classifier



“cat”

# Traditional Recognition



**SVM**



“cat”



Edges



**SVM**



“cat”



Edges



Histogram



**SVM**



“cat”



Edges



Histogram



K-means  
Sparse  
code

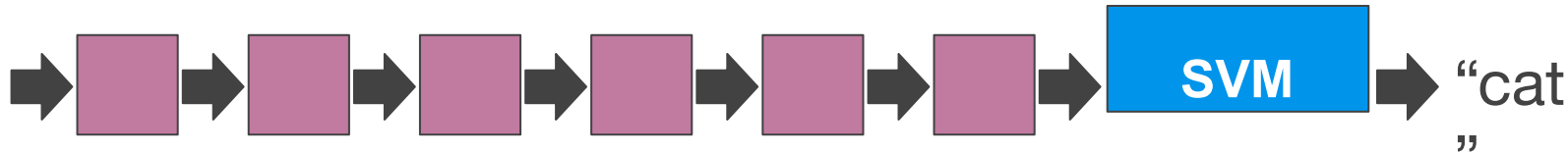
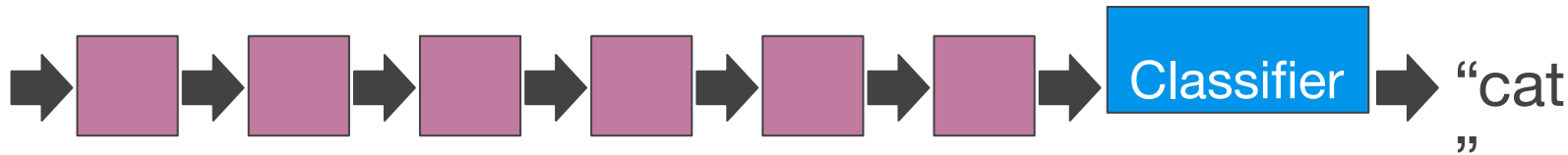
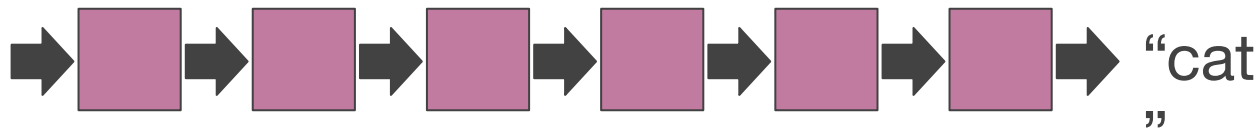


**SVM**



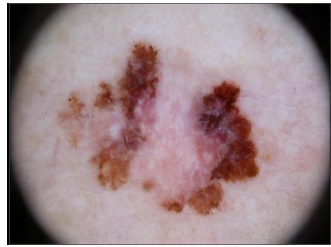
“cat”

# Deep Learning





# Transfer Learning with CNNs



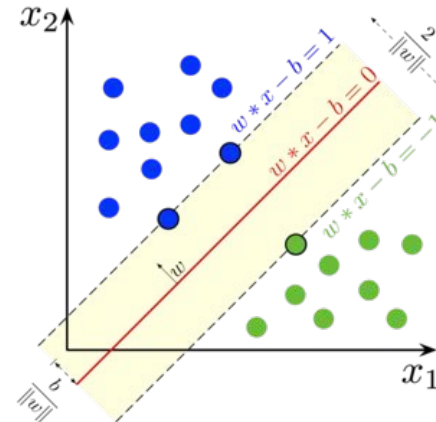
VGG as **Feature Extractor**



$[0.01 \ 0.8 \ 1 \ 0.5 \ \dots \ 0.3 \ 0.07 \ 0 \ 0.4 \ 0.6 \ 0 \ 0]$   
4096-d



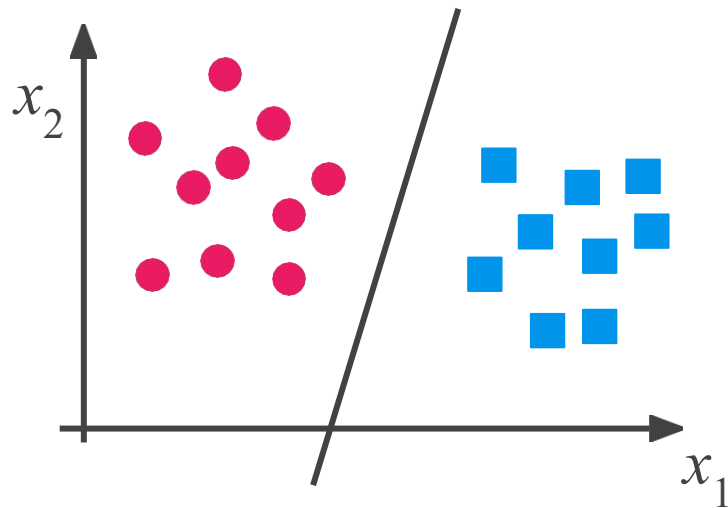
Train a classifier (e.g., SVM)



# Support Vector Machine (SVM)

[Vapnik and Chervonenkis, 1964; Vapnik, 1982; Vapnik, 1995]

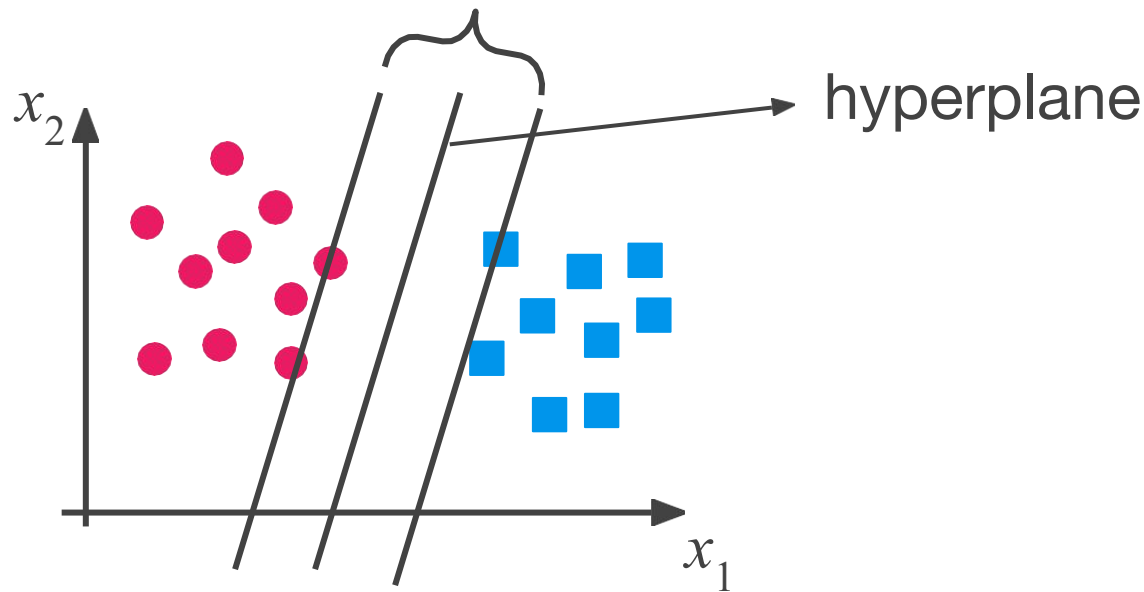
Idea of separating data with a **large** “gap”.



# Support Vector Machine (SVM)

[Vapnik and Chervonenkis, 1964; Vapnik, 1982; Vapnik, 1995]

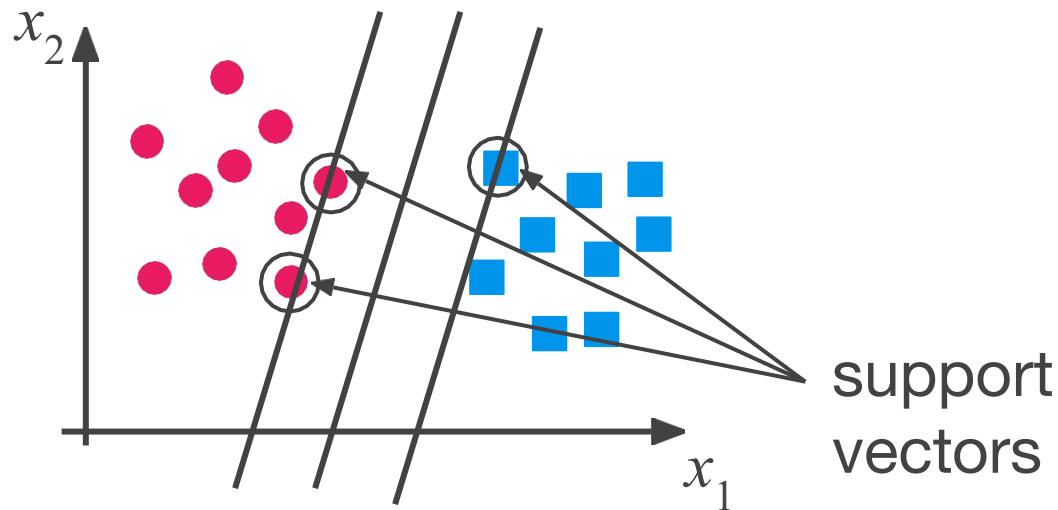
Idea of separating data with a **large** “gap”.



# Support Vector Machine (SVM)

[Vapnik and Chervonenkis, 1964; Vapnik, 1982; Vapnik, 1995]

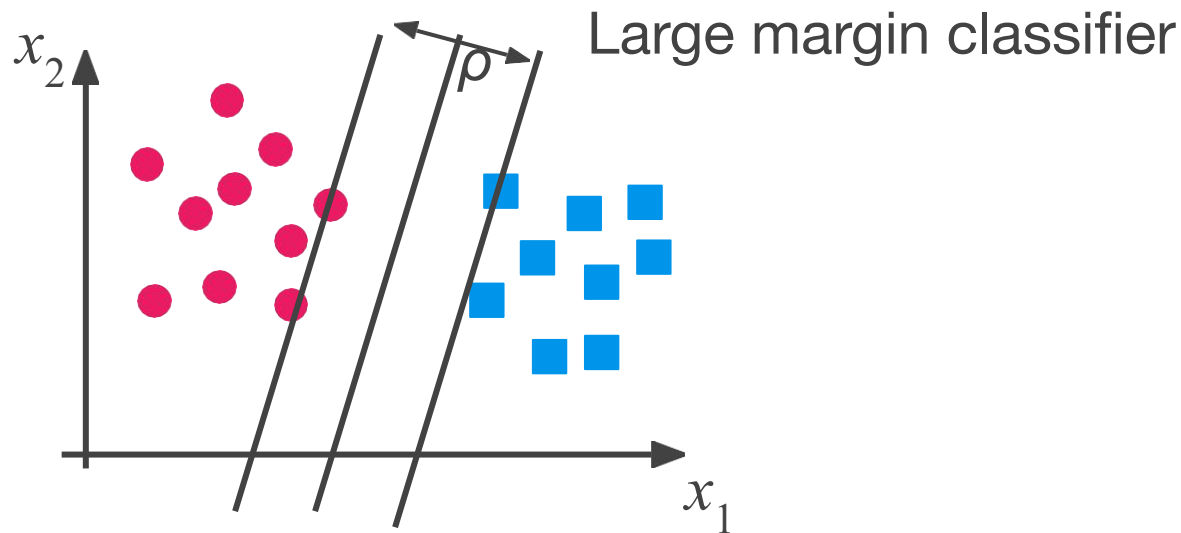
Examples closest to the hyperplane are support vectors.



# Support Vector Machine (SVM)

[Vapnik and Chervonenkis, 1964; Vapnik, 1982; Vapnik, 1995]

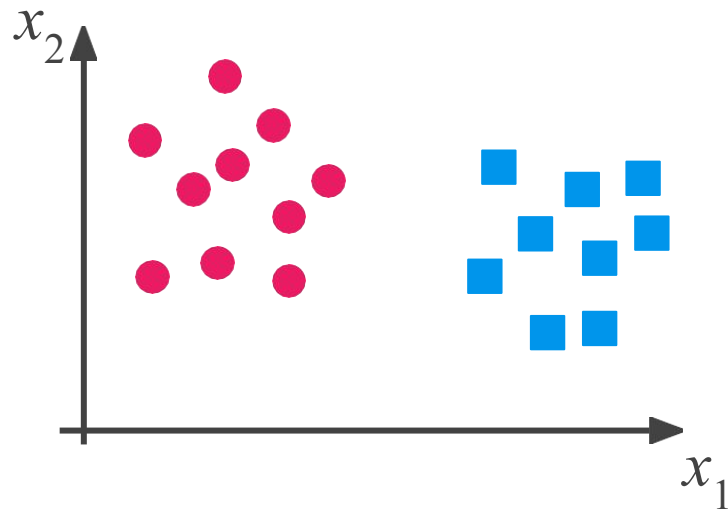
Margin  $\rho$  of the separator is the distance between support vectors.



# How does SVM work?

How can we identify the right hyperplane?

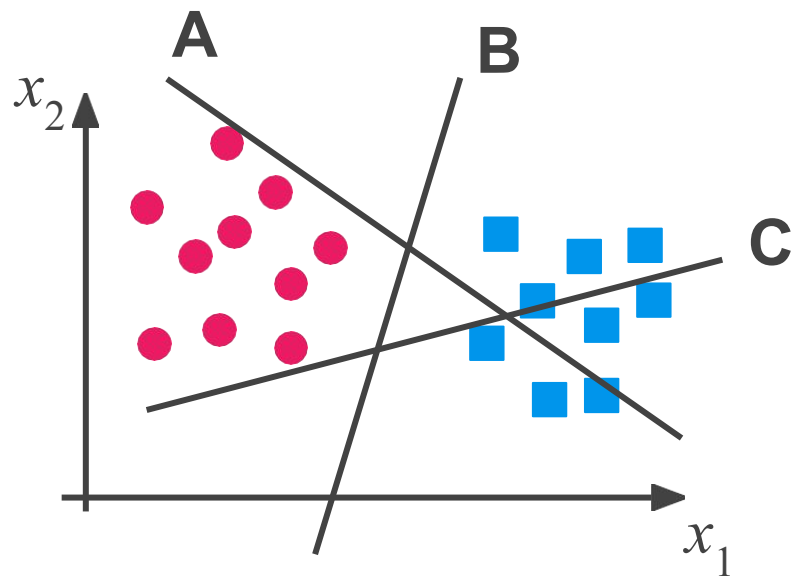
**Scenario 1**



# How can we identify the right hyperplane?

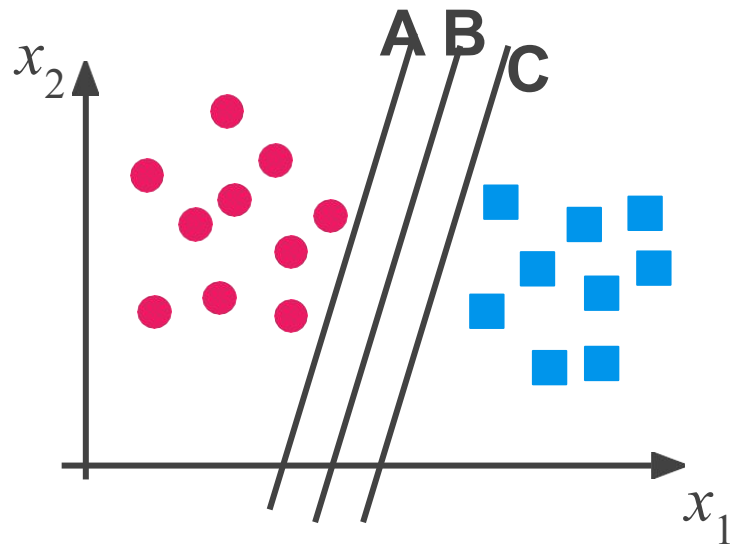
## Scenario

1



# How can we identify the right hyperplane?

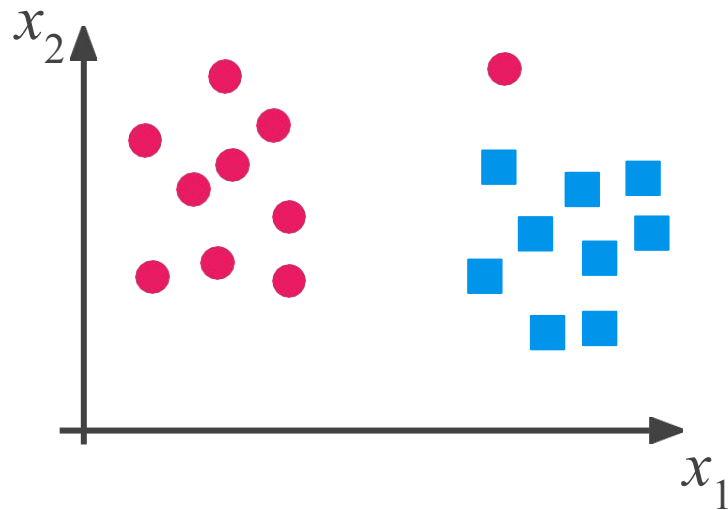
**Scenario  
2**





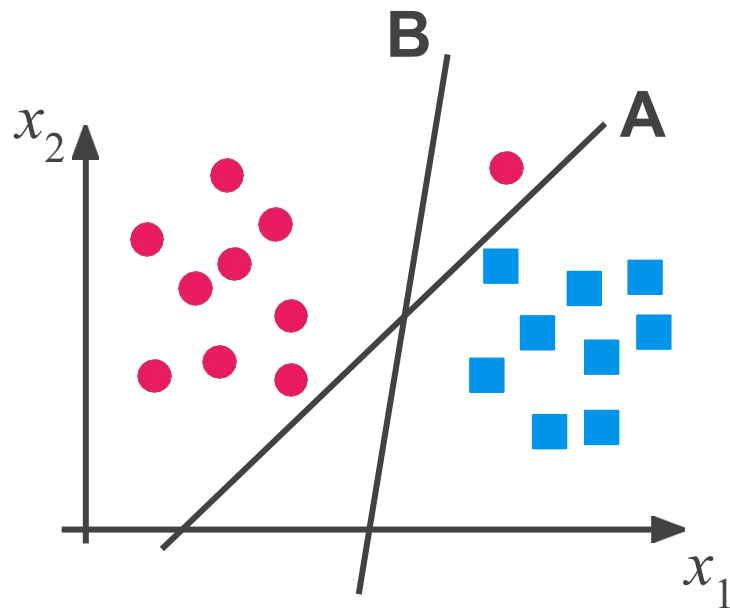
# How can we identify the right hyperplane?

## Scenario 3



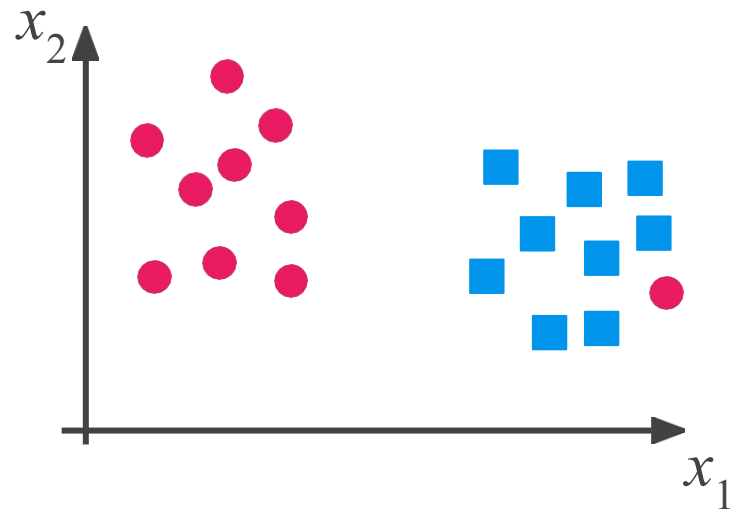
# How can we identify the right hyperplane?

## Scenario 3



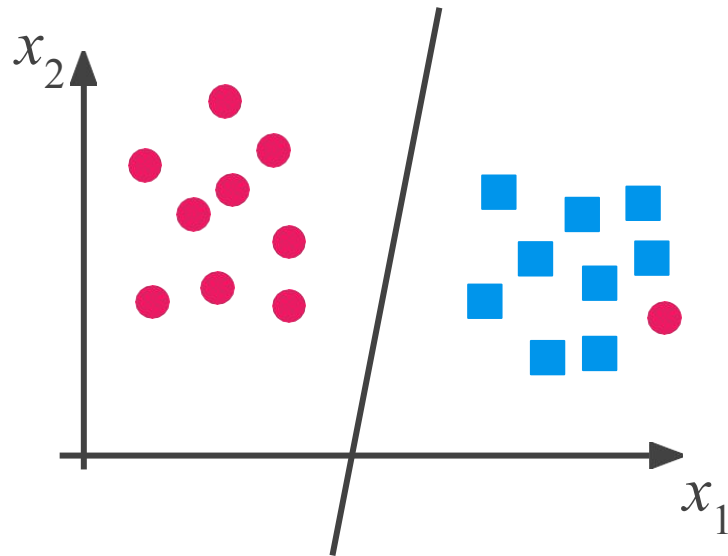
# How can we identify the right hyperplane?

## Scenario 4



# How can we identify the right hyperplane?

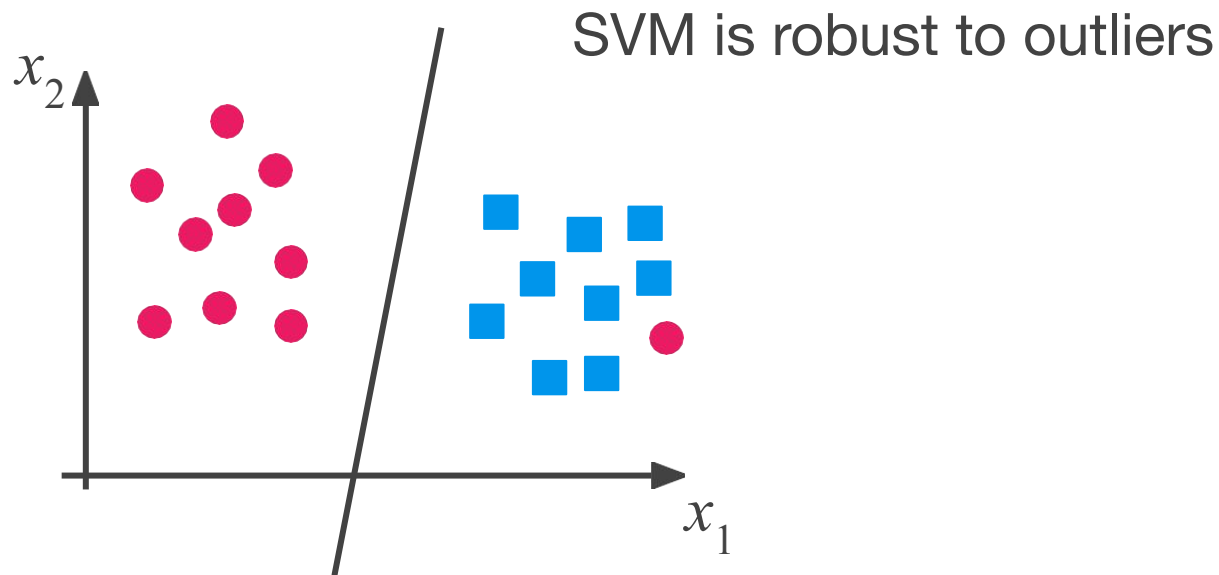
## Scenario 4



## How can we identify the right hyperplane?

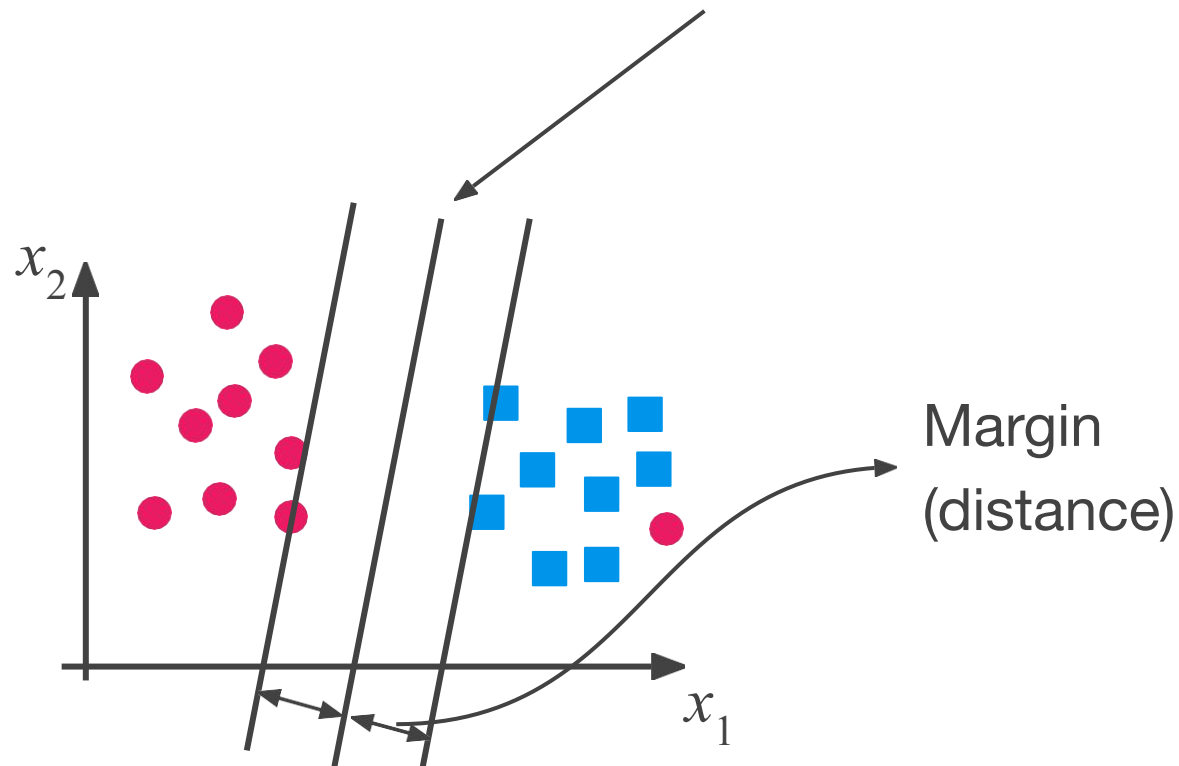
## Scenario

4



# How can we identify the right hyperplane?

## Scenario 4



# SVM: Notation

We will be considering a **linear classifier for a binary classification** problem with labels  $y$  and features  $x$ .

# SVM: Notation

We will be considering a **linear classifier for a binary classification** problem with labels  $y$  and features  $x$ .

- Class labels:  $y \in \{-1, 1\}$  (instead of  $\{0, 1\}$ )
- Parameters:  $w, b$  (instead of vector  $\theta$ )



# SVM: Notation

We will be considering a **linear classifier for a binary classification** problem with labels  $y$  and features  $x$ .

- Class labels:  $y \in \{-1, 1\}$  (instead of  $\{0, 1\}$ )
- Parameters:  $w, b$  (instead of vector  $\theta$ )
- Classifier:  $h_{w,b}(x) = g(w^T x + b)$ 
  - $g(z) = 1$  if  $z \geq 0$ , and  $g(z) = -1$  otherwise

# SVM: The Optimal Hyperplane

Given a training example  $(x^{(i)}, y^{(i)})$ , we define the margin of  $(w, b)$  with respect to the training example:

$$y^{(i)}(w^T x + b) \geq 1, i = \{1, \dots, m\}.$$

# SVM: The Optimal Hyperplane

Let  $P(x^{(1)}, y^{(1)})$  be a point and  $l$  be a line defined by  $ax + by + c = 0$ . The distance  $d$  from  $P$  to  $l$  is defined by:

$$d(l, P) = \frac{|ax^{(1)} + by^{(1)} + c|}{\sqrt{a^2 + b^2}}$$

# SVM: The Optimal Hyperplane

Let  $P(x^{(1)}, y^{(1)})$  be a point and  $l$  be a line defined by  $ax + by + c = 0$ . The distance  $d$  from  $P$  to  $l$  is defined by:

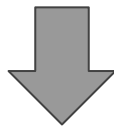
$$d(l, P) = \frac{|ax^{(1)} + by^{(1)} + c|}{\sqrt{a^2 + b^2}}$$



$$d(w, b, x) = \frac{|w^T x + b|}{||w||}$$

# SVM: The Optimal Hyperplane

$$d(w, b, x) = \frac{|w^T x + b|}{||w||}$$



$$\begin{aligned} & \min_{w, b} \frac{1}{2} ||w||^2 \\ \text{s.t. } & y^{(i)}(w^T x + b) \geq 1, i = \{1, \dots, m\} \end{aligned}$$

# SVM: The Optimal Hyperplane

$$d(w, b, x) = \frac{|w^T x + b|}{\|w\|}$$

<http://cs229.stanford.edu/notes2019fall/cs229-notes3.pdf>

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x + b) \geq 1, i = \{1, \dots, m\} \end{aligned}$$

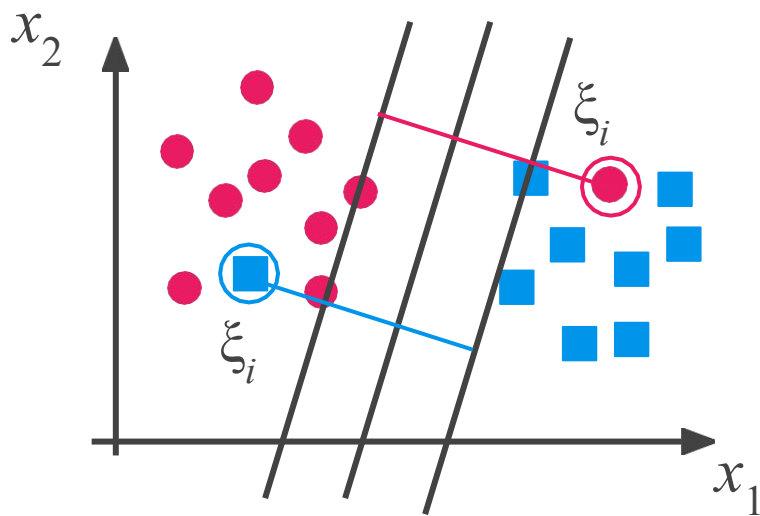
Need to optimize a quadratic function subject to linear

# Soft Margin Classification

What if the training set is not linearly separable?

# Soft Margin Classification

**Slack variables**  $\xi_i$  can be added to allow misclassification of difficult or noisy examples, resulting margin called **soft**.





# Soft Margin Classification

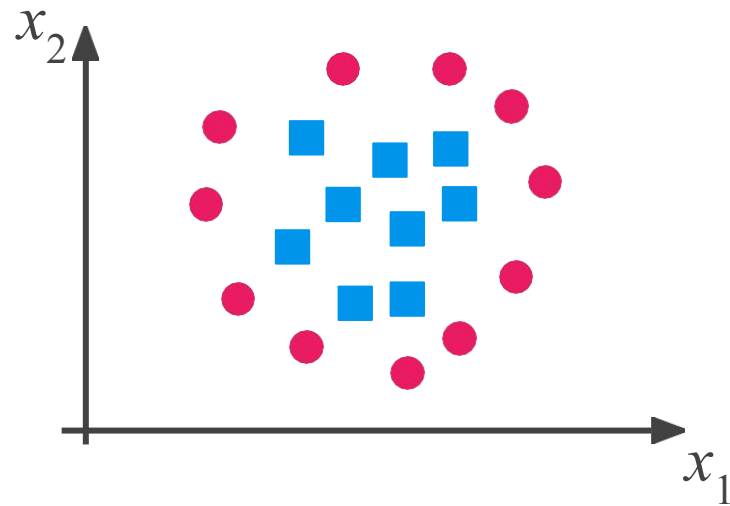
Modified formulation incorporates slack variables:

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum \xi_i \\ \text{s.t.} \quad & y_i(w^T x + b) \geq 1 - \xi_i, \xi_i \geq 0, i = \{1, \dots, m\} \end{aligned}$$

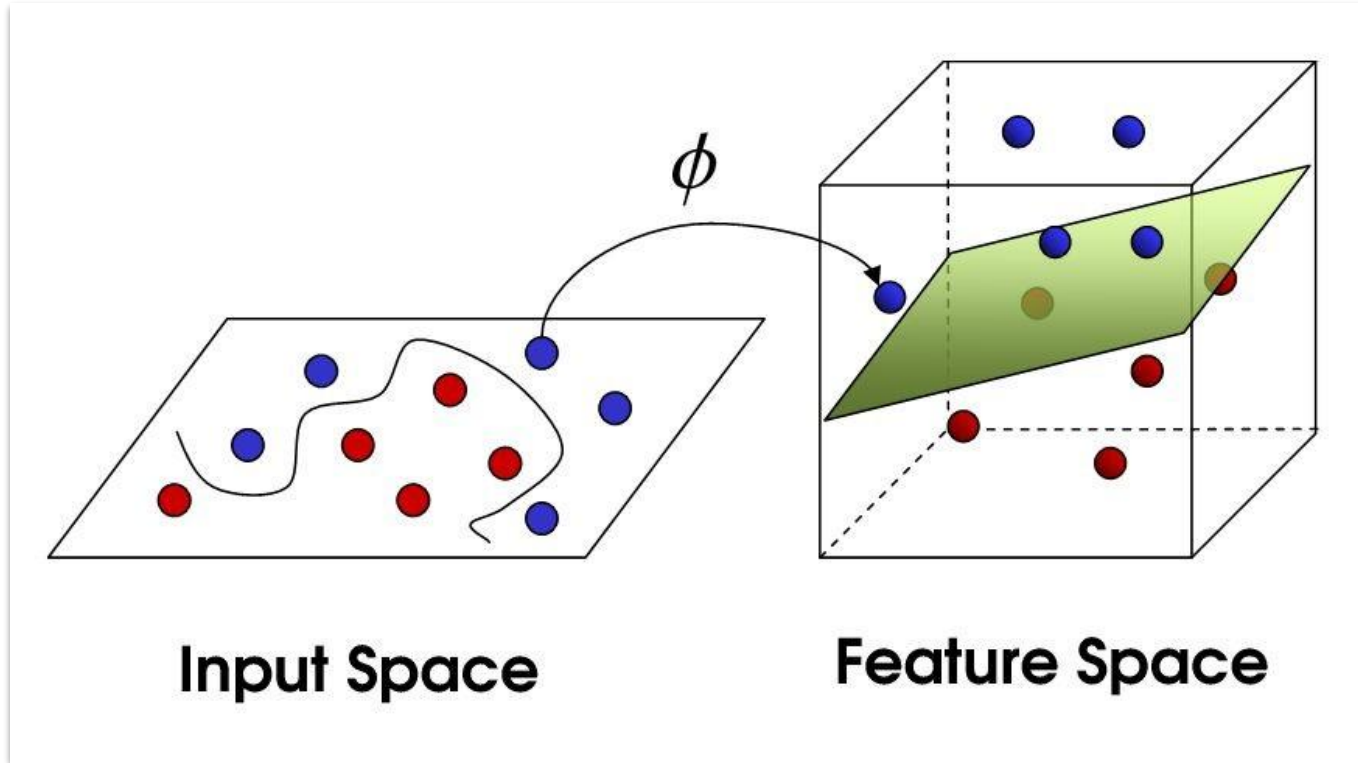
Parameter  $C$  can be viewed as a way to control overfitting: it “trades off” the relative importance of maximizing the margin and fitting the training data.

## How can we identify the right hyperplane?

## Scenario 5



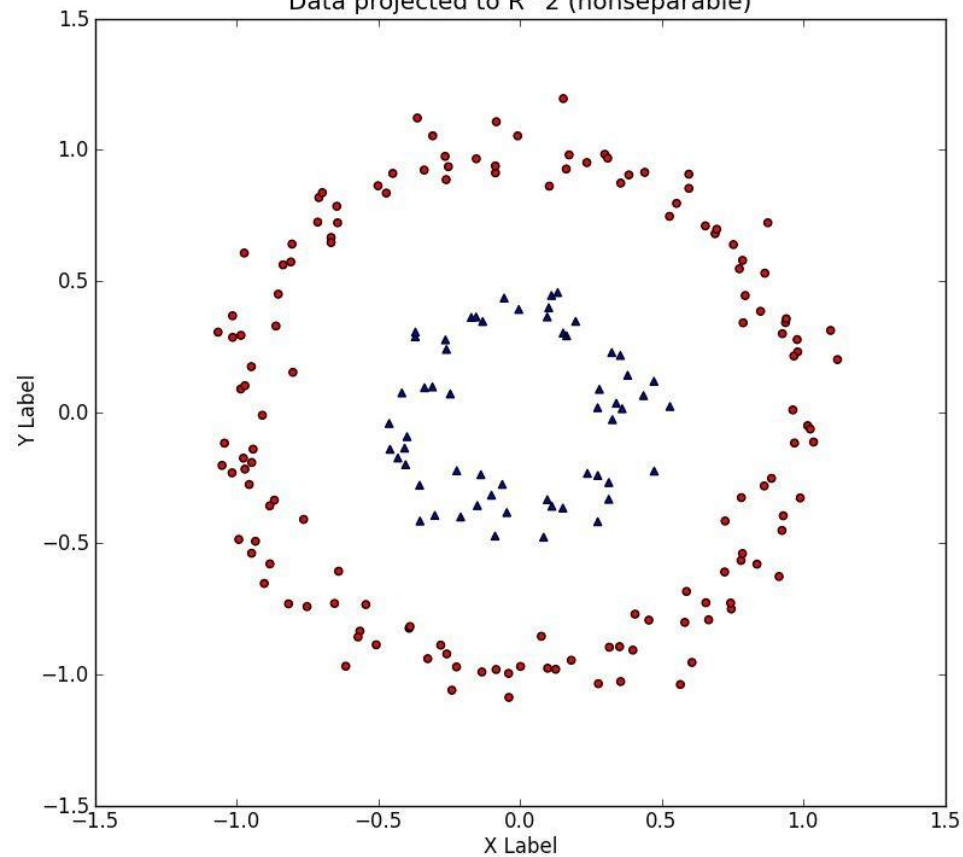
# Kernel Trick



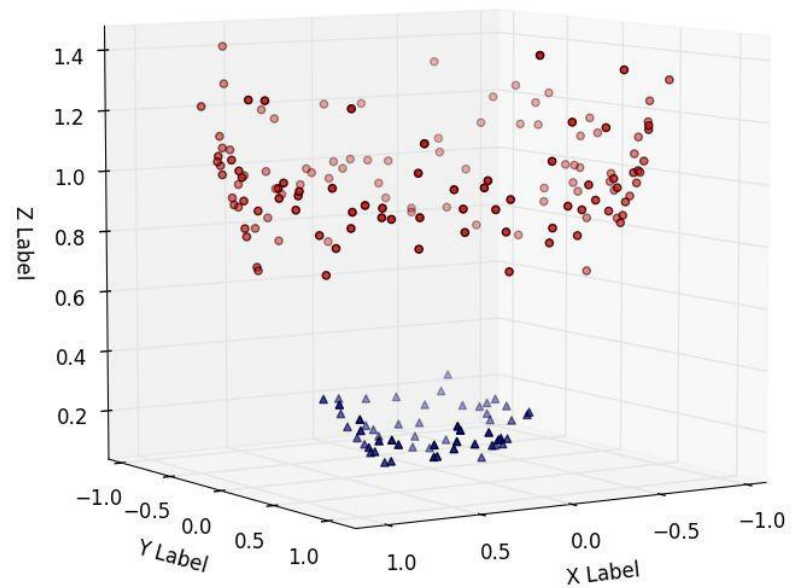
# Kernel Trick

- Linear SVM:  $x_i \cdot x_j$
- Nonlinear SVM:  $K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$ , feature mapping  $\phi$
- Kernel matrix  $K_{ij} = K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j) = \phi(x_j) \cdot \phi(x_i) = K_{ji}$
- Radial Basis Function (RBF) kernel
- Gaussian kernel
- Polynomial kernel
- Chi-square kernel, histogram intersection kernel, string kernel, ....

Data projected to  $R^2$  (nonseparable)



Data in  $R^3$  (separable)



# Important Parameters

Important parameters having higher impact on model performance, “kernel”, “gamma” and “C”.

# Important Parameters

Important parameters having higher impact on model performance, “kernel”, “gamma” and “C”.

C: Penalty parameter C of the error term. It also controls the trade off between smooth decision boundary and classifying the training points correctly.

# Important Parameters

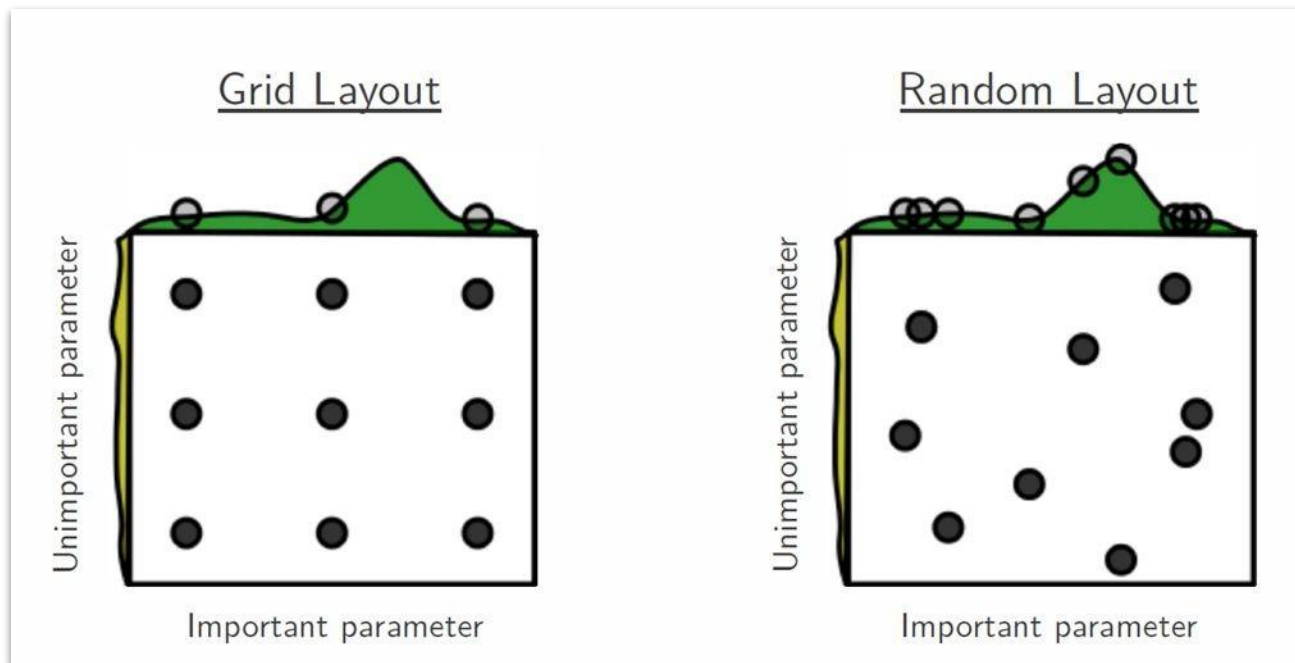
Important parameters having higher impact on model performance, “kernel”, “gamma” and “C”.

C: Penalty parameter C of the error term. It also controls the trade off between smooth decision boundary and classifying the training points correctly.

The parameters can be tuned using grid-search. 



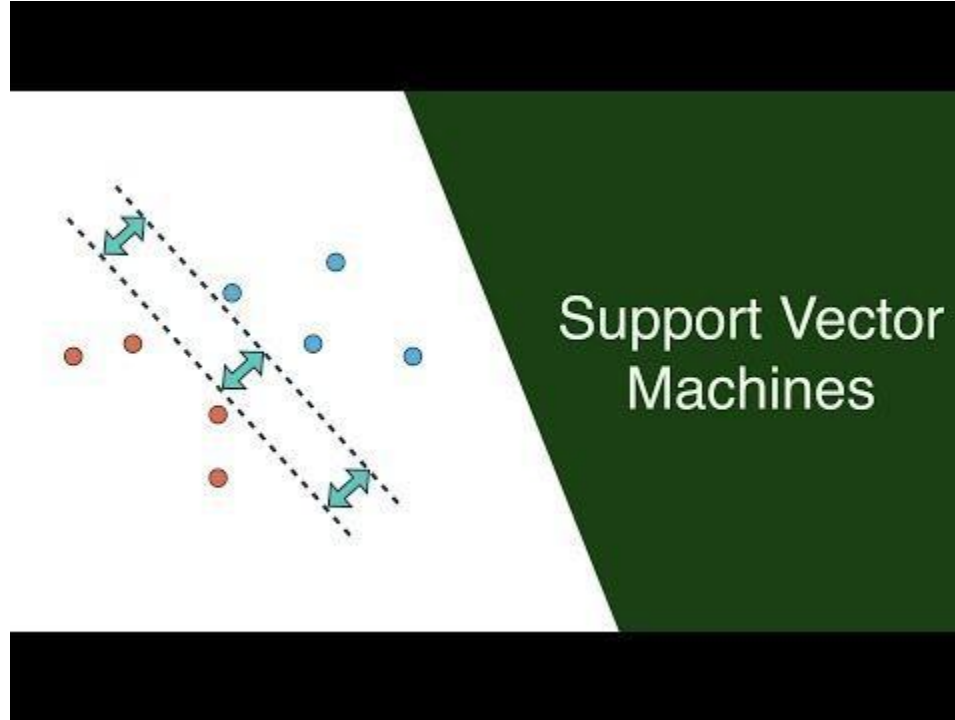
# Grid Search



# Libraries

- Scikit-learn: <https://scikit-learn.org/stable/modules/svm.html>
- LIBSVM: <https://www.csie.ntu.edu.tw/~cjlin/libsvm>
- LIBLINEAR: <https://www.csie.ntu.edu.tw/~cjlin/liblinear>
- PmSVM:  
<https://sites.google.com/site/wujx2001/home/power-mean-svm>

# Support Vector Machines (SVMs): A friendly introduction



[https://www.youtube.com/watch?v=Lpr\\_X8zuE8](https://www.youtube.com/watch?v=Lpr_X8zuE8)

# References

## Machine Learning Books

- Hands-On Machine Learning with Scikit-Learn and TensorFlow
- Pattern Recognition and Machine Learning, Chap. 6 & 7

## Machine Learning Courses

- <https://www.coursera.org/learn/machine-learning>, Week 7
- <http://cs229.stanford.edu/notes2019fall/cs229-notes3.pdf>

# Decision Trees

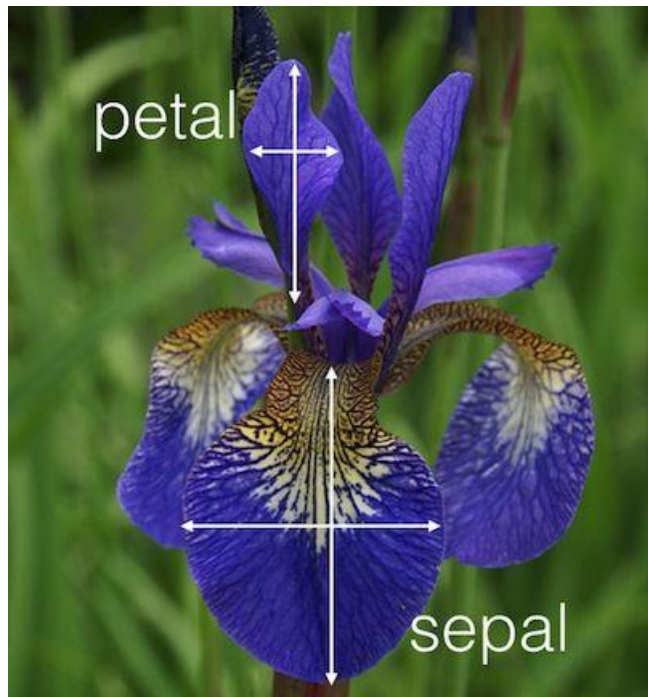
# Decision Tree & Random Forest

- **Decision Trees** are versatile Machine Learning algorithms that can perform both classification and regression tasks, and even multi-output tasks.

# Decision Tree & Random Forest

- Decision Trees are versatile Machine Learning algorithms that can perform both classification and regression tasks, and even multi-output tasks.
- **Random Forest is an ensemble of Decision Trees**, generally trained using the Bagging method (or sometimes Pasting).

# Decision Tree: Iris Dataset



[http://sebastianraschka.com/Articles/2014\\_python\\_lda.html](http://sebastianraschka.com/Articles/2014_python_lda.html)

150 iris flowers from three different species.

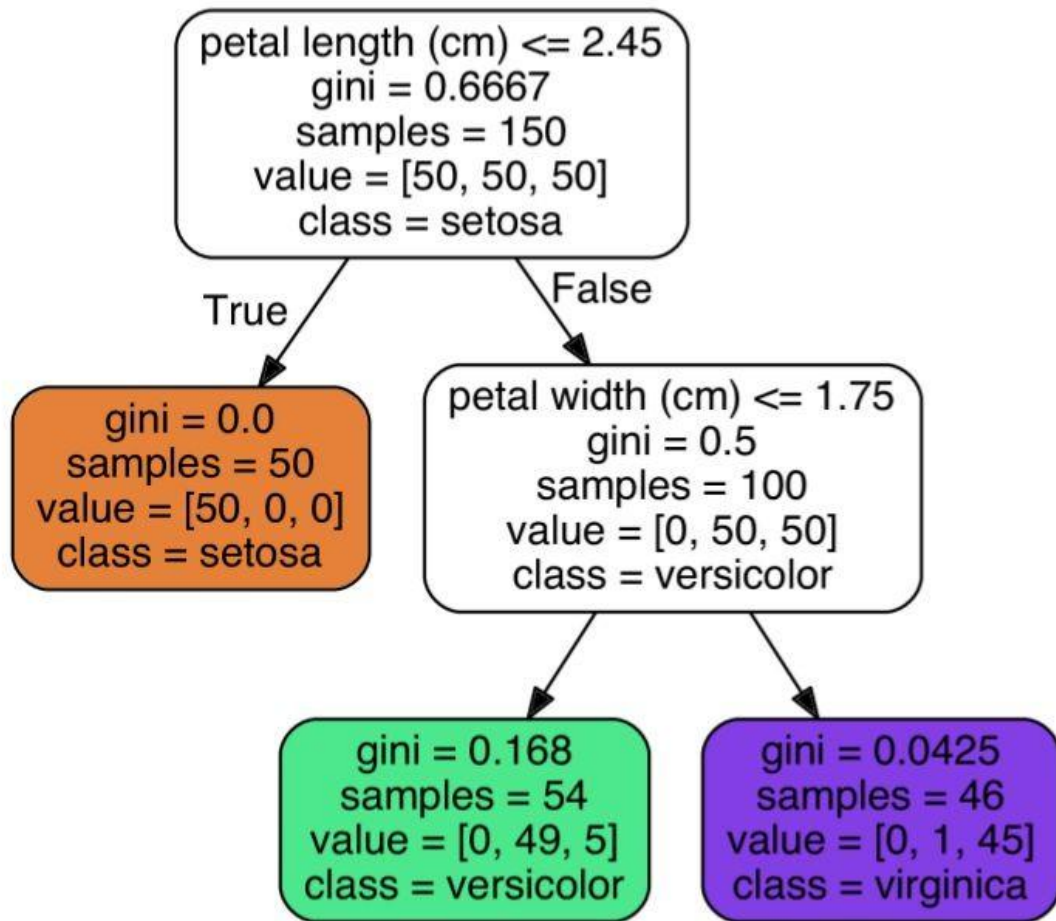
The three classes in the Iris dataset:

1. Iris-setosa ( $n=50$ )
2. Iris-versicolor ( $n=50$ )
3. Iris-virginica ( $n=50$ )

The four features of the Iris dataset:

1. sepal length in cm
2. sepal width in cm
3. petal length in cm
4. petal width in cm





```
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier
```

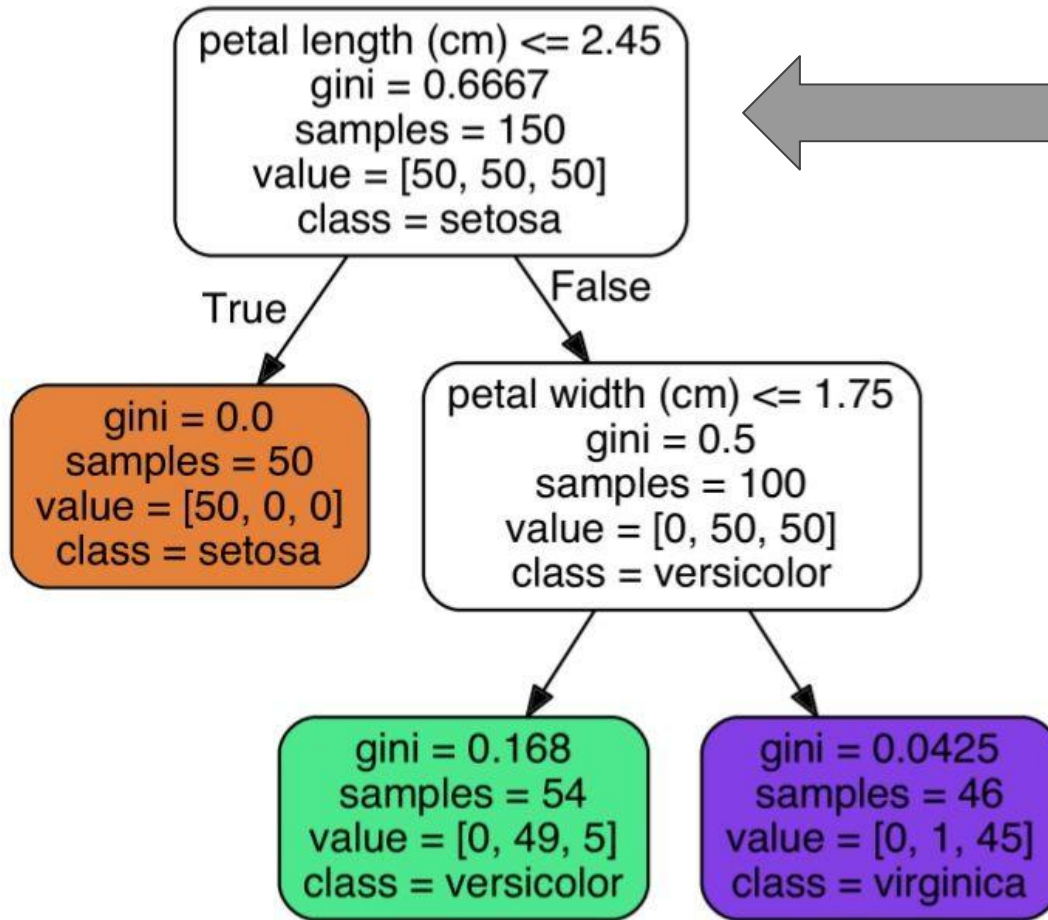
```
iris = load_iris()
X = iris.data[:, 2:] # petal length and width
y = iris.target
tree_clf = DecisionTreeClassifier(max_depth=2)
tree_clf.fit(X, y)
```

#You can visualize the trained Decision Tree by first using the `export_graphviz()` method to output a graph definition file called `iris_tree.dot`:

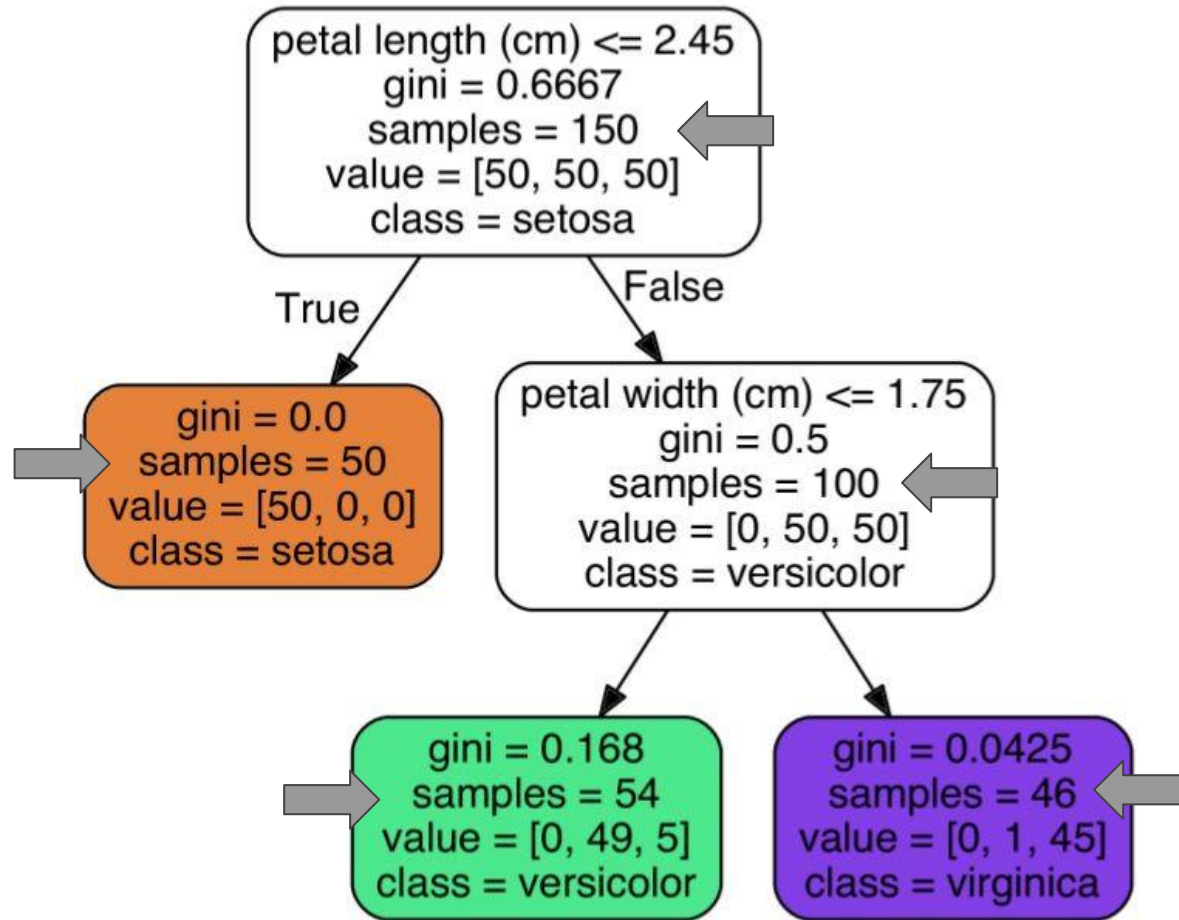
```
from sklearn.tree import export_graphviz
```

```
export_graphviz(
    tree_clf,
    out_file="iris_tree.dot",
    feature_names=iris.feature_names[2:],
    class_names=iris.target_names,
    rounded=True,
    filled=True
)
```

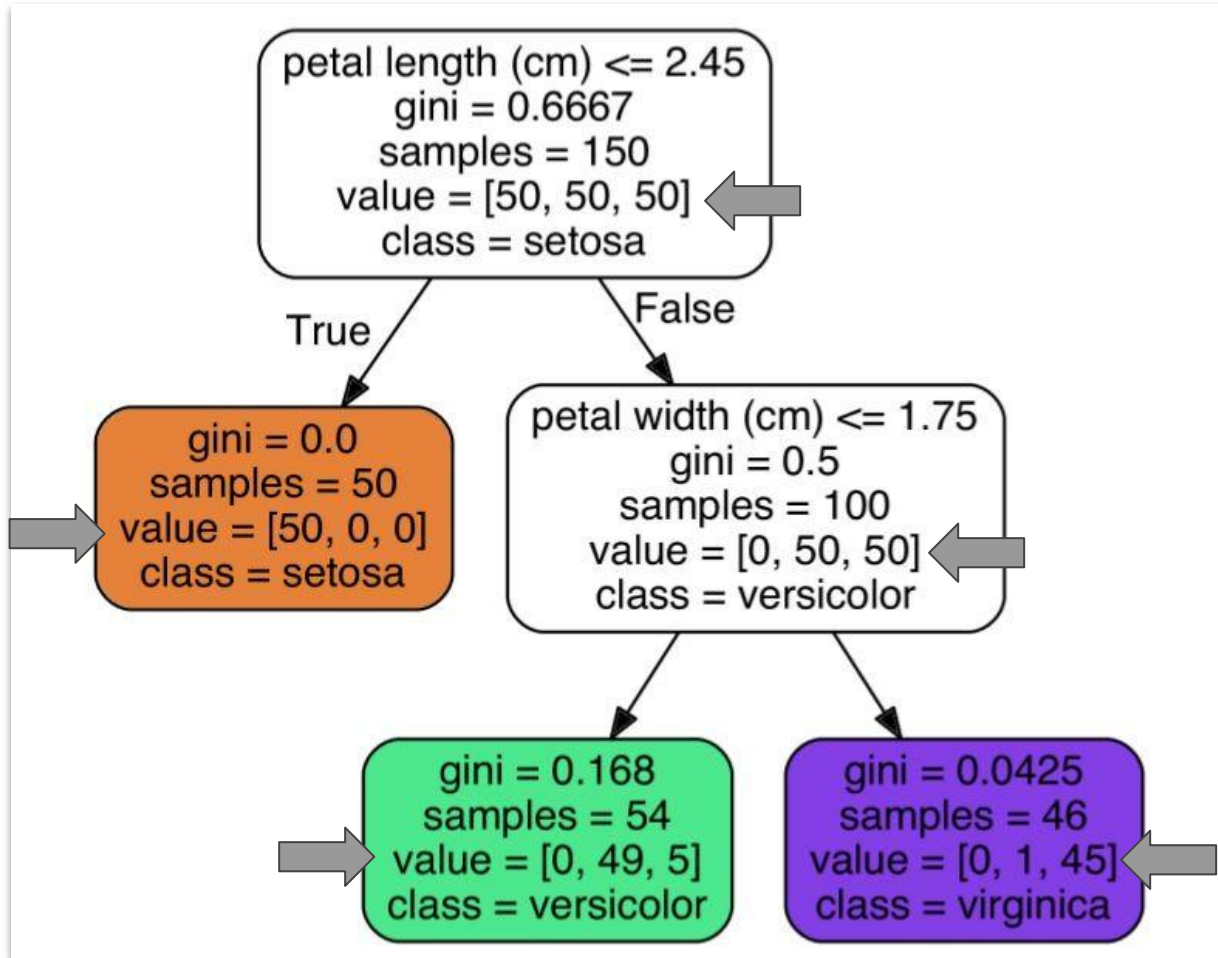
\$ `dot -Tpng iris_tree.dot -o iris_tree.png` #convert this .dot file to a variety of formats such as PDF or PNG using the dot command-line tool from the graphviz package



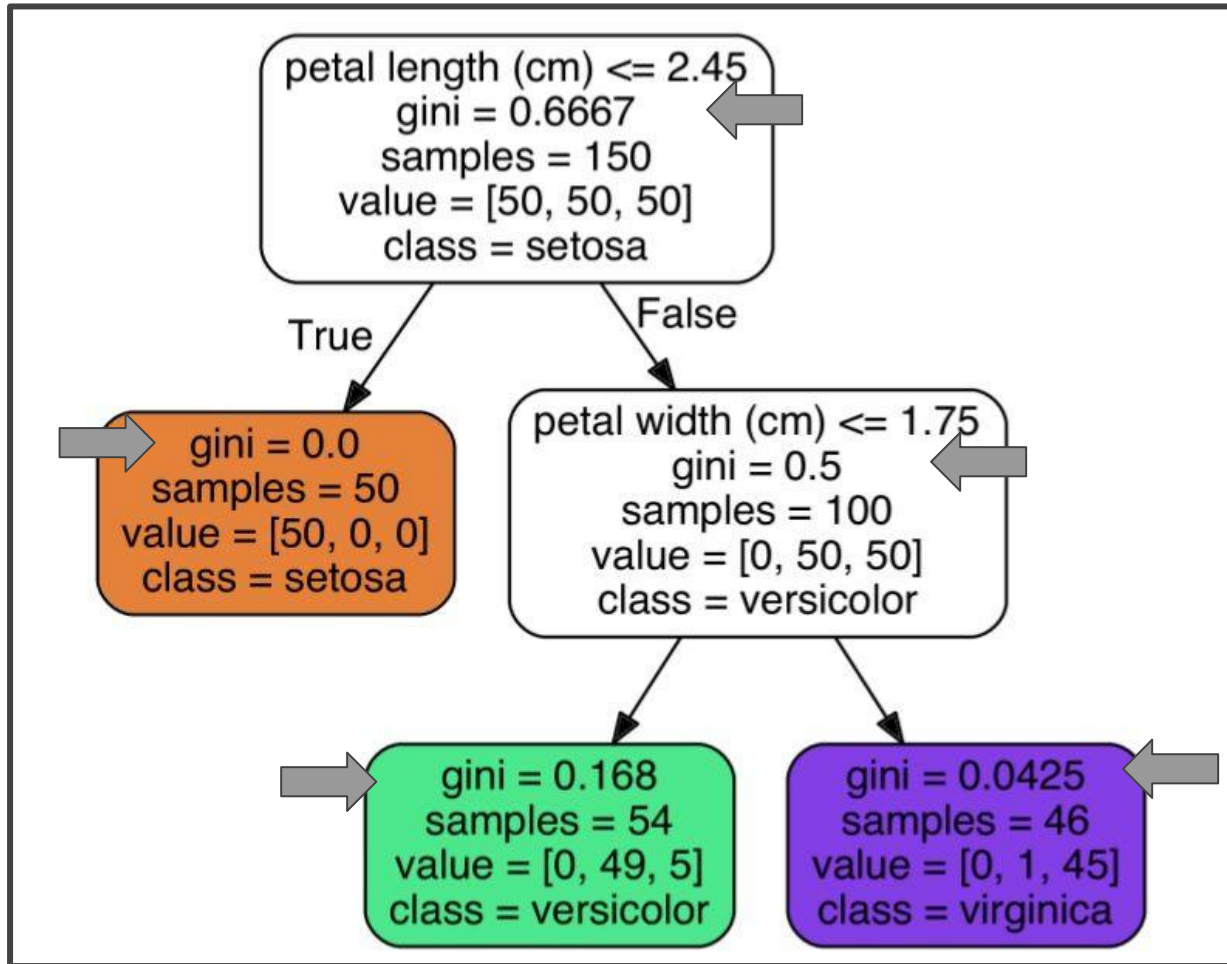
This node asks whether the flower's petal length is smaller than 2.45 cm



A **node's samples** attribute counts how many training instances it applies to.



**A node's value** attribute tells you how many training instances of each class this node applies to.



A **node's gini** attribute measures its impurity.

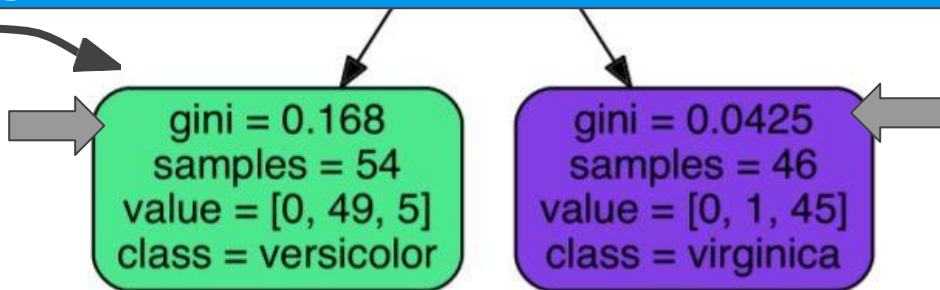
“pure” (gini=0): all training instances belong to the same class.

petal length (cm) <= 2.45

For example, the depth 2 left node has a gini score equal to  $1 - (0/54)^2 - (49/54)^2 - (5/54)^2 \approx 0.168$ .

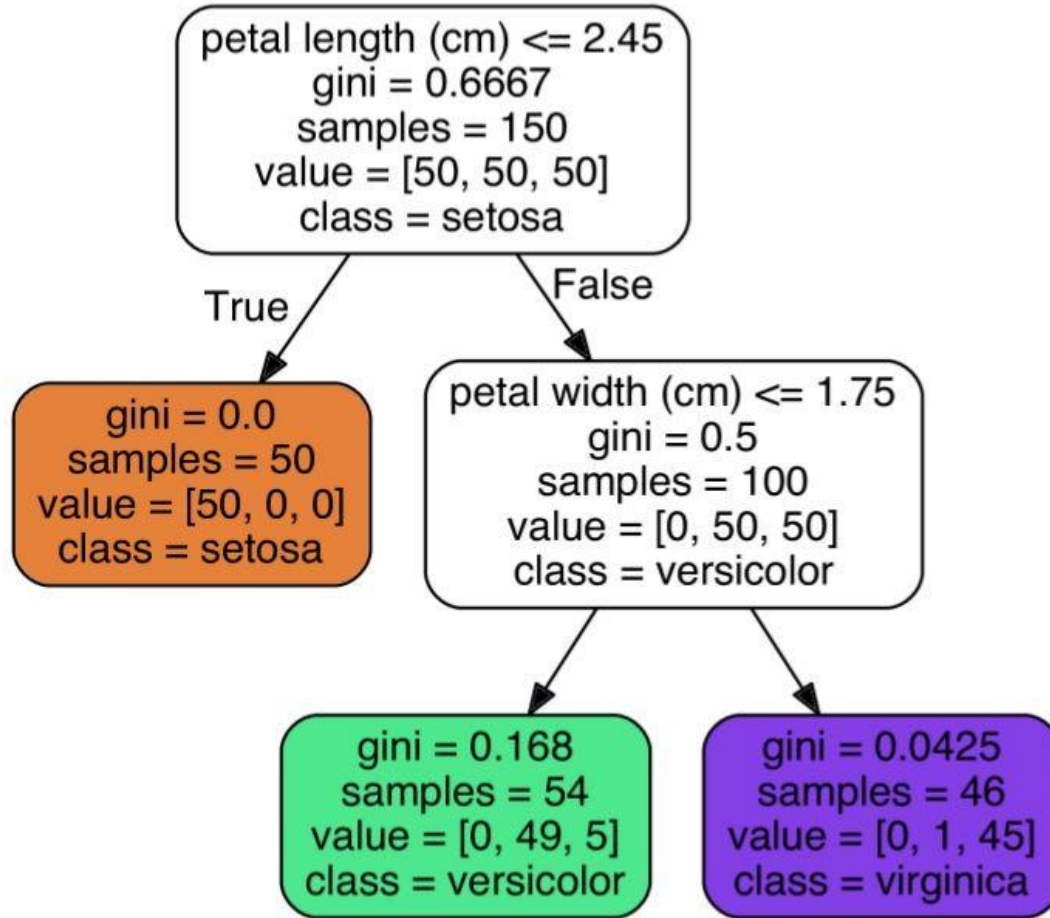
$$G_i = 1 - \sum p_{i,k}^2$$

$p_{i,k}$  is the ratio of class  $k$  instances among the training instances in the  $i^{th}$  node



belong to the same class.

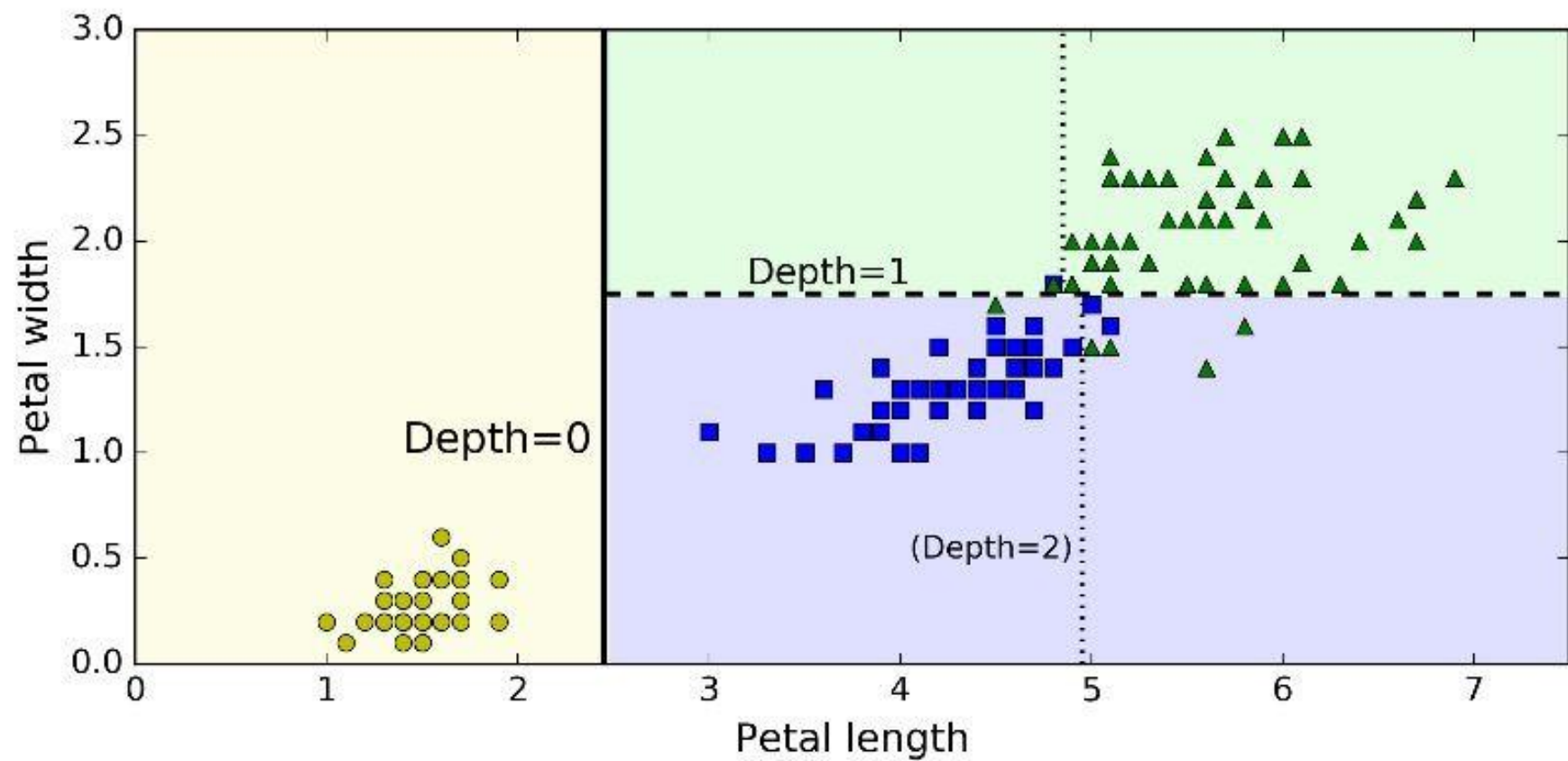




A **node's gini** attribute measures its impurity.

“pure” (gini=0): all training instances belong to the same class.





# Model Interpretation

- White Box vs Black Box
- **Decision Trees** are fairly intuitive and their decisions are easy to interpret (white box models).
- **Random Forests** or **neural networks** are generally considered black box models.

# The CART Algorithm (Scikit-Learn default)

- Classification And Regression Tree (CART) algorithm.

# The CART Algorithm (Scikit-Learn default)

- Classification And Regression Tree (CART) algorithm.
- The idea is really quite simple: the algorithm first splits the training set in two subsets using a single feature  $k$  and a threshold  $t_k$  (e.g. “petal length  $\leq 2.45$  cm”).

# The CART Algorithm (Scikit-Learn default)

- Classification And Regression Tree (CART) algorithm.
- The idea is really quite simple: the algorithm first splits the training set in two subsets using a single feature  $k$  and a threshold  $t_k$  (e.g. “petal length  $\leq 2.45$  cm”).
- How does it choose  $k$  and  $t_k$ ?

# The CART Algorithm (Scikit-Learn default)

- Classification And Regression Tree (CART) algorithm.
- The idea is really quite simple: the algorithm first splits the training set in two subsets using a single feature  $k$  and a threshold  $t_k$  (e.g. “petal length  $\leq 2.45$  cm”).
- How does it choose  $k$  and  $t_k$ ?  
It searches for the pair  $(k, t_k)$  that produces the purest subsets (weighted by their size).

# The CART Algorithm (Scikit-Learn default)

$$J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$

where  $\begin{cases} G_{\text{left/right}} \text{ measures the impurity of the left/right subset,} \\ m_{\text{left/right}} \text{ is the number of instances in the left/right subset.} \end{cases}$

CART cost function for classification

It stops recursing once it reaches the maximum depth (hyperparameter), or if it cannot find a split that will reduce impurity

# The CART Algorithm (Scikit-Learn default)

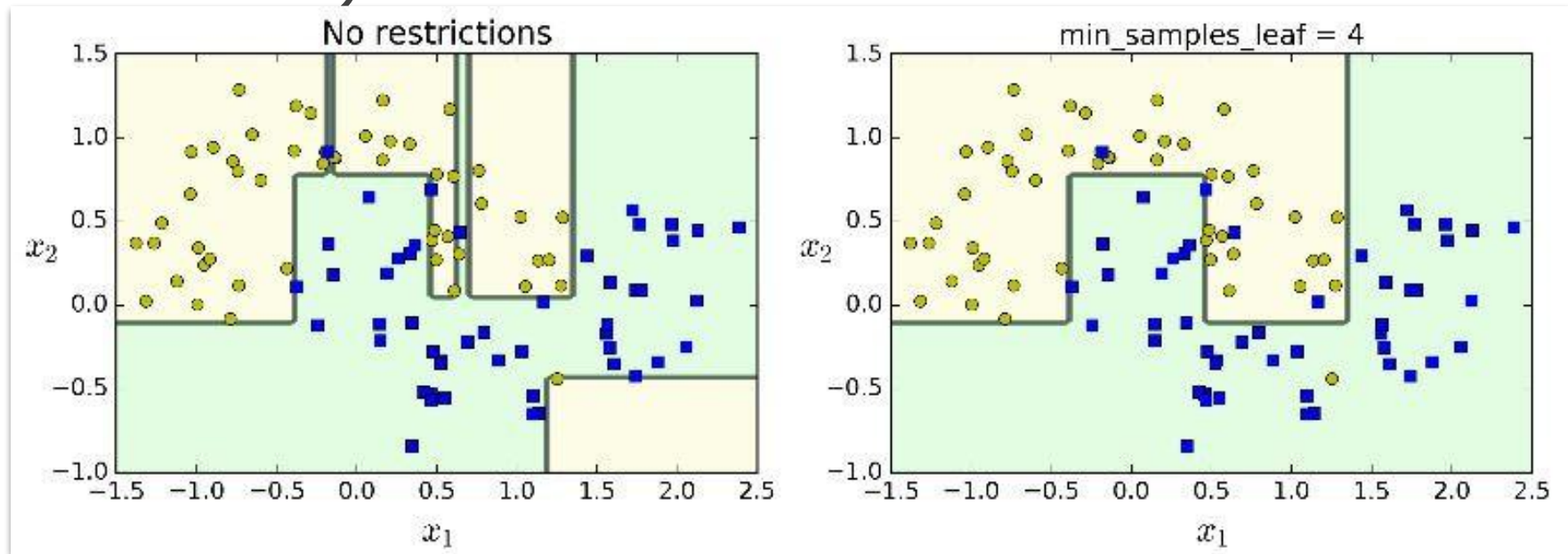
- CART algorithm is a greedy algorithm.
- Training:  $O(n \times m \log(m))$ ,  $m$  samples,  $n$  features.
- Overall prediction:  $O(\log^2(m))$ .



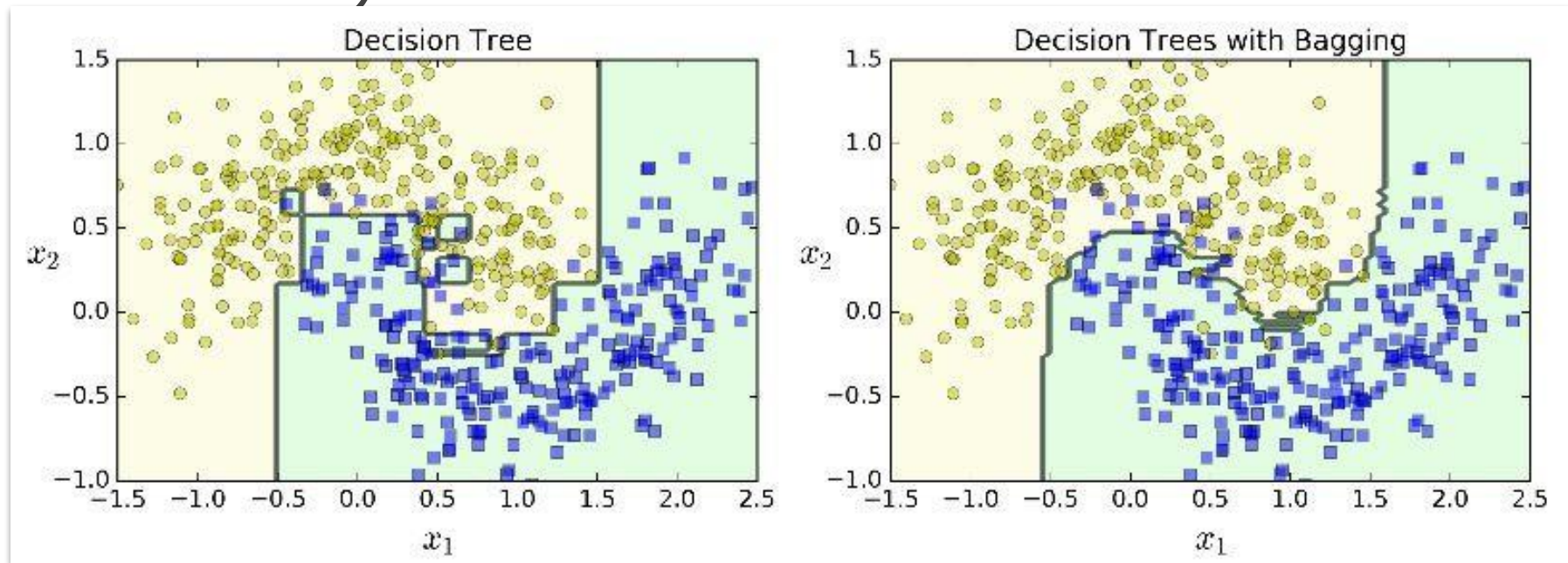
# Regularization (in Scikit-Learn)

- `max_depth` hyperparameter: the default value is `None`, which means unlimited.
- `min_samples_split`: the minimum number of samples a node must have before it can be split.
- `min_samples_leaf`: the minimum number of samples a leaf node must have.

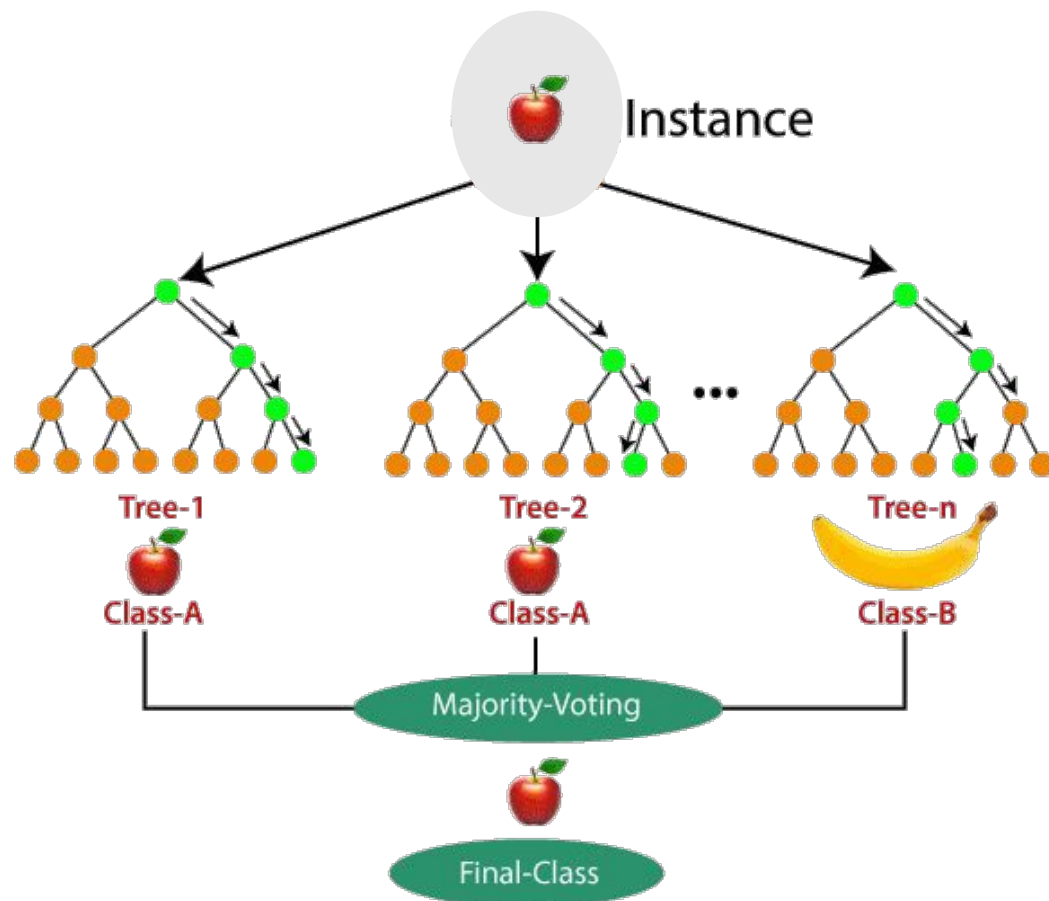
# Regularization (in Scikit-Learn)



# Regularization (in Scikit-Learn)



# Random Forest



# Random Forest [Ho, 1995]

- Random Forest is an ensemble of Decision Trees, generally trained using the Bagging method.

# Random Forest [Ho, 1995]

- Random Forest is an ensemble of Decision Trees, generally trained using the Bagging method.
- **Extra randomness when growing trees:**
  - Instead of searching for the very best feature when splitting a node, it searches for the best feature among a **random subset of features**.

“Random Decision Forests”, Tin Kam Ho (1995): <http://goo.gl/zVOGQ1>

# Random Forest [Ho, 1995]

1. Assume number of cases in the training set is  $N$ . Then, sample of these  $N$  cases is taken at random but with replacement.



# Random Forest [Ho, 1995]

2. If there are  $M$  input variables, a number  $m < M$  is specified such that at each node,  $m$  variables are selected at random out of the  $M$ .

The best split on these  $m$  is used to split the node. The value of  $m$  is held constant while we grow the forest.

# Random Forest [Ho, 1995]

3. Each tree is grown to the largest extent possible and there is no pruning.
4. Predict new data by aggregating the predictions of the ntree trees (i.e., majority votes for classification, average for regression).

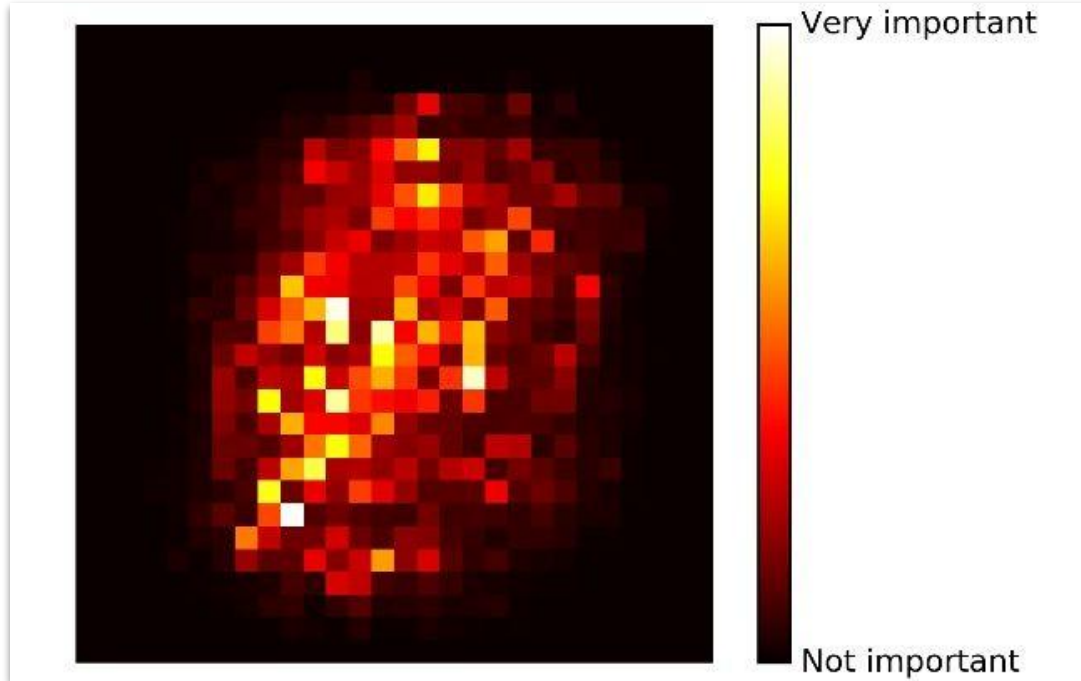
# Random Forest: Feature Importance

```
from sklearn.datasets import load_iris
from sklearn.ensemble import RandomForestClassifier
iris = load_iris()
rnd_clf = RandomForestClassifier(n_estimators= 500, n_jobs= -1)
rnd_clf.fit(iris[ "data"], iris["target"])

for name, score in zip(iris["feature_names" ],
rnd_clf.feature_importances_):
    print(name, score)
```

```
sepal          length          (cm)
0.112492250999 sepal width (cm)
0.0231192882825 petal   length
(cm) 0.441030464364 petal width
(cm) 0.423357996355
```

# Random Forest: Feature Importance



MNIST pixel importance  
(according to a Random  
Forest classifier)

# Forward Thinking: Building Deep Random Forests

Kevin Miller, Chris Hettinger, Jeffrey Humpherys, Tyler Jarvis, and David K

Department of Mathematics  
Brigham Young University  
Provo, Utah 84602

millerk5@byu.edu, hettinger@math.byu.edu, jeffh@math.byu.edu  
jarvis@math.byu.edu, david.kartchner@math.byu.edu

## Abstract

# Training Big Random Forests with Little Resources

Fabian Gieseke

Department of Computer Science  
University of Copenhagen  
Copenhagen, Denmark  
fabian.gieseke@di.ku.dk

## ABSTRACT

Without access to large compute clusters, building random forests on large datasets is still a challenging problem. This is, in particular, the case if fully-grown trees are desired. We propose a simple yet effective framework that allows to efficiently construct ensembles of huge trees for hundreds of millions or even billions of training instances using a cheap desktop computer with commodity hardware. The basic idea is to consider a multi-level construction scheme, which builds top trees for small random subsets of the available data and which subsequently distributes all training instances to the top trees' leaves for further processing. While being conceptually simple, the overall efficiency crucially depends on the particular implementation of the different phases. The practical merits of our

ensembles in a parallel or distributed manner (e.g., by node). While this can significantly reduce the training time in such environments, further, the efficiency of these frameworks might cause problems in case the data does not fit into the main memory.

In this work, we propose a scheme for building random forests on large datasets. The main idea is to build the forests in three phases: Starting with a top tree for the whole data, one subsequently distributes the data to the leaves of that tree. For each leaf

Christian Fasel  
Department of Computer Science  
University of Copenhagen  
Copenhagen, Denmark  
fasel@di.ku.dk

# Distributed Deep Forest and its Application to Automatic Detection of Cash-out Fraud

Ya-Lin Zhang<sup>†</sup>, Jun Zhou<sup>‡</sup>, Wenhao Zheng<sup>†</sup>, Ji Feng<sup>†</sup>, Longfei Li<sup>†</sup>, Ziqi Liu<sup>†</sup>, Ming Li<sup>†</sup>, Zhiqiang Zhang<sup>†</sup>, Chaochao Chen<sup>‡</sup>, Xiaolong Li<sup>‡</sup>, Zhi-Hua Zhou<sup>†</sup>

<sup>†</sup>National Key Lab for Novel Software Technology, Nanjing University, China

<sup>†</sup>{zhangyl, zhengwh, fengj, lim, zhoush}@lamda.nju.edu.cn

<sup>‡</sup>Ant Financial Services Group, China

<sup>‡</sup>{jun.zhoujun, longyao.llf, ziqiliu, lingyao.zzq, chaochao.ccc, xl.li}@antfin.com

# Deep Forest: Towards an Alternative to Deep Neural Networks\*

Zhi-Hua Zhou and Ji Feng

National Key Lab for Novel Software Technology, Nanjing University, Nanjing 210023, China  
{zhouzh, fengj}@lamda.nju.edu.cn

## Abstract

In this paper, we propose gcForest, a decision tree

ensemble, even when several authors all use convolutional neural networks [LeCun *et al.*, 1998; Krizhevsky *et al.*, 2012; Simonyan and Zisserman, 2014], they are actually using dif-

# References

## Machine Learning Books

- Hands-On Machine Learning with Scikit-Learn and TensorFlow, Chap. 6 & 7
- Pattern Recognition and Machine Learning, Chap. 14
- Pattern Classification, Chap 8 & 9 (Sec. 9.5)
- <https://towardsdatascience.com/random-forest-in-python-24d0893d51c0>