

# Ensemble Methods in Machine Learning

LO: Be able to

Describe the difference between strong and weak learnings.

Describe the techniques of bootstrapping, gradient boosting, and Ada boosting.

Implement the random forest machine learning model.

# Boosting

- **Boosting** is
  - a machine learning ensemble meta-algorithm for reducing bias and variance in supervised learning,.
  - a family of machine learning algorithms that convert **weak learners** to **strong learners**.

# Boosting methods

- Most boosting methods consist of iteratively adding the results of a weak learner to form a strong learner.
  - When adding the results of the weak learners, they are typically weighted by a factor related to the weak learner's accuracy.
  - Each time the results are added, the sample data are reweighted before the next weak learner is applied to them.

STOP

Algorithm

Gradient Boosting

**AdaBoost**

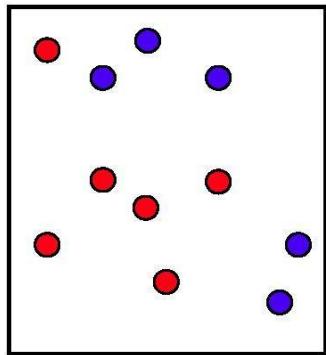
Tree Bagging and Random Forest

# AdaBoost (short for Adaptive Boosting)

- Formulated by Yoav Freund and Robert Schapire
- Awarded the Gödel Prize in 2003 for this work.

# Initial data

$m=0$



$m=1$

$m=2$

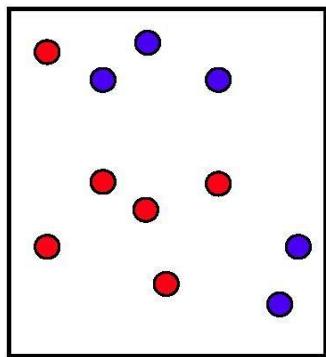
$m=3$

$k_m$

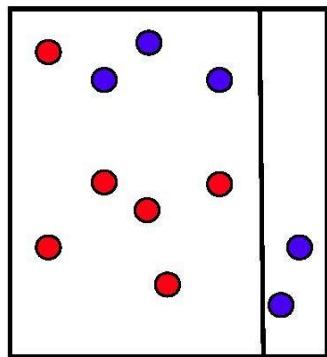
$w_{m+1}$

$k_1$

$m=0$



$m=1$



$m=2$

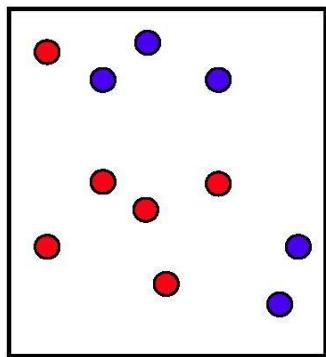
$m=3$

$k_m$

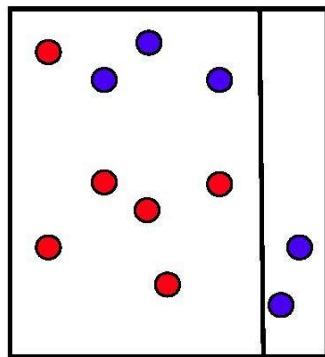
$w_{m+1}$

$w_2$

$m=0$



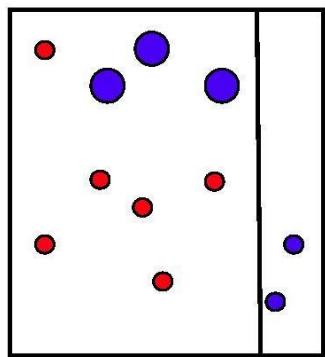
$m=1$



$m=2$

$m=3$

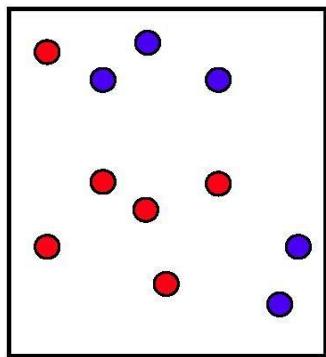
$k_m$



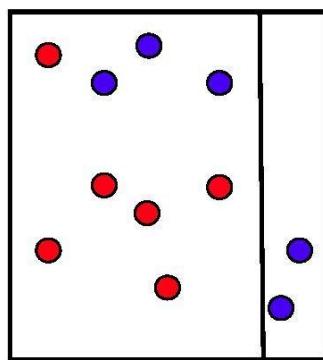
$w_{m+1}$

$k_2$

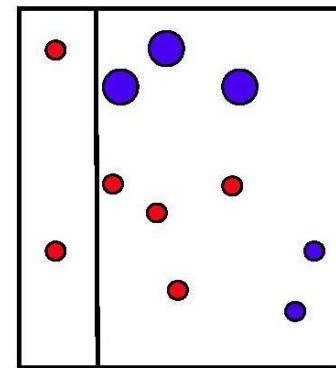
$m=0$



$m=1$

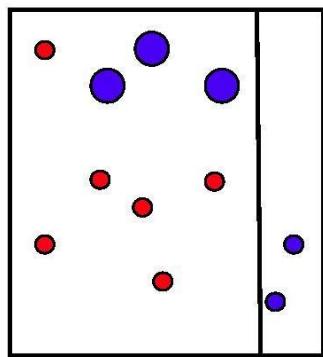


$m=2$



$m=3$

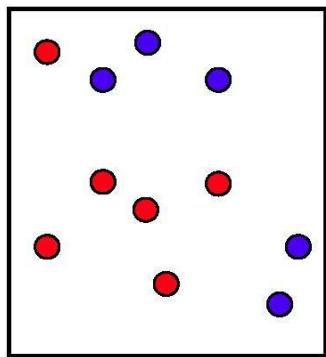
$k_m$



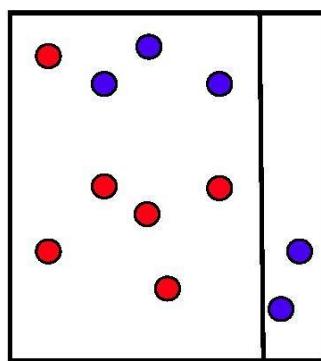
$w_{m+1}$

$w_3$

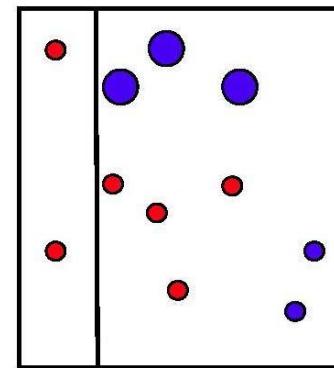
$m=0$



$m=1$

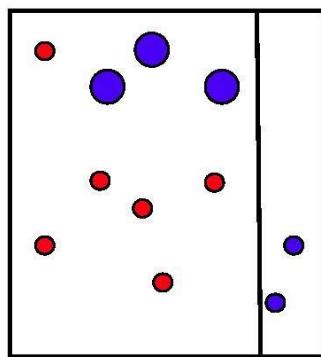


$m=2$



$m=3$

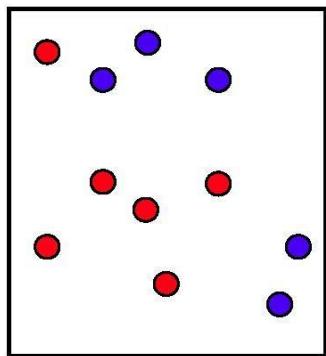
$k_m$



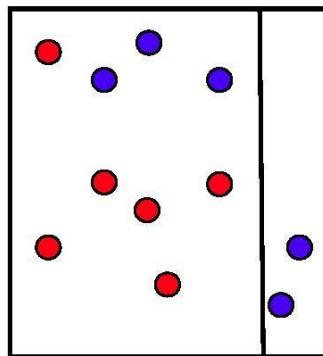
$w_{m+1}$

$k_3$

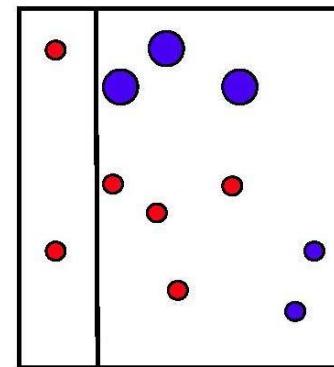
$m=0$



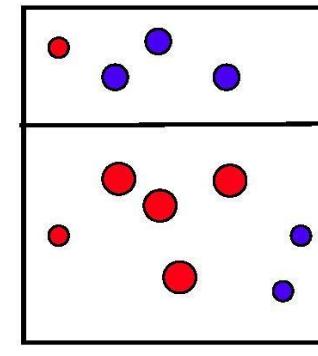
$m=1$



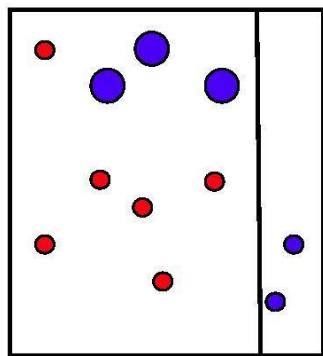
$m=2$



$m=3$



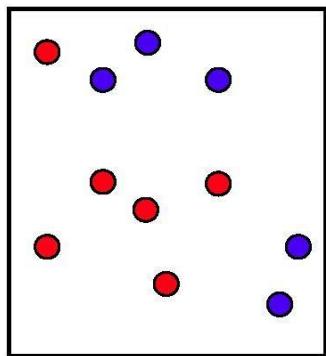
$k_m$



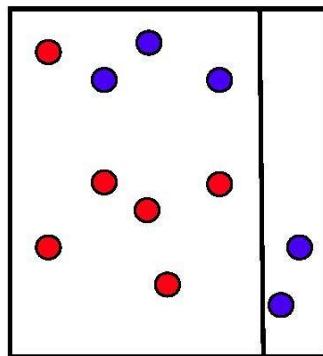
$w_{m+1}$

$w_4$

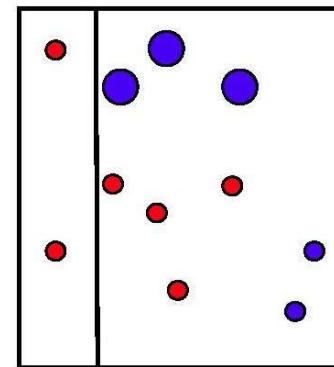
$m=0$



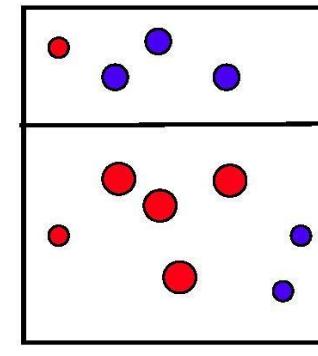
$m=1$



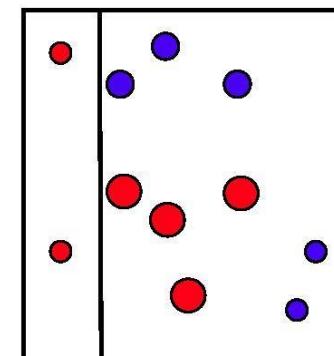
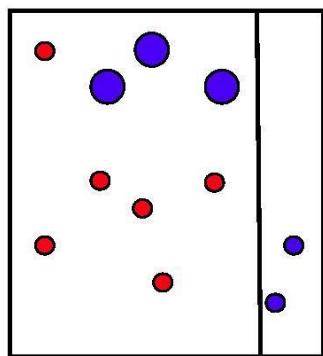
$m=2$



$m=3$



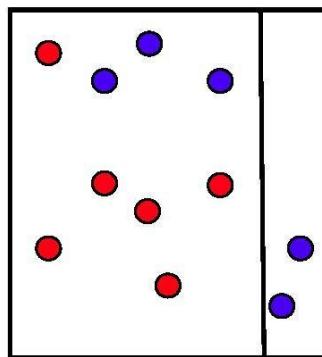
$k_m$



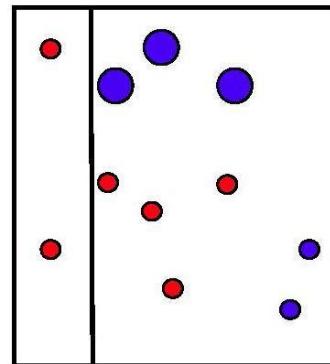
$w_{m+1}$

# Final classifier components

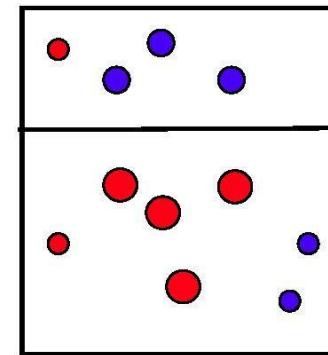
$m=1$



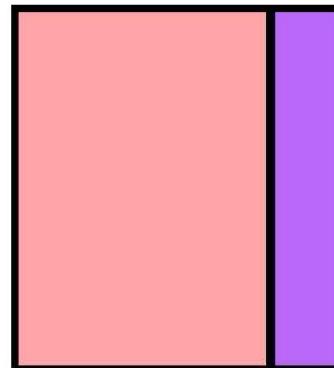
$m=2$



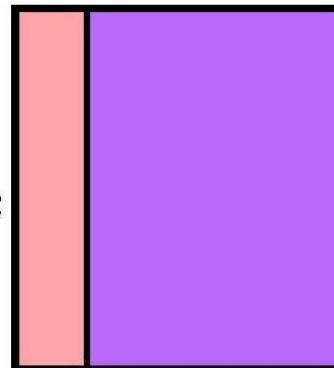
$m=3$



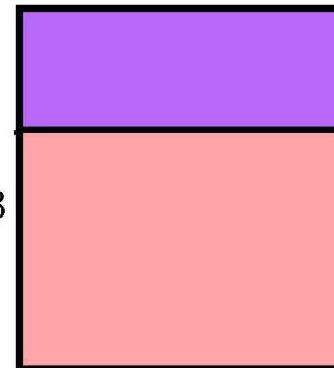
$\alpha_1$



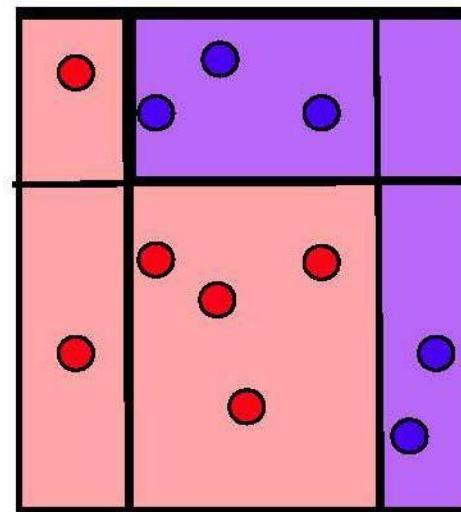
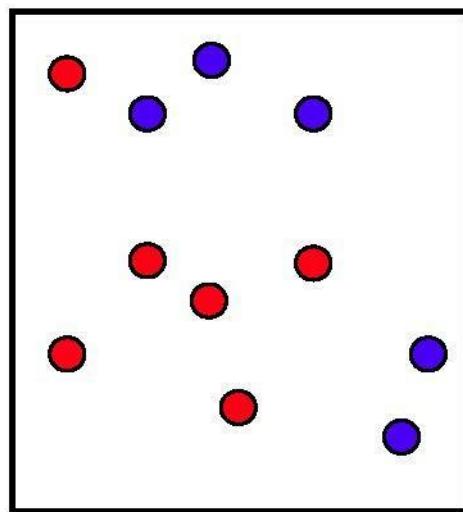
$+ \alpha_2$



$+ \alpha_3$



# Final classifier



STOP

# Bootstrapping

# Bootstrapping

- **Bootstrapping** is any test or metric that relies on random sampling **with replacement**.
- Bootstrapping allows assigning measures of accuracy (defined in terms of bias, variance, confidence intervals, prediction error or some other such measure) to sample estimates.

# Bootstrapping

- Bootstrapping can be used to generate multiple training dataset from a sample dataset.
- If a set of observations can be assumed to be from an I.I.D. population, an approximating distribution of the population can be implemented by constructing a number of resamples (with replacement) of the observed dataset and **of equal size to the observed dataset**.

STOP

Algorithm

Gradient Boosting

AdaBoost

**Tree Bagging and Random Forest**

# Tree Bagging algorithm

- 1) Randomly choose  $n$  samples from the training set with replacement to create a **bootstrap sample**.
- 2) Grow a decision tree from the bootstrap sample. At each node:
  - 1) Split the node using the feature for the best split according to the objective function (e.g., maximizing information gain).
- 3) Repeat steps 1) and 2)  $k$  times.
- 4) Aggregate the prediction by each tree to assign a classification by **majority vote**.

# Random Forest algorithm

- 1) Randomly choose  $n$  samples from the training set with replacement to create a **bootstrap sample**.
- 2) Grow a decision tree from the bootstrap sample. At each node:
  - 1) **Randomly select  $d$  features without replacement**
  - 2) Split the node using the feature for the best split according to the objective function (e.g., maximizing information gain).
- 3) Repeat steps 1) and 2)  $k$  times.
- 4) Aggregate the prediction by each tree to assign a classification by **majority vote**.

# What distinguishes the Random Forest algorithm from Bagging

## 1) Randomly select $d$ features without replacement

This added step is sometimes call **feature bagging**.

## original training data

Sample ID	1	2	3	4	...	N
Feature 1	$F_1^{(1)}$	$F_1^{(2)}$	$F_1^{(3)}$	$F_1^{(4)}$	...	$F_1^{(N)}$
Feature 2	$F_2^{(1)}$	$F_2^{(2)}$	$F_2^{(3)}$	$F_2^{(4)}$	...	$F_2^{(N)}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Feature M	$F_M^{(1)}$	$F_M^{(2)}$	$F_M^{(3)}$	$F_M^{(4)}$	...	$F_M^{(N)}$
Class	①	②	②	①	...	②

1	2	3	4	...	N
①	②	②	①	...	②

original training data

1	2	3	4	...	N
①	①	①	①	...	①

omitted data

bootstrap samples

$S_1$

1	2	2	4	...
①	①	①	①	...

3	...
①	...

$S_2$

1	1	3	4	...
①	①	①	①	...

2	...
①	...

$S_3$

1	1	2	4	...
①	①	①	①	...

3	...
①	...

$S_4$

2	3	3	4	...
①	①	①	①	...

1	...
①	...

⋮

⋮

$S_K$

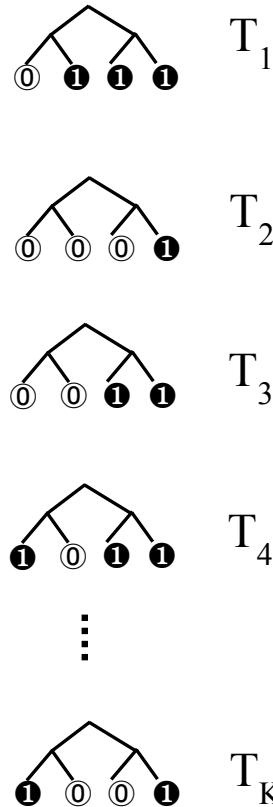
1	1	3	3	...
①	①	①	①	...

2	4	...
①	①	...

training  
data

$S_1$	<table border="1"><tr><td>1</td><td>2</td><td>2</td><td>4</td><td>...</td></tr><tr><td>0</td><td>1</td><td>1</td><td>0</td><td>...</td><td>1</td></tr></table>	1	2	2	4	...	0	1	1	0	...	1
1	2	2	4	...								
0	1	1	0	...	1							
$S_2$	<table border="1"><tr><td>1</td><td>1</td><td>3</td><td>4</td><td>...</td></tr><tr><td>0</td><td>0</td><td>1</td><td>0</td><td>...</td><td>0</td></tr></table>	1	1	3	4	...	0	0	1	0	...	0
1	1	3	4	...								
0	0	1	0	...	0							
$S_3$	<table border="1"><tr><td>1</td><td>1</td><td>2</td><td>4</td><td>...</td></tr><tr><td>0</td><td>0</td><td>1</td><td>0</td><td>...</td><td>1</td></tr></table>	1	1	2	4	...	0	0	1	0	...	1
1	1	2	4	...								
0	0	1	0	...	1							
$S_4$	<table border="1"><tr><td>2</td><td>3</td><td>3</td><td>4</td><td>...</td></tr><tr><td>1</td><td>1</td><td>1</td><td>0</td><td>...</td><td>1</td></tr></table>	2	3	3	4	...	1	1	1	0	...	1
2	3	3	4	...								
1	1	1	0	...	1							
⋮												
$S_K$	<table border="1"><tr><td>1</td><td>1</td><td>3</td><td>3</td><td>...</td></tr><tr><td>0</td><td>0</td><td>1</td><td>1</td><td>...</td><td>0</td></tr></table>	1	1	3	3	...	0	0	1	1	...	0
1	1	3	3	...								
0	0	1	1	...	0							

resulting  
tree

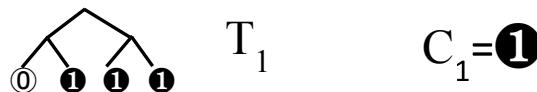


new  
instance

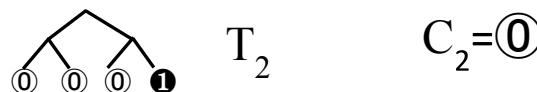
tree

classification  
using tree  $T_i$

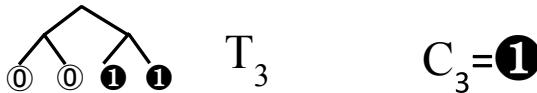
majority vote  
classification



$$C_1 = \mathbf{1}$$

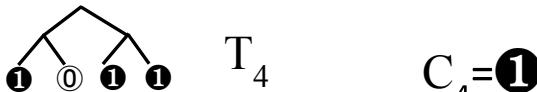


$$C_2 = \mathbf{0}$$



$$C_3 = \mathbf{1}$$

$$C = \mathbf{1}$$



$$C_4 = \mathbf{1}$$

⋮

⋮



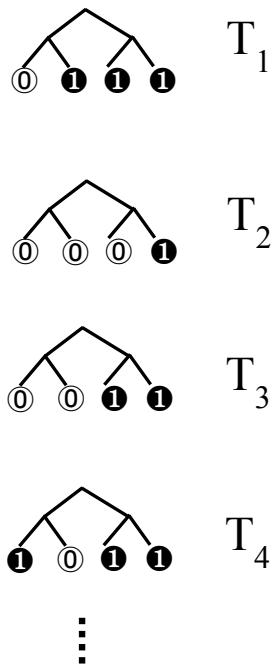
$$C_K = \mathbf{0}$$

Sample ID	j
Feature 1	$F_1^{(j)}$
Feature 2	$F_2^{(j)}$
⋮	⋮
Feature M	$F_M^{(j)}$
Class	?

## training data

Sample ID	1	2	...	N
Feature 1	$F_1^{(1)}$	$F_1^{(2)}$	...	$F_1^{(N)}$
Feature 2	$F_2^{(1)}$	$F_2^{(2)}$	...	$F_2^{(N)}$
⋮	⋮	⋮	⋮	⋮
Feature M	$F_M^{(1)}$	$F_M^{(2)}$	...	$F_M^{(N)}$
Class	①	②	...	②

## random forest

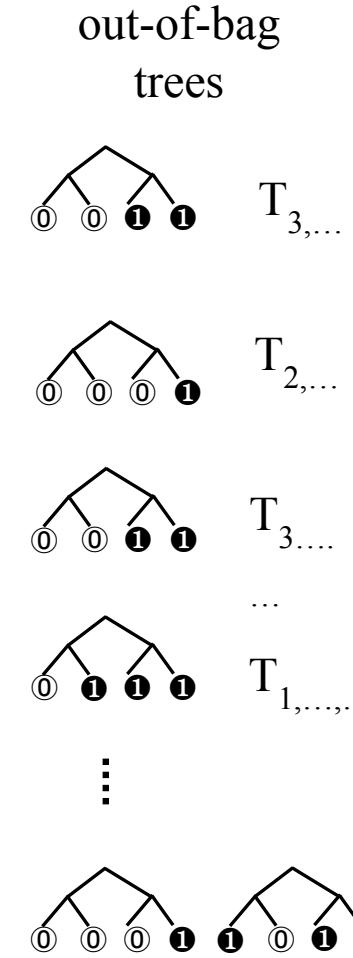


## majority –vote predictions

Sample ID	1	2	...	N
Feature 1	$F_1^{(1)}$	$F_1^{(2)}$	...	$F_1^{(N)}$
Feature 2	$F_2^{(1)}$	$F_2^{(2)}$	...	$F_2^{(N)}$
⋮	⋮	⋮	⋮	⋮
Feature M	$F_M^{(1)}$	$F_M^{(2)}$	...	$F_M^{(N)}$
Class	①	②	...	②
Predicted Class	①	②	...	①
Error	0	0	...	1

$$Error_{training} = \frac{1}{N} \sum_{i=1}^N Error^{(i)}$$

	training data					
$S_1$	1	2	2	4	...	
	0	1	1	0	...	1
$S_2$	1	1	3	4	...	
	0	0	1	0	...	0
$S_3$	1	1	2	4	...	
	0	0	1	0	...	1
$S_4$	2	3	3	4	...	
	1	1	1	0	...	1
	⋮					
$S_K$	1	1	3	3	...	
	0	0	1	1	...	0



Classify each  $x_i$  using only the  $S_k$  that do not contain  $x_i$  to obtain  $\hat{F}_{oob,k}(x_i)$ .

$$Error_{oob} = \frac{1}{N} \sum_{i=1}^N 1_{\hat{F}_{oob}(x_i) \neq y_i}$$

STOP

# Out-of-bag samples

When choosing a bootstrap sample  $S$  of size  $N$ , we select sample points without replacement; some points can (and likely will) be repeated in the bootstrap sample. Other points can (and likely will) be omitted.

The probability that a particular sample point is omitted from  $S$  is

$$\Pr(N) = \left(1 - \frac{1}{N}\right)^N.$$

$$\lim_{N \rightarrow \infty} \Pr(N) = \lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^N = e^{-1} \approx 0.3678. N \text{ does not have to}$$

be very large for  $\Pr(N)$  to be close to its asymptotic value.

For example,  $\Pr(5) \approx 0.328$  and  $\Pr(50) \approx 0.364$ .

# Out-of-bag sample advantages

- The fact any given data point,  $(x_i, y_i)$  is missing from a significant fraction ( $\sim 1/3$ ) of bootstrap samples,  $\hat{S}_k$ , drawn from even a small sample set,  $S$ , creates an interesting opportunity.
- We can train our model on all of the bootstrap samples  $\{\hat{S}_1, \hat{S}_2, \dots, \hat{S}_K\}$  and then test how well the model does in predicting  $y_i$  given  $x_i$  using only those bootstrap samples that do not contain  $x_i$ .
- A data point is said to be **out of bag** for the bootstrap samples that do not contain it.
- Using out of bag samples for testing allows us to use all of our data for training and still have data for testing that has not been used in training.

# Training a majority voting classifier using out-of-bag samples

Let  $F(\mathbf{x}_i) = \text{MajorityVoting}(F_1(\mathbf{x}_i), F_2(\mathbf{x}_i), \dots, F_K(\mathbf{x}_i))$ ,

The training error is given by is

$$E[F(\mathbf{x})] = \frac{1}{N} \sum_{i=1}^N 1_{F(\mathbf{x}_i) \neq y_i},$$