



# Programming for Data Science

## – Combining Models

Henrik Boström

Prof. of Computer Science - Data Science Systems

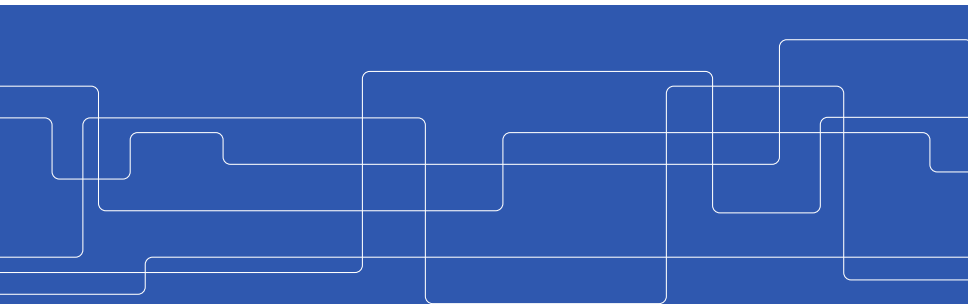
Division of Software and Computer Systems

Department of Computer Science

School of Electrical Engineering and Computer Science

KTH Royal Institute of Technology

[bostromh@kth.se](mailto:bostromh@kth.se)





# Outline

Condorcet's jury theorem

Bagging

Randomization

Random forests

Boosting

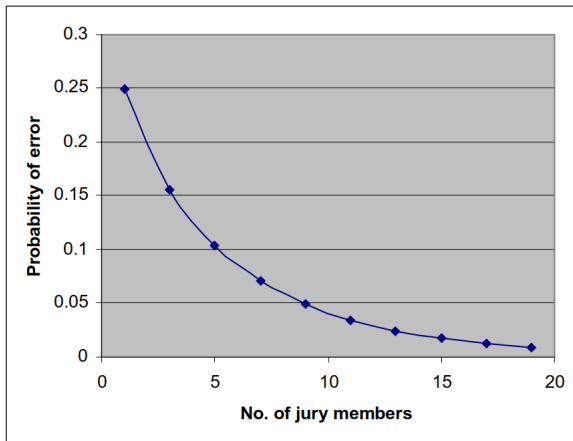
Stacking

## Condorcet's jury theorem

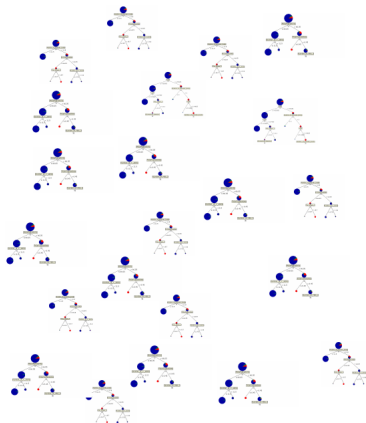
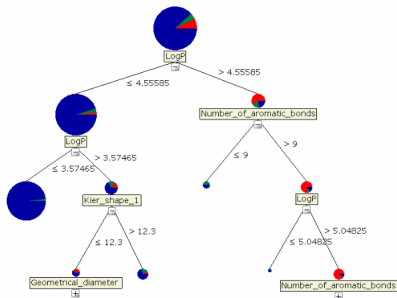
*If each member of a jury is more likely to be right than wrong, then the majority of the jury, too, is more likely to be right than wrong; and the probability that the right outcome is supported by a majority of the jury is a (swiftly) increasing function of the size of the jury, converging to 1 as the size of the jury tends to infinity.*

Condorcet, 1785

# Condorcet's jury theorem (example)



# Single trees vs. forests



- ▶ A bootstrap sample  $B$  of a set of instances  $I$  is created by randomly selecting  $n = |I|$  instances from  $I$  *with replacement*
- ▶ The probability of an instance in  $I$  appearing in  $B$  is

$$1 - \left(1 - \frac{1}{n}\right)^n = 1 - \frac{1}{e} \approx 0.632$$

- ▶ Instances in  $I \setminus B$  are said to be *out-of-bag*

L. Breiman. 1996. Bagging predictors. *Machine Learning*, 24(2):123–140

# Bootstrap sample (example)

Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e1	yes	no	no	yes	some	yes
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e5	yes	no	yes	no	none	no
e6	no	yes	no	yes	some	yes

## Bootstrap sample (example)

Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e2	yes	no	no	yes	full	no
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e4	yes	no	yes	yes	full	yes
e6	no	yes	no	yes	some	yes



# Bagging

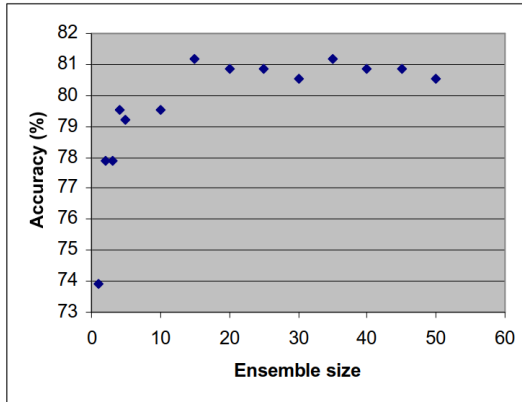
Input: instances  $I$ , learner  $L$ , iterations  $m$

Output: a combined model  $M$

```
for i = 1 to m:  
    B = bootstrap sample of I  
    M_i = L(B)  
M = average({M_1, ..., M_m})
```

For classification, [average](#) can be majority vote (the mode) or the mean of the class probability distributions, while for regression it could be the mean of the predicted values.

## Bagging (example)



- ▶ Heart-disease dataset from UCI repository
- ▶ Stratified 10-fold cross-validation
- ▶ Bagged decision trees

## Bagging in practice

- ▶ Any type of base learner may be used, but better results can often be obtained with more brittle methods, e.g., decision trees, neural networks
- ▶ Works in the same way for classification and regression
- ▶ For decision trees, pruning often has a detrimental effect, i.e., the more variance of each individual model the better
- ▶ Predictive performance can be measured without a separate validation/test set; each training instance may be predicted using ensemble members for which the instance is out-of-bag. This leaves more instances for model generation and all instances will be used for testing.

- ▶ Repeated application of a base learner that contains some stochastic element (or that can be adjusted to contain stochastic elements) such as
  - ▶ the setting of initial weights in a neural network
  - ▶ sampling of grow and prune data
  - ▶ choice of attribute to split on in decision trees
- ▶ Randomization may be combined with other techniques

Random forests (Breiman 2001) are generated by combining two techniques:

- ▶ bagging (Breiman 1996)
- ▶ the random subspace method (Ho 1998)
  - consider only a random sample of the features; either one sample for each tree or one sample for each split. The latter is employed in random forests.

L. Breiman. 2001. Random forests. *Machine Learning*, 45(1):532

T. K. Ho. 1998. The random subspace method for constructing decision forests, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(8):832-844

# Random Forest (example)

Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e2	yes	no	no	yes	full	no
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e4	yes	no	yes	yes	full	yes
e6	no	yes	no	yes	some	yes

e2 e2 e3 e4 e4 e6

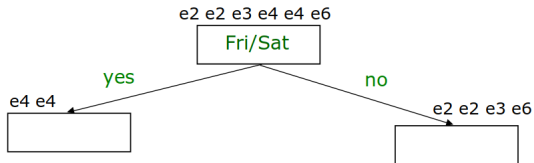
# Random Forest (example)

Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e2	yes	no	no	yes	full	no
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e4	yes	no	yes	yes	full	yes
e6	no	yes	no	yes	some	yes

e2 e2 e3 e4 e4 e6

# Random Forest (example)

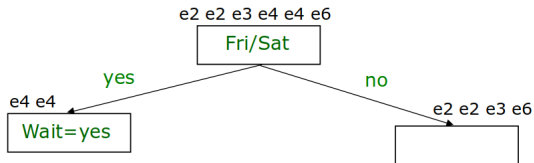
Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e2	yes	no	no	yes	full	no
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e4	yes	no	yes	yes	full	yes
e6	no	yes	no	yes	some	yes





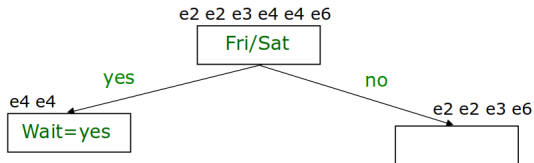
# Random Forest (example)

Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e2	yes	no	no	yes	full	no
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e4	yes	no	yes	yes	full	yes
e6	no	yes	no	yes	some	yes



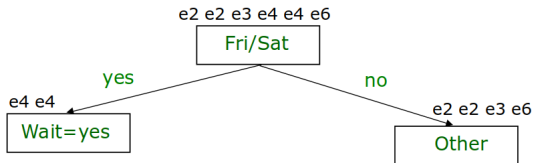
# Random Forest (example)

Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e2	yes	no	no	yes	full	no
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e4	yes	no	yes	yes	full	yes
e6	no	yes	no	yes	some	yes



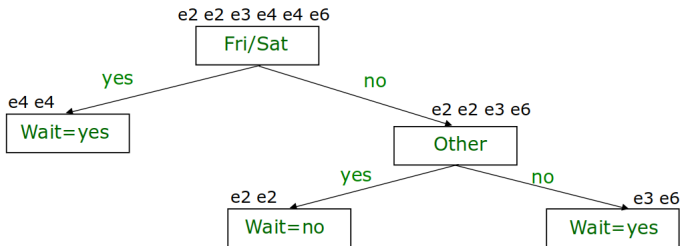
# Random Forest (example)

Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e2	yes	no	no	yes	full	no
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e4	yes	no	yes	yes	full	yes
e6	no	yes	no	yes	some	yes



# Random Forest (example)

Ex.	Other	Bar	Fri/Sat	Hungry	Guests	Wait
e2	yes	no	no	yes	full	no
e2	yes	no	no	yes	full	no
e3	no	yes	no	no	some	yes
e4	yes	no	yes	yes	full	yes
e4	yes	no	yes	yes	full	yes
e6	no	yes	no	yes	some	yes



## Random forests = robust performance

In (Fernandez-Delgado et al 2014), an empirical investigation was presented using:

- ▶ 179 classifiers from 17 families, incl. all standard approaches
- ▶ 121 datasets, incl. all of UCI except the largest ones
- ▶ The random forest versions ranked the highest and they obtained near to the best accuracy for almost all the data sets

It was concluded that:

*The classifiers most likely to be the best are the random forest (RF) versions*

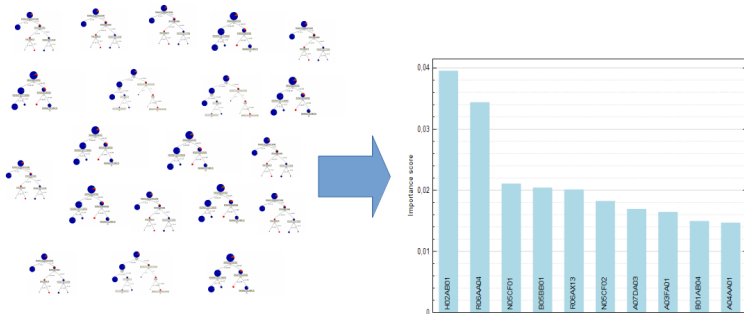
M. Fernandez-Delgado, E. Cernadas, S. Barro, and D. Amorim. Do we Need Hundreds of Classifiers to Solve Real World Classification Problems? *Journal of Machine Learning Research*, 15(1), pp. 3133-3181, 2014.



# Interpretability

- ▶ Can we interpret the model/understand the predictions?
  - ▶ To be discussed ...

# Variable Importance



- Breiman (2001) proposed an approach to estimating variable importance by measuring the relative performance degradation (on OOB predictions) when permuting the values of each feature in turn.

# Boosting: AdaBoost<sup>1</sup>

Input: instances  $\{Z_1, \dots, Z_n\}$ , learner  $L$ , iterations  $m$

Output: a set of model-weight pairs  $M$

$w_1, \dots, w_n = 1$

$M = \{\}$

for  $i = 1$  to  $m$ :

$M_i = L(\{(Z_1, w_1), \dots, (Z_n, w_n)\})$

$\text{Err} = (w_1 \cdot \text{err}(M_i, z_1) + \dots + w_n \cdot \text{err}(M_i, z_n)) /$   
         $(w_1 + \dots + w_n)$

    if  $\text{Err} = 0$  or  $\text{Err} > 0.5$  then break

    for  $j = 1$  to  $n$ :

        if  $\text{err}(M_i, Z_j) = 0$  then  $w_j = w_j \cdot \text{Err} / (1 - \text{Err})$

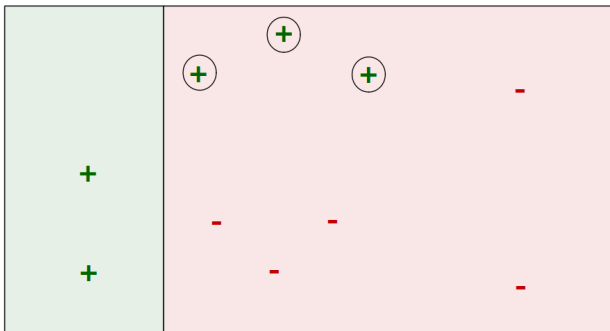
$M = M + \{(M_i, -\log \text{Err} / (1 - \text{Err}))\}$

---

<sup>1</sup>Y. Freund and R. Schapire. 1997. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. *J. Comput. Syst. Sci.* 55(1): 119-139



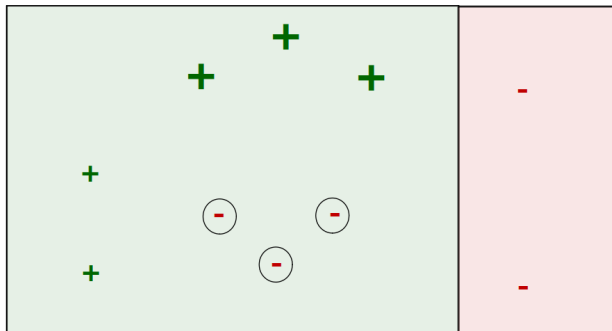
# AdaBoost (example)



M1

$$Err = 0.3 \quad \frac{Err}{1 - Err} \approx 0.43$$

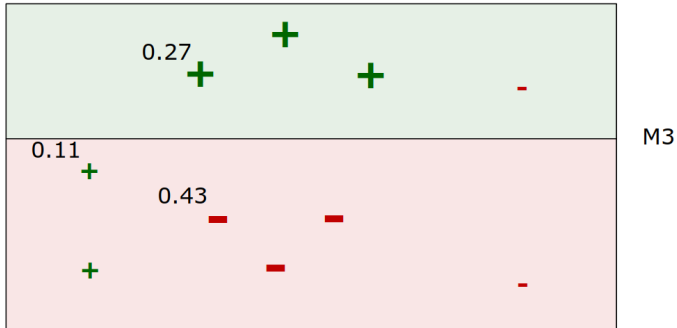
## AdaBoost (example, cont.)



M2

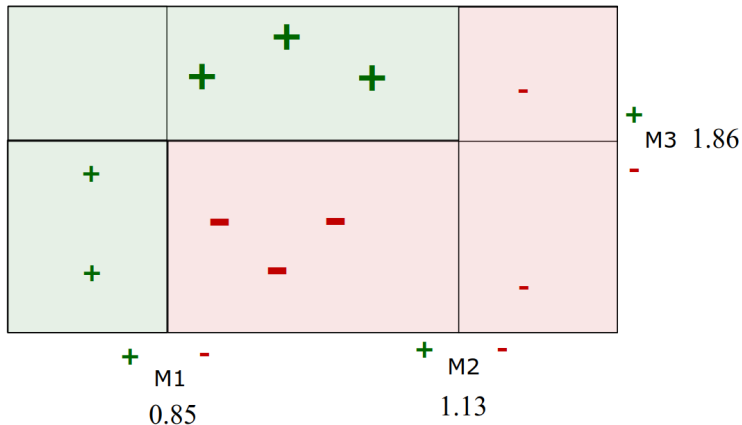
$$Err = 0.21 \quad \frac{Err}{1 - Err} \approx 0.27$$

## AdaBoost (example, cont.)



$$Err \approx 0.13 \quad \frac{Err}{1 - Err} \approx 0.16$$

# AdaBoost (example, cont.)



# Gradient boosting machine (regression)<sup>1</sup>

Input: instances  $\{(X_1, y_1), \dots, (X_n, y_n)\}$ ,  
learner  $L$ , iterations  $m$

Output: a sequence of models  $M_0, \dots, M_m$

$M_0 = (y_1 + \dots + y_n)/n$

for  $i = 1$  to  $m$ :

    for  $j = 1$  to  $n$ :

$y_j = y_j - M_{i-1}(X_j)$

$M_i = L(\{(X_1, y_1), \dots, (X_n, y_n)\})$

---

<sup>1</sup>J. Friedman. 2001. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pp.1189-1232

# Gradient boosting machine (example)

Input:  $\{(0, y_1=1), (1, y_2=1), (2, y_3=2),$   
 $(3, y_4=2), (4, y_5=3), (5, y_6=3)\},$   
 $L = \text{decision stump}, m = 3$

$$M_0 = (y_1 + \dots + y_6)/6 = 2$$

Iteration 1:

$$y_1 \text{ --} 2 = -1$$

$$y_2 \text{ --} 2 = -1$$

$$y_3 \text{ --} 2 = 0$$

$$y_4 \text{ --} 2 = 0$$

$$y_5 \text{ --} 2 = 1$$

$$y_6 \text{ --} 2 = 1$$

Iteration 2:

$$y_1 \text{ --} -0.5 = -0.5$$

$$y_2 \text{ --} -0.5 = -0.5$$

$$y_3 \text{ --} -0.5 = 0.5$$

$$y_4 \text{ --} -0.5 = 0.5$$

$$y_5 \text{ --} 1 = 0$$

$$y_6 \text{ --} 1 = 0$$

Iteration 3:

$$y_1 \text{ --} -0.5 = 0$$

$$y_2 \text{ --} -0.5 = 0$$

$$y_3 \text{ --} 0.25 = 0.25$$

$$y_4 \text{ --} 0.25 = 0.25$$

$$y_5 \text{ --} 0.25 = -0.25$$

$$y_6 \text{ --} 0.25 = -0.25$$

$$M_1: \begin{array}{l} 1 \text{ if } x > 3 \\ \text{else } -0.5 \end{array}$$

$$M_2: \begin{array}{l} -0.5 \text{ if } x < 2 \\ \text{else } 0.25 \end{array}$$

$$M_3: \begin{array}{l} -0.25 \text{ if } x > 3 \\ \text{else } 0.125 \end{array}$$

## Boosting in practice

- ▶ For multi-class problems, it may be difficult for AdaBoost to reduce the error below 50% with weak base learners, e.g., decision stumps
  - ▶ More powerful base learners may be employed
  - ▶ The problem may be transformed into multiple binary classification problems
  - ▶ Specific multi-class versions of AdaBoost have been developed
- ▶ Various loss functions may be used for the Gradient boosting machine (GBM), including log likelihood for classification
- ▶ GBM requires careful tuning, in particular of tree depth and *learning rate*
- ▶ XGBoost, well-known for winning many competitions, see e.g., [www.kaggle.com](http://www.kaggle.com), implements GBM
- ▶ Boosting can be sensitive to label noise, e.g., erroneously labeled training examples

# Stacking

Input: instances  $I$ , learner  $L$ , learners  $L_1, \dots, L_m$

Output: a model  $M$ , a set of models  $M_1, \dots, M_m$

divide  $I$  into a training and test set

for each learner  $L_i$ :

    generate a model  $M_i$  from the training set

    and apply it to the test set

$I' = \{(y_1, x_{11}, \dots, x_{1m}), \dots, (y_n, x_{n1}, \dots, x_{nm})\}$ ,

    where  $y_i$  is the true label for the  $i$ th

    test instance, and  $x_{ij}$  is the prediction

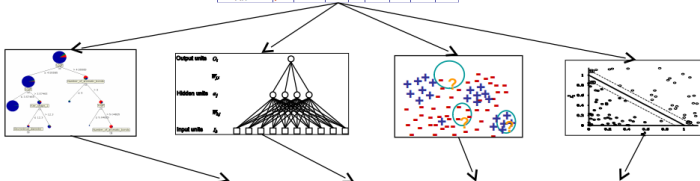
    of model  $M_j$  for that instance

$M = L(I')$



# Stacking (example)

Name	Stability	Frac. of molecules	Gen. diam.	Log P	Zagreb index	Topol. diam.	No. C atoms	No. heavy atoms
methylpentane	good	0.40	3.45	2.44	19	4	6	5
methylcyclohexane	good	0	3.00	2.51	31	4	7	7
hexane	med.	0.75	6.93	3.53	28	8	9	8
heptadecane	good	0.90	4.30	3.19	19	3	6	5
butadiene	good	0.33	2.05	1.36	8	3	4	3
naphthalene	good	0	3.61	2.84	57	3	10	11
nonaphthalene	good	0	3.56	3.32	69	3	12	14
pyrene	poor	0	5.00	4.59	117	7	16	19
dimethylanthracene	poor	0	5.29	4.61	108	7	16	18
hexadecaprene	med.	0	5.00	3.82	117	7	16	19
naphthalene	poor	0	5.00	5.13	128	7	18	21
benzodicyclopent	poor	0	5.29	5.66	132	7	20	24



Class	M1	M2	M3	M4
good	good	good	good	med.
good	med.	good	good	good
med.	med.	good	med.	poor
poor	good	poor	good	good

IF M2 = good & M3 = good THEN Class = good  
 IF M2 = poor THEN Class = poor

...



## Stacking in practice

- ▶ Linear models have been shown to be effective when learning the combination function
- ▶ Predictive performance can be improved by combining class probability estimates rather than class labels generated by the base models

## Summary

- ▶ We have considered various ways of combining models; bagging, randomization, boosting and stacking (as well as some combinations of these)
- ▶ Compared to the base learning algorithms, the predictive performance can be substantially improved; sometimes leading to state-of-the-art performance
- ▶ The increased predictive power comes, however, at the cost of reduced interpretability
- ▶ The combination strategies may not necessarily lead to increased computational cost; and this cost can often be further reduced by parallelization