

Committees, ensembles.

Petr Pošík
petr.posik@fel.cvut.cz

Czech Technical University in Prague
Faculty of Electrical Engineering
Dept. of Cybernetics

Introduction	2
Quiz	3
Committee	4
Examples	5
Aggregation.....	6
Stacking	7
Bagging	8
Quiz	9
Bootstrapping	10
Bootstrap sample	11
Bagging	12
Features	13
Random forests	14
RF	15
Features	16
Boosting	17
Boosting	18
Quiz	19
AdaBoost	20
Graphically	21
Algorithm	22
AdaBoost: remarks	23
Another view	24
L2Boosting	25
GBM	26
Further considerations	27
Summary	28
Competencies	29

Quiz

Ensemble methods use essentially *wisdom of crowd*, i.e. combining the opinions of many entities to get a single „average“ opinion.

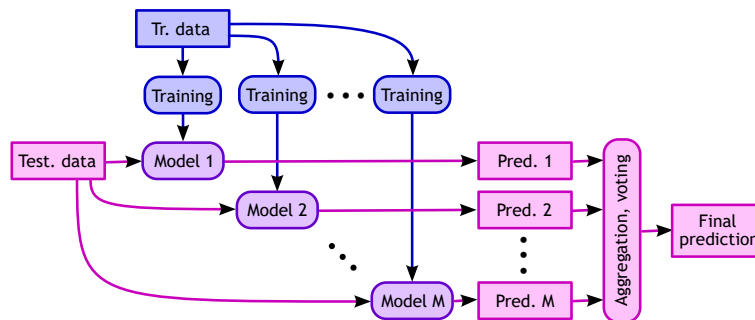
What configuration of individual entities is the best for the ensemble?

- A** All entities shall be of a similar type.
- B** The entities shall be as diverse as possible.
- C** The diversity of the entities does not matter.
- D** The diversity matters but we cannot control it.

Note: Similarity/diversity of entities relates to the ways (internal processes) they arrive at their opinion. It is not about similarity/diversity of their opinions.

Ensemble a.k.a committee

- ML model composing multiple different models to obtain better predictive performance than could be obtained from any of the constituent models.
- A way to compensate for poor learning algorithms by performing a lot of extra computations.
- Ensembles tend to yield better results when there is a significant diversity among the models (Intuition: averaging reduces variance).
- Individual ensemble/committee methods differ in the way they create individual *models different from each other*.
- Use different kinds of models, or models unstable w.r.t. a change in the training data.



Ensemble examples

Some examples of committee/ensemble methods:

- Stacking
- Bagging
- Random forests
- Boosting
- ...

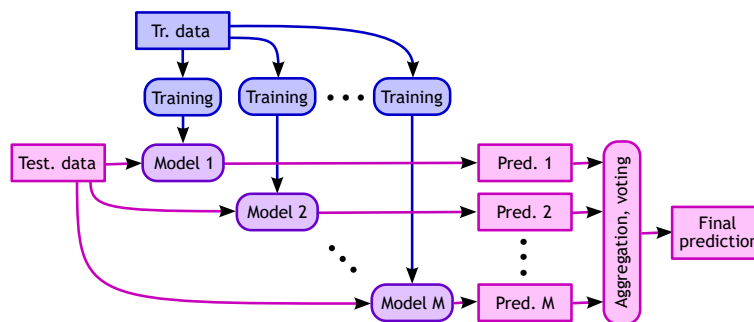
Decision trees (classification and regression) are used most often as the base models because

- they are relatively fast to learn,
- they are unstable w.r.t. the changes in the training dataset, and thus
- it is quite easy to make a lot of trees which are very diverse.

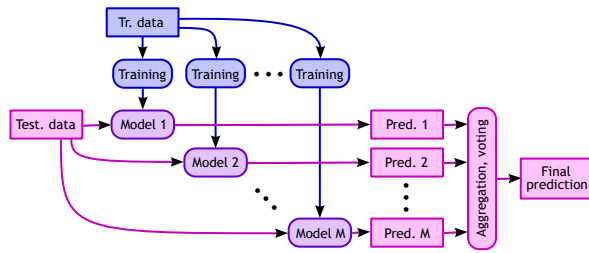
Aggregation

The final aggregation of results of individual models is usually done by

- (weighted) *voting* of individual models for *classification* problems,
- (weighted) *averaging* of individual models for *regression* problems,
- or by other techniques.



Stacking



- Assume we have M different models h_m created for the same modeling task, each being a function $h_m(x)$ of the input features x .
- The predictions of these models, $h(x) = (h_1(x), \dots, h_M(x))$, may be considered new features extracted from the data set (similar to basis expansion, but new features are learned, not predefined).
- We can thus train a higher-level classification/regression model h_{stack} as a function of these new features, i.e. $h_{\text{stack}}(h)$ (sometimes together with the original features, i.e. $h_{\text{stack}}(x, h)$).

- For classification, logistic regression is often used as h_{stack} .
- For regression, multiple linear regression is often used as h_{stack} with the constraint on the weights w_i such that $\sum w_i = 1$ and $w_i > 0 \forall i$.

- An obvious way to estimate the weights w as

$$w^* = \arg \min_w \sum_{i=1}^{|T|} L \left(y_i, \sum_{m=1}^M w_m h_m(x_i) \right),$$

however, can result in overfitting; this is solved by LOO cross-validation, i.e. using the estimate

$$w^* = \arg \min_w \sum_{i=1}^{|T|} L \left(y_i, \sum_{m=1}^M w_m \hat{h}_m^{-i}(x_i) \right),$$

where \hat{h}_m^{-i} is a predictor obtained by training on data excluding (x_i, y_i) , i.e. at the price of high-computational demands.

Bagging

Quiz

Assume we have an urn with N numbered balls, $1, \dots, N$. We repeat the following process N times:

1. Shake the urn.
2. Pick one ball randomly and note its number.
3. Return the ball to the urn.

Some of the balls were picked more than once, some of the balls were never picked.

What is the expected number of balls never picked from the urn for $N = 1,000$?

- A 5
- B 43
- C 136
- D 368

Bootstrapping

- A general statistical technique for assessing the accuracy of parameter estimates and for hypotheses testing.
- It relies on *many repetitions* and *random sampling with replacement*.

Example: Assume we want to estimate the average height of all the people in the world. How to do that?

- Cannot measure the whole population, measure just a sample of N people.
- Using this sample, we can obtain a (single) point estimate of the average population height: $\hat{h} = \frac{1}{N} \sum_{i=1}^N h_i$.
- We also need some measure of uncertainty/variability of this estimate. How to do that?
- Either use “classic” statistics: compute the sample variance \hat{s}_h^2 and compute the variance of the estimate as $\hat{s}_{\hat{h}}^2 = \frac{\hat{s}_h^2}{N-1}$,
- or use *bootstrapping*:
 1. For b in $1, \dots, M$ where $M \in (10^2, \dots, 10^6)$:
 - Create a bootstrap sample T^b from the original dataset T .
 - Compute b th estimate of the statistic (here average) from the bootstrap sample.
 2. Now you have a histogram of the estimates (here averages), from which you can estimate the mean, variance, ... of the sampling distribution.

Similar process works for many other estimators.

Bootstrap sample

Assume we have a dataset T with N items. What is the **bootstrap sample** T^b ?

- A perturbed version of the original dataset T .
- Each item of T^b was chosen uniformly with replacement from the original dataset T . Usually, $|T| = N = |T^b|$.
- Some items of T are copied to T^b more than once. Some items are not copied at all.

How many unique elements of T are present in T^b (on average)?

- Probability that a particular item will not be chosen in one particular pick: $1 - \frac{1}{N}$
- Probability that a particular item will not be chosen in any of N picks: $\left(1 - \frac{1}{N}\right)^N$
- The expected number of items that will not be copied to a bootstrap sample: $N \left(1 - \frac{1}{N}\right)^N \approx Ne^{-1} = N \cdot 0.368$
- The expected number of unique elements copied from T : $N \left(1 - \left(1 - \frac{1}{N}\right)^N\right) \approx N(1 - e^{-1}) = N \cdot 0.632$

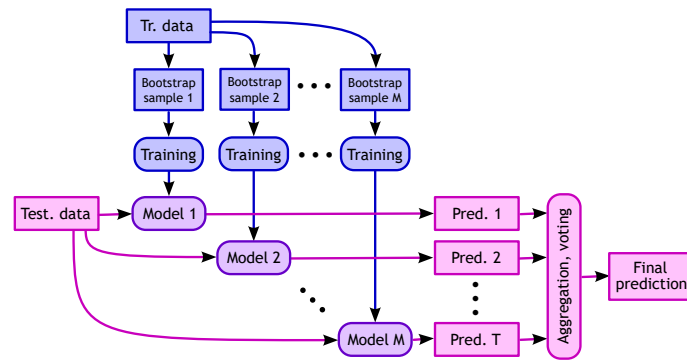
Bagging a.k.a. Bootstrap aggregation

- Uses bootstrap to improve the estimate or the prediction itself.
- Aggregating results of several models reduces variance and prevents overfitting.
- Algorithm:

1. Create M bootstrap samples T^b from training data T ($b = 1, \dots, M$).
2. Build a model h_b on each bootstrap sample T^b .

3. Construct final model by averaging/voting the predictions of individual models: $\hat{y} = h_{\text{bag}}(x) = \frac{1}{M} \sum_{b=1}^M h_b(x)$, resp.

$$\hat{y} = h_{\text{bag}}(x) = \arg \max_{y \in C} \sum_{b=1}^M I(y = h_b(x))$$



Features

Bagging

- leads to improvements for unstable procedures (artificial neural networks, classification and regression trees, etc.), but
- it can mildly degrade the performance of stable methods such as K-nearest neighbors.
- Thanks to bootstrapping, it can provide not only predictions, but also estimates of uncertainty of those predictions.

Estimate of prediction error (out-of-bag error):

- Around 37 % of training examples are not part of a bootstrap sample; they are called OOB (out of bag).
- We can predict the model response for each training sample x_i using only the models that did not have x_i in their bootstrap sample.
- We can average these predicted responses (regression) or can take a majority vote (classification) to get a single “OOB prediction” for each observation.
- OOB predictions then can be used to compute OOB estimate of the error.
- With M sufficiently large, OOB error is virtually equivalent to leave-one-out cross-validation error.

Random forest (RF)

An ensemble method using set of decision trees (i.e. forest):

- Trees that are grown very deep tend to learn highly irregular patterns: they *overfit* their training sets, i.e. have *low bias*, but *very high variance*.
- RF perform averaging of multiple deep decision trees, trained on different parts of the same training set, with the goal of reducing the variance.

RF combine

- bagging, and
- random subspace method (see below).

Predictions are computed using voting/averaging.

To train a single tree, RF algorithm

- creates a bootstrap sample of the training data (bagging), and
- uses a modified tree-learning algorithm which considers only *a random subset of input features* at each candidate split in the learning process ("feature bagging"; this further decorrelates the resulting trees). Suggestions:
 - Classification: consider \sqrt{D} features at each split.
 - Regression: consider $D/3$ features at each split, use minimum node size of 5.
- In *ExtRaTrees* (extremely randomized trees), instead of searching for the locally optimal split for each variable, a random value is used for the split.

RF features

Estimate of prediction uncertainty and OOB error:

- See bagging.

Variable importance:

1. Grow the forest. Compute OOB error for each data point averaged over the whole forest.
2. To measure the importance of j th variable, permute its values, and compute OOB error on this perturbed dataset. Compute the difference of the estimates before and after permutation.
3. The larger the difference, the larger the importance of variable j .

Boosting

Hypothesis Boosting Problem

- If there exists an efficient algorithm able to create *weak classifiers* (i.e. classifiers only slightly better than random guessing), does it also mean that there is an efficient algorithm able to build *strong classifiers* (i.e. classifiers with an arbitrary precision)?
- No constraint on the algorithm.

Most (not all) **Boosting algorithms**

- sequentially learn weak classifiers using weighted training set (based on information from previously learnt weak classifiers),
- construct the final strong classifier as a weighted sum of the weak classifiers,
- assign the weights to individual weak learners depending on their accuracy,
- re-weight the training data for another round of the weak learner,
- differ in the way how they weight the training data and/or the individual weak classifiers.

Quiz

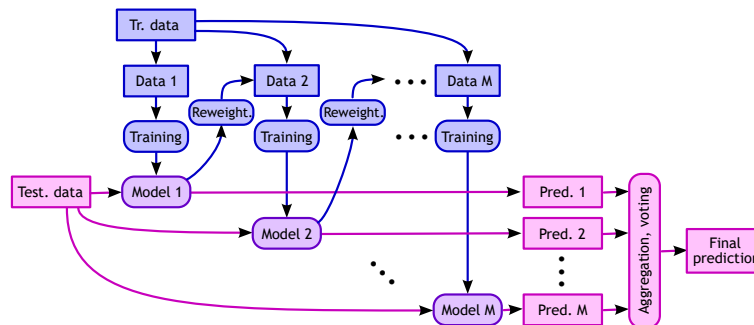
How should the weight of a weak classifier in the final strong classifier be set up?
 How should the weight of the training examples be set up for the next learning iteration?

- A** Accurate weak classifiers get *larger* weights.
 Misclassified examples get *larger* weights.
- B** Accurate weak classifiers get *smaller* weights.
 Misclassified examples get *larger* weights.
- C** Accurate weak classifiers get *larger* weights.
 Misclassified examples get *smaller* weights.
- D** Accurate weak classifiers get *smaller* weights.
 Misclassified examples get *smaller* weights.

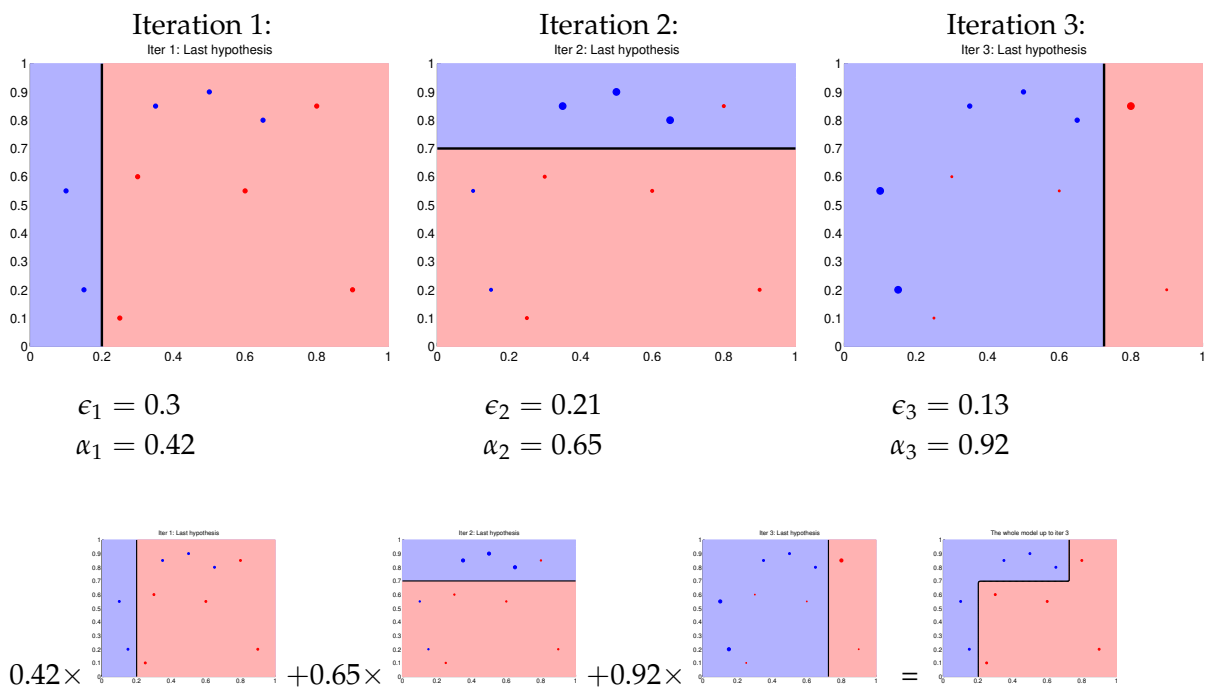
AdaBoost (informally)

AdaBoost

- Training data:
 - In each iteration $t = 1, \dots, M$, it uses different weights $w_t(i)$ of the training examples x_i .
 - *Misclassified examples get a larger weight* for the next iteration.
- The resulting classifier:
 - Weighted voting.
 - More accurate models get larger weight.



AdaBoost.M1 graphically



AdaBoost.M1

- AdaBoost for classification problem and weak learners with class label as output.
- A slightly different version exists for weak learners with output in the form of class probabilities.

Algorithm 1: AdaBoost.M1

Input: Training set of labeled examples: $T = \{\mathbf{x}_i, y_i\}, \mathbf{x}_i \in \mathcal{R}^D, y_i \in \{+1, -1\}, i = 1, \dots, |T|$

Output: Final classifier $H_{\text{final}}(x) = \text{sign} \left(\sum_{m=1}^M \alpha_m h_m(x) \right)$

```
1 begin
2   Initialize the weights of training examples:  $w_1(i) = \frac{1}{|T|}$ .
3   for  $m = 1, \dots, M$  do
4     Train a weak classifier  $h_m$  using  $T$  with weights  $w_m$ .
5     Compute the weighted error:  $\epsilon_m = \frac{\sum_{m=1}^{|T|} w_m(i) I(y_i \neq h_m(\mathbf{x}_i))}{\sum_{m=1}^{|T|} w_m(i)}$ 
6     Compute the weight of classifier  $h_m$ :  $\alpha_m = \ln \left( \frac{1 - \epsilon_m}{\epsilon_m} \right) > 0$ 
7     Update the weights of the training examples:  $w_{m+1}(i) = w_m(i) \cdot \exp[\alpha_m I(y_i \neq h_m(\mathbf{x}_i))]$ .
```

AdaBoost: remarks

The training error:

- Let $\gamma_t = 0.5 - \epsilon_t$ be the improvement of the t -th model over a random guess.
- Let $\gamma = \min_t \gamma_t$ be the minimal improvement, i.e. the difference of error of all models h_t compared to the error of random guessing is at least γ , i.e.

$$\forall t : \gamma_t \geq \gamma > 0.$$

- It can be shown that the training error

$$\text{Err}_{\text{Tr}}(H_{\text{final}}) \leq e^{-2\gamma^2 M}$$

Forward stagewise additive modeling

Boosting tries to solve the following optimization problem: $f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^{|T|} L(y_i, f(x_i))$ Finding the optimal f^* is hard; we shall tackle it sequentially:

Algorithm 2: Forward stagewise additive modeling (FSAM)

```
1 begin
2   Initialize  $f_0(x) = 0$ .
3   for  $m = 1, \dots, M$  do
4     Compute  $(\alpha_m, \theta_m) = \arg \min_{\alpha, \theta} \sum_{i=1}^{|T|} L(y_i, f_{m-1}(x_i) + \alpha h(x_i; \theta))$ .
5     Set  $f_m(x) = f_{m-1}(x) + \alpha_m h(x; \theta_m)$ .
```

AdaBoost.M1 is equivalent to FSAM using the **exponential loss function** $L(y, f(x)) = \exp(-y \cdot f(x))$.

$$\begin{aligned}(\alpha_m, \theta_m) &= \arg \min_{\alpha, \theta} \sum_{i=1}^{|T|} \exp[-y_i (f_{m-1}(x_i) + \alpha h(x_i; \theta))] \\ &= \arg \min_{\alpha, \theta} \sum_{i=1}^{|T|} w_m(i) \exp[-y_i \alpha h(x_i; \theta)],\end{aligned}$$

where $w_m(i) = \exp(-y_i f_{m-1}(x_i))$ depend neither on α_m nor θ_m and can be regarded as weights of training examples, which change each iteration. AdaBoost.M1 then follows from minimization of the last expression.

L2Boosting

Suppose we need to solve *regression problem* with squared error loss (L2).

- Then at step m we have:

$$\begin{aligned}L(y_i, f_m(x_i)) &= L(y_i, f_{m-1}(x_i) + \alpha_m h(x_i; \theta_m)) = \\ &= (y_i - f_{m-1}(x_i) - \alpha_m h(x_i; \theta_m))^2 = \\ &= (r_{im} - \alpha_m h(x_i; \theta_m))^2,\end{aligned}$$

where we define $r_{im} = y_i - f_{m-1}(x_i)$ to be the current **residual** of the model for i th data point.

- By fitting each weak model h_m to the residuals r_{im} , the m th model f_m learns to correct its predecessor f_{m-1} .
- Observation: the residuals $r_{im} = y_i - f_{m-1}(x_i)$ are negative gradients of the squared error loss function $\frac{1}{2}(y - f(x))^2$.
- The algorithm can be viewed as a *gradient descent in the space of functions*.
- The generalization of
 - FSAM using exponential loss (AdaBoost.M1) and
 - FSAM using L2 loss (L2Boosting)

for a general differentiable loss function L is called **Gradient Boosting Machine**.

Gradient Boosting Algorithms

Algorithm 3: Gradient boosting

Input: Training set of labeled examples: $T = \{x_i, y_i\}, i = 1, \dots, |T|$, a differentiable loss function $L(y, f(x))$, number of iterations M .

Output: Final model $f_M(x)$.

```
1 begin
2   Initialize model with constant value:  $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^{|T|} L(y_i, \gamma)$ 
3   for  $m = 1, \dots, M$  do
4     Compute pseudo-residuals  $r_{im} = -\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \Big|_{f(x)=f_{m-1}(x)}$  for all  $i = 1, \dots, |T|$ .
5     Fit model  $h_m(x)$  to pseudo-residuals, i.e. use training set  $\{(x_i, y_i)\}_{i=1}^{|T|}$ .
6     Compute multiplier  $\alpha_m$  by solving the following 1D opt. problem:  $\alpha_m = \arg \min_{\alpha} \sum_{i=1}^{|T|} L(y_i, f_{m-1}(x_i) + \alpha h_m(x_i))$ .
7     Update the model:  $f_m(x) = f_{m-1}(x) + \alpha_m h_m(x)$ 
```

By plugging in different loss functions, we can construct different boosting variants like

- AdaBoost,
- L2Boost,
- LogitBoost,
- etc.

Further considerations

Choosing the number of models M :

- The optimal value usually found by tracking the error on validation set.
- Often, we do not bother; we just set it sufficiently high (several hundreds). Boosting can overfit, but is quite resistant to it.

Shrinkage:

- Often, the so-called shrinkage is applied, i.e. only a small part of the m th model is used:

$$f_m(x) = f_{m-1}(x) + \nu \alpha_m h(x; \theta_m),$$

where $\nu \in (0, 1)$, often $\nu \approx 0.1$, is the so-called *learning rate*.

- Learning is slowed down; it requires more models to be added to the model, providing a configuration trade-off between the number of trees and learning rate.

Stochastic gradient boosting

- It is possible to subsample the training data set and use only a subset of it to train each model.
- Subsample examples as in bagging (but without replacement).
- Subsample features as in random forests.
- It further *prevents overfitting*, *speeds up learning* of individual models, and gives chance to *compute out-of-bag error* estimates.

Competencies

After this lecture, a student shall be able to ...

1. describe the basic principle behind all committee/ensemble methods;
2. list and conceptually compare several methods to achieve diversity among models trained on the same data, and know which of these methods are used in which ensemble algorithms;
3. explain the purpose and the basic principle of stacking;
4. explain how a bootstrap sample is created from the available data, and describe its properties;
5. describe features of bagging;
6. explain how to compute out-of-bag error estimate when using bagging;
7. explain the principle of random forests and describe their difference to bagging with trees;
8. explain how to compute a score of variable importance using random forest;
9. explain the hypothesis boosting problem, and define a weak and a strong classifier in this context;
10. explain the basic principle of AdaBoost.M1 algorithm;
11. relate the training error of the AdaBoost algorithm to the number of constituent models and to the errors of individual models;
12. describe the relations of AdaBoost.M1, L2Boost, and Gradient Boosting.