

Ensemble Methods: Types

- **Bagging:** train learners in parallel on different samples of the data, then combine by voting (discrete output) or by averaging (continuous output).
- **Stacking:** combine model outputs using a second-stage learner like linear regression (different types of learners).
- **Boosting:** train learners on the filtered output of other learners (sequential).

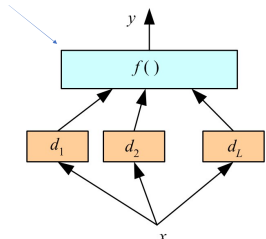
Bagging Issues

1. **Homogeneous Learners:** Similar classifiers usually make similar errors - so forming an ensemble with similar classifiers may not improve the classification accuracy.
2. At the presence of a classifier that performs much better than all available base classifiers, may cause degradation in the overall performance.
3. A poorly performing classifier may cause performance deterioration of the ensemble.
4. **Amount of correlation** among the incorrect classifications made by the base classifiers.
5. **Dominance of a classifier:** If the consistent classifiers tend to misclassify the same instances, then combining their results will have no benefit.
6. **Diversity of classifiers – Heterogeneous Learners:** In contrast, a greater degree of independence among the classifiers can result in errors made by individual classifiers being overlooked when the results of the ensemble are combined.

Stacking

Combiner/Meta-Learner $f()$ is another learner.

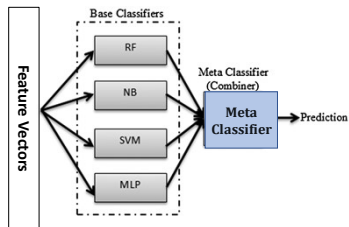
- Uses *meta learner (level-1)* instead of voting to combine predictions of base (level-0) learners
- Predictions of base learners (*level-0 models*) are used as input for meta learner (*level-1 model*).
- Base learners usually different learning schemes (heterogeneous models).
- Hard to analyze theoretically: “black magic”



The diagram illustrates the stacking architecture. At the bottom, an input x is fed into multiple base learners, represented by boxes labeled d_1 , d_2 , and d_L . Each base learner produces a prediction. These predictions are then fed into a meta-learner, represented by a box labeled $f()$. The meta-learner combines these predictions to produce the final output y . A blue arrow points from the text 'Combiner/Meta-Learner $f()$ is another learner.' to the $f()$ box.

1

Stacking



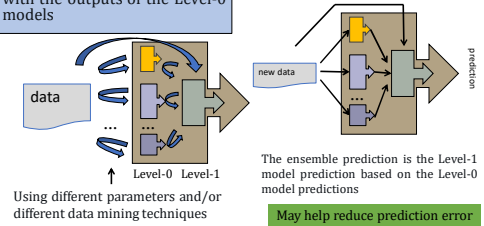
1. Train **several different models** (base/level-0 models) on the original data
2. Train a new model (meta classifier/ level-1) based on the outputs of level-0 models.

5

Stacking

1. Train different models on the same data ("Level-0 models")
2. Train a new ("Level-1") model with the outputs of the Level-0 models

2. Given a new unlabeled data instance, input it to each Level-0 model:



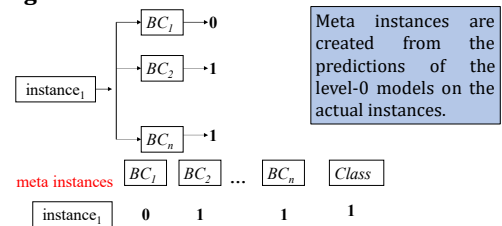
6

Stacking

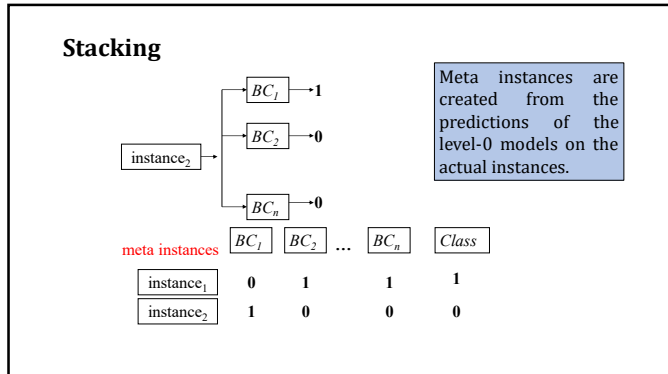
Algorithm	Stacking
1:	Input: training data $D = \{x_i, y_i\}_{i=1}^m$
2:	Output: ensemble classifier H
3:	Step 1: learn base-level classifiers
4:	for $t = 1$ to T do
5:	learn h_t based on D
6:	end for
7:	Step 2: construct new data set of predictions
8:	for $i = 1$ to m do
9:	$D_h = \{x'_i, y_i\}$, where $x'_i = \{h_1(x_i), \dots, h_T(x_i)\}$
10:	end for
11:	Step 3: learn a meta-classifier
12:	learn H based on D_h
13:	return H

7

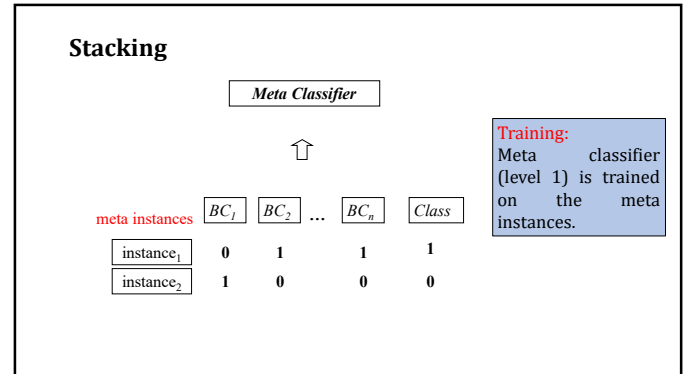
Stacking



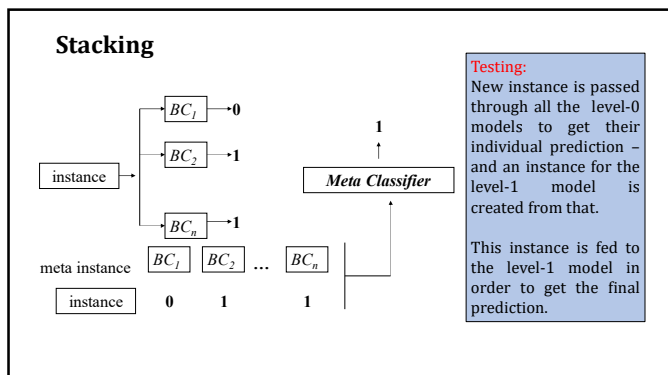
8



9



10



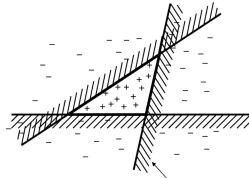
11

Boosting Intuition

- Train a number of weak classifiers (e.g. decision trees) **in a sequence**.
- We **adaptively weigh each data** case.
- Data cases which are **wrongly classified** get **high weight**.
- A new classifier should focus on those cases which were incorrectly classified in the last round.
- Each boosting round learns a new (**simple/weak**) classifier on the **weighted dataset**.
- These **weak** classifiers are **weighted** to **combine** them into a single **powerful/strong classifier**. Combine the classifiers by letting them vote on the final prediction (like bagging).
- Classifiers that obtain **low training error rate** have **high weight**.
- We **stop** by using monitoring a **hold out set** (**cross-validation**).

12

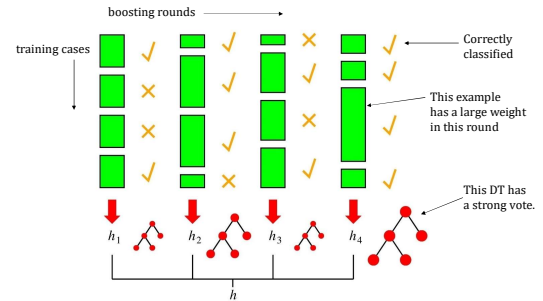
Example



This line is one simple classifier saying that everything to the left is + and everything to the right is -

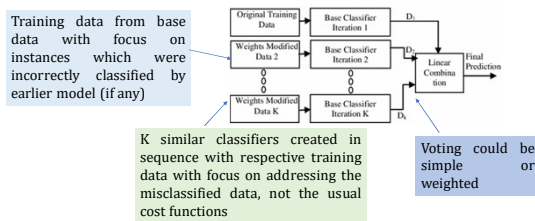
13

Boosting Process



14

Adaboost



It is called **Adaptive Boosting** as the **weights are re-assigned** to each instance, with higher weights to incorrectly classified instance.

15

Adaboost

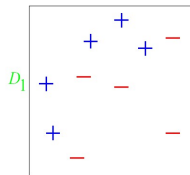
Algorithm AdaBoost

- 1: Init data weights $\{w_n\}$ to $1/N$
- 2: for $m = 1$ to M do
- 3: fit a classifier $y_m(x)$ by minimizing weighted error function J_m :
- 4: $J_m = \sum_{n=1}^N w_n^{(m)} 1[y_m(x_n) \neq t_n]$
- 5: compute $\epsilon_m = \sum_{n=1}^N w_n^{(m)} 1[y_m(x_n) \neq t_n] / \sum_{n=1}^N w_n^{(m)}$
- 6: evaluate $\alpha_m = \log \left(\frac{1-\epsilon_m}{\epsilon_m} \right)$
- 7: update the data weights: $w_n^{(m+1)} = w_n^{(m)} \exp\{\alpha_m 1[y_m(x_n) \neq t_n]\}$
- 8: end for
- 9: Make predictions using the final model: $Y_M(x) = \text{sign} \left(\sum_{m=1}^M \alpha_m y_m(x) \right)$

AdaBoost, short for *Adaptive Boosting*, is a machine learning meta-algorithm formulated by [Yoav Freund](#) and [Robert Schapire](#), who won the 2003 [Gödel Prize](#) for their work.

16

Adaboost

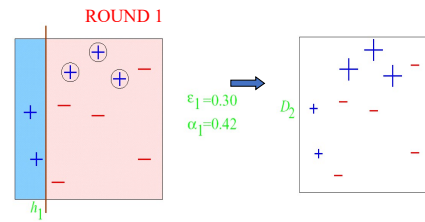


Original training set: equal weights to all training samples

17

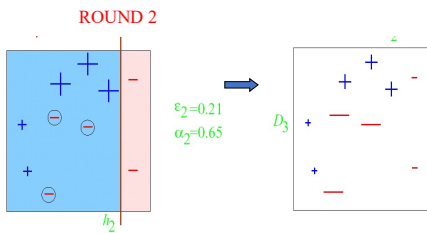
Adaboost

ϵ = error rate of classifier
 α = weight of classifier



18

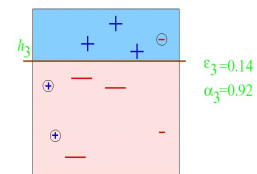
Adaboost



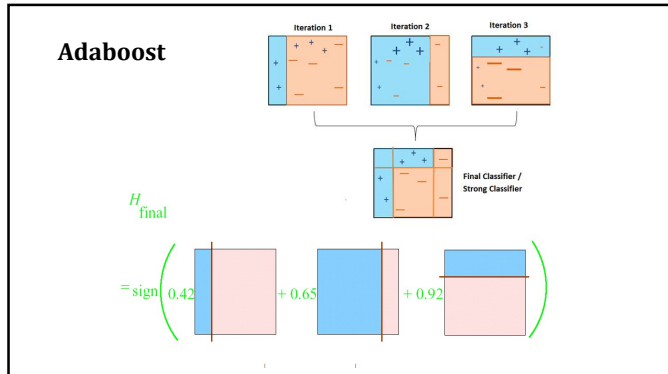
19

Adaboost

ROUND 3



20



21

How to Train each Classifier?

input: \mathbf{x} , output: $y(\mathbf{x}) \in \{1, -1\}$
 target: $t \in \{1, -1\}$,
 weight on case n for classifier m : w_n^m
 Cost function for classifier m :

$$J_m = \sum_{n=1}^N w_n^m [y_m(\mathbf{x}_n) \neq t_n] = \sum \text{weighted errors}$$

1 if error,
0 if correct

22

How to weigh each training case for Classifier m ?

Let $\varepsilon_m = \frac{J_m}{\sum_n w_n^m}$ ← weighted error rate of classifier

Let $\alpha_m = \ln \left(\frac{1 - \varepsilon_m}{\varepsilon_m} \right)$ ← This is the quality of the classifier. It is zero if the classifier has weighted error rate of 0.5 and infinity if the classifier is perfect

$$w_n^{m+1} = w_n^m \exp \{ \alpha_m [y_m(\mathbf{x}_n) \neq t_n] \}$$

23

How to weigh each training case for Classifier m ?

- Weight the binary prediction of each classifier by the quality of that classifier:

$$Y_M(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right)$$

24

Adaboost Algorithm

- Initialize the data weights $w_n = 1/N$.
- For $m=1, \dots, M$:
 - Fit a classifier $y_m(\mathbf{x})$ to the training data by minimizing the weighted error function:

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n),$$

where $I(y_m(\mathbf{x}_n) \neq t_n)$ is the **indicator function** and equals to one when $y_m(\mathbf{x}_n) \neq t_n$ and zero otherwise.

- Evaluate:

$$\alpha_m = \ln \frac{1 - \epsilon_m}{\epsilon_m}, \quad \epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}.$$

Weighting coefficients.

weighted measures of the error rates.

Adaboost Algorithm

- Initialize the data weights $w_n = 1/N$.
- For $m=1, \dots, M$:
 - Fit a classifier $y_m(\mathbf{x})$ to the training data by minimizing:

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n),$$

- Evaluate:

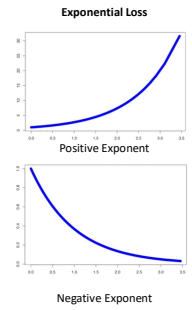
$$\alpha_m = \ln \frac{1 - \epsilon_m}{\epsilon_m}, \quad \epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}.$$

- Update the data weights:

$$w_n^{(m+1)} = w_n^{(m)} \exp(\alpha_m I(y_m(\mathbf{x}_n) \neq t_n)).$$

- Make predictions using the final model:

$$Y_M(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right).$$



25

26



Case-study

27

28

Classification of Credit Risk - Predicting Possible Credit Defaulter based on German Credit Dataset

Context

The original dataset contains 1000 entries with 20 categorial/symbolic attributes prepared by Prof. Hofmann.

In this dataset, each entry represents a person who takes a credit by a bank. Each person is classified as good or bad credit risks according to the set of attributes.

Description of the Attributes

- > Age (numeric)
 - > Sex (text: male, female)
 - > Job (numeric: 0 - unskilled and non-resident, 1 - unskilled and resident, 2 - skilled, 3 - highly skilled)
 - > Housing (text: own, rent, or free)
 - > Saving accounts (text - little, moderate, quite rich, rich)
 - > Checking account (numeric, in DM - Deutsch Mark)
 - > Credit amount (numeric, in DM)
 - > Duration (numeric, in month)
 - > Purpose (text: car, furniture/equipment, radio/TV, domestic appliances, repairs, education, business, vacation/others) etc....
- Total of 17 features (16 input features and 1 target variable ("default" column))

29

Let's go to the Coding Demo...

30

To be continued in the next session.....

31