

# Multiple Models

João Gama

LIAAD-INESC Porto, University of Porto, Portugal

[jgama@fep.up.pt](mailto:jgama@fep.up.pt)

November 2020

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography













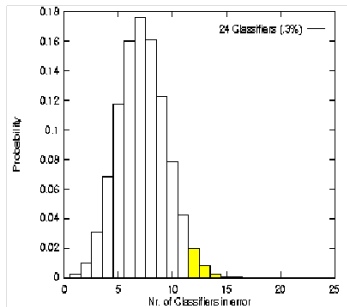








# Multiple Models: Simulation



Assume an ensemble of 23 classifiers:

- probability of error of each classifier: 30%;
- aggregation by uniform vote.

Given a test example:

- the ensemble will be in error iif 12 or more classifiers are in error.
- The probability of error in the ensemble is given by the area under the curve of a binomial distribution;
- In this case this area is 0.026.
- Much less than each individual classifier

# Necessary Conditions

*To achieve higher accuracy the models should be diverse and each model must be quite accurate Ali & Pazzani 96*

## Necessary Conditions

Classifiers in the ensemble, should have:

- performance better than random guess;
- non-correlated errors;
- errors in different regions of the instance space.

# Multiple Models

- Combining Outputs
  - Voting Methods
  - Fusion of Classifiers
  - Model Applicability
- Perturbing the set of training examples
  - Homogeneous Classifiers
    - Bagging
    - Boosting
  - Heterogeneous Classifiers
    - Cascading
    - Stacking
- Perturbing the set of attributes
- Perturbing test examples

# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography

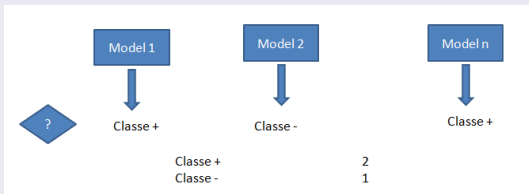




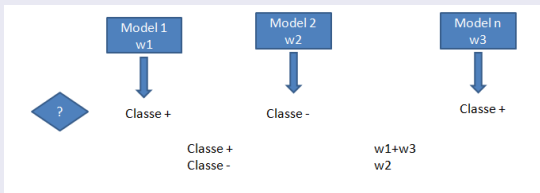


# Combining Outputs: voting

## Uniform Voting



## Weighted Voting



# Pos and Cons of Voting Methods:

## Advantages

- Simplicity;
- Applicable everywhere.

## Disadvantages

- Does not take into account the example to classify;
- Does not make classifiers selection.



# Fusion of Classifiers: Agregation Functions

*J. Kittler; Combining Classifiers: A theoretical framework, Pattern Analysis and Applications, Vol. 1, No. 1,*

Problem: Fusion of  $m$  probabilistic classifiers in a problem with  $j$  classes

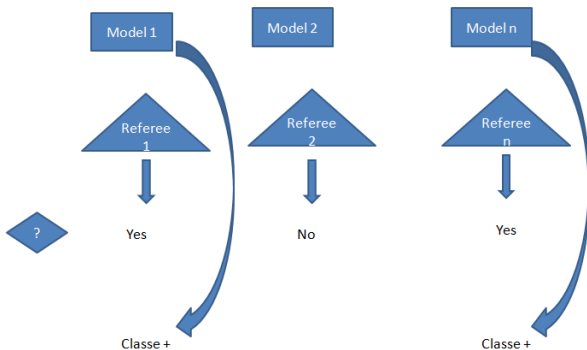
- Sum rule:  $S_j = \sum_{k=1}^m P_{kj}$
- Mean rule:  $S_j = \sum_{k=1}^m \frac{P_{kj}}{m}$
- Geometric mean rule:  $S_j = \sqrt[m]{\prod_{k=1}^m P_{kj}}$
- Product rule:  $S_j = \prod_{k=1}^m P_{kj}$
- Maximum rule:  $S_j = \max_k P_{kj}$
- Minimum rule:  $S_j = \min_k P_{kj}$

Classify the example in the class that maximizes  $S_j$ .





# Model Applicability Induction





# Model Applicability Induction

## Learning the meta classifier

- Input: Base Classifier  $\phi$ , Training Set  $T_0$ , Meta Classifier  $\Phi$ .
- Output: A meta-model for  $\phi$

## Algorithm

- Let  $T_1 = \{\}$
- For each example  $\{\vec{x}, y\}$  of the training set  $T_0$ 
  - Let  $T' = T_0 - \{\vec{x}, y\}$
  - Learn a model  $M = \phi(T')$
  - $\hat{y} = M(\vec{x})$
  - If  $(\hat{y} == y)$  Then  $T_1 = T_1 \cup \{\vec{x}, +\}$  Else  $T_1 = T_1 \cup \{\vec{x}, -\}$
- Output  $\Phi(T_1)$





# MAI: Example

```
sunny, 85, 85, false, +
sunny, 80, 90, true, +
overcast, 83, 78, false, +
rain, 70, 96, false, -
rain, 68, 80, false, +
rain, 65, 70, true, -
overcast, 64, 65, true, +
sunny, 72, 95, false, -
sunny, 69, 70, false, +
rain, 75, 80, false, -
sunny, 75, 70, true, +
overcast, 72, 90, true, +
overcast, 81, 75, false, -
rain, 71, 80, true, +
```

Decision Tree:

```
windy = true: + (6.0/1.0)
```

```
windy = false:
```

```
| humidity <= 90 : + (6.0/2.0)
```

```
| humidity > 90 : - (2.0)
```

## Meta Data ( $T_1$ )

- Positive Examples: those correctly predicted by the base classifier.
- Negative Examples: those wrongly predicted by the base classifier.

Meta Model in the form of a decision tree (can be any other classifier)



# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples**
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography

# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples**
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography





# Bootstrap Aggregation - Bagging

- Learning:
  - Obtain  $N$  replicas of the training set, with reposition;
  - All the samples with the same number of examples of the training set;
  - Learn a classifier for each sample.
- Testing
  - For each test example;
  - All classifiers classify the test example;
  - Predictions are aggregated by uniform vote.











# Random Forests

Breiman, *Random Forests*, MLJ 2001;

## A variant of Bagging;

- Repeat  $k$  times
  - Training set = Draw with replacement  $N$  examples;
  - Built a decision tree
    - Choose (without replacement)  $i$  features
    - Choose best of these  $i$  as the root of this (sub)tree
  - Do NOT prune

where  $N$  is the nr. of examples,  $F$  nr. of features, and  $i$  some number  $\ll F$ .

# Boosting

- Can a set of *weak learners* create a single *strong learner*?
- A weak learner is defined to be a classifier which is only slightly correlated with the true classification.
- A strong learner is a classifier that is arbitrarily well-correlated with the true classification.

Rob Schapire, *Strength of Weak Learnability* Journal of Machine Learning Vol. 5, pages 197-227. 1990





# Boosting

## Characteristics

- Boosting is an iterative algorithm;
- Associates a weight with each example;
- The weight indicates the probability of the example being select in a uniform sampling;









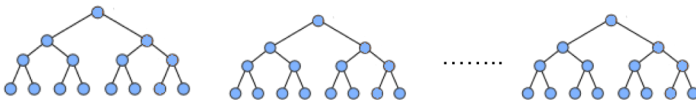




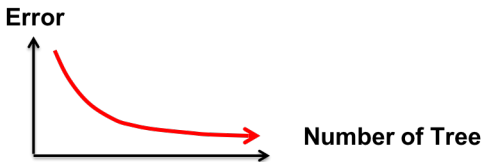


# XGBoost - Extreme Gradient Boosting Tree

- Additive tree model: add new trees that complement the already-built ones



## Greedy Algorithm



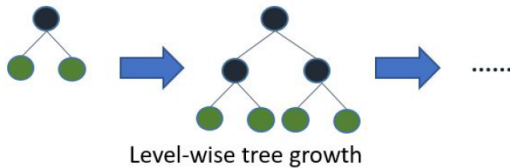




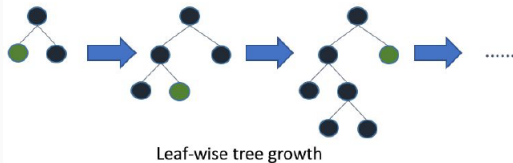


# XGBoost vs lightGBM

XGBoost:



LightGBM:





# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples**
  - Homogeneous Classifiers
  - **Heterogeneous Classifiers**
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography

# Stacking Generalization

Wolpert, *Stacking Generalization*, Neural Networks, Nr. 5, 1992

## Layered Learning

The output of an ensemble of trained classifiers is used as input to the next-layer of classifiers.

## Stacked Generalization with 2 layers

- $Layer_0$

**Data** : is original training set;

**Models** : classifiers trained from the  $layer_0$  data;

- $Layer_1$

**Data** : the predictions of  $layer_0$  classifiers on  $layer_0$  data using cross-validation;

**Models** : classifier trained from the  $layer_1$  data;

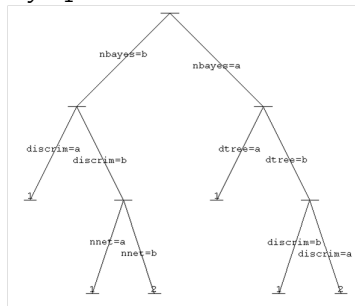




# Stacking Generalization: Example

Base models: naive Bayes, neural net, decision tree, linear discriminant (LDA);

*layer*<sub>1</sub> model: decision tree



*layer*<sub>1</sub> model: LDA

```
> lda(observed~.,df)
Call:
lda.formula(observed ~ ., data = df)

Prior probabilities of groups:
  1  2
0.51 0.49

Group means:
  discrim2  nbayes2  nnet2  dtree2
1 0.4509804 0.5686275 0.5098039 0.3725490
2 0.4285714 0.4693878 0.4489796 0.4693878

Coefficients of linear discriminants:
      LD1
discrim2 -0.3424098
nbayes2  -1.3950698
nnet2    -1.0185849
dtree2    1.0547212
```





# Cascade Generalization

Sequential composition of a naive-Bayes and a Decision Tree:

Dataset Original

```
3,4,3,4,B
4,1,4,1,B
4,2,2,1,L
5,2,5,3,R
2,5,4,4,R
2,3,4,3,R
5,1,4,5,R
4,3,2,5,L
3,3,2,5,R
1,3,4,5,R
```

Dataset Extendido

```
3,4,3,4,0.461183,0.077635,0.461183,B
4,1,4,1,0.413818,0.172365,0.413818,B
4,2,2,1,0.838750,0.089446,0.071804,L
5,2,5,3,0.307441,0.089143,0.603416,R
2,5,4,4,0.283686,0.104362,0.611952,R
2,3,4,3,0.213796,0.070892,0.715312,R
5,1,4,5,0.072916,0.075340,0.851744,R
4,3,2,5,0.505602,0.094848,0.399550,L
3,3,2,5,0.391624,0.080813,0.527563,R
1,3,4,5,0.030005,0.043305,0.926691,R
```

Arvore de Decisao  
(dataset extendido)

```
File stem <balnew>
Read 625 cases (7 attributes) from balnew.data
Decision Tree:
p3 > 0.471812 : R (288.0)
p3 <= 0.471812 :
| p1 <= 0.471812 : B (49.0)
| p1 > 0.471812 : L (288.0)
Tree saved
```

# Cascade Generalization in KNIME

The diagram illustrates a KNIME workflow for Cascade Generalization. It starts with a **File Reader** (Node 1) feeding into a **Naive Bayes Learner** (Node 10) and a **Naive Bayes Predictor** (Node 11). The output of the Naive Bayes Predictor is then processed by a **Column Filter** (Node 12), which is configured to exclude the 'Winner(Naive Bayes)' column. The filtered data is then fed into a **Decision Tree Learner** (Node 13).

The **Column Filter** dialog (Node 12) shows the following configuration:

- Exclude:** Winner(Naive Bayes)
- Include:** Col0, Col1, Col2, Col3, Col4, setosa, versicolor, virginica

The **Decision Tree View** (Node 13) displays the resulting decision tree structure:

- [root]: class 'virginica' (50 of 150)
  - [virginica <= 0]: class 'setosa' (50 of 50)
    - [setosa <= 0]: class 'virginica' (49 of 52)
    - [setosa > 0]: class 'versicolor' (47 of 48)
  - [virginica > 0]: class 'virginica' (50 of 100)



# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set**
- 5 Perturbing Test Examples
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography







# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples**
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography

# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples**
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography

# Dual Perturb & Combine

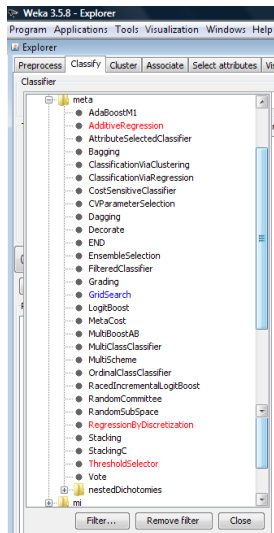
Approaches that use only a single model and delays at the prediction stage the generation of multiple predictions by perturbing the attribute vector corresponding to a test case.

# Dual Perturb & Combine

Geurts & Wehenkel *Closed-form dual perturb and combine for tree-based models*. In Proc. of the 22nd international Conference on Machine Learning, 2005

- Only a single model is generated from the training set.
- In the prediction phase, each test example is perturbed several times.
  - To perturb a test example, white noise is added to the attribute-values.
  - The predictive model makes a prediction for each perturbed version of the test example.
  - The final prediction is obtained by aggregating the different predictions.
- Geurts presents evidence that this method is efficient in variance reduction.

# Multiple Models in Weka



# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography



# Summary

Well designed ensembles of classifiers allow improve performance over their individual elements.

## Necessary Conditions

- Variability between elements;
- Low Error correlation;
- Each individual classifier must be better than a random choice.

# Summary

## General Methods

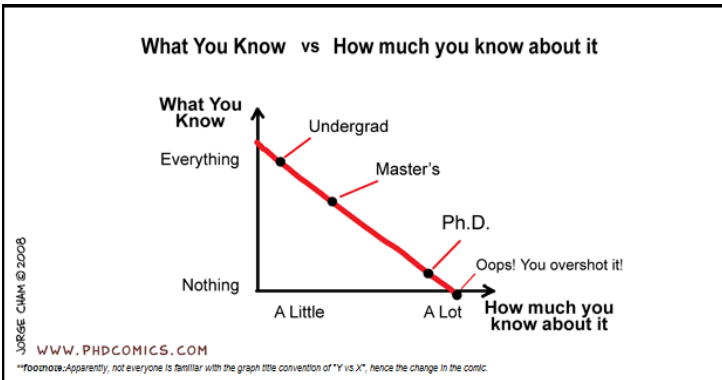
- Voting Methods;
- Fusion of Classifiers (probabilistic classifiers)
- Perturbing the training examples:
  - Bootstrap Aggregation (Bagging)
  - Adaptive Boosting (AdaBoosting)
- Perturbing the set of attributes
- Perturbing the test examples
- Using different classifiers:
  - Stacking Generalization
  - Cascade Generalization

# Outline

- 1 Motivation
- 2 Combining Outputs
- 3 Perturbing Training Examples
  - Homogeneous Classifiers
  - Heterogeneous Classifiers
- 4 Perturbing the attribute set
- 5 Perturbing Test Examples
  - Dual Perturb & Combine
- 6 Summary
- 7 Bibliography

# Bibliography

- Ali and Pazzani, *Error Reduction through learning multiple descriptions*, Machine Learning, 23, 1996
- Breiman, L. Stacking Predictors, Machine Learning, 25, 1997
- Breiman, L. Bagging Predictors, Machine Learning, 24, 1997
- Bauer, Kohavi An empirical comparison of voting classification algorithms: Bagging, Boosting and Variants, Machine Learning, 36, 1999
- Dietterich, T., Machine Learning Research-Four current directions, AI Magazine, 98
- Freund, Y. and Schapire Experiments with a new boosting algorithm, ICML96
- Gama, J. Combining Classifiers by Constructive Induction, ECML98
- Gama, J. Cascade Generalization, Machine Learning, 2000
- Kohavi, R., Kunz, C., Option Trees with majority votes, ICML97
- Quinlan, R., Bagging, Boosting and C4.5, AAAI96
- Ting, K.; Witten, I. Stacked Generalization: when it works, IJCAI, 1997
- Wolpert, D. Stacked Generalization, Neural Networks, N.5



Would you like to learn more? Wait for ECDII ...