

Introduction to Machine Learning

Combining Models, Bayesian Average and Boosting

Andres Mendez-Vazquez

July 31, 2018

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

Outline

1 Combining Models

- Introduction
- Average for Committee
- Beyond Simple Averaging
- Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
- The Differences

3 Committees

- Introduction
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

Introduction

Observation

- It is often found that improved performance can be obtained by combining multiple classifiers together in some way.

Introduction

Observation

- It is often found that improved performance can be obtained by combining multiple classifiers together in some way.

Example, Committees

- We might train L different classifiers and then make predictions:
 - ▶ by using the average of the predictions made by each classifier.

Introduction

Observation

- It is often found that improved performance can be obtained by combining multiple classifiers together in some way.

Example, Committees

- We might train L different classifiers and then make predictions:
 - ▶ by using the average of the predictions made by each classifier.

Example, Boosting

- It involves training multiple models in sequence:
 - ▶ A error function used to train a particular model depends on the performance of the previous models.

Introduction

Observation

- It is often found that improved performance can be obtained by combining multiple classifiers together in some way.

Example, Committees

- We might train L different classifiers and then make predictions:
 - ▶ by using the average of the predictions made by each classifier.

Example, Boosting

- It involves training multiple models in sequence:
 - ▶ A error function used to train a particular model depends on the performance of the previous models

Introduction

Observation

- It is often found that improved performance can be obtained by combining multiple classifiers together in some way.

Example, Committees

- We might train L different classifiers and then make predictions:
 - ▶ by using the average of the predictions made by each classifier.

Example, Boosting

- It involves training multiple models in sequence:
 - ▶ A error function used to train a particular model depends on the performance of the previous models.

Outline

1 Combining Models

- Introduction
- **Average for Committee**
- Beyond Simple Averaging
- Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
- The Differences

3 Committees

- Introduction
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

We could use simple averaging

Given a series of observed samples $\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_N\}$ with noise $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \epsilon$$

We can use our knowledge of probability to remove such noise

$$E[\hat{\mathbf{x}}_i] = E[\mathbf{x}_i + \epsilon] = E[\mathbf{x}_i] + E[\epsilon]$$

Then, because $E[\epsilon] = 0$

$$E[\mathbf{x}_i] = E[\hat{\mathbf{x}}_i] \approx \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{x}}_i$$

We could use simple averaging

Given a series of observed samples $\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_N\}$ with noise $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \epsilon$$

We can use our knowledge of probability to remove such noise

$$E[\hat{\mathbf{x}}_i] = E[\mathbf{x}_i + \epsilon] = E[\mathbf{x}_i] + E[\epsilon]$$

Then, because $E[\epsilon] = 0$

$$E[\mathbf{x}_i] = E[\hat{\mathbf{x}}_i] \approx \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{x}}_i$$

We could use simple averaging

Given a series of observed samples $\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_N\}$ with noise $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \epsilon$$

We can use our knowledge of probability to remove such noise

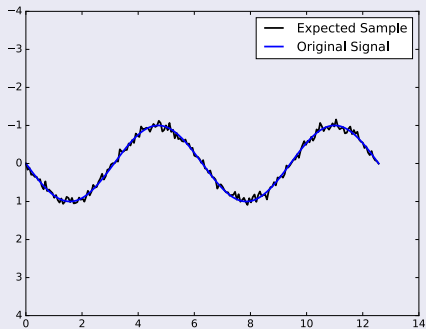
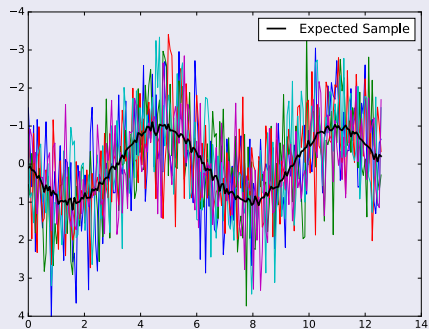
$$E[\hat{\mathbf{x}}_i] = E[\mathbf{x}_i + \epsilon] = E[\mathbf{x}_i] + E[\epsilon]$$

Then, because $E[\epsilon] = 0$

$$E[\mathbf{x}_i] = E[\hat{\mathbf{x}}_i] \approx \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{x}}_i$$

For Example

We have a nice result



Outline

1 Combining Models

- Introduction
- Average for Committee
- **Beyond Simple Averaging**
- Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
- The Differences

3 Committees

- Introduction
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

Beyond Simple Averaging

Instead of averaging the predictions of a set of models

- You can use an alternative form of combination that selects one of the models to make the prediction.

Where

- The choice of model is a function of the input variables.

How

- Different Models become responsible for making decisions in different regions of the input space.

Beyond Simple Averaging

Instead of averaging the predictions of a set of models

- You can use an alternative form of combination that selects one of the models to make the prediction.

Where

- The choice of model is a function of the input variables.

How

- Different Models become responsible for making decisions in different regions of the input space.

Beyond Simple Averaging

Instead of averaging the predictions of a set of models

- You can use an alternative form of combination that selects one of the models to make the prediction.

Where

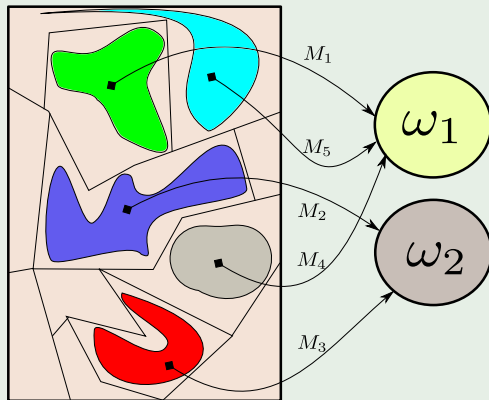
- The choice of model is a function of the input variables.

Thus

- Different Models become responsible for making decisions in different regions of the input space.

Something like this

Models in charge of different set of inputs



Outline

1 Combining Models

- Introduction
- Average for Committee
- **Beyond Simple Averaging**
 - **Example**

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
 - The Differences

3 Committees

- Introduction
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

Example, Decision Trees

We can have the decision trees on top of the models

Given a set of models, a model is chosen to take a decision in certain area of the input.

Limitation: It is based on hard splits in which only one model is responsible for making predictions for any given value.

Example, Decision Trees

We can have the decision trees on top of the models

Given a set of models, a model is chosen to take a decision in certain area of the input.

Limitation: It is based on hard splits in which only one model is responsible for making predictions for any given value.

Thus it is hard to utilize the combination by using

- If we have M classifier for a conditional distribution $p(t|x, k)$.
 - ▶ x is the input variable.
 - ▶ t is the target variable.
 - ▶ $k = 1, 2, \dots, M$ indexes the classifiers.

Example, Decision Trees

We can have the decision trees on top of the models

Given a set of models, a model is chosen to take a decision in certain area of the input.

Limitation: It is based on hard splits in which only one model is responsible for making predictions for any given value.

Thus it is better to soften the combination by using

- If we have M classifier for a conditional distribution $p(t|\mathbf{x}, k)$.

- ▶ \mathbf{x} is the input variable.
- ▶ t is the target variable.
- ▶ $k = 1, 2, \dots, M$ indexes the classifiers.

Example, Decision Trees

We can have the decision trees on top of the models

Given a set of models, a model is chosen to take a decision in certain area of the input.

Limitation: It is based on hard splits in which only one model is responsible for making predictions for any given value.

Thus it is better to soften the combination by using

- If we have M classifier for a conditional distribution $p(t|\mathbf{x}, k)$.
 - ▶ \mathbf{x} is the input variable.
 - ▶ t is the target variable.
 - ▶ $k = 1, 2, \dots, M$ indexes the classifiers.

Example, Decision Trees

We can have the decision trees on top of the models

Given a set of models, a model is chosen to take a decision in certain area of the input.

Limitation: It is based on hard splits in which only one model is responsible for making predictions for any given value.

Thus it is better to soften the combination by using

- If we have M classifier for a conditional distribution $p(t|\mathbf{x}, k)$.
 - ▶ \mathbf{x} is the input variable.
 - ▶ t is the target variable.
 - ▶ $k = 1, 2, \dots, M$ indexes the classifiers.

Example, Decision Trees

We can have the decision trees on top of the models

Given a set of models, a model is chosen to take a decision in certain area of the input.

Limitation: It is based on hard splits in which only one model is responsible for making predictions for any given value.

Thus it is better to soften the combination by using

- If we have M classifier for a conditional distribution $p(t|\mathbf{x}, k)$.
 - ▶ \mathbf{x} is the input variable.
 - ▶ t is the target variable.
 - ▶ $k = 1, 2, \dots, M$ indexes the classifiers.

Example, Decision Trees

We can have the decision trees on top of the models

Given a set of models, a model is chosen to take a decision in certain area of the input.

Limitation: It is based on hard splits in which only one model is responsible for making predictions for any given value.

Thus it is better to soften the combination by using

- If we have M classifier for a conditional distribution $p(t|\mathbf{x}, k)$.
 - ▶ \mathbf{x} is the input variable.
 - ▶ t is the target variable.
 - ▶ $k = 1, 2, \dots, M$ indexes the classifiers.

This is used in the mixture of distributions

Thus (Mixture of Experts)

$$p(t|\mathbf{x}) = \sum_{k=1}^M \pi_k(\mathbf{x}) p(t|\mathbf{x}, k) \quad (1)$$

where $\pi_k(\mathbf{x}) = p(k|\mathbf{x})$ represent the input-dependent mixing coefficients.

This type of models

They can be viewed as mixture distribution in which the component densities and the mixing coefficients are conditioned on the input variables and are known as mixture experts.

This is used in the mixture of distributions

Thus (Mixture of Experts)

$$p(t|\mathbf{x}) = \sum_{k=1}^M \pi_k(\mathbf{x}) p(t|\mathbf{x}, k) \quad (1)$$

where $\pi_k(\mathbf{x}) = p(k|\mathbf{x})$ represent the input-dependent mixing coefficients.

This type of models

They can be viewed as mixture distribution in which the component densities and the mixing coefficients are conditioned on the input variables and are known as mixture experts.

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

It is important to differentiate between them

Although

- Model Combinations and Bayesian Model Averaging look similar.
 - ▶ However, they are actually different

For this

We have the following example.

It is important to differentiate between them

Although

- Model Combinations and Bayesian Model Averaging look similar.
 - ▶ However, they are actually different

For this

We have the following example.

Example of the Differences

For this consider the following

- Mixture of Gaussians with a binary latent variable z indicating to which component a point belongs to.

This the model is specified in terms a joint distribution

$$p(\mathbf{x}, z)$$

Corresponding density over the observed variable x using marginalization

$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, z)$$

Example of the Differences

For this consider the following

- Mixture of Gaussians with a binary latent variable z indicating to which component a point belongs to.

Thus the model is specified in terms a joint distribution

$$p(\mathbf{x}, z)$$

Corresponding density over the observed variable \mathbf{x} using marginalization

$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, z)$$

Example of the Differences

For this consider the following

- Mixture of Gaussians with a binary latent variable z indicating to which component a point belongs to.

Thus the model is specified in terms a joint distribution

$$p(\mathbf{x}, z)$$

Corresponding density over the observed variable \mathbf{x} using marginalization

$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, z)$$

Example

In the case of Mixture of Gaussian's

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

This is an example of model combination

- What about other Models

Example

In the case of Mixture of Gaussian's

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

This is an example of model combination.

- What about other Models

More Models

Now, for independent, identically distributed data

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

$$p(\mathbf{X}) = \prod_{n=1}^N p(\mathbf{x}_n) = \prod_{n=1}^N \left[\sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n) \right]$$

Therefore

Something Notable

- Each observed data point x_n has a corresponding latent variable z_n .

Here, we are doing a Combination of Worlds

- Each Gaussian indexed by z_n is in charge of generating one section of the sample space

Therefore

Something Notable

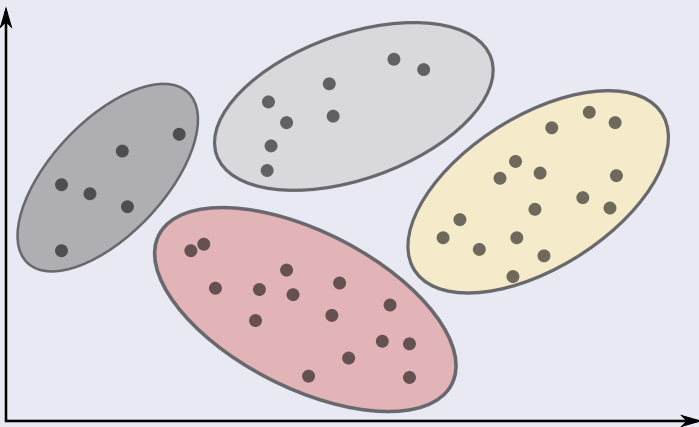
- Each observed data point x_n has a corresponding latent variable z_n .

Here, we are doing a Combination of Models

- Each Gaussian indexed by z_n is in charge of generating one section of the sample space

Example

We have



Outline

1 Combining Models

- Introduction
- Average for Committee
- Beyond Simple Averaging
 - Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- **Now Model Averaging**
 - The Differences

3 Committees

- Introduction
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

Now, suppose

We have several different models indexed by $h = 1, \dots, H$ with prior probabilities

- One model might be a mixture of Gaussians and another model might be a mixture of Cauchy distributions

The Marginal Distribution is

$$p(X) = \sum_{h=1}^H p(X, h) = \sum_{h=1}^H \underbrace{p(X|h) p(h)}_{\approx p(h|X)}$$

- This is an example of Bayesian model averaging

Now, suppose

We have several different models indexed by $h = 1, \dots, H$ with prior probabilities

- One model might be a mixture of Gaussians and another model might be a mixture of Cauchy distributions

The Marginal Distribution is

$$p(X) = \sum_{h=1}^H p(X, h) = \sum_{h=1}^H \underbrace{p(X|h) p(h)}_{\approx p(h|X)}$$

- This is an example of **Bayesian model averaging**

Bayesian Model Averaging

Remark

- The summation over h means that just one model is responsible for generating the whole data set.

Observation

- The probability over h simply reflects our uncertainty of which is the correct model to use.

Bayesian Model Averaging

Remark

- The summation over h means that just one model is responsible for generating the whole data set.

Observation

- The probability over h simply reflects our uncertainty of which is the correct model to use.

This is the sum of all the probabilities

- This uncertainty reduces
 - ▶ Posterior probabilities $p(h|X)$ become increasingly focused on just one of the models.

Bayesian Model Averaging

Remark

- The summation over h means that just one model is responsible for generating the whole data set.

Observation

- The probability over h simply reflects our uncertainty of which is the correct model to use.

Thus, as the size of the data set increases

- This uncertainty reduces
 - ▶ Posterior probabilities $p(h|X)$ become increasingly focused on just one of the models.

Example

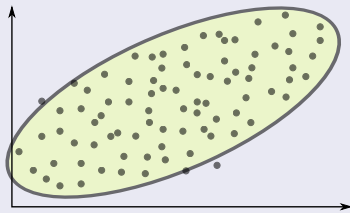
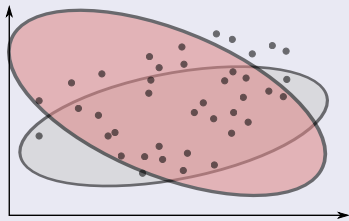
We have

● h_1

○ h_3

● h_2

INCREASING THE
NUMBER OF SAMPLES



Outline

1

Combining Models

- Introduction
- Average for Committee
- Beyond Simple Averaging
- Example

2

Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- **Now Model Averaging**
- **The Differences**

3

Committees

- Introduction
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4

Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

The Differences

Bayesian model averaging

- The whole data set is generated by a single model h .

Model combination

- Different data points within the data set can potentially be generated from different by different components.

The Differences

Bayesian model averaging

- The whole data set is generated by a single model h .

Model combination

- Different data points within the data set can potentially be generated from different by different components.

Outline

1 Combining Models

- Introduction
- Average for Committee
- Beyond Simple Averaging
 - Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
 - The Differences

3 Committees

- **Introduction**
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

Committees

Idea, the simplest way to construct a committee

- It is to average the predictions of a set of individual models.

Committees

Idea, the simplest way to construct a committee

- It is to average the predictions of a set of individual models.

Thinking as a frequentist

- This is coming from taking in consideration the trade-off between bias and variance.

Committees

Idea, the simplest way to construct a committee

- It is to average the predictions of a set of individual models.

Thinking as a frequentist

- This is coming from taking in consideration the trade-off between bias and variance.

Where the error in the model into

- The bias component that arises from differences **between the model and the true function to be predicted.**
- The variance component that represents the sensitivity of the model to the individual data points.

Committees

Idea, the simplest way to construct a committee

- It is to average the predictions of a set of individual models.

Thinking as a frequentist

- This is coming from taking in consideration the trade-off between bias and variance.

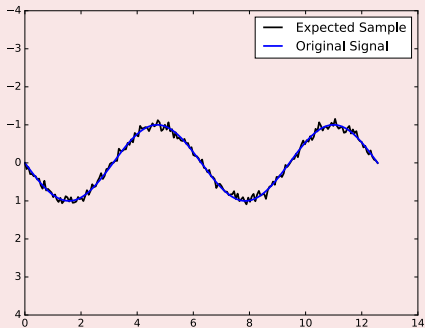
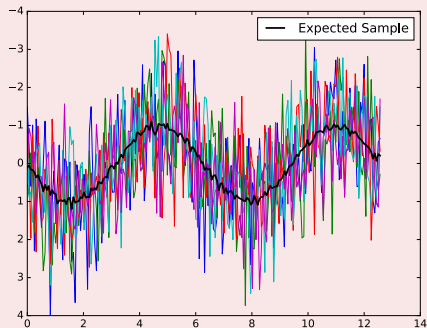
Where the error in the model into

- The bias component that arises from differences **between the model and the true function to be predicted.**
- The variance component that represents **the sensitivity of the model to the individual data points.**

For example

When we averaged a set of low-bias models

- We obtained accurate predictions of the underlying sinusoidal function from which the data were generated.



However

Big Problem

- We have normally a single data set

This

- We need to introduce certain variability between the different committee members.

One approach

- You can use bootstrap data sets.

However

Big Problem

- We have normally a single data set

Thus

- We need to introduce certain variability between the different committee members.

One approach

- You can use bootstrap data sets.

Outline

1 Combining Models

- Introduction
- Average for Committee
- Beyond Simple Averaging
 - Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
 - The Differences

3 Committees

- Introduction
- **Bootstrap Data Sets**
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

The Idea of Bootstrap

We denote the training set by $Z = \{z_1, z_2, \dots, z_N\}$

- Where $z_i = (\mathbf{x}_i, y_i)$

The basic idea is to randomly draw datasets with replacement from the training data.

- Each sample the same size as the original training set.

This is done B times.

- Producing B bootstrap datasets.

The Idea of Bootstrap

We denote the training set by $Z = \{z_1, z_2, \dots, z_N\}$

- Where $z_i = (\mathbf{x}_i, y_i)$

The basic idea is to randomly draw datasets with replacement from the training data

- Each sample the same size as the original training set.

This is done B times

- Producing B bootstrap datasets.

The Idea of Bootstrap

We denote the training set by $Z = \{z_1, z_2, \dots, z_N\}$

- Where $z_i = (\mathbf{x}_i, y_i)$

The basic idea is to randomly draw datasets with replacement from the training data

- Each sample the same size as the original training set.

This is done B times

- Producing B bootstrap datasets.

Then

Then a quantity is computed

- $S(Z)$ is any quantity computed from the data Z

From the bootstrap sampling

- We can estimate any aspect of the distribution of $S(Z)$.

Then

Then a quantity is computed

- $S(Z)$ is any quantity computed from the data Z

From the bootstrap sampling

- We can estimate any aspect of the distribution of $S(Z)$.

Then

Then a quantity is computed

- $S(Z)$ is any quantity computed from the data Z

From the bootstrap sampling

- We can estimate any aspect of the distribution of $S(Z)$.

Then

we refit the model to each of the bootstrap datasets

- You generate $S(Z^{*b})$ to refit the model to this dataset.

Then

- You examine the behavior of the fits over the B replications.

Then

we refit the model to each of the bootstrap datasets

- You generate $S(Z^{*b})$ to refit the model to this dataset.

Then

- You examine the behavior of the fits over the B replications.

For Example

Its variance

$$\widehat{Var} [S(Z)] = \frac{1}{B-1} \sum_{b=1}^B (S(Z^{*b}) - \bar{S}^*)^2$$

Where

$$\bar{S}^* = \frac{1}{B} \sum_{b=1}^B S(Z^{*b})$$

For Example

Its variance

$$\widehat{Var} [S (Z)] = \frac{1}{B-1} \sum_{b=1}^B (S (Z^{*b}) - \bar{S}^*)^2$$

Where

$$\bar{S}^* = \frac{1}{B} \sum_{b=1}^B S (Z^{*b})$$

Outline

1 Combining Models

- Introduction
- Average for Committee
- Beyond Simple Averaging
 - Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
 - The Differences

3 Committees

- Introduction
- Bootstrap Data Sets
- **Relation with Monte-Carlo Estimation**

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

Relation with Monte-Carlo Estimation

Note that $\widehat{Var}[S(Z)]$

- It can be thought of as a Monte-Carlo estimate of the variance of $S(Z)$ under sampling.

This is similar

- From the empirical distribution function \hat{F} for the data $Z = \{z_1, z_2, \dots, z_N\}$

Relation with Monte-Carlo Estimation

Note that $\widehat{Var}[S(Z)]$

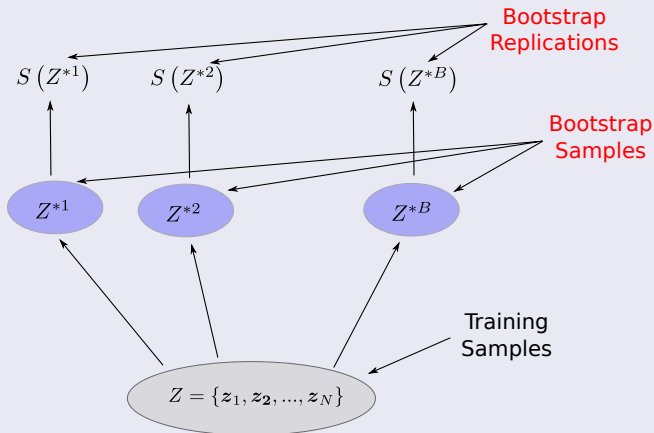
- It can be thought of as a Monte-Carlo estimate of the variance of $S(Z)$ under sampling.

This is coming

- From the empirical distribution function \widehat{F} for the data $Z = \{z_1, z_2, \dots, z_N\}$

For Example

Schematic of the bootstrap process



Thus

Use each of them to train a copy $y_b(\mathbf{x})$ of a predictive regression model to predict a single continuous variable

Then,

$$y_{com}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B y_b(\mathbf{x}) \quad (2)$$

This is also known as **Bootstrap Aggregation or Bagging**.

What do we do with this samples?

Now, assume a true regression function $h(\mathbf{x})$ and a estimation $y_b(\mathbf{x})$

$$y_b(\mathbf{x}) = h(\mathbf{x}) + \epsilon_b(\mathbf{x}) \quad (3)$$

The average squared error over the data takes the form

$$E_{\mathbf{x}} \left[(y_b(\mathbf{x}) - h(\mathbf{x}))^2 \right] = E_{\mathbf{x}} \left[\epsilon_b^2(\mathbf{x}) \right] \quad (4)$$

What is $E_{\mathbf{x}}$?

It denotes a frequentist expectation with respect to the distribution of the input vector \mathbf{x} .

What do we do with this samples?

Now, assume a true regression function $h(\mathbf{x})$ and a estimation $y_b(\mathbf{x})$

$$y_b(\mathbf{x}) = h(\mathbf{x}) + \epsilon_b(\mathbf{x}) \quad (3)$$

The average sum-of-squares error over the data takes the form

$$E_{\mathbf{x}} \left[(y_b(\mathbf{x}) - h(\mathbf{x}))^2 \right] = E_{\mathbf{x}} \left[\epsilon_b^2(\mathbf{x}) \right] \quad (4)$$

What is $E_{\mathbf{x}}$?

It denotes a frequentest expectation with respect to the distribution of the input vector \mathbf{x} .

What do we do with this samples?

Now, assume a true regression function $h(\mathbf{x})$ and a estimation $y_b(\mathbf{x})$

$$y_b(\mathbf{x}) = h(\mathbf{x}) + \epsilon_b(\mathbf{x}) \quad (3)$$

The average sum-of-squares error over the data takes the form

$$E_{\mathbf{x}} \left[(y_b(\mathbf{x}) - h(\mathbf{x}))^2 \right] = E_{\mathbf{x}} \left[\epsilon_b^2(\mathbf{x}) \right] \quad (4)$$

What is $E_{\mathbf{x}}$?

It denotes a frequentest expectation with respect to the distribution of the input vector \mathbf{x} .

Meaning

Thus, the average error is

$$E_{AV} = \frac{1}{B} \sum_{b=1}^B E_x \left[\{\epsilon_b(\mathbf{x})\}^2 \right] \quad (5)$$

Similarly the Expected error over the committee

$$E_{COM} = E_x \left[\left\{ \frac{1}{B} \sum_{b=1}^B (y_m(\mathbf{x}) - h(\mathbf{x})) \right\}^2 \right] = E_x \left[\left\{ \frac{1}{B} \sum_{b=1}^B \epsilon_b(\mathbf{x}) \right\}^2 \right] \quad (6)$$

Meaning

Thus, the average error is

$$E_{AV} = \frac{1}{B} \sum_{b=1}^b E_{\mathbf{x}} \left[\{\epsilon_b(\mathbf{x})\}^2 \right] \quad (5)$$

Similarly the Expected error over the committee

$$E_{COM} = E_{\mathbf{x}} \left[\left\{ \frac{1}{B} \sum_{b=1}^B (y_m(\mathbf{x}) - h(\mathbf{x})) \right\}^2 \right] = E_{\mathbf{x}} \left[\left\{ \frac{1}{B} \sum_{b=1}^B \epsilon_b(\mathbf{x}) \right\}^2 \right] \quad (6)$$

Assume that the errors have zero mean and are uncorrelated

Assume that the errors have zero mean and are uncorrelated

- Something Reasonable to assume given the way we produce the Bootstrap Samples

$$E_{\mathbf{x}} [\epsilon_b(\mathbf{x})] = 0$$

$$E_{\mathbf{x}} [\epsilon_b(\mathbf{x}) \epsilon_l(\mathbf{x})] = 0, \text{ for } b \neq l$$

Then

We have that

$$\begin{aligned} E_{COM} &= \frac{1}{b^2} E_{\mathbf{x}} \left[\left\{ \sum_{b=1}^B (\epsilon_b(\mathbf{x})) \right\}^2 \right] \\ &= \frac{1}{B^2} E_{\mathbf{x}} \left[\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) + \sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right] \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + E_{\mathbf{x}} \left(\sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right) \right\} \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + \sum_{\substack{h=1 \\ h \neq k}}^M \sum_{k=1}^M E_{\mathbf{x}} (\epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x})) \right\} \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} = \frac{1}{B} \left\{ \frac{1}{B} E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} \end{aligned}$$

Then

We have that

$$\begin{aligned} E_{COM} &= \frac{1}{b^2} E_{\mathbf{x}} \left[\left\{ \sum_{b=1}^B (\epsilon_b(\mathbf{x})) \right\}^2 \right] \\ &= \frac{1}{B^2} E_{\mathbf{x}} \left[\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) + \sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right] \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + E_{\mathbf{x}} \left(\sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right) \right\} \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + \sum_{\substack{h=1 \\ h \neq k}}^M \sum_{k=1}^M E_{\mathbf{x}}(\epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x})) \right\} \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} = \frac{1}{B} \left\{ \frac{1}{B} E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} \end{aligned}$$

Then

We have that

$$\begin{aligned} E_{COM} &= \frac{1}{b^2} E_{\mathbf{x}} \left[\left\{ \sum_{b=1}^B (\epsilon_b(\mathbf{x})) \right\}^2 \right] \\ &= \frac{1}{B^2} E_{\mathbf{x}} \left[\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) + \sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right] \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + E_{\mathbf{x}} \left(\sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right) \right\} \end{aligned}$$

$$= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + \sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B E_{\mathbf{x}}(\epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x})) \right\}$$

$$= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} = \frac{1}{B} \left\{ \frac{1}{B} E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\}$$

Then

We have that

$$\begin{aligned} E_{COM} &= \frac{1}{b^2} E_{\mathbf{x}} \left[\left\{ \sum_{b=1}^B (\epsilon_b(\mathbf{x})) \right\}^2 \right] \\ &= \frac{1}{B^2} E_{\mathbf{x}} \left[\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) + \sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right] \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + E_{\mathbf{x}} \left(\sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right) \right\} \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + \sum_{\substack{h=1 \\ h \neq k}}^M \sum_{k=1}^M E_{\mathbf{x}} (\epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x})) \right\} \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} = \frac{1}{B} \left\{ \frac{1}{B} E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} \end{aligned}$$

Then

We have that

$$\begin{aligned} E_{COM} &= \frac{1}{b^2} E_{\mathbf{x}} \left[\left\{ \sum_{b=1}^B (\epsilon_b(\mathbf{x})) \right\}^2 \right] \\ &= \frac{1}{B^2} E_{\mathbf{x}} \left[\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) + \sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right] \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + E_{\mathbf{x}} \left(\sum_{\substack{h=1 \\ h \neq k}}^B \sum_{k=1}^B \epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x}) \right) \right\} \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) + \sum_{\substack{h=1 \\ h \neq k}}^M \sum_{k=1}^M E_{\mathbf{x}} (\epsilon_h(\mathbf{x}) \epsilon_k(\mathbf{x})) \right\} \\ &= \frac{1}{B^2} \left\{ E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} = \frac{1}{B} \left\{ \frac{1}{B} E_{\mathbf{x}} \left(\sum_{b=1}^B \epsilon_b^2(\mathbf{x}) \right) \right\} \end{aligned}$$

We finally obtain

We obtain

$$E_{COM} = \frac{1}{B} E_{AV} \quad (7)$$

Looks great BUT!!

Unfortunately, it depends on the key assumption that the errors at the individual Bootstrap Models are uncorrelated.

We finally obtain

We obtain

$$E_{COM} = \frac{1}{B} E_{AV} \quad (7)$$

Looks great BUT!!!

Unfortunately, it depends on the key assumption that the errors at the individual Bootstrap Models are uncorrelated.

Thus

The Reality!!!

The errors are typically highly correlated, and the reduction in overall error is generally small.

Something Notable

However, it can be shown that the expected committee error will not exceed the expected error of the constituent models, so

$$E_{COM} \leq E_{AV} \quad (8)$$

However, we need something better

A more sophisticated technique known as **boosting**.

Thus

The Reality!!!

The errors are typically highly correlated, and the reduction in overall error is generally small.

Something Notable

However, It can be shown that the expected committee error will not exceed the expected error of the constituent models, so

$$E_{COM} \leq E_{AV} \quad (8)$$

However, we need something better.

A more sophisticated technique known as **boosting**.

Thus

The Reality!!!

The errors are typically highly correlated, and the reduction in overall error is generally small.

Something Notable

However, It can be shown that the expected committee error will not exceed the expected error of the constituent models, so

$$E_{COM} \leq E_{AV} \quad (8)$$

However, we need something better

A more sophisticated technique known as **boosting**.

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 **Boosting**
 - **AdaBoost Development**
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

Boosting

What Boosting does?

It combines several classifiers to produce a form of a committee.

We will describe AdaBoost

“Adaptive Boosting” developed by Freund and Schapire (1995).

Boosting

What Boosting does?

It combines several classifiers to produce a form of a committee.

We will describe AdaBoost

“Adaptive Boosting” developed by Freund and Schapire (1995).

Sequential Training

Main difference between boosting and committee methods

The base classifiers are trained in sequence.

Explanation

Consider a two-class classification problem:

- Samples x_1, x_2, \dots, x_N
- Binary labels $(-1, 1)$ t_1, t_2, \dots, t_N

Sequential Training

Main difference between boosting and committee methods

The base classifiers are trained in sequence.

Explanation

Consider a two-class classification problem:

- 1 Samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$
- 2 Binary labels $(-1, 1)$ t_1, t_2, \dots, t_N

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 **Boosting**
 - **AdaBoost Development**
 - **Cost Function**
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

Cost Function

Now

You want to put together a set of M experts able to recognize the most difficult inputs in an accurate way!!!

This

For each pattern x_i each expert classifier outputs a classification $y_j(x_i) \in \{-1, 1\}$

The final decision of the committee of M experts is $\sum_{j=1}^M \alpha_j y_j(x_i)$

$$C(x_i) = \alpha_1 y_1(x_i) + \alpha_2 y_2(x_i) + \dots + \alpha_M y_M(x_i) \quad (9)$$

Cost Function

Now

You want to put together a set of M experts able to recognize the most difficult inputs in an accurate way!!!

Thus

For each pattern \mathbf{x}_i each expert classifier outputs a classification $y_j(\mathbf{x}_i) \in \{-1, 1\}$

The final decision of the committee of M experts is $\sum_{j=1}^M \alpha_j y_j(\mathbf{x}_i)$

$$C(\mathbf{x}_i) = \alpha_1 y_1(\mathbf{x}_i) + \alpha_2 y_2(\mathbf{x}_i) + \dots + \alpha_M y_M(\mathbf{x}_i) \quad (9)$$

Cost Function

Now

You want to put together a set of M experts able to recognize the most difficult inputs in an accurate way!!!

Thus

For each pattern \mathbf{x}_i each expert classifier outputs a classification $y_j(\mathbf{x}_i) \in \{-1, 1\}$

The final decision of the committee of M experts is $\text{sign}(C(\mathbf{x}_i))$

$$C(\mathbf{x}_i) = \alpha_1 y_1(\mathbf{x}_i) + \alpha_2 y_2(\mathbf{x}_i) + \dots + \alpha_M y_M(\mathbf{x}_i) \quad (9)$$

Adaptive Boosting

It works even with a continuum of classifiers.

However

For the sake of simplicity, we will assume that the set of expert is finite.

Now

Adaptive Boosting

It works even with a continuum of classifiers.

However

For the sake of simplicity, we will assume that the set of expert is finite.

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 **Boosting**
 - **AdaBoost Development**
 - Cost Function
 - **Selection Process**
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

Getting the correct classifiers

We want the following

- We want to review possible element members.
- Select them, if they have certain properties.
- Assigning a weight to their contribution to the set of experts.

Getting the correct classifiers

We want the following

- We want to review possible element members.
- Select them, if they have certain properties.
- Assigning a weight to their contribution to the set of experts.

Getting the correct classifiers

We want the following

- We want to review possible element members.
- Select them, if they have certain properties.
- Assigning a weight to their contribution to the set of experts.

Now

Selection is done the following way

Testing the classifiers in the pool using a training set T of N multidimensional data points x_i :

- For each point x_i we have a label $t_i = 1$ or $t_i = -1$.

Now

Selection is done the following way

Testing the classifiers in the pool using a training set T of N multidimensional data points x_i :

- For each point x_i we have a label $t_i = 1$ or $t_i = -1$.

Assigning a cost for actions

We test and rank all classifiers in the expert pool by

- Charging a cost $\exp\{\beta\}$ any time a classifier fails (a miss).
- Charging a cost $\exp\{-\beta\}$ any time a classifier provides the right label (a hit).

Now

Selection is done the following way

Testing the classifiers in the pool using a training set T of N multidimensional data points x_i :

- For each point x_i we have a label $t_i = 1$ or $t_i = -1$.

Assigning a cost for actions

We test and rank all classifiers in the expert pool by

- Charging a cost $\exp\{\beta\}$ any time a classifier fails (a miss).
- Charging a cost $\exp\{-\beta\}$ any time a classifier provides the right label (a hit).

Now

Selection is done the following way

Testing the classifiers in the pool using a training set T of N multidimensional data points x_i :

- For each point x_i we have a label $t_i = 1$ or $t_i = -1$.

Assigning a cost for actions

We test and rank all classifiers in the expert pool by

- Charging a cost $\exp\{\beta\}$ any time a classifier fails (a miss).
- Charging a cost $\exp\{-\beta\}$ any time a classifier provides the right label (a hit).

Now

Selection is done the following way

Testing the classifiers in the pool using a training set T of N multidimensional data points x_i :

- For each point x_i we have a label $t_i = 1$ or $t_i = -1$.

Assigning a cost for actions

We test and rank all classifiers in the expert pool by

- Charging a cost $\exp\{\beta\}$ any time a classifier fails (a miss).
- Charging a cost $\exp\{-\beta\}$ any time a classifier provides the right label (a hit).

Remarks about β

We require $\beta > 0$

- Thus misses are penalized more heavily than hits

Remarks about β

We require $\beta > 0$

- Thus misses are penalized more heavily than hits

Although

- It looks strange to penalize hits,
- However, as long as the penalty of a success is smaller than the penalty for a miss:

$$\exp\{-\beta\} < \exp\{\beta\}$$

Remarks about β

We require $\beta > 0$

- Thus misses are penalized more heavily than hits

Although

- It looks strange to penalize hits,
- However, as long as the penalty of a success is smaller than the penalty for a miss:

$$\exp\{-\beta\} < \exp\{\beta\}$$

- if we assign cost a to misses and cost b to hits, where $a > b > 0$.
- We can rewrite such costs as $a = e^d$ and $b = e^{-d}$ for constants c and d .
 - ▶ It does not compromise generality.

Remarks about β

We require $\beta > 0$

- Thus misses are penalized more heavily than hits

Although

- It looks strange to penalize hits,
- However, as long as the penalty of a success is smaller than the penalty for a miss:

$$\exp\{-\beta\} < \exp\{\beta\}$$

Why?

- if we assign cost a to misses and cost b to hits, where $a > b > 0$.
- We can rewrite such costs as $a = c^d$ and $b = c^{-d}$ for constants c and d .
- It does not compromise generality.

Remarks about β

We require $\beta > 0$

- Thus misses are penalized more heavily than hits

Although

- It looks strange to penalize hits,
- However, as long as the penalty of a success is smaller than the penalty for a miss:

$$\exp\{-\beta\} < \exp\{\beta\}$$

Why?

- if we assign cost a to misses and cost b to hits, where $a > b > 0$.
- We can rewrite such costs as $a = c^d$ and $b = c^{-d}$ for constants c and d .

• It does not compromise generality.

Remarks about β

We require $\beta > 0$

- Thus misses are penalized more heavily than hits

Although

- It looks strange to penalize hits,
- However, as long as the penalty of a success is smaller than the penalty for a miss:

$$\exp\{-\beta\} < \exp\{\beta\}$$

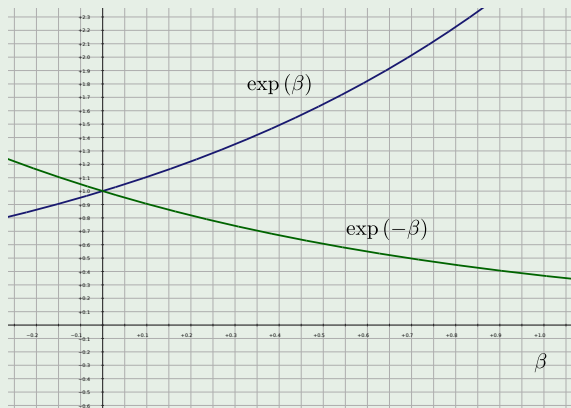
Why?

- if we assign cost a to misses and cost b to hits, where $a > b > 0$.
- We can rewrite such costs as $a = c^d$ and $b = c^{-d}$ for constants c and d .
 - ▶ It does not compromise generality.

Exponential Loss Function

This kind of error function is different from Squared Euclidean distance

- The classification target is called an exponential loss function.
- AdaBoost uses exponential error loss as error criterion.



Outline

1 Combining Models

- Introduction
- Average for Committee
- Beyond Simple Averaging
 - Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
 - The Differences

3 Committees

- Introduction
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- **How do we select classifiers?**
 - Selecting New Classifiers
 - Deriving against the weight α_m
- AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
- Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
- Example using an Infinitude of Perceptrons

Selection of the Classifier

We need to have a way to select the best Classifier in the Pool

- When we test the M classifiers in the pool, we build a matrix S

Then:

- We record the misses (with a ONE) and hits (with a ZERO) of each classifiers.

Selection of the Classifier

We need to have a way to select the best Classifier in the Pool

- When we test the M classifiers in the pool, we build a matrix S

Then

- We record the misses (with a ONE) and hits (with a ZERO) of each classifiers.

The Matrix S

Row i in the matrix is reserved for the data point \mathbf{x}_i

- Column m is reserved for the m th classifier in the pool.

Classifiers

	1	2	...	M
\mathbf{x}_1	0	1	...	1
\mathbf{x}_2	0	0	...	1
\mathbf{x}_3	1	1	...	0
\vdots	\vdots	\vdots		\vdots
\mathbf{x}_N	0	0	...	0

Something interesting about the S

The sum along the rows is the sum at the empirical risk

$$\text{ER}(y_j) = \frac{1}{N} \sum_{i=1}^N S_{ij} \text{ with } j = 1, \dots, M$$

Therefore, the candidate to be used at each iteration

- It is the classifier y_j with the smallest empirical risk!!!

Something interesting about the S

The sum along the rows is the sum at the empirical risk

$$\text{ER}(y_j) = \frac{1}{N} \sum_{i=1}^N S_{ij} \text{ with } j = 1, \dots, M$$

Therefore, the candidate to be used at certain iteration

- It is the classifier y_j with the smallest empirical risk!!!

Main Idea

What does AdaBoost want to do?

The main idea of AdaBoost is to proceed systematically by extracting one classifier from the pool in each of M iterations.

Main Idea

What does AdaBoost want to do?

The main idea of AdaBoost is to proceed systematically by extracting one classifier from the pool in each of M iterations.

Thus

The elements in the data set are weighted according to their current relevance (or urgency) at each iteration.

Main Idea

What does AdaBoost want to do?

The main idea of AdaBoost is to proceed systematically by extracting one classifier from the pool in each of M iterations.

Thus

The elements in the data set are weighted according to their current relevance (or urgency) at each iteration.

Thus at the beginning of the iterations

All data samples are assigned the same weight:

• Just 1, or $\frac{1}{N}$, if we want to have a total sum of 1 for all weights.

Main Idea

What does AdaBoost want to do?

The main idea of AdaBoost is to proceed systematically by extracting one classifier from the pool in each of M iterations.

Thus

The elements in the data set are weighted according to their current relevance (or urgency) at each iteration.

Thus at the beginning of the iterations

All data samples are assigned the same weight:

- Just 1, or $\frac{1}{N}$, if we want to have a total sum of 1 for all weights.

The process of the weights

As the selection progresses

- The more difficult samples, those where the committee still performs badly, are assigned larger and larger weights.

The process of the weights

As the selection progresses

- The more difficult samples, those where the committee still performs badly, are assigned larger and larger weights.

The selection process concentrates in selecting new classifiers

- For the committee by focusing on those which can help with the still misclassified examples.

The process of the weights

As the selection progresses

- The more difficult samples, those where the committee still performs badly, are assigned larger and larger weights.

The selection process concentrates in selecting new classifiers

- For the committee by focusing on those which can help with the still misclassified examples.

Then

- The best classifiers are those which can provide new insights to the committee.
- Classifiers being selected should complement each other in an optimal way.

The process of the weights

As the selection progresses

- The more difficult samples, those where the committee still performs badly, are assigned larger and larger weights.

The selection process concentrates in selecting new classifiers

- For the committee by focusing on those which can help with the still misclassified examples.

Then

- The best classifiers are those which can provide new insights to the committee.
- Classifiers being selected should complement each other in an optimal way.

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - **How do we select classifiers?**
 - **Selecting New Classifiers**
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

Selecting New Classifiers

What we want

In each iteration, we rank all classifiers, so that we can select the current best out of the pool.

At m th iteration

We have already included $m - 1$ classifiers in the committee and we want to select the next one.

Thus, we have the following cost function which is actually the output of the committee

$$C_{(m-1)}(x_i) = \alpha_1 y_1(x_i) + \alpha_2 y_2(x_i) + \dots + \alpha_{m-1} y_{m-1}(x_i) \quad (10)$$

Selecting New Classifiers

What we want

In each iteration, we rank all classifiers, so that we can select the current best out of the pool.

At m th iteration

We have already included $m - 1$ classifiers in the committee and we want to select the next one.

Thus, we have the following cost function which is actually the output of the committee:

$$C_{(m-1)}(x_i) = \alpha_1 y_1(x_i) + \alpha_2 y_2(x_i) + \dots + \alpha_{m-1} y_{m-1}(x_i) \quad (10)$$

Selecting New Classifiers

What we want

In each iteration, we rank all classifiers, so that we can select the current best out of the pool.

At m th iteration

We have already included $m - 1$ classifiers in the committee and we want to select the next one.

Thus, we have the following cost function which is actually the output of the committee

$$C_{(m-1)}(\mathbf{x}_i) = \alpha_1 y_1(\mathbf{x}_i) + \alpha_2 y_2(\mathbf{x}_i) + \dots + \alpha_{m-1} y_{m-1}(\mathbf{x}_i) \quad (10)$$

Thus, we have that

Extending the cost function by the new regression y_m

$$C_{(m)}(\mathbf{x}_i) = C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \quad (11)$$

At the first iteration $m = 1$:

- $C_{(0)}$ is the zero function.

Thus, the total cost or total error is defined as the exponential error

$$E = \sum_{i=1}^N \exp \left\{ -t_i \left(C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \right) \right\} \quad (12)$$

Thus, we have that

Extending the cost function by the new regression y_m

$$C_{(m)}(\mathbf{x}_i) = C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \quad (11)$$

At the first iteration $m = 1$

- $C_{(0)}$ is the zero function.

Thus, the total cost or total error is defined as the exponential error

$$E = \sum_{i=1}^N \exp \left\{ -t_i \left(C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \right) \right\} \quad (12)$$

Thus, we have that

Extending the cost function by the new regression y_m

$$C_{(m)}(\mathbf{x}_i) = C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \quad (11)$$

At the first iteration $m = 1$

- $C_{(0)}$ is the zero function.

Thus, the total cost or total error is defined as the exponential error

$$E = \sum_{i=1}^N \exp \left\{ -t_i \left(C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \right) \right\} \quad (12)$$

Thus

We want to determine

α_m and y_m in optimal way

Thus, rewriting

$$E = \sum_{i=1}^N w_i^{(m)} \exp \{-t_i \alpha_m y_m(\mathbf{x}_i)\} \quad (13)$$

where, for $i = 1, 2, \dots, N$

$$w_i^{(m)} = \exp \{-t_i C_{(m-1)}(\mathbf{x}_i)\} \quad (14)$$

Thus

We want to determine

α_m and y_m in optimal way

Thus, rewriting

$$E = \sum_{i=1}^N w_i^{(m)} \exp \{-t_i \alpha_m y_m(\mathbf{x}_i)\} \quad (13)$$

where, for $m = 1, 2, \dots, M$

$$w_i^{(m)} = \exp \{-t_i C_{(m-1)}(\mathbf{x}_i)\} \quad (14)$$

Thus

We want to determine

α_m and y_m in optimal way

Thus, rewriting

$$E = \sum_{i=1}^N w_i^{(m)} \exp \{-t_i \alpha_m y_m(\mathbf{x}_i)\} \quad (13)$$

Where, for $i = 1, 2, \dots, N$

$$w_i^{(m)} = \exp \{-t_i C_{(m-1)}(\mathbf{x}_i)\} \quad (14)$$

Remark

We have that the weight

$$w_i^{(m)} = \exp \left\{ -t_i C_{(m-1)}(\mathbf{x}_i) \right\}$$

Needs to be used in some way for the training of the new classifier

- This is of the out most importance!!!

Remark

We have that the weight

$$w_i^{(m)} = \exp \left\{ -t_i C_{(m-1)}(\mathbf{x}_i) \right\}$$

Needs to be used in some way for the training of the new classifier

- This is of the out most importance!!!

Therefore

You could use such weight

- As a output in the estimator function when applied to the loss function

$$\sum_{i=1}^N \left(y_i - w_i^{(m)} f(\mathbf{x}_i) \right)^2$$

Therefore

You could use such weight

- As a output in the estimator function when applied to the loss function

$$\sum_{i=1}^N \left(y_i - w_i^{(m)} f(\mathbf{x}_i) \right)^2$$

You could use such weight

- You could sub-sample with substitution by using the distribution $D_m \left\{ w_i^{(m)} \right\}$ of \mathbf{x}_i

▶ The train using that sub-sample

Therefore

You could use such weight

- As a output in the estimator function when applied to the loss function

$$\sum_{i=1}^N \left(y_i - w_i^{(m)} f(\mathbf{x}_i) \right)^2$$

You could use such weight

- You could sub-sample with substitution by using the distribution $D_m \left\{ w_i^{(m)} \right\}$ of \mathbf{x}_i
 - ▶ The train using that sub-sample

You could apply the weight function to the loss function itself used for training

$$\sum_{i=1}^N w_i^{(m)} (y_i - w_i f(\mathbf{x}_i))^2$$

Therefore

You could use such weight

- As a output in the estimator function when applied to the loss function

$$\sum_{i=1}^N \left(y_i - w_i^{(m)} f(\mathbf{x}_i) \right)^2$$

You could use such weight

- You could sub-sample with substitution by using the distribution $D_m \left\{ w_i^{(m)} \right\}$ of \mathbf{x}_i
 - ▶ The train using that sub-sample

You could apply the weight function to the loss function itself used for training

$$\sum_{i=1}^N w_i^{(m)} (y_i - w_i f(\mathbf{x}_i))^2$$

Thus

In the first iteration $w_i^{(1)} = 1$ for $i = 1, \dots, N$

- Meaning all the points have the same importance.

During later iterations, the vector w

- It represents the weight assigned to each data point in the training set at iteration m .

Thus

In the first iteration $w_i^{(1)} = 1$ for $i = 1, \dots, N$

- Meaning all the points have the same importance.

During later iterations, the vector $\mathbf{w}^{(m)}$

- It represents the weight assigned to each data point in the training set at iteration m .

Rewriting the Cost Equation

We can split (Eq. 13)

$$E = \sum_{t_i=y_m(\mathbf{x}_i)} w_i^{(m)} \exp\{-\alpha_m\} + \sum_{t_i \neq y_m(\mathbf{x}_i)} w_i^{(m)} \exp\{\alpha_m\} \quad (15)$$

Meaning:

The total cost is the weighted cost of all hits plus the weighted cost of all misses.

Rewriting the Cost Equation

We can split (Eq. 13)

$$E = \sum_{t_i=y_m(\mathbf{x}_i)} w_i^{(m)} \exp\{-\alpha_m\} + \sum_{t_i \neq y_m(\mathbf{x}_i)} w_i^{(m)} \exp\{\alpha_m\} \quad (15)$$

Meaning

The total cost is the weighted cost of all hits plus the weighted cost of all misses.

Therefore

Writing the first summand as $W_c \exp \{-\alpha_m\}$ and the second as $W_e \exp \{\alpha_m\}$

$$E = W_c \exp \{-\alpha_m\} + W_e \exp \{\alpha_m\} \quad (16)$$

Empty

Now, for the selection of y_m

- The exact value of $\alpha_m > 0$ is irrelevant

Since arrived at minimizing E

- It is equivalent to minimizing $\exp\{\alpha_m\} E$

Or in other words

$$\exp\{\alpha_m\} E = W_c + W_e \exp\{2\alpha_m\} \quad (17)$$

Empty

Now, for the selection of y_m

- The exact value of $\alpha_m > 0$ is irrelevant

Since a fixed α_m minimizing E

- It is equivalent to minimizing $\exp\{\alpha_m\} E$

Or in other words

$$\exp\{\alpha_m\} E = W_e + W_e \exp\{2\alpha_m\} \quad (17)$$

Empty

Now, for the selection of y_m

- The exact value of $\alpha_m > 0$ is irrelevant

Since a fixed α_m minimizing E

- It is equivalent to minimizing $\exp\{\alpha_m\} E$

Or in other words

$$\exp\{\alpha_m\} E = W_c + W_e \exp\{2\alpha_m\} \quad (17)$$

Now, we have

Given that $\alpha_m > 0$

$$2\alpha_m > 0$$

We have

$$\exp\{2\alpha_m\} > \exp\{0\} = 1$$

Now, we have

Given that $\alpha_m > 0$

$$2\alpha_m > 0$$

We have

$$\exp\{2\alpha_m\} > \exp\{0\} = 1$$

Then

We can rewrite (Eq. 17)

$$\exp\{\alpha_m\} E = W_c + W_e - W_e + W_e \exp\{2\alpha_m\} \quad (18)$$

Thus

$$\exp\{\alpha_m\} E = (W_c + W_e) + W_e (\exp\{2\alpha_m\} - 1) \quad (19)$$

Now $(W_c + W_e)$ is the total gain W of the weights.

- Of all data points which is constant in the current iteration.

Then

We can rewrite (Eq. 17)

$$\exp\{\alpha_m\} E = W_c + W_e - W_e + W_e \exp\{2\alpha_m\} \quad (18)$$

Thus

$$\exp\{\alpha_m\} E = (W_c + W_e) + W_e (\exp\{2\alpha_m\} - 1) \quad (19)$$

Now $W_c + W_e$ is the total sum W of the weights

- Of all data points which is constant in the current iteration.

Then

We can rewrite (Eq. 17)

$$\exp\{\alpha_m\} E = W_c + W_e - W_e + W_e \exp\{2\alpha_m\} \quad (18)$$

Thus

$$\exp\{\alpha_m\} E = (W_c + W_e) + W_e (\exp\{2\alpha_m\} - 1) \quad (19)$$

Now, $W_c + W_e$ is the total sum W of the weights

- Of all data points which is constant in **the current iteration**.

Thus

The right hand side of the equation is minimized

- When at the m -th iteration, we pick the classifier with the lowest total cost W_e
 - ▶ That is the lowest rate of weighted error.

Intuitively,

The next selected y_m should be the one with the lowest penalty given the current set of weights.

Thus

The right hand side of the equation is minimized

- When at the m -th iteration, we pick the classifier with the lowest total cost W_e
 - ▶ That is the lowest rate of weighted error.

Intuitively

The next selected y_m should be the one with the lowest penalty given the current set of weights.

Do you remember?

The Matrix S

- We pick the classifier with the lowest total cost W_e

Now, we need to do some updates

- Specifically the value α_m .

Do you remember?

The Matrix S

- We pick the classifier with the lowest total cost W_e

Now, we need to do some updates

- Specifically the value α_m .

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - **How do we select classifiers?**
 - Selecting New Classifiers
 - **Deriving against the weight α_m**
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

Deriving against the weight α_m

Going back to the original E , we can use the derivative trick

$$\frac{\partial E}{\partial \alpha_m} = -W_c \exp\{-\alpha_m\} + W_e \exp\{\alpha_m\} \quad (20)$$

Making the equation equal to zero and multiplying by $\exp\{\alpha_m\}$

$$-W_c + W_e \exp\{2\alpha_m\} = 0 \quad (21)$$

The optimal value is thus

$$\alpha_m = \frac{1}{2} \ln \left(\frac{W_e}{W_c} \right) \quad (22)$$

Deriving against the weight α_m

Going back to the original E , we can use the derivative trick

$$\frac{\partial E}{\partial \alpha_m} = -W_c \exp\{-\alpha_m\} + W_e \exp\{\alpha_m\} \quad (20)$$

Making the equation equal to zero and multiplying by $\exp\{\alpha_m\}$

$$-W_c + W_e \exp\{2\alpha_m\} = 0 \quad (21)$$

The optimal value is thus

$$\alpha_m = \frac{1}{2} \ln \left(\frac{W_e}{W_c} \right) \quad (22)$$

Deriving against the weight α_m

Going back to the original E , we can use the derivative trick

$$\frac{\partial E}{\partial \alpha_m} = -W_c \exp \{-\alpha_m\} + W_e \exp \{\alpha_m\} \quad (20)$$

Making the equation equal to zero and multiplying by $\exp \{\alpha_m\}$

$$-W_c + W_e \exp \{2\alpha_m\} = 0 \quad (21)$$

The optimal value is thus

$$\alpha_m = \frac{1}{2} \ln \left(\frac{W_c}{W_e} \right) \quad (22)$$

Now

Making the total sum of all weights

$$W = W_c + W_e \quad (23)$$

We can rewrite the previous equation as

$$\alpha_m = \frac{1}{2} \ln \left(\frac{W - W_e}{W_e} \right) = \frac{1}{2} \ln \left(\frac{1 - e_m}{e_m} \right) \quad (24)$$

With the percentage rate of error given the weights of the data points

$$e_m = \frac{W_e}{W} \quad (25)$$

Now

Making the total sum of all weights

$$W = W_c + W_e \quad (23)$$

We can rewrite the previous equation as

$$\alpha_m = \frac{1}{2} \ln \left(\frac{W - W_e}{W_e} \right) = \frac{1}{2} \ln \left(\frac{1 - e_m}{e_m} \right) \quad (24)$$

With the percentage rate of error given the weights of the data points

$$e_m = \frac{W_e}{W} \quad (25)$$

Now

Making the total sum of all weights

$$W = W_c + W_e \quad (23)$$

We can rewrite the previous equation as

$$\alpha_m = \frac{1}{2} \ln \left(\frac{W - W_e}{W_e} \right) = \frac{1}{2} \ln \left(\frac{1 - e_m}{e_m} \right) \quad (24)$$

With the percentage rate of error given the weights of the data points

$$e_m = \frac{W_e}{W} \quad (25)$$

What about the weights?

Using the equation

$$w_i^{(m)} = \exp \left\{ -t_i C_{(m-1)}(\mathbf{x}_i) \right\} \quad (26)$$

And because we have $w_i^{(m)}$ and $w_i^{(m-1)}$

What about the weights?

Using the equation

$$w_i^{(m)} = \exp \left\{ -t_i C_{(m-1)}(\mathbf{x}_i) \right\} \quad (26)$$

And because we have α_m and $y_m(\mathbf{x}_i)$

$$\begin{aligned} w_i^{(m+1)} &= \exp \left\{ -t_i C_{(m)}(\mathbf{x}_i) \right\} \\ &= \exp \left\{ -t_i \left[C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \right] \right\} \\ &= w_i^{(m)} \exp \left\{ -t_i \alpha_m y_m(\mathbf{x}_i) \right\} \end{aligned}$$

What about the weights?

Using the equation

$$w_i^{(m)} = \exp \left\{ -t_i C_{(m-1)}(\mathbf{x}_i) \right\} \quad (26)$$

And because we have α_m and $y_m(\mathbf{x}_i)$

$$\begin{aligned} w_i^{(m+1)} &= \exp \left\{ -t_i C_{(m)}(\mathbf{x}_i) \right\} \\ &= \exp \left\{ -t_i \left[C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \right] \right\} \\ &= w_i^{(m)} \exp \left\{ -t_i \alpha_m y_m(\mathbf{x}_i) \right\} \end{aligned}$$

What about the weights?

Using the equation

$$w_i^{(m)} = \exp \left\{ -t_i C_{(m-1)}(\mathbf{x}_i) \right\} \quad (26)$$

And because we have α_m and $y_m(\mathbf{x}_i)$

$$\begin{aligned} w_i^{(m+1)} &= \exp \left\{ -t_i C_{(m)}(\mathbf{x}_i) \right\} \\ &= \exp \left\{ -t_i \left[C_{(m-1)}(\mathbf{x}_i) + \alpha_m y_m(\mathbf{x}_i) \right] \right\} \\ &= w_i^{(m)} \exp \left\{ -t_i \alpha_m y_m(\mathbf{x}_i) \right\} \end{aligned}$$

Sequential Training

Thus

- AdaBoost trains a new classifier using a data set
- There the weighting coefficients are adjusted according to the performance of the previously trained classifier
- To give greater weight to the misclassified data points.

Sequential Training

Thus

- AdaBoost trains a new classifier using a data set
- There the weighting coefficients are adjusted according to the performance of the previously trained classifier
- To give greater weight to the misclassified data points.

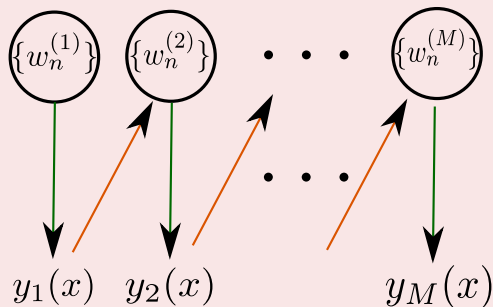
Sequential Training

Thus

- AdaBoost trains a new classifier using a data set
- There the weighting coefficients are adjusted according to the performance of the previously trained classifier
- To give greater weight to the misclassified data points.

Illustration

Schematic Illustration



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

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 **Boosting**
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - **AdaBoost Algorithm**
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

AdaBoost Algorithm

Step 1

Initialize $\{w_i^{(1)}\}$ to $\frac{1}{N}$

AdaBoost Algorithm

Step 1

Initialize $\{w_i^{(1)}\}$ to $\frac{1}{N}$

Step 2

For $m = 1, 2, \dots, M$

- Select a weak classifier $y_m(x)$ to the training data by minimizing the weighted error function or

$$\arg \min_{y_m} \sum_{i=1}^N w_i^{(m)} I(y_m(x_i) \neq t_n) = \arg \min_{y_m} \sum_{t_i \neq y_m(x_i)} w_i^{(m)} = \arg \min_{y_m} W_e \quad (27)$$

Where I is an indicator function.

AdaBoost Algorithm

Step 1

Initialize $\{w_i^{(1)}\}$ to $\frac{1}{N}$

Step 2

For $m = 1, 2, \dots, M$

- Select a weak classifier $y_m(\mathbf{x})$ to the training data by minimizing the weighted error function or

$$\arg \min_{y_m} \sum_{i=1}^N w_i^{(m)} I(y_m(\mathbf{x}_i) \neq t_i) = \arg \min_{y_m} \sum_{t_i \neq y_m(\mathbf{x}_i)} w_i^{(m)} = \arg \min_{y_m} W_e \quad (27)$$

Where I is an indicator function.

AdaBoost Algorithm

Step 1

Initialize $\{w_i^{(1)}\}$ to $\frac{1}{N}$

Step 2

For $m = 1, 2, \dots, M$

- Select a weak classifier $y_m(\mathbf{x})$ to the training data by minimizing the weighted error function or

$$\arg \min_{y_m} \sum_{i=1}^N w_i^{(m)} I(y_m(\mathbf{x}_i) \neq t_n) = \arg \min_{y_m} \sum_{t_i \neq y_m(\mathbf{x}_i)} w_i^{(m)} = \arg \min_{y_m} W_e \quad (27)$$

Where I is an indicator function.

AdaBoost Algorithm

Step 2

- Evaluate

$$e_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (28)$$

Where I is an indicator function

AdaBoost Algorithm

Step 3

Set the α_m weight to

$$\alpha_m = \frac{1}{2} \ln \left\{ \frac{1 - e_m}{e_m} \right\} \quad (29)$$

AdaBoost Algorithm

Step 3

Set the α_m weight to

$$\alpha_m = \frac{1}{2} \ln \left\{ \frac{1 - e_m}{e_m} \right\} \quad (29)$$

Now update the weights of the data for the next iteration

- If $t_i \neq y_m(\mathbf{x}_i)$ i.e. a miss

$$w_i^{(m+1)} = w_i^{(m)} \exp \{ \alpha_m \} = w_i^{(m)} \sqrt{\frac{1 - e_m}{e_m}} \quad (30)$$

- If $t_i = y_m(\mathbf{x}_i)$ i.e. a hit

$$w_i^{(m+1)} = w_i^{(m)} \exp \{ -\alpha_m \} = w_i^{(m)} \sqrt{\frac{e_m}{1 - e_m}} \quad (31)$$

AdaBoost Algorithm

Step 3

Set the α_m weight to

$$\alpha_m = \frac{1}{2} \ln \left\{ \frac{1 - e_m}{e_m} \right\} \quad (29)$$

Now update the weights of the data for the next iteration

- If $t_i \neq y_m(\mathbf{x}_i)$ i.e. a miss

$$w_i^{(m+1)} = w_i^{(m)} \exp \{ \alpha_m \} = w_i^{(m)} \sqrt{\frac{1 - e_m}{e_m}} \quad (30)$$

- If $t_i = y_m(\mathbf{x}_i)$ i.e. a hit

$$w_i^{(m+1)} = w_i^{(m)} \exp \{ -\alpha_m \} = w_i^{(m)} \sqrt{\frac{e_m}{1 - e_m}} \quad (31)$$

Finally, make predictions

For this use

$$Y_M(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right) \quad (32)$$

Outline

1 Combining Models

- Introduction
- Average for Committee
- Beyond Simple Averaging
 - Example

2 Bayesian Model Averaging

- Model Combination Vs. Bayesian Model Averaging
- Now Model Averaging
 - The Differences

3 Committees

- Introduction
- Bootstrap Data Sets
- Relation with Monte-Carlo Estimation

4 Boosting

- AdaBoost Development
 - Cost Function
 - Selection Process
- How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
- **AdaBoost Algorithm**
 - **Some Remarks**
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

Observations

First

The first base classifier is the usual procedure of training a single classifier.

Second

From (Eq. 30) and (Eq. 31), we can see that the weighting coefficient are increased for data points that are misclassified.

Third

- The quantity e_m represent weighted measures of the error rate.
- Thus α_m gives more weight to the more accurate classifiers.

Observations

First

The first base classifier is the usual procedure of training a single classifier.

Second

From (Eq. 30) and (Eq. 31), we can see that the weighting coefficient are increased for data points that are misclassified.

- The quantity e_m represent weighted measures of the error rate.
- Thus α_m gives more weight to the more accurate classifiers.

Observations

First

The first base classifier is the usual procedure of training a single classifier.

Second

From (Eq. 30) and (Eq. 31), we can see that the weighting coefficient are increased for data points that are misclassified.

Third

- The quantity e_m represent weighted measures of the error rate.
- Thus α_m gives more weight to the more accurate classifiers.

In addition

The pool of classifiers in Step 1 can be substituted by a family of classifiers

One whose members are trained to minimize the error function given the current weights

Now

If indeed a finite set of classifiers is given, we only need to test the classifiers once for each data point

The Scoring Matrix

It can be reused at each iteration by multiplying the transposed vector of weights $w^{(m)}$ with S to obtain W_c of each machine

In addition

The pool of classifiers in Step 1 can be substituted by a family of classifiers

One whose members are trained to minimize the error function given the current weights

Now

If indeed a finite set of classifiers is given, we only need to test the classifiers once for each data point

The scoring matrix

It can be reused at each iteration by multiplying the transposed vector of weights $w^{(m)}$ with S to obtain W_c of each machine

In addition

The pool of classifiers in Step 1 can be substituted by a family of classifiers

One whose members are trained to minimize the error function given the current weights

Now

If indeed a finite set of classifiers is given, we only need to test the classifiers once for each data point

The Scouting Matrix S

It can be reused at each iteration by multiplying the transposed vector of weights $w^{(m)}$ with S to obtain W_e of each machine

We have then

The following

$$\left[W_e^{(1)} \ W_e^{(2)} \ \dots \ W_e^{(M)} \right] = \left(\mathbf{w}^{(m)} \right)^T S \quad (33)$$

We have then

The following

$$\left[W_e^{(1)} \ W_e^{(2)} \ \dots \ W_e^{(M)} \right] = \left(\mathbf{w}^{(m)} \right)^T S \quad (33)$$

This allows to reformulate the weight update step such that

It only misses lead to weight modification.

We have then

The following

$$\left[W_e^{(1)} \ W_e^{(2)} \ \dots \ W_e^{(M)} \right] = \left(\mathbf{w}^{(m)} \right)^T S \quad (33)$$

This allows to reformulate the weight update step such that

It only misses lead to weight modification.

Note

- Note that the weight vector $\mathbf{w}^{(m)}$ is constructed iteratively.
- It could be recomputed completely at every iteration, but the iterative construction is more efficient and simple to implement.

We have then

The following

$$\left[W_e^{(1)} \ W_e^{(2)} \ \dots \ W_e^M \right] = \left(\mathbf{w}^{(m)} \right)^T S \quad (33)$$

This allows to reformulate the weight update step such that

It only misses lead to weight modification.

Note

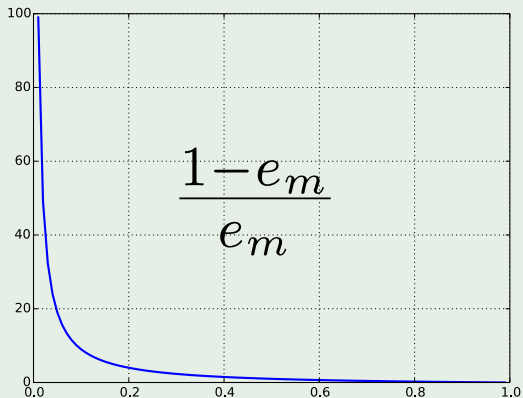
- Note that the weight vector $\mathbf{w}^{(m)}$ is constructed iteratively.
- It could be recomputed completely at every iteration, but the iterative construction is more efficient and simple to implement.

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - **AdaBoost Algorithm**
 - Some Remarks
 - **Explanation about AdaBoost's behavior**
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

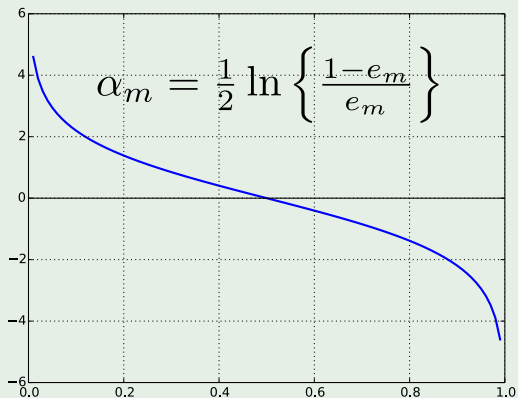
Explanation

Graph for $\frac{1-e_m}{e_m}$ in the range $0 \leq e_m \leq 1$



So we have

Graph for α_m



We have the following cases

We have the following

If $e_m \rightarrow 1$, we have that all the samples were not correctly classified!!!

This

We get that for all miss-classified sample $\lim_{c_{m+1}} \frac{1-e_m}{e_m} \rightarrow 0$, then
 $\alpha_m \rightarrow -\infty$

We have the following cases

We have the following

If $e_m \rightarrow 1$, we have that all the samples were not correctly classified!!!

Thus

We get that for all miss-classified sample $\lim_{e_m \rightarrow 1} \frac{1-e_m}{e_m} \rightarrow 0$, then

$\alpha_m \rightarrow -\infty$

Now

We get that for all miss-classified sample

$$w_i^{(m+1)} = w_i^{(m)} \exp\{\alpha_m\} \rightarrow 0$$

Therefore

- We only need to reverse the answers to get the perfect classifier and select it as the only committee member.

Now

We get that for all miss-classified sample

$$w_i^{(m+1)} = w_i^{(m)} \exp\{\alpha_m\} \rightarrow 0$$

Therefore

- We only need to reverse the answers to get the perfect classifier and select it as the only committee member.

Now, the Last Case

If $e_m \rightarrow 1/2$

- We have $\alpha_m \rightarrow 0$

Thus we have that if the sample is well or bad classified

$$\exp \{-\alpha_m t_i y_m(x_i)\} \rightarrow 1 \quad (34)$$

Therefore

- The weight does not change at all.

Now, the Last Case

If $e_m \rightarrow 1/2$

- We have $\alpha_m \rightarrow 0$

Thus we have that if the sample is well or bad classified

$$\exp \{ -\alpha_m t_i y_m(\mathbf{x}_i) \} \rightarrow 1 \quad (34)$$

Therefore

- The weight does not change at all.

Now, the Last Case

If $e_m \rightarrow 1/2$

- We have $\alpha_m \rightarrow 0$

Thus we have that if the sample is well or bad classified

$$\exp \{ -\alpha_m t_i y_m (\mathbf{x}_i) \} \rightarrow 1 \quad (34)$$

Therefore

- The weight does not change at all.

Thus, we have

What about $e_m \rightarrow 0$

- We have that $\alpha_m \rightarrow +\infty$

Thus, we have

What about $e_m \rightarrow 0$

- We have that $\alpha_m \rightarrow +\infty$

Thus, we have

Samples always correctly classified

$$w_i^{(m+1)} = w_i^{(m)} \exp \{-\alpha_m t_i y_m(\mathbf{x}_i)\} \rightarrow 0$$

- Thus, the only need m committee members, we do not need another $m+1$ member.

Thus, we have

What about $e_m \rightarrow 0$

- We have that $\alpha_m \rightarrow +\infty$

Thus, we have

Samples always correctly classified

$$w_i^{(m+1)} = w_i^{(m)} \exp \{-\alpha_m t_i y_m(\mathbf{x}_i)\} \rightarrow 0$$

- Thus, the only need m committee members, we do not need another $m + 1$ member.

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - **Statistical Analysis of the Exponential Loss**
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

This comes from

The paper

- “Additive Logistic Regression: A Statistical View of Boosting” by Friedman, Hastie and Tibshirani

Something to Note

- In this paper, a proof exists to show that boosting algorithms are procedures to fit an additive logistic regression model.

$$E[y|\mathbf{x}] = F(\mathbf{x}) \text{ with } F(\mathbf{x}) = \sum_{m=1}^M f_m(\mathbf{x})$$

This comes from

The paper

- “Additive Logistic Regression: A Statistical View of Boosting” by Friedman, Hastie and Tibshirani

Something Notable

- In this paper, a proof exists to show that boosting algorithms are procedures to fit an additive logistic regression model.

$$E[y|\mathbf{x}] = F(\mathbf{x}) \text{ with } F(\mathbf{x}) = \sum_{m=1}^M f_m(\mathbf{x})$$

Consider the Additive Regression Model

We are interested in modeling the mean $E[y|\mathbf{x}] = F(\mathbf{x})$

- With Additive Model

$$F(\mathbf{x}) = \sum_{i=1}^d f_i(x_i)$$

Where each f_i is a function for each feature input.

- A convenient algorithm for updating these models is the **backfitting algorithm** with update:

$$f_i(x_i) = E \left[y - \sum_{k \neq i} f_k(x_k) \mid x_i \right]$$

Consider the Additive Regression Model

We are interested in modeling the mean $E[y|\mathbf{x}] = F(\mathbf{x})$

- With Additive Model

$$F(\mathbf{x}) = \sum_{i=1}^d f_i(x_i)$$

Where each $f_i(x_i)$ is a function for each feature input x_i

- A convenient algorithm for updating these models is the **backfitting algorithm** with update:

$$f_i(x_i) = E \left[y - \sum_{k \neq i} f_k(x_k) \mid x_i \right]$$

Remarks

An example of these additive models is the matching pursuit

$$f(t) = \sum_{n=-\infty}^{+\infty} a_n g_{\gamma_n}(t)$$

Backfitting ensures that under fairly general conditions

- Backfitting converges to the minimizer of $E[(y - f(x))^2]$

Remarks

An example of these additive models is the matching pursuit

$$f(t) = \sum_{n=-\infty}^{+\infty} a_n g_{\gamma_n}(t)$$

Backfitting ensures that under fairly general conditions

- Backfitting converges to the minimizer of $E[(y - f(\mathbf{x}))^2]$

In the case of AdaBoost

We have an additive model

- Which considers functions $\{f_m(\mathbf{x})\}_{m=1}^M$ that take in account all the features - Perceptron, Decision Trees, etc

Each of these functions is characterized by a set of parameters and multiplier α_m

$$f_m(\mathbf{x}) = \alpha_m y_m(\mathbf{x}|\gamma_m)$$

With additive model

$$F_M(\mathbf{x}) = \alpha_1 y_1(\mathbf{x}|\gamma_1) + \dots + \alpha_M y_M(\mathbf{x}|\gamma_M)$$

In the case of AdaBoost

We have an additive model

- Which considers functions $\{f_m(\mathbf{x})\}_{m=1}^M$ that take in account all the features - Perceptron, Decision Trees, etc

Each of these functions is characterized by a set of parameters γ_m and multiplier α_m

$$f_m(\mathbf{x}) = \alpha_m y_m(\mathbf{x}|\gamma_m)$$

With additive model

$$F_M(\mathbf{x}) = \alpha_1 y_1(\mathbf{x}|\gamma_1) + \dots + \alpha_M y_M(\mathbf{x}|\gamma_M)$$

In the case of AdaBoost

We have an additive model

- Which considers functions $\{f_m(\mathbf{x})\}_{m=1}^M$ that take in account all the features - Perceptron, Decision Trees, etc

Each of these functions is characterized by a set of parameters γ_m and multiplier α_m

$$f_m(\mathbf{x}) = \alpha_m y_m(\mathbf{x}|\gamma_m)$$

With additive model

$$F_M(\mathbf{x}) = \alpha_1 y_1(\mathbf{x}|\gamma_1) + \dots + \alpha_M y_M(\mathbf{x}|\gamma_M)$$

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - **Statistical Analysis of the Exponential Loss**
 - **Moving from Regression to Classification**
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

Remark - Moving from Regression to Classification

Given that Regression have wide ranges of outputs

- Logistic Regression is widely used to move Regression to Classification

$$\log \frac{P(Y = 1|\mathbf{x})}{P(Y = -1|\mathbf{x})} = \sum_{m=1}^M f_m(\mathbf{x})$$

Another property: the probability estimates lie in $[0, 1]$

- Now, solving by assuming $P(Y = 1|\mathbf{x}) + P(Y = -1|\mathbf{x}) = 1$

$$P(Y = 1|\mathbf{x}) = \frac{e^{F(\mathbf{x})}}{1 + e^{F(\mathbf{x})}}$$

Remark - Moving from Regression to Classification

Given that Regression have wide ranges of outputs

- Logistic Regression is widely used to move Regression to Classification

$$\log \frac{P(Y = 1|\mathbf{x})}{P(Y = -1|\mathbf{x})} = \sum_{m=1}^M f_m(\mathbf{x})$$

A nice property, the probability estimates lie in $[0, 1]$

- Now, solving by assuming $P(Y = 1|\mathbf{x}) + P(Y = -1|\mathbf{x}) = 1$

$$P(Y = 1|\mathbf{x}) = \frac{e^{F(\mathbf{x})}}{1 + e^{F(\mathbf{x})}}$$

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - **Statistical Analysis of the Exponential Loss**
 - Moving from Regression to Classification
 - **Minimization of the Exponential Criterion**
 - Finally, The Additive Logistic Regression
 - Example using an Infinitude of Perceptrons

The Exponential Criterion

We have our exponential Criterion under an Expected Value with $y \in \{1, -1\}$

$$J(F) = E \left[e^{-yF(x)} \right]$$

Lemma

- $E \left[e^{-yF(x)} \right]$ is minimized at

$$F(x) = \frac{1}{2} \log \frac{P(Y = 1|x)}{P(Y = -1|x)}$$

Hence:

$$P(Y = 1|x) = \frac{e^{F(x)}}{e^{-F(x)} + e^{F(x)}}$$
$$P(Y = -1|x) = \frac{e^{-F(x)}}{e^{-F(x)} + e^{F(x)}}$$

The Exponential Criterion

We have our exponential Criterion under an Expected Value with $y \in \{1, -1\}$

$$J(F) = E \left[e^{-yF(\mathbf{x})} \right]$$

Lemma

- $E \left[e^{-yF(\mathbf{x})} \right]$ is minimized at

$$F(\mathbf{x}) = \frac{1}{2} \log \frac{P(Y = 1|\mathbf{x})}{P(Y = -1|\mathbf{x})}$$

Hence:

$$P(Y = 1|\mathbf{x}) = \frac{e^{F(\mathbf{x})}}{e^{-F(\mathbf{x})} + e^{F(\mathbf{x})}}$$
$$P(Y = -1|\mathbf{x}) = \frac{e^{-F(\mathbf{x})}}{e^{-F(\mathbf{x})} + e^{F(\mathbf{x})}}$$

Proof

Given the discrete nature of $y \in \{1, -1\}$

$$\frac{\partial E \left[e^{-yF(\mathbf{x})} \right]}{\partial F(\mathbf{x})} = -P(Y = 1|\mathbf{x}) e^{-F(\mathbf{x})} + P(Y = -1|\mathbf{x}) e^{F(\mathbf{x})}$$

Therefore

$$-P(Y = 1|\mathbf{x}) e^{-F(\mathbf{x})} + P(Y = -1|\mathbf{x}) e^{F(\mathbf{x})} = 0$$

Proof

Given the discrete nature of $y \in \{1, -1\}$

$$\frac{\partial E \left[e^{-yF(\mathbf{x})} \right]}{\partial F(\mathbf{x})} = -P(Y = 1|\mathbf{x}) e^{-F(\mathbf{x})} + P(Y = -1|\mathbf{x}) e^{F(\mathbf{x})}$$

Therefore

$$-P(Y = 1|\mathbf{x}) e^{-F(\mathbf{x})} + P(Y = -1|\mathbf{x}) e^{F(\mathbf{x})} = 0$$

Then

We have that

$$\begin{aligned}P(Y = 1|\mathbf{x}) e^{-F(\mathbf{x})} &= P(Y = -1|\mathbf{x}) e^{F(\mathbf{x})} \\ &= [1 - P(Y = 1|\mathbf{x})] e^{F(\mathbf{x})}\end{aligned}$$

Solving

$$e^{F(\mathbf{x})} = [e^{-F(\mathbf{x})} + e^{F(\mathbf{x})}] P(Y = 1|\mathbf{x})$$

Then

We have that

$$\begin{aligned} P(Y = 1|\mathbf{x}) e^{-F(\mathbf{x})} &= P(Y = -1|\mathbf{x}) e^{F(\mathbf{x})} \\ &= [1 - P(Y = 1|\mathbf{x})] e^{F(\mathbf{x})} \end{aligned}$$

Solving

$$e^{F(\mathbf{x})} = [e^{-F(\mathbf{x})} + e^{F(\mathbf{x})}] P(Y = 1|\mathbf{x})$$

Finally, we have

The first equation

$$P(Y = 1|\mathbf{x}) = \frac{e^{F(\mathbf{x})}}{e^{-F(\mathbf{x})} + e^{F(\mathbf{x})}}$$

Similarly

$$P(Y = -1|\mathbf{x}) = \frac{e^{-F(\mathbf{x})}}{e^{-F(\mathbf{x})} + e^{F(\mathbf{x})}}$$

Finally, we have

The first equation

$$P(Y = 1|\mathbf{x}) = \frac{e^{F(\mathbf{x})}}{e^{-F(\mathbf{x})} + e^{F(\mathbf{x})}}$$

Similarly

$$P(Y = -1|\mathbf{x}) = \frac{e^{-F(\mathbf{x})}}{e^{-F(\mathbf{x})} + e^{F(\mathbf{x})}}$$

Basically

We have that the $E \left[e^{-yF(x)} \right]$

- When you minimize the cost function

Then at the optimal you have the Binary Classification

- Of the Logistic Regression

Basically

We have that the $E \left[e^{-yF(x)} \right]$

- When you minimize the cost function

Then at the optimal you have the Binary Classification

- Of the Logistic Regression

Furthermore

Corollary

- If E is replaced by averages over regions of \mathbf{x} where $F(\mathbf{x})$ is constant (Similar to a decision tree),

► The same result applies to the sample proportions of $y = 1$ and $y = -1$

Furthermore

Corollary

- If E is replaced by averages over regions of \mathbf{x} where $F(\mathbf{x})$ is constant (Similar to a decision tree),
 - ▶ The same result applies to the sample proportions of $y = 1$ and $y = -1$

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 Boosting
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - **Statistical Analysis of the Exponential Loss**
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - **Finally, The Additive Logistic Regression**
 - Example using an Infinitude of Perceptrons

Finally, The Additive Logistic Regression

Proposition

- The AdaBoost algorithm fits an additive logistic regression model by stage-wise optimization of

$$J(F) = E \left[e^{-yF(\mathbf{x})} \right]$$

Proof

- Imagine you have an estimate $F(\mathbf{x})$ then we seek an improved estimate:

$$F(\mathbf{x}) + f(\mathbf{x})$$

Finally, The Additive Logistic Regression

Proposition

- The AdaBoost algorithm fits an additive logistic regression model by stage-wise optimization of

$$J(F) = E \left[e^{-yF(\mathbf{x})} \right]$$

Proof

- Imagine you have an estimate $F(\mathbf{x})$ then we seek an improved estimate:

$$F(\mathbf{x}) + f(\mathbf{x})$$

For This

We minimize at each x

$$J(F(x) + f(x))$$

This can be expanded

$$\begin{aligned} J(F(x) + f(x)) &= E \left[e^{-y(F(x)+f(x))} | x \right] \\ &= e^{-f(x)} E \left[e^{-yF(x)} I(y = 1) | x \right] + \\ &\dots e^{f(x)} E \left[e^{-yF(x)} I(y = -1) | x \right] \end{aligned}$$

For This

We minimize at each \mathbf{x}

$$J(F(\mathbf{x}) + f(\mathbf{x}))$$

This can be expanded

$$\begin{aligned} J(F(\mathbf{x}) + f(\mathbf{x})) &= E \left[e^{-y(F(\mathbf{x}) + f(\mathbf{x}))} | \mathbf{x} \right] \\ &= e^{-f(\mathbf{x})} E \left[e^{-yF(\mathbf{x})} I(y = 1) | \mathbf{x} \right] + \\ &\dots e^{f(\mathbf{x})} E \left[e^{-yF(\mathbf{x})} I(y = -1) | \mathbf{x} \right] \end{aligned}$$

Deriving w.r.t. $f(\mathbf{x})$

We get

$$-e^{-f(\mathbf{x})} E \left[e^{-yF(\mathbf{x})} I(y = 1) | \mathbf{x} \right] + e^{f(\mathbf{x})} E \left[e^{-yF(\mathbf{x})} I(y = -1) | \mathbf{x} \right] = 0$$

We have the following

If we divide by $E [e^{-yF(\mathbf{x})} | \mathbf{x}]$, the first term

$$\frac{E [e^{-yF(\mathbf{x})} I(y = 1) | \mathbf{x}]}{E [e^{-yF(\mathbf{x})} | \mathbf{x}]} = E_w [I(y = 1) | \mathbf{x}]$$

Also

$$\frac{E [e^{-yF(\mathbf{x})} I(y = -1) | \mathbf{x}]}{E [e^{-yF(\mathbf{x})} | \mathbf{x}]} = E_w [I(y = -1) | \mathbf{x}]$$

We have the following

If we divide by $E \left[e^{-yF(\mathbf{x})} | \mathbf{x} \right]$, the first term

$$\frac{E \left[e^{-yF(\mathbf{x})} I(y = 1) | \mathbf{x} \right]}{E \left[e^{-yF(\mathbf{x})} | \mathbf{x} \right]} = E_w \left[I(y = 1) | \mathbf{x} \right]$$

Also

$$\frac{E \left[e^{-yF(\mathbf{x})} I(y = -1) | \mathbf{x} \right]}{E \left[e^{-yF(\mathbf{x})} | \mathbf{x} \right]} = E_w \left[I(y = -1) | \mathbf{x} \right]$$

Thus, we have

We apply the natural log to both sides

$$\log e^{-f(\mathbf{x})} + \log E_w [I(y = 1) | \mathbf{x}] = \log e^{f(\mathbf{x})} + \log E_w [I(y = -1) | \mathbf{x}]$$

Then

$$2f(\mathbf{x}) = \log E_w [I(y = 1) | \mathbf{x}] - \log E_w [I(y = -1) | \mathbf{x}]$$

Thus, we have

We apply the natural log to both sides

$$\log e^{-f(\mathbf{x})} + \log E_w [I(y = 1) | \mathbf{x}] = \log e^{f(\mathbf{x})} + \log E_w [I(y = -1) | \mathbf{x}]$$

Then

$$2f(\mathbf{x}) = \log E_w [I(y = 1) | \mathbf{x}] - \log E_w [I(y = -1) | \mathbf{x}]$$

Finally

We have that

$$\hat{f}(\mathbf{x}) = \frac{1}{2} \log \frac{E_w [I(y = 1) | \mathbf{x}]}{E_w [I(y = -1) | \mathbf{x}]}$$

in term of probabilities

$$\hat{f}(\mathbf{x}) = \frac{1}{2} \log \frac{P_w(y = 1 | \mathbf{x})}{P_w(y = -1 | \mathbf{x})}$$

Finally

We have that

$$\hat{f}(\mathbf{x}) = \frac{1}{2} \log \frac{E_w [I(y = 1) | \mathbf{x}]}{E_w [I(y = -1) | \mathbf{x}]}$$

In term of probabilities

$$\hat{f}(\mathbf{x}) = \frac{1}{2} \log \frac{P_w(y = 1 | \mathbf{x})}{P_w(y = -1 | \mathbf{x})}$$

The Weight Update

Finally, we have a way to update the weights by setting

$$w_t(\mathbf{x}, y) = e^{-yF(\mathbf{x})}$$

$$w_{t+1}(\mathbf{x}, y) = w_t(\mathbf{x}, y) e^{-y\hat{f}(\mathbf{x})}$$

Additionally, the weighted conditional mean

Corollary

- At the Optimal $F(\mathbf{x})$, the weighted conditional mean of y is 0.

Proof

- When $F(\mathbf{x})$ is optimal

$$\frac{\partial J(F(\mathbf{x}))}{\partial F(\mathbf{x})} = \frac{\partial \{ P(Y=1|\mathbf{x}) e^{-yF(\mathbf{x})} + P(Y=-1|\mathbf{x}) e^{yF(\mathbf{x})} \}}{\partial F(\mathbf{x})}$$

Additionally, the weighted conditional mean

Corollary

- At the Optimal $F(\mathbf{x})$, the weighted conditional mean of y is 0.

Proof

- When $F(\mathbf{x})$ is optimal

$$\frac{\partial J(F(\mathbf{x}))}{\partial F(\mathbf{x})} = \frac{\partial \{P(Y = 1|\mathbf{x}) e^{-yF(\mathbf{x})} + P(Y = -1|\mathbf{x}) e^{yF(\mathbf{x})}\}}{\partial F(\mathbf{x})}$$

Therefore

We have

$$\frac{\partial J(F(\mathbf{x}))}{\partial F(\mathbf{x})} = [P(Y = 1|\mathbf{x}) e^{-yF(\mathbf{x})}] \{-y\} + [P(Y = -1|\mathbf{x}) e^{-yF(\mathbf{x})}] \{-y\}$$

Therefore

$$E \left[e^{yF(\mathbf{x})} y \right] = 0$$

Therefore

We have

$$\frac{\partial J(F(\mathbf{x}))}{\partial F(\mathbf{x})} = [P(Y = 1|\mathbf{x}) e^{-yF(\mathbf{x})}] \{-y\} + [P(Y = -1|\mathbf{x}) e^{-yF(\mathbf{x})}] \{-y\}$$

Therefore

$$E \left[e^{yF(\mathbf{x})} y \right] = 0$$

Outline

- 1 Combining Models
 - Introduction
 - Average for Committee
 - Beyond Simple Averaging
 - Example
- 2 Bayesian Model Averaging
 - Model Combination Vs. Bayesian Model Averaging
 - Now Model Averaging
 - The Differences
- 3 Committees
 - Introduction
 - Bootstrap Data Sets
 - Relation with Monte-Carlo Estimation
- 4 **Boosting**
 - AdaBoost Development
 - Cost Function
 - Selection Process
 - How do we select classifiers?
 - Selecting New Classifiers
 - Deriving against the weight α_m
 - AdaBoost Algorithm
 - Some Remarks
 - Explanation about AdaBoost's behavior
 - Statistical Analysis of the Exponential Loss
 - Moving from Regression to Classification
 - Minimization of the Exponential Criterion
 - Finally, The Additive Logistic Regression
 - **Example using an Infinitude of Perceptrons**

Here, we decide to use Perceptrons

As Weak Learners

- We could be using a finite number of Perceptrons
- But we want to have a infinitude of possible weak learners
 - ▶ Thus avoiding the need of a matrix S

Here, we decide to use Perceptrons

As Weak Learners

- We could be using a finite number of Perceptrons
- But we want to have a infinitude of possible weak learners
 - ▶ Thus avoiding the need of a matrix S

Remark

- We need to use a Gradient Based Learner for this

Here, we decide to use Perceptrons

As Weak Learners

- We could be using a finite number of Perceptrons
- But we want to have a infinitude of possible weak learners
 - ▶ Thus avoiding the need of a matrix S

Remark

- We need to use a Gradient Based Learner for this

Perceptron

We use the following formula of error per sample

$$E(\mathbf{w}) = \frac{1}{2} \sum_{j=1}^N (w_j(t) y_j(t) - d_j)^2$$

- With $y_j(t) = \varphi(\mathbf{w}^T(t) \mathbf{x}_j)$

Deriving against w_i

$$\frac{\partial E(\mathbf{w})}{\partial w_i} = \sum_{j=1}^N (w_j(t) y_j(t) - d_j) \varphi'(\mathbf{w}^T(t) \mathbf{x}_j) w_i^j x_{ij}$$

Then, using gradient descent, we have the following update

$$w_i(n+1) = w_i(n) - \eta \left[\sum_{j=1}^N (w_j(t) y_j(t) - d_j) \varphi'(\mathbf{w}^T(t) \mathbf{x}_j) w_i^j x_{ij} \right]$$

Perceptron

We use the following formula of error per sample

$$E(\mathbf{w}) = \frac{1}{2} \sum_{j=1}^N (w_j(t) y_j(t) - d_j)^2$$

- With $y_j(t) = \varphi(\mathbf{w}^T(t) \mathbf{x}_j)$

Deriving against w_i

$$\frac{\partial E(\mathbf{w})}{\partial w_i} = \sum_{j=1}^N (w_j(t) y_j(t) - d_j) \varphi'(\mathbf{w}^T(t) \mathbf{x}_j) w_j^b x_{ij}$$

Then, using gradient descent, we have the following update

$$w_i(n+1) = w_i(n) - \eta \left[\sum_{j=1}^N (w_j(t) y_j(t) - d_j) \varphi'(\mathbf{w}^T(t) \mathbf{x}_j) w_j^b x_{ij} \right]$$

Perceptron

We use the following formula of error per sample

$$E(\mathbf{w}) = \frac{1}{2} \sum_{j=1}^N (w_j(t) y_j(t) - d_j)^2$$

- With $y_j(t) = \varphi(\mathbf{w}^T(t) \mathbf{x}_j)$

Deriving against w_i

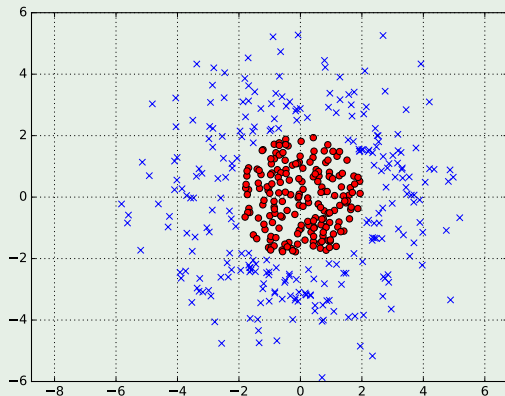
$$\frac{\partial E(\mathbf{w})}{\partial w_i} = \sum_{j=1}^N (w_j(t) y_j(t) - d_j) \varphi'(\mathbf{w}^T(t) \mathbf{x}_j) w_j^b x_{ij}$$

Then, using gradient descent, we have the following update

$$w_i(n+1) = w_i(n) - \eta \left[\sum_{j=1}^N (w_j(t) y_j(t) - d_j) \varphi'(\mathbf{w}^T(t) \mathbf{x}_j) w_j x_{ij} \right]$$

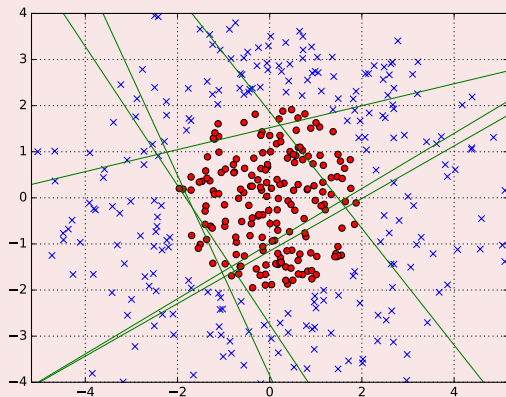
Data Set

Training set with classes $\omega_1 = N(0, 1)$ and $\omega_2 = N(0, \sigma^2) - N(0, 1)$



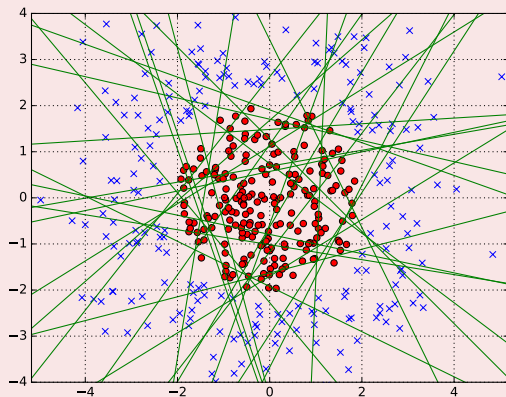
Example

For $m = 10$



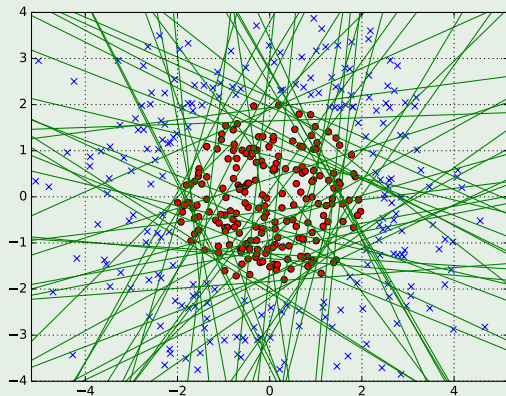
Example

For 40



At the end of the process

For $m = 80$



Final Confusion Matrix

When $m = 80$

	C_1	C_2
C_1	1.0	0.0
C_2	0.0	1.0

However

There are other versions to the Cryptic Phrase

- At “Boosting: Foundation and Algorithms” by Schaphire and Freund
 - ▶ “Train weak learner using distribution D_t ”

We could re-sample using the distribution D_t

- Basically using sampling with substitution over the data set $\{x_1, x_2, \dots, x_N\}$

However

There are other versions to the Cryptic Phrase

- At “Boosting: Foundation and Algorithms” by Schaphire and Freund
 - ▶ “Train weak learner using distribution D_t ”

We could re-sample using the distribution w_t

- Basically using sampling with substitution over the data set $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$

Other Interpretations exist

But you can use a weighted version of the cost function

$$\frac{1}{2} \sum_j w_j(t) (y_j(t) - d_j)^2$$

For More, Take a look

- “Boosting Neural Networks” by Holger Schwenk and Yoshua Bengio