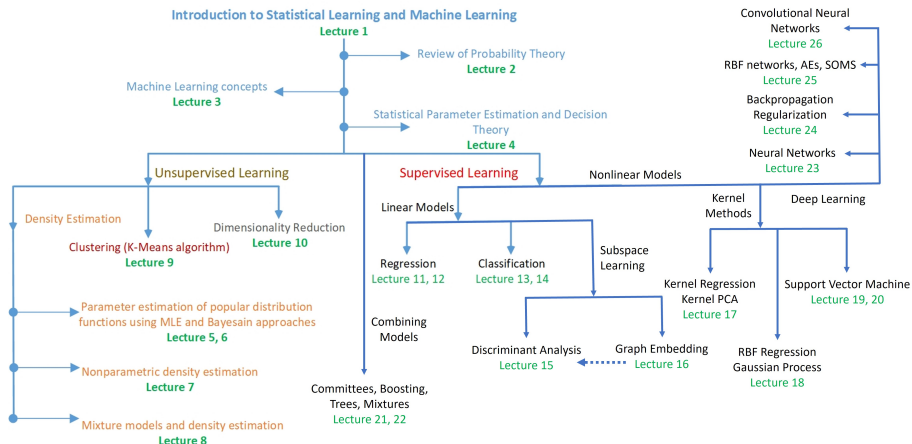


# Statistical Learning and Machine Learning

## Lecture 21 - Combining Models 1

October 13, 2021

# Course overview and where do we stand



# Why to combine multiple models?

Sometimes combining multiple models can lead to better performance compared to using only one:

- *Committees*: Use  $L$  different models and then make predictions using the average of the predictions of these  $L$  models
- *Boosting*: Use multiple models *in a sequence* in which each model's training is depending on the models preceding it in the sequence
- *Decision trees*: Divide the space in multiple regions and train one model for each region.
- *Mixture of experts*: Train  $K$  models and combine them based on a probabilistic mixture of the form:

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

Given a set of data points and their labels:

- We sample  $M$  subsets independently and we train a model on each of these subsets  $y_m(x)$
- The committee prediction is:

$$y_{COM} = \frac{1}{M} \sum_{m=1}^M y_m(x) \quad (2)$$

This procedure is called bootstrap aggregation or *bagging*.

# Committees

Suppose the function we want to predict is  $h(w)$ , then:

$$y_m(x) = h(x) + \epsilon_m(x) \quad (3)$$

The average sum-of-squares error has the form:

$$\mathbb{E}_x[(y_m(x) - h(x))^2] = \mathbb{E}_x[\epsilon_m(x)^2] \quad (4)$$

So, the average error by the models acting individually is:

$$E_{AV} = \frac{1}{M} \sum_{m=1}^M \mathbb{E}_x[\epsilon_m(x)^2] \quad (5)$$

The expected error of the committee is:

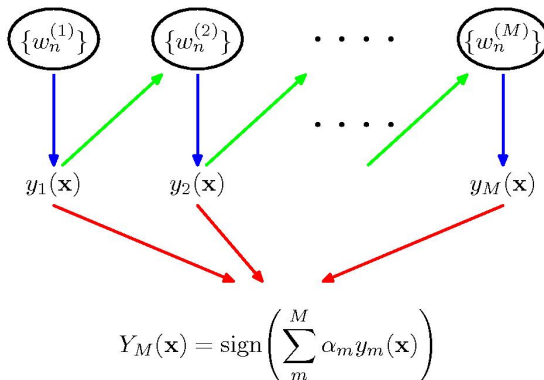
$$\mathbb{E}_x \left[ \left( \frac{1}{M} \sum_{m=1}^M y_m(x) - h(x) \right)^2 \right] = \mathbb{E}_x \left[ \left( \frac{1}{M} \sum_{m=1}^M \epsilon_m(x) \right)^2 \right] \quad (6)$$

For zero-mean and uncorrelated errors, we obtain:  $E_{COM} = \frac{1}{M} E_{AV}$ .

# AdaBoost

## Boosting:

- Combination of multiple classifiers in a sequence
- The *base* classifiers do not need to be highly performing (*weak classifiers*)
- The combination of multiple weak classifiers leads to a high-performing committee of classifiers



## AdaBoost

1. Initialize the data weighting coefficients  $\{w_n\}$  by setting  $w_n^{(1)} = 1/N$  for  $n = 1, \dots, N$ .
2. For  $m = 1, \dots, M$ :

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

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

where  $I(y_m(\mathbf{x}_n) \neq t_n)$  is the indicator function and equals 1 when  $y_m(\mathbf{x}_n) \neq t_n$  and 0 otherwise.

- (b) Evaluate the quantities

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

and then use these to evaluate

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}.$$

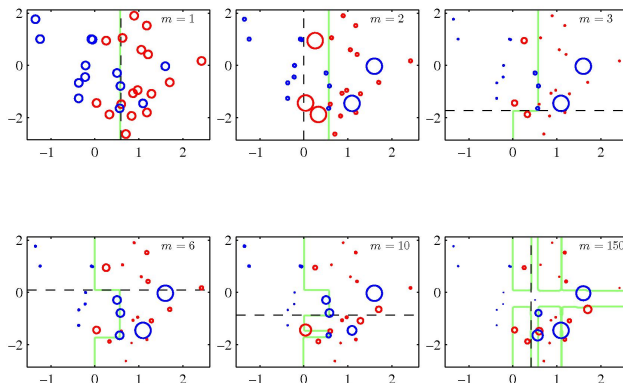
- (c) Update the data weighting coefficients

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

3. Make predictions using the final model, which is given by

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

# AdaBoost



Dashed black line: decision line of the most recent learner

Solid green line: combined decision line of the ensemble

Radius of points: weight assigned for training the most recent learner



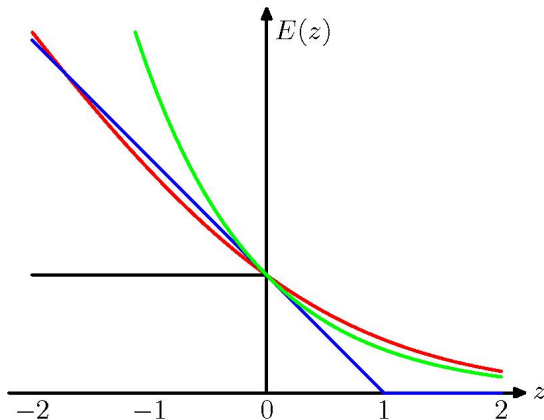
Error function for AdaBoost:

$$E = \sum_{n=1}^N \exp(-t_n f_m(x_n)) \quad (7)$$

where  $t_n \in \{-1, 1\}$  and:

$$f_m(x) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(x) \quad (8)$$

# AdaBoost



Green: Exponential function  
Red: Cross-entropy function  
Blue: Hinge error (used in SVM)  
Black: misclassification error