

Ensemble Learning

CS4780/5780 – Machine Learning
Fall 2013

Igor Labutov
Cornell University

Ensemble Learning

- A class of “meta” learning algorithms
- Combining multiple classifiers to increase performance
- Very effective in practice
- Good theoretical guarantees
- Easy to implement!

Ensemble

Problem : given T binary classification hypotheses (h_1, \dots, h_T) , **find** a combined classifier:

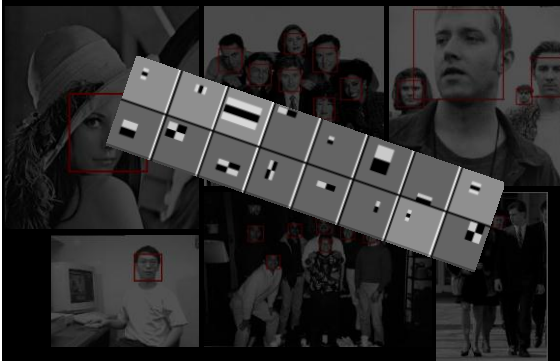
$$h_S(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

with better performance.

Teaser



Teaser



BAGGING



Bagging

Bagging (Bootstrap aggregating).

(Breiman, 1996)

```

BAGGING( $S = ((x_1, y_1), \dots, (x_m, y_m))$ )
1 for  $t \leftarrow 1$  to  $T$  do
2    $S_t \leftarrow \text{BOOTSTRAP}(S)$  > i.i.d. sampling with replacement from  $S$ .
3    $h_t \leftarrow \text{TRAINCLASSIFIER}(S_t)$ 
4 return  $h_S = x \mapsto \text{MAJORITYVOTE}((h_1(x), \dots, h_T(x)))$ 
    
```

Bagging

Ensemble :

$$h_S(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

Bagging : Special case where we fix:

$$\alpha_t = 1 \quad \text{and} \quad h_t = \mathbb{L}(S_t)^*$$

* \mathbb{L} is some learning algorithm

S_t is a training set drawn from distribution $P(\langle x, y \rangle)$

Bias-Variance Tradeoff

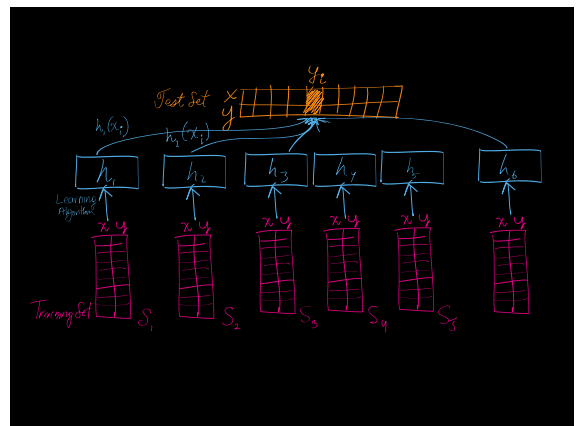
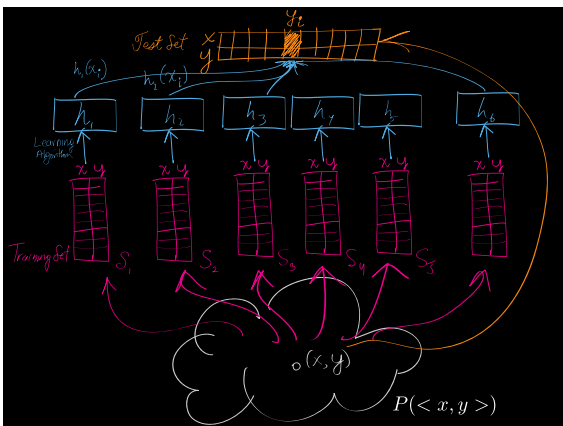
Generalization Error

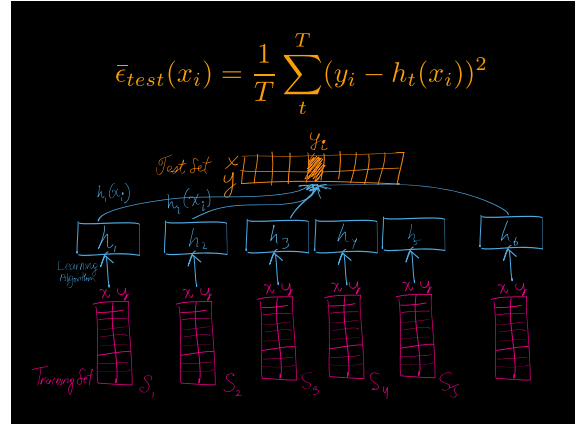
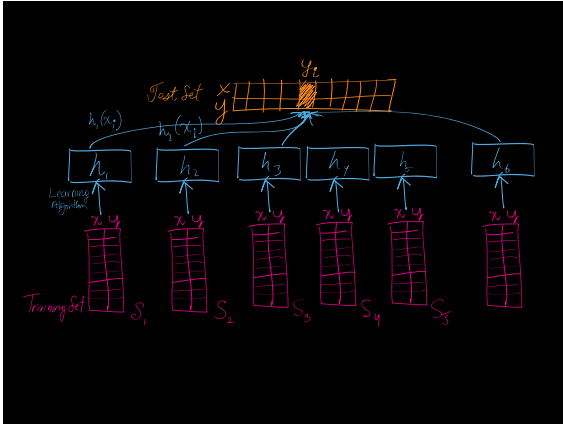
Classification :

$$\epsilon_{test} = \frac{1}{n} \sum_i \text{Zero-One-Loss}(y_i, h(x_i))$$

Regression :

$$\epsilon_{test} = \frac{1}{n} \sum_i (y_i - h(x_i))^2$$





$$\bar{\epsilon}_{test}(x_i) = \frac{1}{T} \sum_t (y_i - h_t(x_i))^2$$

OR, as an expectation:

$$\mathbb{E}_S [(y_i - h_S(x_i))^2]$$

For the entire test set:

$$\mathbb{E}_{X,Y} \mathbb{E}_S [(y_i - h_S(x_i))^2]$$

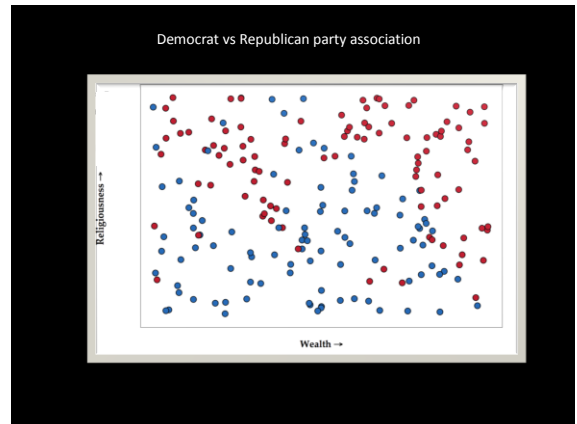
CLAIM:

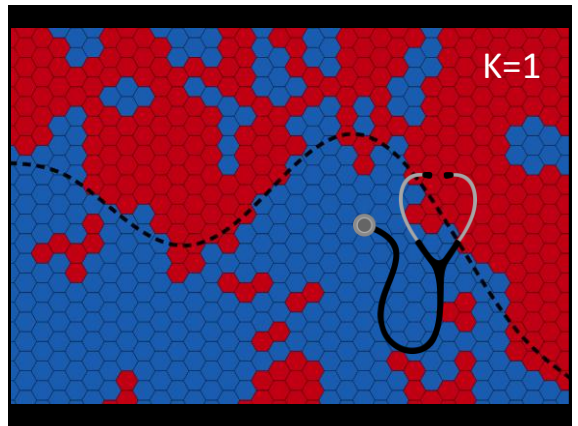
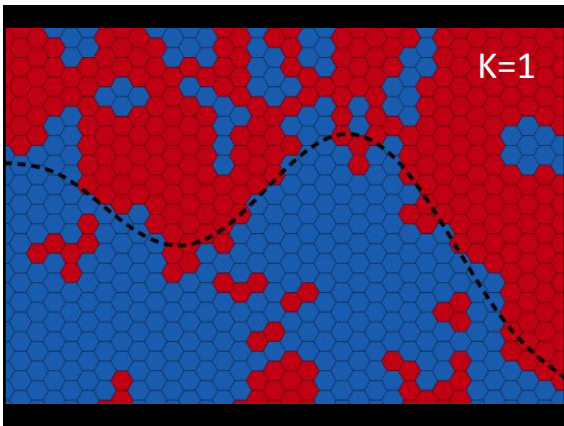
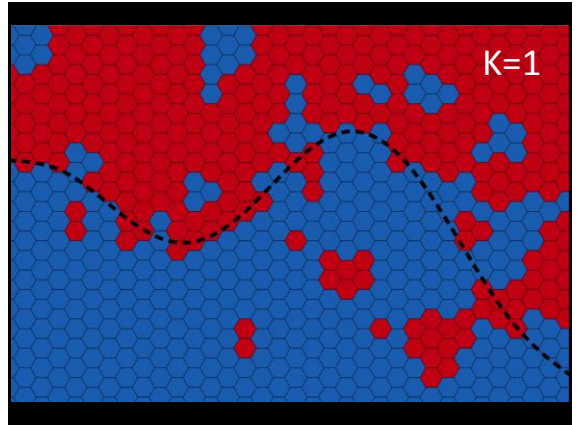
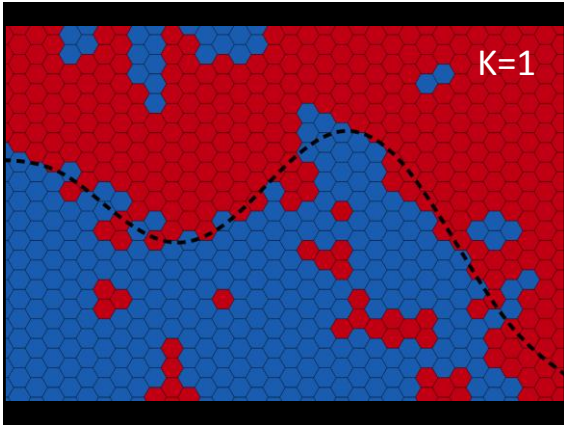
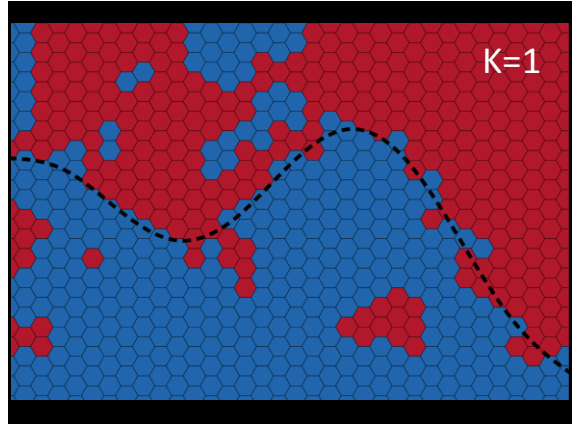
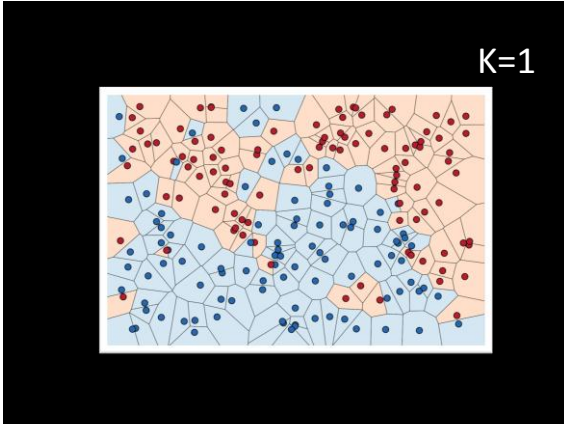
$$\mathbb{E}_S [(y_i - h_S(x_i))^2] =$$

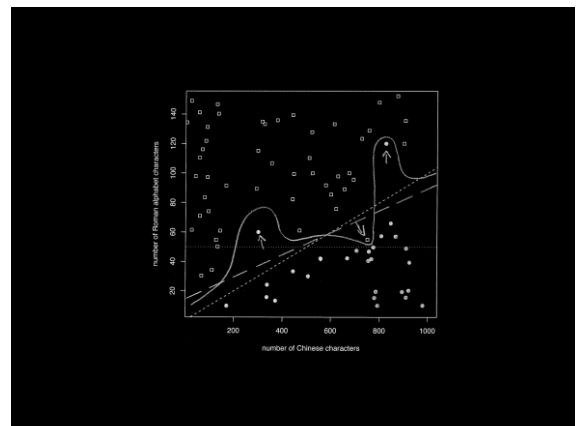
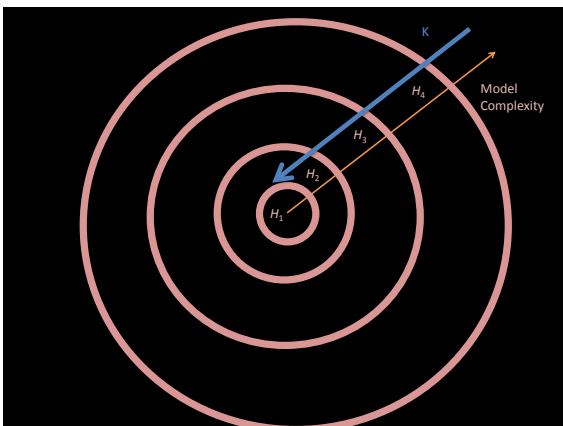
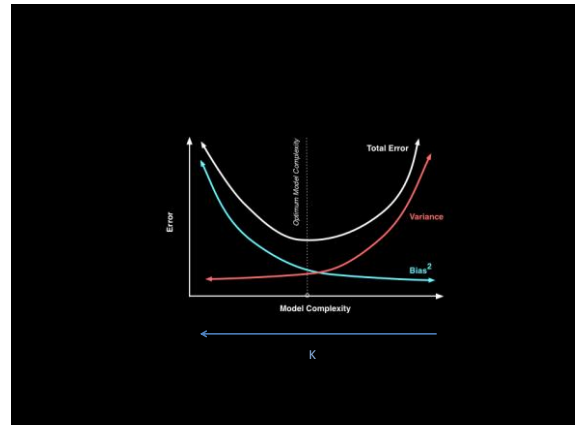
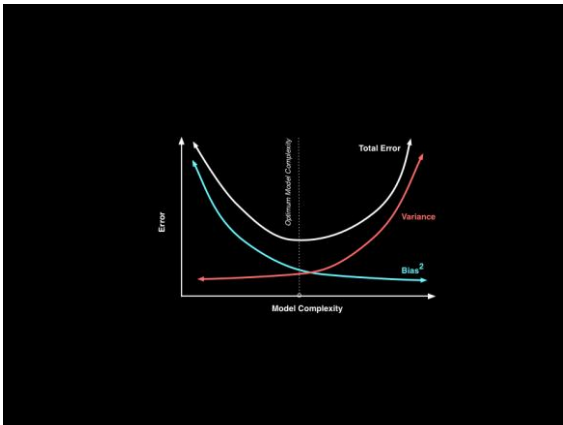
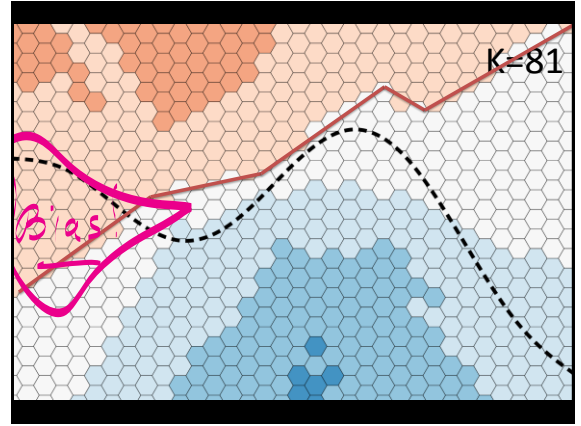
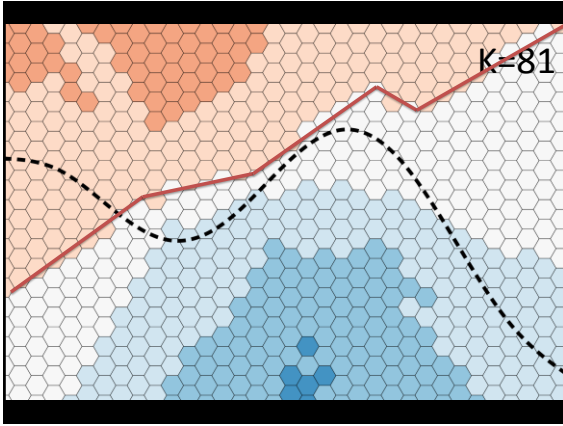
bias² $(y_i - \mathbb{E}_S[h_S(x_i)])^2 +$

variance $+ \mathbb{E}_S[(h_S(x_i) - \mathbb{E}_S[h_S(x_i)])^2]$

Example
(kNN)







CLAIM:

$$\begin{aligned} \mathbb{E}_S [(y_i - h_S(x_i))^2] = \\ \text{bias}^2 & (y_i - \mathbb{E}_S[h_S(x_i)])^2 + \\ \text{variance} & + \mathbb{E}_S[(h_S(x_i) - \mathbb{E}_S[h_S(x_i)])^2] \end{aligned}$$

USEFUL LEMMA:

$$\mathbb{E}[(\alpha - \mathbb{E}[\alpha])^2] = \mathbb{E}[\alpha^2] - \mathbb{E}[\alpha]^2$$

$$y_i = f(x_i) + \mathcal{N}(0, \sigma^2)$$

$$\begin{aligned} \mathbb{E}_S [(y_i - h_S(x_i))^2] = \\ \text{bias}^2 & (y_i - \mathbb{E}_S[h_S(x_i)])^2 + \\ \text{variance} & + \mathbb{E}_S[(h_S(x_i) - \mathbb{E}_S[h_S(x_i)])^2] \end{aligned}$$

$$y_i = f(x_i) + \mathcal{N}(0, \sigma^2)$$

$$\begin{aligned} \mathbb{E}_S [(y_i - h_S(x_i))^2] = \\ \text{bias}^2 & (f(x_i) - \mathbb{E}_S[h_S(x_i)])^2 + \\ \text{variance} & + \mathbb{E}_S[(h_S(x_i) - \mathbb{E}_S[h_S(x_i)])^2] \\ \text{noise} & + \mathbb{E}_S[(f(x_i) - y_i)^2] \end{aligned}$$

$$y_i = f(x_i) + \mathcal{N}(0, \sigma^2)$$

$$\begin{aligned} \mathbb{E}_S [(y_i - h_S(x_i))^2] = \\ \text{bias}^2 & (f(x_i) - \mathbb{E}_S[h_S(x_i)])^2 + \\ \text{variance} & + \mathbb{E}_S[(h_S(x_i) - \mathbb{E}_S[h_S(x_i)])^2] \\ \text{noise} & + \sigma^2 \end{aligned}$$

$$\mathbb{E}_S [(y_i - h_S(x_i))^2] =$$

$$\begin{aligned} \text{bias}^2 & (y_i - \mathbb{E}_S[h_S(x_i)])^2 + \\ \text{variance} & + \mathbb{E}_S[(h_S(x_i) - \mathbb{E}_S[h_S(x_i)])^2] \end{aligned}$$

BAGGING

revisited



Bagging

Bagging (Bootstrap aggregating).

(Breiman, 1996)

```
BAGGING( $S = ((x_1, y_1), \dots, (x_m, y_m))$ )
1 for  $t \leftarrow 1$  to  $T$  do
2    $S_t \leftarrow \text{BOOTSTRAP}(S)$  > i.i.d. sampling with replacement from  $S$ .
3    $h_t \leftarrow \text{TRAINCLASSIFIER}(S_t)$ 
4 return  $h_S = x \mapsto \text{MAJORITYVOTE}((h_1(x), \dots, h_T(x)))$ 
```

Why does it work?

Bagging

Ensemble :

$$h_S(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

Bagging : Special case where we fix:

$$\alpha_t = 1 \quad \text{and} \quad h_t = \mathbb{L}(S_t)^*$$

* \mathbb{L} is some learning algorithm

S_t is a training set drawn from distribution $P(\langle x, y \rangle)$

Bagging

Bagging Ensemble :

$$h_S(x) = \text{sign} \left(\sum_{t=1}^T h_t(x) \right)$$

What happens to *bias* and *variance*?

Bagging

Bagging Ensemble (regression) :

$$h_S(x) = \frac{1}{T} \sum_{t=1}^T h_t(x)$$

*bias*²

$$(y_i - \mathbb{E}_S[h_S(x_i)])^2$$

variance

$$\mathbb{E}_S[(h_S(x_i) - \mathbb{E}_S[h_S(x_i)])^2]$$

Bagging

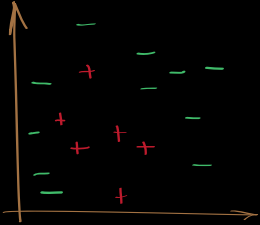
What happens to *bias* and *variance*?

$$\text{Bias}(h_S, x_i) = \frac{1}{T} \sum_{t=1}^T \text{Bias}(h_t, x_i)$$

$$\text{Var}(h_S, x_i) \approx \frac{1}{T} \text{Var}(h_1, x_i)$$

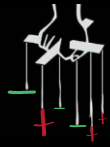
Bagging has approximately the same bias, but reduces variance of individual classifiers!

Bagging

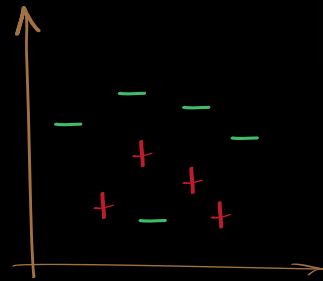


Bagging as a "Training set manipulator"

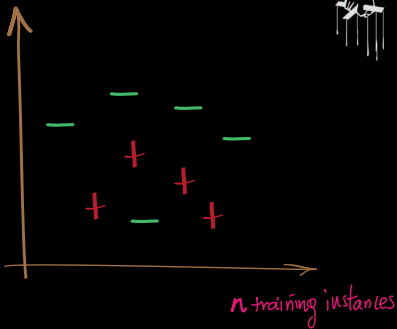
Bagging as a "Training set manipulator"



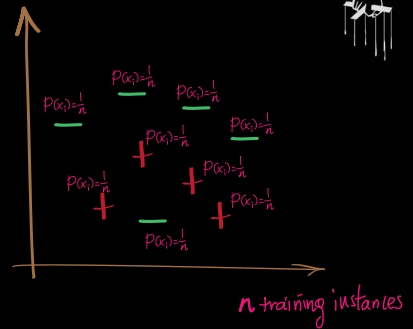
Bagging as a "Training set manipulator"



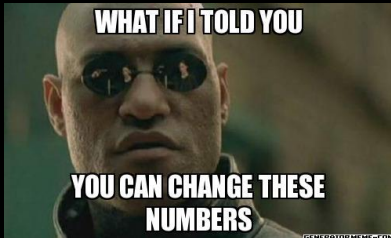
Bagging as a "Training set manipulator"



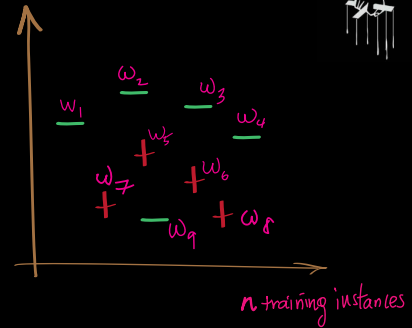
Bagging as a "Training set manipulator"



Bagging as a "Training set manipulator"



Bagging as a "Training set manipulator"



Ensemble

Problem : given T binary classification hypotheses (h_1, \dots, h_T) , **find** a combined classifier:

$$h_S(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

with better performance.

Teaser



Hypothetical Algorithm

Hypothetical Algorithm

(incomplete)

Given $x_i \in X, y_i \in Y = \{-1, 1\}$

where $(x_1, y_1), \dots, (x_n, y_n)$

Initialize $W_1(i) = 1/n$

Initialize set $H = \{h_1, \dots, h_T\}$

For $t = 1, \dots, T$:

- Pick hypothesis h_t out of the set H
- Compute error rate ϵ_t of h_t
- Assign new weights W_t to X
- Compute new weight α_t for h_t

Output $h_S(x) = \sum_{t=1}^T \alpha_t h_t(x)$

Hypothetical Algorithm

(incomplete)

Given $x_i \in X, y_i \in Y = \{-1, 1\}$

where $(x_1, y_1), \dots, (x_n, y_n)$

Initialize $W_1(i) = 1/n$

Learning algorithm \mathbb{L}

For $t = 1, \dots, T$:

- Generate hypothesis h_t with \mathbb{L}
- Compute error rate ϵ_t of h_t
- Assign new weights W_t to X
- Compute new weight α_t for h_t

Output $h_S(x) = \sum_{t=1}^T \alpha_t h_t(x)$

Hypothetical Algorithm

(incomplete)

Given $x_i \in X, y_i \in Y = \{-1, 1\}$

where $(x_1, y_1), \dots, (x_n, y_n)$

Initialize $W_1(i) = 1/n$

Initialize set $H = \{h_1, \dots, h_T\}$

For $t = 1, \dots, T$:

- Pick hypothesis h_t out of the set H
- Compute error rate ϵ_t of h_t
- Assign new weights W_t to X
- Compute new weight α_t for h_t

Output $h_S(x) = \sum_{t=1}^T \alpha_t h_t(x)$

Hypothetical Algorithm

(incomplete)

Given $x_i \in X, y_i \in Y = \{-1, 1\}$

where $(x_1, y_1), \dots, (x_n, y_n)$

Initialize $W_1(i) = 1/n$

Initialize set $H = \{h_1, \dots, h_T\}$

For $t = 1, \dots, T$:

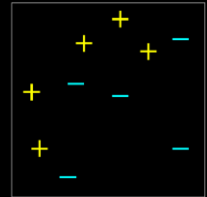
- Pick hypothesis h_t out of the set H
- Compute error rate ϵ_t of h_t
- Assign new weights W_t to X
- Compute new weight α_t for h_t

Output $h_S(x) = \sum_{t=1}^T \alpha_t h_t(x)$

Hypothetical Algorithm

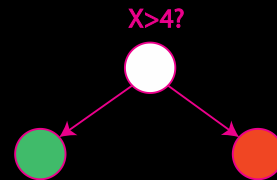
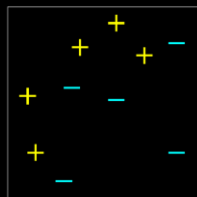
Toy Example

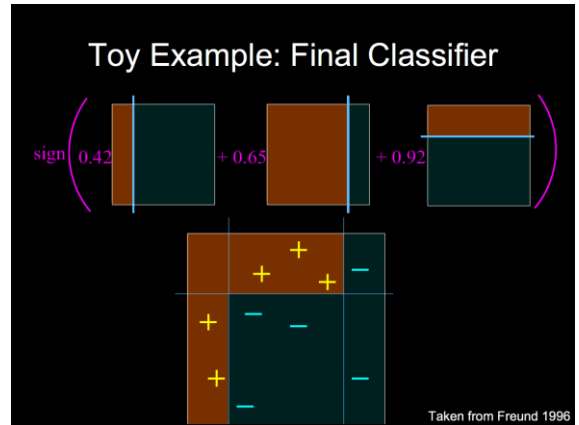
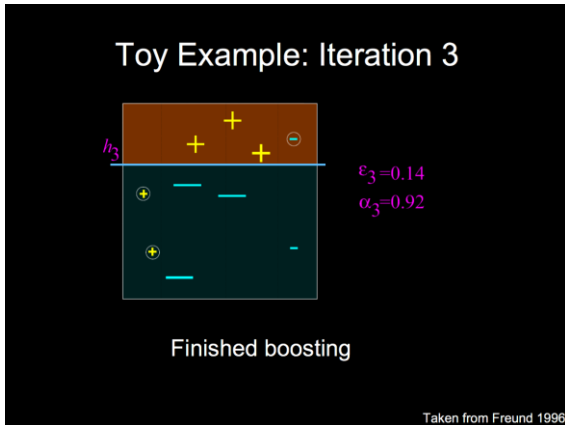
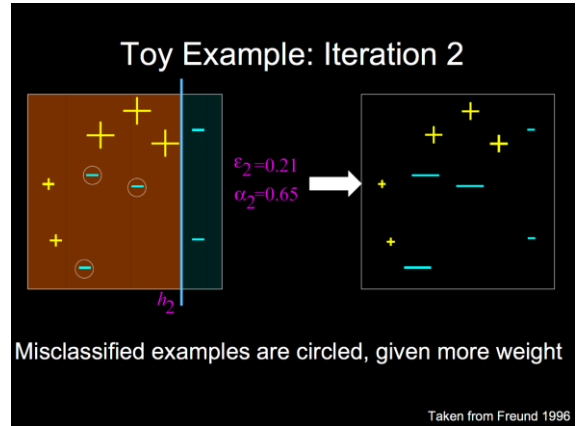
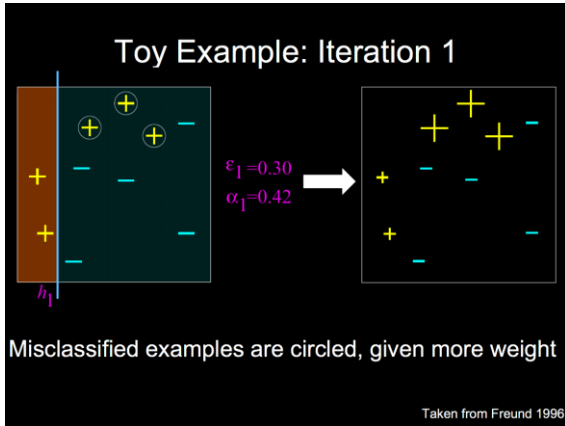
- Positive examples
- Negative examples
- 2-Dimensional plane
- Weak hyps: linear separators
- 3 iterations



Toy Example

- Positive examples
- Negative examples
- 2-Dimensional plane
- Weak hyps: linear separators
- 3 iterations





Questions

- Which hypothesis do we choose at every iteration?
- How should we weight the hypotheses?
- How should we weight the examples?

Answers

Choose h_t that maximizes $W_{correct}$ (minimizes ϵ_t)

Choose α_t according to:

$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Update the weight of instance i as follows:

$$w_t(i) = w_{t-1}(i) * e^{-\alpha_t} \quad \text{if } y_i = h_t(x_i)$$

$$w_t(i) = w_{t-1}(i) * e^{\alpha_t} \quad \text{if } y_i \neq h_t(x_i)$$

AdaBoost

Hypothetical Algorithm

(incomplete)

Given $x_i \in X, y_i \in Y = \{-1, 1\}$

where $(x_1, y_1), \dots, (x_n, y_n)$

Initialize $W_1(i) = 1/n$

Initialize set $H = \{h_1, \dots, h_T\}$

For $t = 1, \dots, T$:

- Pick hypothesis h_t out of the set H
- Compute error rate ϵ_t of h_t
- Assign new weights W_t for X
- Compute new weight α_t for h_t

Output $h_S(x) = \sum_{t=1}^T \alpha_t h_t(x)$

Hypothetical Algorithm

(incomplete)

Given $x_i \in X, y_i \in Y = \{-1, 1\}$

where $(x_1, y_1), \dots, (x_n, y_n)$

Initialize $W_1(i) = 1/n$

Initialize set $H = \{h_1, \dots, h_T\}$

For $t = 1, \dots, T$:

- Pick hypothesis h_t with smallest ϵ_t
- Compute weight α_t on h_t
- Update weights W_t for X

Output $h_S(x) = \sum_{t=1}^T \alpha_t h_t(x)$

Training Error for AdaBoost

Write for some h_t weighted error ϵ_t :

$$\epsilon_t = \frac{1}{2} - \gamma_t$$

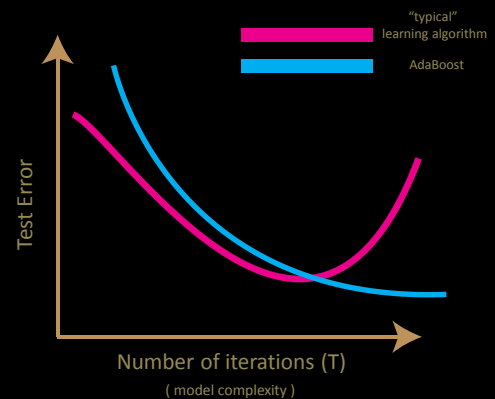
We can then bound the training error:

$$\text{training error} \leq \exp(-2T\gamma^2)$$

For some γ such that:

$$\gamma_t \geq \gamma > 0$$

What about
Generalization Error?



Why?

Margin

$$\text{margin}_f(x, y) =$$

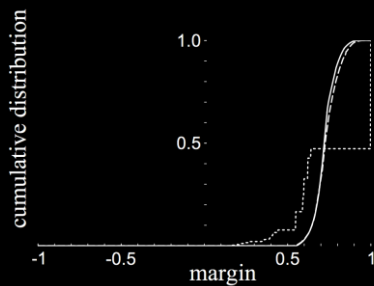
Margin

$$\text{margin}_f(x, y) = \frac{yf(x)}{\sum_t |\alpha_t|} =$$

Margin

$$\text{margin}_f(x, y) = \frac{yf(x)}{\sum_t |\alpha_t|} = \frac{y \sum_t \alpha_t h_t(x)}{\sum_t |\alpha_t|}$$

Margin



Viola Jones Classifier



Image Features

“Rectangle filters”

Value =

$$\frac{\sum (\text{pixels in white area}) - \sum (\text{pixels in black area})}{\sum (\text{pixels in white area}) - \sum (\text{pixels in black area})}$$

Source

Result

Fast computation with integral images

- The *integral image* computes a value at each pixel (x,y) that is the sum of the pixel values above and to the left of (x,y) , inclusive
- This can quickly be computed in one pass through the image

Computing sum within a rectangle

- Let A,B,C,D be the values of the integral image at the corners of a rectangle
- Then the sum of original image values within the rectangle can be computed as:

$$\text{sum} = A - B - C + D$$
- Only 3 additions are required for any size of rectangle!
 - This is now used in many areas of computer vision

Example

“Rectangle filters”



A

B

Similar to Haar wavelets

Papageorgiou, et al.

$$h_i(x_i) = \begin{cases} \alpha_i & \text{if } f_i(x_i) > \theta_i \\ \beta_i & \text{otherwise} \end{cases}$$

C



D

$$C(x) = \theta \left(\sum_i h_i(x) + b \right)$$

60,000 features to choose from