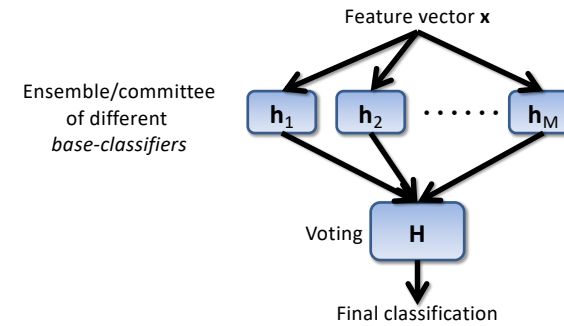


Neural Networks and Learning Systems
TBMI26 / 732A55
2019

Lecture 4
Ensemble Learning

Magnus Borga
magnus.borga@liu.se

Classifier ensemble

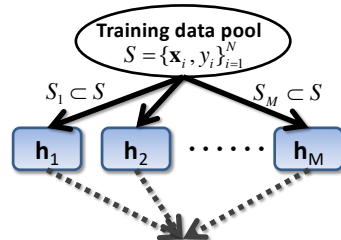


2

Bootstrap Aggregating (Bagging)

Breiman, 1994

Train each base-classifier using a subset of the training data



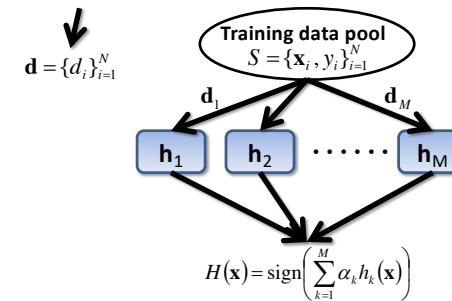
Reduced variance – reduced risk of overfitting

3

Boosting

Schapire and others, (1989-1990)

Train each base-classifier using all training data but with weights indicating how important each training sample is.



4

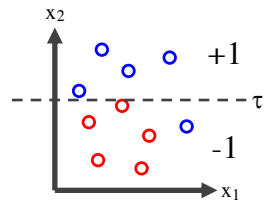
Simple/Weak classifiers

cf. "a rule of thumb"

Example: Threshold **one** feature

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix}$$

$$h(x_2) = \begin{cases} +1 & x_2 \geq \tau \\ -1 & x_2 < \tau \end{cases}$$



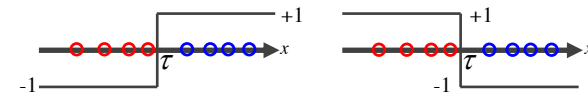
5

Decision stump

Polarity $p = \{-1, 1\}$

$$h(x) = \begin{cases} +1 & x \geq \tau \\ -1 & x < \tau \end{cases} \quad \text{for } p = 1$$

$$h(x) = \begin{cases} +1 & x \leq \tau \\ -1 & x > \tau \end{cases} \quad \text{for } p = -1$$

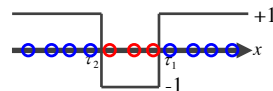


6

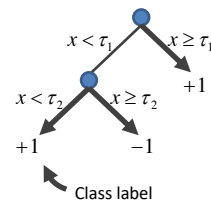
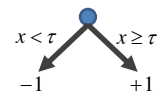
Classification and Regression Trees (CART)

$$h(x) = \begin{cases} +1 & x \geq \tau \\ -1 & x < \tau \end{cases}$$

$$h(x) = \begin{cases} +1 & x \geq \tau_1 \\ -1 & x < \tau_1 \text{ and } x \geq \tau_2 \\ +1 & x < \tau_2 \end{cases}$$

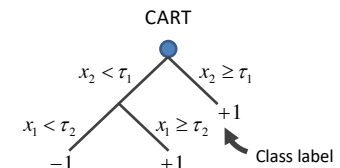
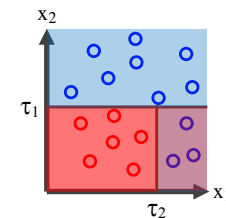


Decision stump



7

CART – 2D example



Piecewise flat classification function

$$f(\mathbf{x}; \mathbf{w}_1, \dots, \mathbf{w}_k) \rightarrow \Omega$$

Feature index and thresholds

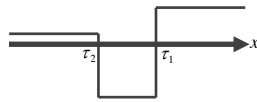
8

Regression Tree

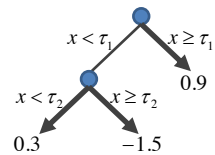
Same, but with real-valued output!

$$h(x) = \begin{cases} 0.75 & x \geq \tau \\ 0.25 & x < \tau \end{cases}$$

$$h(x) = \begin{cases} 0.9 & x \geq \tau_1 \\ -1.5 & x < \tau_1 \text{ and } x \geq \tau_2 \\ 0.3 & x < \tau_2 \end{cases}$$

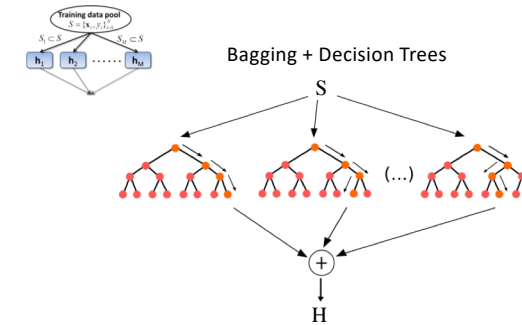


Decision stump



9

Random Forest



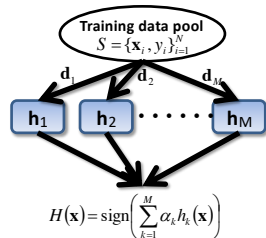
1. For each tree, select a random subset of the training samples
 2. For each split, select a random subset of the features
- Make the trees "independent"

Image from <https://towardsdatascience.com/random-forest-learning-essential-understanding-1ca586963cb>

10

General boosting algorithm

Train weak classifiers sequentially!



1. Set weights $d_1 = 1/N$
2. Train weak classifier $h_1(x)$ using weights d_1
3. Increase and decrease weight for wrongly and correctly classified training examples respectively $\rightarrow d_2$
4. Train weak classifier $h_2(x)$ using weights d_2
5. Repeat until $h_M(x)$

11

Training a decision stump

Find best split threshold τ !

Training input: $\{x_i, y_i, d_i\}_{i=1}^M$ Normalized weights: $\sum_{i=1}^M d_i = 1$

Class label $\{-1, +1\}$ Consider only one feature Weight

Threshold function: $h(x; \tau, p) = \begin{cases} +1 & p x \geq p \tau \\ -1 & p x < p \tau \end{cases}$

Polarity $\{-1, 1\}$

0-1 cost function: $\min_{\tau, p} \varepsilon(\tau, p) = \sum_{i=1}^M d_i I(y_i \neq h(x_i; \tau, p))$

1 for false classifications

12

Training a decision stump, cont.

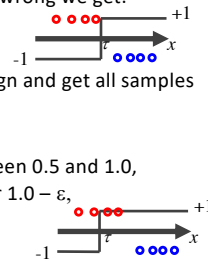
$$\min_{\tau, p} \varepsilon(\tau, p) = \sum_{i=1}^M d_i I(y_i \neq h(x_i; \tau, p)) \text{ is always } \leq 0.5!$$

Why? If we classify all training samples wrong we get:

$$\varepsilon(\tau, p) = \sum_{i=1}^M d_i = 1$$

But we can then just change polarity/sign and get all samples correct, i.e., $\varepsilon = 0$!

In general, if we obtain an error ε between 0.5 and 1.0, we can switch polarity and get the error $1.0 - \varepsilon$, which is smaller than 0.5.

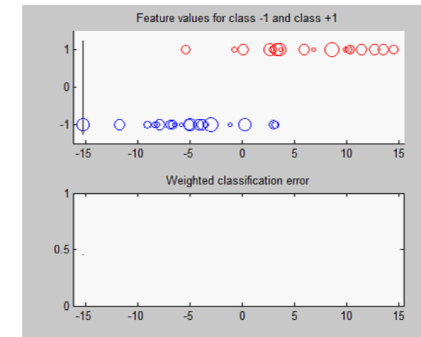


13

Brute force optimization

$$\min_{\tau, p} \varepsilon(\tau, p) = \sum_{i=1}^M d_i I(y_i \neq h(x_i; \tau, p))$$

Movie!

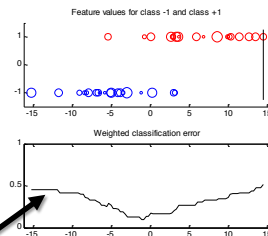


14

Brute force optimization

$$\min_{\tau, p} \varepsilon(\tau, p) = \sum_{i=1}^M d_i I(y_i \neq h(x_i; \tau, p))$$

Cost-function jumps at the x_i 's. Enough to test all thresholds in the set $\tau \in \{x_i\}_{i=1}^N$ and see which one gives the smallest error.



Cost function is piece-wise constant!

15

Brute force optimization

$$\text{Training samples } \mathbf{x}_i = \begin{pmatrix} x_{1,i} \\ \vdots \\ x_{N,i} \end{pmatrix}, i = 1 \cdots M$$

training samples

feature components

Pseudo code:

```

εmin = inf
for all feature components k = 1:N
  for all thresholds τ ∈ {xk,i}i=1M
    ε(τ, p = 1) = ∑i=1M I(yi ≠ h(xk,i; τ, p = 1))
    if ε > 0.5
      p = -1
      ε = 1 - ε
    end
    if ε < εmin ..... end
  end
end
end
  
```

16

Discrete AdaBoost

Freund & Schapire, 1995

Training data $\{\mathbf{x}_i, y_i\}_{i=1}^M$, $y_i \in \{-1, +1\}$

Initialization: $d_1(i) = \frac{1}{M}$, $T = \#$ base classifiers

for $t = 1$ to T

Find weak classifier $h_t(\mathbf{x}) \in \{-1, +1\}$ that minimizes the weighted classification error:

$$\varepsilon_t = \sum_{i=1}^M d_t(i) I(y_i \neq h_t(\mathbf{x}_i))$$

Update weights:

$$d_{t+1}(i) = d_t(i) e^{-\alpha_t y_i h_t(\mathbf{x}_i)}, \text{ where } \alpha_t = \frac{1}{2} \ln \frac{1 - \varepsilon_t}{\varepsilon_t}$$

and renormalize so that $\sum_{i=1}^M d_{t+1}(i) = 1$

end

Final strong classifier: $H(\mathbf{x}) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x})\right)$

17

Discrete AdaBoost

- Discrete output from the weak classifier

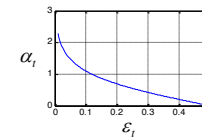
$$h_t(\mathbf{x}) \in \{-1, +1\}$$

- Weight update

$$d_{t+1}(i) = d_t(i) e^{-\alpha_t y_i h_t(\mathbf{x}_i)} = \begin{cases} d_t(i) e^{-\alpha_t} & \text{If } \mathbf{x}_i \text{ correctly classified} \\ d_t(i) e^{\alpha_t} & \text{If } \mathbf{x}_i \text{ wrongly classified} \end{cases}$$

- Performance of weak classifier:

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \varepsilon_t}{\varepsilon_t}$$

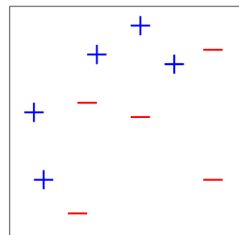


$$\text{Final strong classifier: } H(\mathbf{x}) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x})\right)$$

18

Discrete AdaBoost – Toy example

Initial weights $d_1(i) = \frac{1}{10}$, $T = 3$

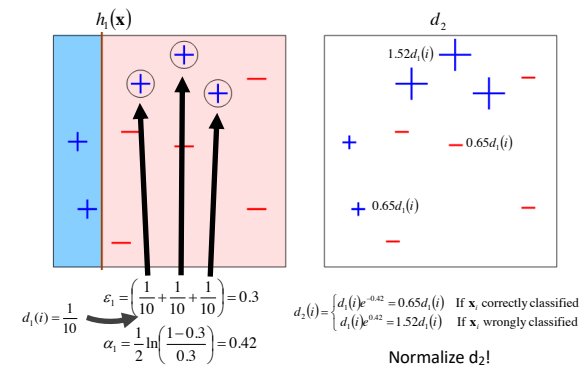


10 training samples $\mathbf{x}_i, i = 1 \dots 10$

Images borrowed from R. Schapire, "A Boosting Tutorial"

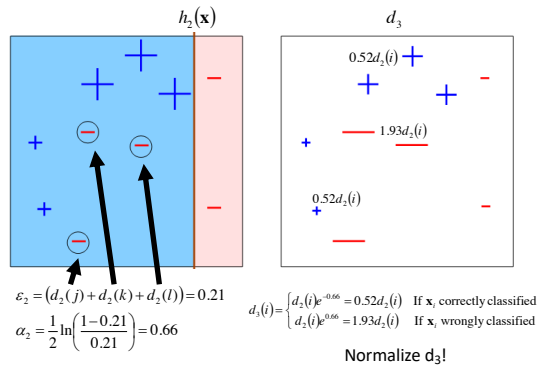
19

Discrete AdaBoost – Toy example



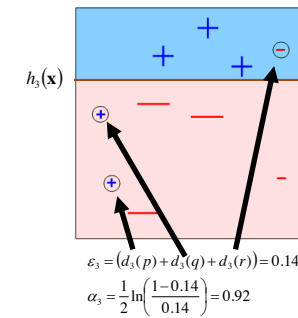
20

Discrete AdaBoost – Toy example



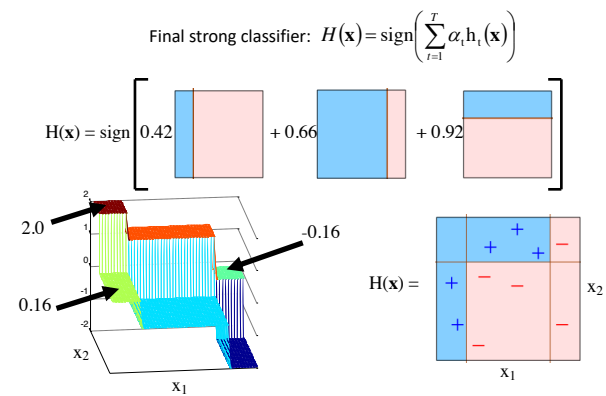
21

Discrete AdaBoost – Toy example



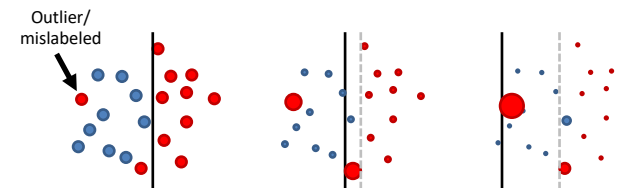
22

Discrete AdaBoost – Toy example



23

Problem - Outliers



- Outliers gain weight exponentially
- Will eventually result in bad weak classifiers
- May ruin the final ensemble classifier

24

Outlier strategies

- Keep an eye on the weights (plot them!)
- Weight trimming
 - Don't allow weights larger than a certain threshold
 - Disregard training samples with too large weights
- Use alternative weight update schemes with less aggressive increases for misclassified training data
 - LogitBoost
 - GentleBoost

25

Summary Ensemble Learning

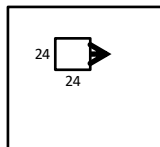
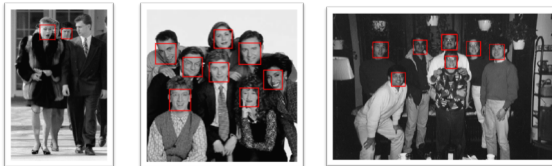
(using decision trees)

- Nonlinear classifiers that are easy to implement
- Easy to use - just one or a few parameters (T)
- Inherent feature selection
- Slow to train – fast to classify (real-time)
- Look out for outlier problems (AdaBoost)

26

Real-time object detection

P. Viola and M. Jones "Rapid Object Detection using a Boosted Cascade of Simple Features", 2001

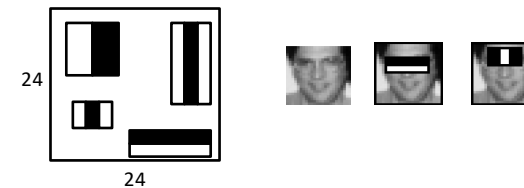


Sweep a sub-window over the image. for each position, determine if the sub-window contains a face or not.

27

Haar-features

Rectangle filters



From the sub-window, calculate contrast features (Haar features) at different locations and scales, and in different orientations. LOTS of different combinations, can be several 100,000 features!

28

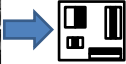
Train detector with AdaBoost

As described previously!

24 x 24 pixels face images



Apply Haar filters
to each image



$$\mathbf{x}_i = \begin{pmatrix} x_{1,i} \\ \vdots \\ x_{N,i} \end{pmatrix}, i=1 \dots M$$

Similar with non-face images to obtain negative examples.

29