

**Filter** - pre classifier.

- information gain  
"how much information about the Classification the feature provides"
- variance of individual features  
""
- correlation among features

**Wrapper** -

#### Filter Methods

Pros: Fast

Cons: Chosen filter might not be relevant for a specific kind of classifier.

Doesn't take into account interactions among features  
Often hard to know how

#### Embedded methods

-Result is that most of the weights go to zero, leaving a small subset of the weights.

information gain. How it's used for feature selection.

$$\text{Log}[a]b = (\log[10]b / \log[10]a)$$

many features to select.

#### Wrapper Methods

Pros: Features are evaluated in context of classification

Wrapper method selects number of features to use

Cons: Slow

## Adaboost algorithm

• Given  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  where  $\mathbf{x} \in X, y_i \in \{+1, -1\}$

• Initialize  $\mathbf{w}_1(i) = 1/N$ . (Uniform distribution over data)

• For  $t = 1, \dots, K$ :

- Select new training set  $S_t$  from  $S$  with replacement, according to  $\mathbf{w}_t$
- Train  $L$  on  $S_t$  to obtain hypothesis  $h_t$
- Compute the training error  $\varepsilon_t$  of  $h_t$  on  $S$ :

$$\varepsilon_t = \sum_{j=1}^N \mathbf{w}_t(j) \delta(y_j \neq h_t(\mathbf{x}_j)), \text{ where}$$

$$\delta(y_j \neq h_t(\mathbf{x}_j)) = \begin{cases} 1 & \text{if } y_j \neq h_t(\mathbf{x}_j) \\ 0 & \text{otherwise} \end{cases}$$

- Compute coefficient

- Compute new weights on data:

For  $i = 1$  to  $N$

$$\mathbf{w}_{t+1}(i) = \frac{\mathbf{w}_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i))}{Z_t}$$

where  $Z_t$  is a normalization factor chosen so that  $\mathbf{w}_{t+1}$  will be a probability distribution:

$$Z_t = \sum_{i=1}^N \mathbf{w}_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i))$$

• At the end of  $K$  iterations of this algorithm, we have

$$h_1, h_2, \dots, h_K$$

We also have

$\alpha_1, \alpha_2, \dots, \alpha_K$ , where

$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$

• Ensemble classifier:

$$H(\mathbf{x}) = \text{sgn} \sum_{t=1}^K \alpha_t h_t(\mathbf{x})$$

• Note that hypotheses with higher accuracy on their training sets are weighted more strongly.

Then define:

$$\text{Entropy}(S_f^{\text{high}}) = -(p_+^{\text{high}} \log_2 p_+^{\text{high}} + p_-^{\text{high}} \log_2 p_-^{\text{high}})$$

$$\text{Entropy}(S_f^{\text{low}}) = -(p_+^{\text{low}} \log_2 p_+^{\text{low}} + p_-^{\text{low}} \log_2 p_-^{\text{low}})$$

$$\text{Entropy}(S_f) = \frac{|S_f^{\text{high}}|}{|S|} \text{Entropy}(S_f^{\text{high}}) + \frac{|S_f^{\text{low}}|}{|S|} \text{Entropy}(S_f^{\text{low}})$$

$$\text{InformationGain}(S_f) = \text{Entropy}(S) - \text{Entropy}(S_f)$$