# Adaboost.M1
exp-loss forward stagewise modelling

$Y = \{-1, 1\}$

$X \in \mathbb{R}^p$

0) Loss-function $L(y, f(x)) = \exp(-yf(x))$ ← exponential loss

$$\| \\ \sum_{i=1}^{N} L(y_i, f(x_i))$$

Forward stepwise modelling:

$$\underset{\beta, \gamma}{\text{argmin}} \sum_{i=1}^{N} L(y_i, \underbrace{\hat{f}^{(m-1)}(x_i)}_{\text{current classifier - not to change}} + \overbrace{\beta b(x_i, \gamma)}^{\text{to be changed}})$$

1) What does this look like at step $m$?

$$(\beta_m, b_m) = \underset{\beta, f}{\text{argmin}} \sum_{i=1}^{N} \exp\left(-y_i \sum_{k=1}^{m} \beta^{(k)} b^{(k)}(x_i)\right)$$

$$= \underset{\beta, b}{\text{argmin}} \sum_{i=1}^{N} \exp\left\{-y_i \left(\sum_{k=1}^{m-1} \beta^{(k)} \hat{f}^{(k-1)}(x_i) + \beta b(x_i)\right)\right\}$$

$$= \underset{\beta, b}{\text{argmin}} \sum_{i=1}^{N} w_i^{(m-1)} \cdot \exp\left\{-y_i \beta b(x_i)\right\}$$

$$w_i^{(m-1)} = \exp\left\{-y_i \hat{f}^{(m-1)}(x_i)\right\}$$

known — from last iteration —

Observe: if $y_i = b(x_i)$ then $y_i \, b(x_i) = +1$

$\quad\quad\quad\quad \neq \quad\quad\quad\quad\quad\quad\quad \div 1$

Next: two step procedure

where a) minimize wrt $b$

b) minimize wrt $\beta$

2)

$$b_m = \underset{b}{\arg\min} \sum_{i=1}^{N} w_i^{(m-1)} \cdot \exp\left(-y_i \, \beta \, b(x_i)\right)$$

$$= \underset{b}{\arg\min} \left\{ \sum_{i: \, y_i = b(x_i)} w_i^{(m-1)} \cdot e^{-\beta} + \sum_{i: \, y_i \neq b(x_i)} w_i^{(m-1)} \, e^{\beta} \right\}$$

$\uparrow$

add and subtract

$$\sum_{i: \, y_i \neq b(x_i)} w_i^{(m-1)} \, e^{-\beta}$$

$$= \underset{b}{\arg\min} \left\{ \underbrace{\sum_{i=1}^{N} w_i^{(m-1)} e^{-\beta}}_{\substack{\text{no } b \\ \text{here}}} + (e^{\beta} - e^{-\beta}) \underbrace{\sum_{i: \, y_i = b(x_i)} w_i^{(m-1)}}_{\mathbb{I}(y_i \neq b(x_i))} \right\}$$

$$\hat{b}_m = \underset{b}{\arg\min} \left\{ \sum_{i=1}^{N} w_i^{(m-1)} \mathbb{I}(y_i \neq b(x_i)) \right\}$$

_____

**3)**

$$\beta^{(m)} = \underset{\beta}{\arg\min} \sum_{i=1}^{N} w_i^{(m-1)} \exp\left\{-y_i \beta\, b(x_i)\right\}$$

$$y_i = b(x_i) \rightarrow \sum_{i:\, y_i = b(x_i)} w_i^{(m-1)} e^{-\beta}$$

$$y_i \neq b(x_i) \rightarrow \sum_{i:\, y_i \neq b(x_i)} w_i^{(m-1)} e^{\beta}$$

$$L = \sum_{i:\, y_i = b(x_i)} w_i^{(m-1)} e^{-\beta} + \sum_{i:\, y_i \neq b(x_i)} w_i\, e^{\beta}$$

$$\frac{\partial L}{\partial \beta} = -\sum_{i:\, y_i = b(x_i)} w_i^{(m-1)} e^{-\beta} + \sum_{i:\, y_i \neq b(x_i)} w_i^{(m-1)} e^{\beta} = 0$$

multiply $e^{\beta}$

$$-\sum_{i:\, y_i = b(x_i)} w_i^{(m-1)} + e^{2\beta} \sum_{i:\, y_i \neq b(x_i)} w_i^{(m-1)} = 0$$

$$e^{2\beta} \sum_{i:\, y_i \neq b(x_i)} w_i^{(m-1)} = \sum_{i=1}^{N} w_i^{(m-1)} - \sum_{i:\, y_i \neq b(x_i)} w_i^{(m-1)}$$

$$e^{2\beta} = \frac{\sum_{i=1}^{N} w_i^{(m-1)} - \sum_{i:\, y_i \neq b(x_i)} w_i^{(m-1)}}{\sum_{i:\, y_i \neq b(x_i)} w_i^{(m-1)}}$$

scale each term with $\sum\limits_{i=1}^{W} w_i^{(M-1)}$ to get

$$e^{2\beta} = \frac{1 - err_m}{err_m}$$

$$2\hat{\beta}_m = \ln\left(\frac{1 - err_m}{err_m}\right)$$

where $err_m = \dfrac{\sum\limits_{i: \, g_i \neq b(x_i)} w_i^{(M-1)}}{\sum\limits_{i=1}^{W} w_i^{(M-1)}}$

In the Adaboost.M1:

$$\alpha_m = 2\beta_m$$

$$\hat{\alpha}_m = \ln\left(\frac{1 - err_m}{err_m}\right)$$

4) In 2) we found that $\hat{b}_m$ is
minimizing the "weighed misclassification error"
and in 3) we found

$$\hat{\beta}_m = \tfrac{1}{2}\ln\left(\frac{1 - err_m}{err_m}\right), \text{ or } \hat{\alpha}_m = \ln\left(\frac{1 - err_m}{err_m}\right)$$

This gives

Our classifier is updated as

$$\hat{f}^{(m)}(x_i) = \hat{f}^{(m-1)}(x_i) + \hat{\beta}_m \cdot \hat{b}_m(x_i)$$

and the weights (from 1))

$$w_i^{(m)} = w_i^{(m-1)} \exp\left(-y_i \,\hat{\beta}_m \,\hat{b}_m(x_i)\right)$$

But in the algorithm the weight update is written differently "$w_i \leftarrow w_i \cdot \exp\left(\alpha_m \cdot I(y_i \neq \hat{b}_m(x_i)\right)$"

Final "adjustment", look at exponent:

$$-y_i \cdot \hat{b}_m(x_i) = -I(y_i = \hat{b}_m(x_i)) + I(y_i \neq \hat{b}_m(x_i))$$

$$= -I(y_i = \hat{b}_m(x_i)) + I(y_i \neq \hat{b}_m(x_i)) - I(y_i \neq \hat{b}(x_i)) + I(y_i \neq \hat{b}_m(x_i))$$

$$= -1 + 2 \cdot I(y_i \neq \hat{b}_m(x_i))$$

Insert info:

$$w_i^{(m)} = w_i^{(m-1)} \cdot \exp\left(\hat{\beta}_m\left(-1 + 2I(y_i \neq \hat{b}_m(x_i)\right)\right)$$

$$= w_i^{(m-1)} \exp\left( \overbrace{2\hat{\beta}_m}^{\hat{\alpha}_m} \cdot \mathbb{I}(y_i \neq \hat{b}_m(x_i) - \hat{\beta}_m \right)$$

$$= w_i^{(m-1)} \exp\left( \hat{\alpha}_m \cdot \mathbb{I}(y_i \neq \hat{b}_m(x_i)) \right) \cdot \exp\left( -\frac{\hat{\alpha}_m}{2} \right)$$

an extra factor = the same value
is multiplied in for each
weight $w_i$ = a scaling