

$$= e^{-2\alpha_{m+1}} \frac{\sum_{i: y_i \neq h_{m+1}(x_i)} w_i^m}{\sum_{i: y_i \neq h_{m+1}(x_i)} w_i^m}$$

but we know that $\alpha_{m+1} = \frac{1}{2} \log \left[\frac{1 - \epsilon_{m+1}}{\epsilon_{m+1}} \right]$

$$\text{So } e^{-2\alpha_{m+1}} = e^{-2 \times \frac{1}{2} \log \left(\frac{1 - \epsilon_{m+1}}{\epsilon_{m+1}} \right)}$$

$$\frac{\epsilon_{m+1}}{1 - \epsilon_{m+1}}$$

But we know that

$$\epsilon_{m+1} = \sum_{i: y_i \neq h_{m+1}(x_i)} w_i^m$$

$$\text{and thus } 1 - \epsilon_{m+1} = \sum_{i: y_i = h_{m+1}(x_i)} w_i^m$$

$$B = \frac{\epsilon_{m+1}}{1 - \epsilon_{m+1}} \times \frac{1 - \epsilon_{m+1}}{\epsilon_{m+1}} = 1$$

$$\therefore A = \frac{1}{1+B} = \frac{1}{1+1} = \underline{\underline{\frac{1}{2}}}$$

□

(2)

①⑥ Proved in 1.a:

$$\sum_{i: h_m(x_i) \neq y_i} w_i^m = \frac{1}{2}$$

Suppose h_1, \dots, h_m models comprise H_M .Adding h_{m+1} , h_{m+1} is determined by

$$\min_{h_{m+1}} \sum_{i: h_{m+1}(x_i) \neq y_i} w_i^m = \epsilon_{m+1}$$

If $h_{m+1} = h_m$, (i.e. optimal $h_{m+1} = h_m$)

$$\epsilon_{m+1} = \sum_{i: h_m(x_i) \neq y_i} w_i^m = \frac{1}{2}$$

We know $\alpha_{m+1} = \log \left(\frac{1 - \epsilon_{m+1}}{\epsilon_{m+1}} \right)$

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

$$= \log(1) = \underline{\underline{0}}$$

Also a classifier h_m is weak only if $\epsilon_m < 1/2$.

If $\epsilon_m \geq 1/2$ it is as good as a random classifier. So $h_{m+1} \neq h_m$.

\therefore if $h_{m+1} = h_m$, it will not get added to the ensemble because $\alpha_{m+1} = 0$, thus we can't have $h_{m+1} = h_m$ for any m .

① ② "weak" h_m
 A classifier \hat{h}_m is useful only if
 error $\epsilon_m = 1/2 - \epsilon$ with $\epsilon > 0$.

if $h_{m+1} = h_m$

then $\epsilon_{m+1} = 1/2$ (as proved before)

but if $h_{m+k} = h_m$

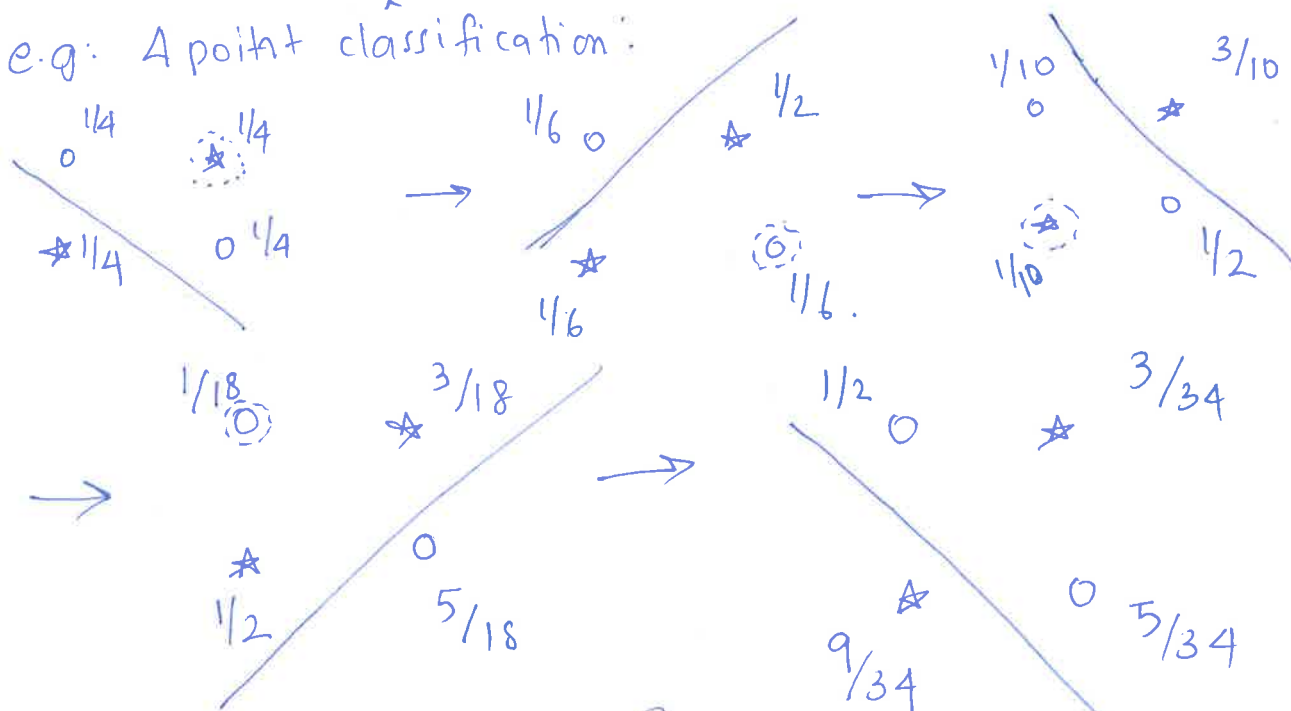
for some $k > 1$,

then $\epsilon_{m+k} = \sum_{i: h_{m+k}(x_i) \neq y_i} W_{m+k-1}$

need not necessarily be $1/2$.

So a classifier can repeat itself in the AdaBoost algo

e.g. A point classification.



{ Note: idea borrowed from Michael Hanna }

classifier is same as the first classifier.

②

Empirical exponential loss,

$$L = \sum_{i=1}^n w_i^{m-1} e^{-\alpha_m y_i h_m(x_i)}$$

Given $h_m(x_i)$, The α_m that minimizes L satisfies:

$$\frac{\partial L}{\partial \alpha_m} = 0$$

$$\Rightarrow \sum_{i=1}^n w_i^{m-1} e^{-\alpha_m y_i h_m(x_i)} (-y_i h_m(x_i)) = 0$$

we know $y_i h_m(x_i) = \begin{cases} 1 & \text{if } y_i = h_m(x_i) \\ -1 & \text{otherwise} \end{cases}$

$$\text{So } \sum_{i: y_i = h_m(x_i)} w_i^{m-1} e^{-\alpha_m} (-1) + \sum_{i: y_i \neq h_m(x_i)} w_i^{m-1} e^{\alpha_m} (1) = 0$$

 α_m is independent of i so

$$-e^{-\alpha_m} \sum_{i: y_i = h_m(x_i)} w_i^{m-1} + e^{\alpha_m} \sum_{i: y_i \neq h_m(x_i)} w_i^{m-1} = 0$$

we know $\epsilon_m = \sum_{i: y_i \neq h_m(x_i)} w_i^{m-1}$

and $1 - \epsilon_m = \sum_{i: y_i = h_m(x_i)} w_i^{m-1}$,

$$\text{So } -e^{-\alpha_m} (1 - \epsilon_m) + e^{\alpha_m} \epsilon_m = 0$$

$$\Rightarrow \frac{e^{2\alpha_m}}{e^{\alpha_m}} = \frac{1 - \epsilon_m}{\epsilon_m} \Rightarrow \alpha_m = \frac{1}{2} \log \left[\frac{1 - \epsilon_m}{\epsilon_m} \right]$$

③

$$(3) \quad p(y=1 | x; w) = \frac{1}{1 + \exp\left(\sum_{j=1}^d w_j \phi_j'(x)\right)} = \frac{1}{1 + \exp(w \cdot \phi'(x))}$$

Define the ^{new} feature space

$\phi'(x_0)$ as:

$$\phi'(x_0) = \begin{bmatrix} K(x_0, x_1) \\ K(x_0, x_2) \\ \vdots \\ K(x_0, x_n) \end{bmatrix}_{n \times 1} \equiv K(x_0, x) \quad \text{and } w \text{ is an } n \times 1 \text{ vector.}$$

where x_1, x_2, \dots, x_n are the training data and $K(x_0, x_i)$ is the kernel defined as

$$K(x_0, x_i) = \phi(x_0)^T \phi(x_i)$$

where ϕ is the original feature space. prediction for x_0 is given as.

$$\hat{y} = \text{sign}(w \cdot \phi'(x_0)) = \text{sign}(w \cdot K(x_0, x))$$

\therefore prediction only depends on the kernel values.

Gradient (as given in notes + adding regularizer term):

$$\text{grad} = \nabla_w [\log p(y_i | x_i; w) + \lambda \|w\|^2] \\ = - [y_i - \sigma(w \cdot K(x_i, x))] K(x_i, x) + 2\lambda w$$

So gradient on a single example also depends on the training data only through the kernel.