

Thursday, 10/29/2020

EE 660

MACHINE LEARNING  
FROM SIGNALS:  
FOUNDATIONS AND METHODS

Prof. B. Keith Jenkins

Lecture 20

---

**Lecture 20**

---

**EE 660****Oct 29, 2020**

---

---

**Announcements**

- Homework 7 is due Friday
- Homework 8 will be posted
- My office hours for today will be: 3:30 - 4:30 PM

---

**Today's topics**

- Boosting
  - Introduction and notation
  - Loss functions
  - Forward stagewise additive modeling: base classifiers and their importance
  - Adaboost
    - Minimization of  $L$
    - Weight of each data point
- $\text{err}_m$
- Adaboost.M1 algorithm
- Example

next  
lecture {

# Boosting [Murphy 16.4.0-16.4.4] [also Hastie et al., Ch.10]

Is also an adaptive basis fcn. model. (ABM)

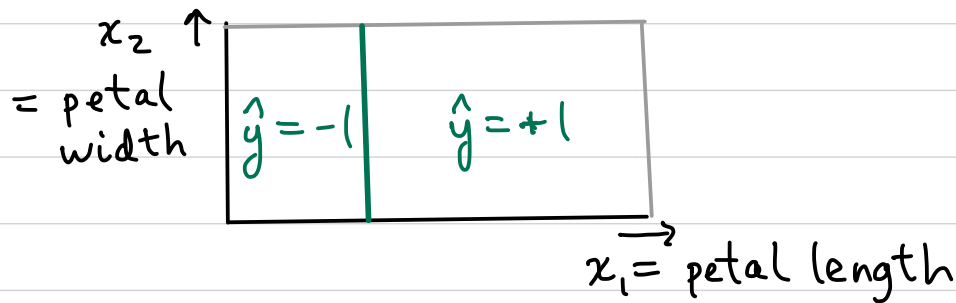
$$(1) \hat{f}(\underline{x}) = w_0 + \sum_{m=1}^M w_m \phi(\underline{x}; \underline{\gamma}_m)$$

$\uparrow$   
 $\triangleq \beta_m$

$\uparrow$   
 parameters of  $\phi_m(\underline{x})$ , learned from data

Each  $\phi(\underline{x}, \underline{\gamma}_m) = \phi_m(\underline{x})$

- is a simple classifier that can classify the entire feature space
- is a "weak learner" that is only required to do better than chance
- is typically a "decision stump"
  - a 1-stage CART resulting in 1 node and 2 leaves (or variants with > 2 leaves)



The simple classifiers  $\phi(\underline{x}; \underline{\gamma}_m)$  are found sequentially:  $m=1, 2, \dots, M$ .

## Notation

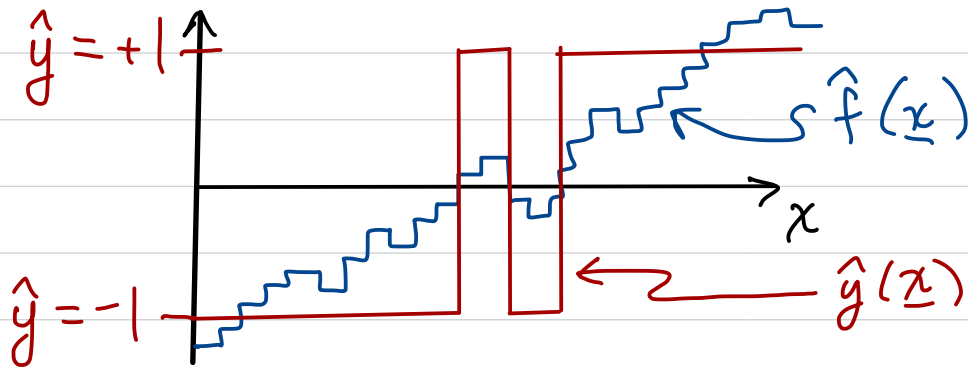
For 2-class problem, let labels be  $\pm 1$ .

$$(2) \hat{f}(x) = f_0 + \sum_{m=1}^M \beta_m \phi(x; \gamma_m)$$

"weight" or importance  
of  $m^{\text{th}}$  base classifier  
 $\beta_m \in \mathbb{R}$  (will have  $\beta_m \geq 0$ )

$\phi(x; \gamma_m) = \phi_m(x)$   
is the  $m^{\text{th}}$  simple "base" classifier  
 $\phi_m \in \{-1, +1\}$

Final classifier:  $\hat{y}(x) = \text{sign}\{\hat{f}(x)\}$ ,  $\hat{y} \in \{-1, +1\}$



Measure of confidence  
in output prediction  $\hat{y}$ ?

$$\rightarrow |\hat{f}(x)|$$

(3) Intermediate expressions for  $f$  ( $m^{\text{th}}$  iteration):

$$\hat{f}_m(x) = f_0 + \sum_{m'=1}^m \beta_{m'} \phi(x, \gamma_{m'})$$

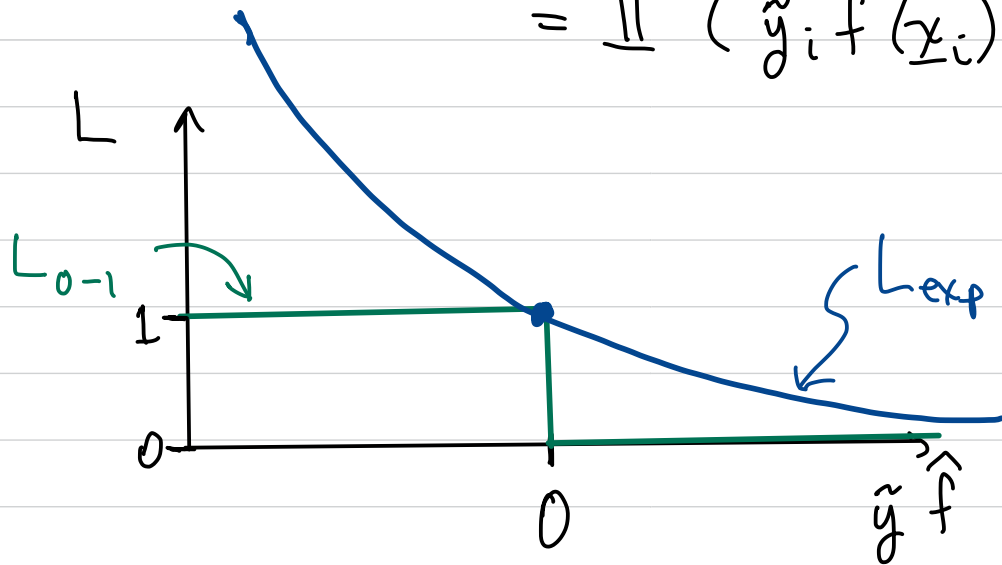
## Loss function

Given  $\mathcal{D} = \{(\underline{x}_i, \tilde{y}_i) \mid \tilde{y}_i \in \{-1, +1\}\}_{i=1}^N$

$$f_{\text{obj}} \left[ \{\beta_m, m=1, 2, \dots, M\}, \{\underline{x}_m, m=1, 2, \dots, M\} \right] = \frac{1}{N} \sum_{i=1}^N L(\tilde{y}_i, \hat{f}(\underline{x}_i))$$

↑ Loss function.

0-1 Loss :  $L_{0-1} = \mathbb{I}(\tilde{y}_i \neq \text{sign}\{\hat{f}(\underline{x}_i)\})$   
 $= \mathbb{I}(\tilde{y}_i \hat{f}(\underline{x}_i) < 0)$



(4) Exponential loss :  $L_{\text{exp}} = \exp[-\tilde{y}_i \hat{f}(\underline{x}_i)]$   
 [other loss functions (e.g., Murphy)]

Why use  $L_{\text{exp}}$  instead of  $L_{0-1}$ ?

- smooth, differentiable
- provides confidence measure
- convex

Want to find:  $\hat{f}(\underline{x}) = \underset{f(\underline{x})}{\operatorname{argmin}} \sum_{i=1}^N L_{\text{exp}}[\tilde{y}_i, f(\underline{x}_i)]$

$$= \underset{f_0, \beta_m, \underline{\gamma}_m, \forall m=1, 2, \dots, M}{\operatorname{argmin}} \sum_{i=1}^N L_{\text{exp}}\left[\tilde{y}_i, f_0 + \sum_{m=1}^M \beta_m \phi_m(\underline{x}_i, \underline{\gamma}_m)\right]$$

Easier to min. each term, sequentially:  
at  $m^{\text{th}}$  iteration:

$$\underset{\beta_m, \underline{\gamma}_m}{\operatorname{argmin}} \sum_{i=1}^N L_{\text{exp}}\left[\tilde{y}_i, \hat{f}_{m-1}(\underline{x}_i) + \beta_m \phi(\underline{x}_i, \underline{\gamma}_m)\right]$$

$$\text{e.g.: } \underline{\gamma}_m = \begin{bmatrix} j' \\ t_{j'} \\ \ell \end{bmatrix} = \begin{bmatrix} \text{Feature to split} \\ \text{threshold value} \\ \text{region labels} \end{bmatrix}$$

## Forward Stagewise Additive Modelling

1. Initialize  $f_0 = 0$

2. For  $m=1$  to  $M$ :

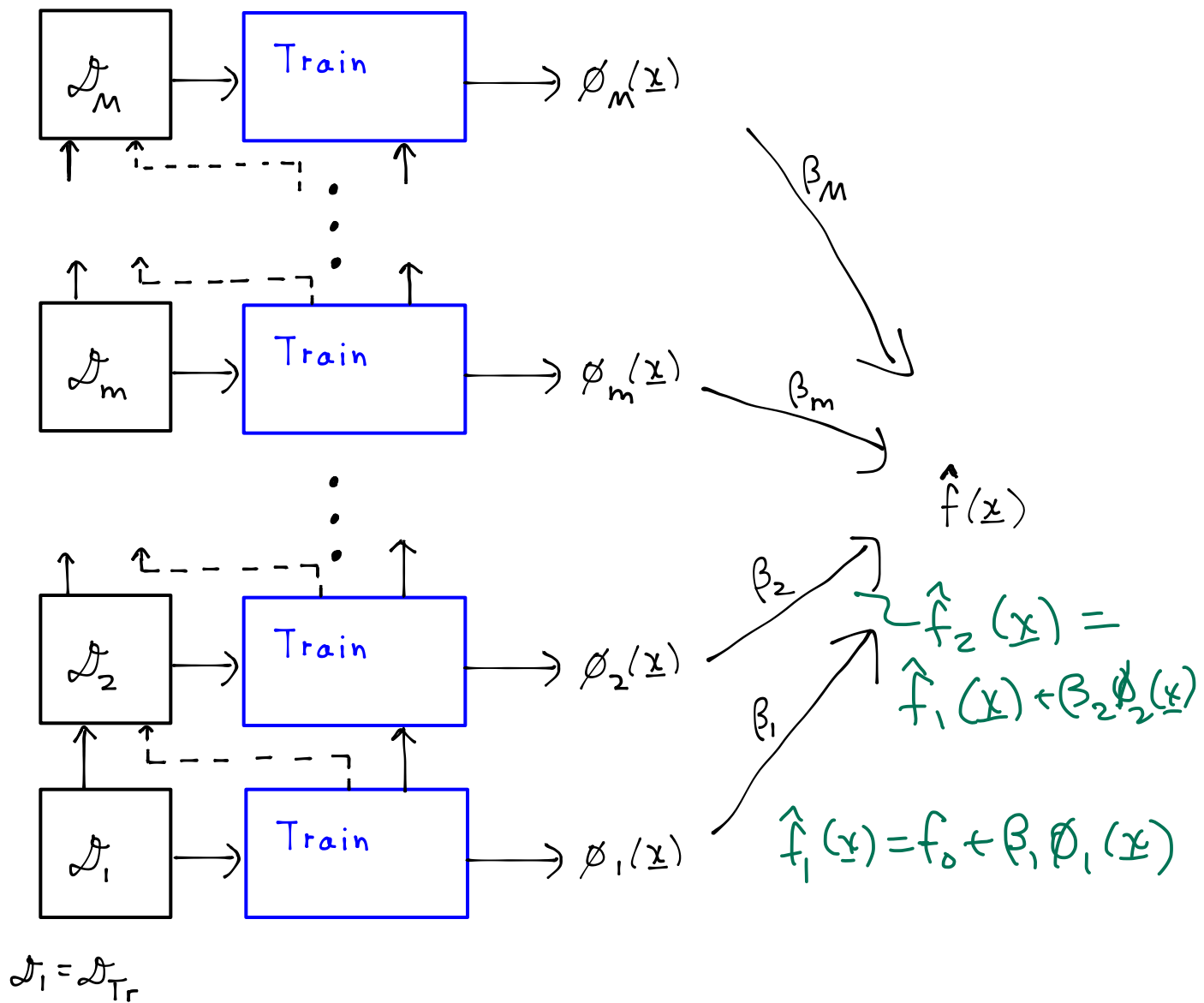
$$(i) \text{ Find } (\beta_m, \underline{\gamma}_m) = \underset{\beta'_m, \underline{\gamma}'_m}{\operatorname{argmin}} \sum_{i=1}^N L[\tilde{y}_i, \hat{f}_{m-1}(x_i) + \beta'_m \phi(x_i, \underline{\gamma}'_m)]$$

$$(ii) \hat{f}_m(\underline{x}) = \hat{f}_{m-1}(\underline{x}) + \beta_m \phi(\underline{x}, \underline{\gamma}_m)$$

3. Final classifier is:  $\hat{y}(\underline{x}) = \operatorname{sign}\{\hat{f}_M(\underline{x})\}$

(For  $L$  could use  $L_{\exp}$  or other loss fcn.)

How to do 2(i) is left unspecified.



$$\hat{y}(\underline{x}) = \text{sign}\{\hat{f}(\underline{x})\} = \text{sign}\left\{\sum_{m=1}^M \beta_m \phi_m(\underline{x})\right\}$$



# Adaboost

→ Use  $L = L_{\text{exp}}$   
 From above 2(i):

$$\text{Find } (\beta_m, \gamma_m) = \underset{\beta'_m, \gamma'_m}{\operatorname{argmin}} \underbrace{\sum_{i=1}^N L_{\text{exp}} \left[ \tilde{y}_i, \hat{f}_{m-1}(x_i) + \beta'_m \phi(x_i, \gamma'_m) \right]}_{L_m}$$

$$(5) \quad L_m = \sum_{i=1}^N \exp \left\{ -\tilde{y}_i \left[ \hat{f}_{m-1}(x_i) + \beta'_m \phi(x_i, \gamma'_m) \right] \right\}$$

$$\underbrace{\quad}_{w_{i,m} = \text{const. of } \beta'_m, \gamma'_m}$$

$$(6) \quad \text{Let } w_{i,m} \triangleq \exp \left\{ -\tilde{y}_i \hat{f}_{m-1}(x_i) \right\}$$

$\nearrow$  data pt.  
 index  
 $\nwarrow$  iteration  
 (base classifier)  
 index

$$L_m = \sum_{i=1}^N w_{i,m} \exp \left[ -\tilde{y}_i \beta'_m \phi(x_i, x'_m) \right], \quad \beta'_m > 0.$$

weight on each  
data pt.  $i$  at  $m^{\text{th}}$   
iteration.

interpret as  $L_{\text{exp}}$  for data pt.  $i$  using base  
classifier  $\phi_m$ .

$L_m$  can be minimized algebraically. [See Murphy].

Can re-arrange  $L_m$  above to get eqns. for Adaboost algorithm [Murphy]:

(\*)  $\phi(x_i, x'_m)$  is chosen to minimize:

$$\phi_m = \underset{\phi}{\operatorname{argmin}} \left\{ \sum_{i=1}^N w_{i,m} \mathbb{I}[\tilde{y}_i \neq \phi(x_i)] \right\}$$

(decision stump  
optimization)

sum of weights of misclassified data pts.

★

Your Murphy may have this omitted (error)  
in Eq. (16.40).