# Lecture November 12

Boosting methods

model $f(x) = \sum\limits_{m=0}^{M} \beta_m b_m(x)$

$b_m(x)$ = weak learner.

- Restrictive model

$f(x_i) = \sum\limits_{j=0}^{P-1} f_j(x_i)$

$= \sum\limits_{j=0}^{P-1} \sum\limits_{m=0}^{M} \beta_{jm} b_{jm}(x_i)$

linear regression and
polynomial expansion.

- selection models
include only those basis
functions $b_m$ that
contribute significantly.
Random forests, Bagging &
boosting models belong
here.

- Regularization models
Ridge is a simple
example. Screen away
selected features.

Boosting :

$$f(x) = f_0(x) + \sum_{m=1}^{M} b_m(x) \beta_m$$

weak learner.

Put emphasis in the iteration on misclassified data.

$$C(f) = \sum_{i=0}^{n-1} L(y_i, f(x_i))$$

$$\hat{f} = \underset{f}{argmin} \; C(f)$$

in linear regression

$$MSE = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - f(x_i))^2$$

L2 boosting for regression

$$f_m(x) = f_{m-1}(x) + \beta_m \underbrace{b(x_i, \gamma_m)}_{b_m(x_i)}$$

(i) Establish a cost function

$$C(f) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - f_m(x_i))^2$$

algo :

   a) initialize $f_0(x)$

   b) for $m = 1 : M$

$$(\hat{\beta}_m, \hat{\gamma}_m) = \underset{\gamma \beta}{argmin}$$

$$\frac{1}{n} \sum_n \left( y_i - f_{m-1}(x_i) - \beta \underline{b(x;\gamma)} \right)^2$$

decision
tree

c) determine

$$f_m(x) = f_{m-1}(x) + \beta_m b(x;\gamma_m)$$

end for

Return $\quad f(x) = \sum_{m=0}^{M} \beta_m b(x;\gamma_m)$

## Adaboost for Regression

Define $\quad Err_m = \dfrac{\sum_{i=0}^{n-1} w_{i,m} \, \mathbb{I}(y_i \neq f_m(x_i))}{\sum_{i=0}^{n-1} w_{i,m}}$

$$f_m(x) = f_{m-1}(x) + \beta_m b_m(x)$$

Cost function

$$C(f) = \sum_{i=0}^{n-1} L_i$$

$$L_i = e^{-y_i f_m(x_i)}$$

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

$$f_m(x_i) = \{-1, 1\}$$

also

a) initialize $w_i = 1/n$

b) for $m = 1:M$   DO

- Fit a classifier $b_m(x_i)$ to the training data using $w_{im}$

- compute $\varepsilon_{mm}$

- compute $\beta_m = \frac{1}{2} \ln\left(\frac{1-\varepsilon_{mm}}{\varepsilon_{mm}}\right)$

- Find new weights $w_{im+1}$

- update
$$f_m(x) = f_{m-1}(x) + \beta_m b_m(x)$$

end for

return $f(x) = \text{sign}\left\{\sum_{m=0}^{M} \beta_m b_m(x)\right\}$

$\rightarrow$ parametrize through minimization $\gamma_m, \beta_m$

<mark>Gradient boosting</mark>

(i) fit $f$ to individual points
$$f : \left\{ f(x_0), f(x_1) -- f(x_{n-1}) \right\}$$

$$g_{im} = \left[ \frac{\partial \mathcal{L}(y_i, f(x_i))}{} \right]$$

$$\left[ \frac{\partial f(x_i)}{\partial f} \right]_{f = f_{m-1}(x_i)}$$

update

$$\underline{f_m = f_{m-1} - \rho_m \, g_m}$$

$\underbrace{\qquad\qquad}$ step length to be optimized

$$\hat{\rho}_m = \arg\min_{\rho} C(f_{m-1} - \rho \, g_m)$$

steepest descent, less use.

(ii) Gradient boosting:

$$f_m = f_{m-1} + \beta_m \, b(x; \gamma_m)$$

$$\underline{f_m - f_{m-1} = \beta_m \, b(x; \gamma_m)}$$

$$\boxed{f_m - f_{m-1} = -g_m \, \beta_m}$$

$$g_m = b(x; \gamma_m) \implies$$

$$-g_m - b(x; \gamma_m) = 0$$

$$\boxed{\gamma_m = \arg\min_{\gamma} \sum_{i=0}^{m-1} \left( -g_{im} - b(x_i; \gamma) \right)^2}$$

## algo

(i) initialize $f_0(x) = \underset{\gamma}{\arg\min} \sum_{i=0}^{n-1} \mathcal{L}(y_i, \gamma)$

(ii) for m = 1:M DO

— compute

$$g_{im} = \frac{\partial \mathcal{L}(y_i, f(x_i))}{\partial f(x_i)} \Big|_{f = f_{m-1}(x_i)}$$

— use weak learner to compute $\gamma_m$

— update

$$f_m(x) = f_{m-1}(x) + b(x; \gamma_m)$$

end for

Return $f(x) = \sum_{m=0}^{M} f_m(x)$