

What is boosting doing?

Additive model:

$$h(x) = \sum_{m=1}^M \beta_m b(x; \delta_m)$$

basis function

weight

parameter for basis

How do we fit these?

For some loss L

$$\{\hat{\beta}_m, \hat{\delta}_m\}_{m=1}^M = \underset{\beta_m, \delta_m}{\operatorname{argmin}} \sum_{n=1}^N L(y_n, \sum_m \beta_m b(x_n, \delta_m))$$

Problem: can be difficult

→ lots of params

→ depending on L may be difficult

Soln: Greedy approach

Forward Stagewise Additive Modeling

Forward Stagewise Additive Modeling

Do one value of m at a time.

$$\hat{G}_0(x) = 0$$

For $m=1, \dots, M$

$$\left[\begin{array}{l} \text{(a) } \hat{\beta}_m, \hat{\gamma}_m = \underset{\beta, \gamma}{\operatorname{argmin}} \sum_{n=1}^N L(y_n, \hat{G}_{m-1}(x) + \beta b(x, \gamma)) \\ \text{(b) } \hat{G}_m(x) = \hat{G}_{m-1}(x) + \hat{\beta}_m b(x, \hat{\gamma}_m) \end{array} \right]$$

$$\text{Return } \hat{G}_M(x) = \sum_{m=1}^M \hat{\beta}_m b(x, \hat{\gamma}_m).$$

Ex. Regression w/ squared error
 $L(y, z) = (y - z)^2$

then

$$L(y_n, \hat{G}_{m-1}(x) + \beta b(x; \gamma))$$

$$= (y_n - \hat{G}_{m-1}(x) - \beta b(x; \gamma))^2$$

r_{nm} = residual of predicting y_n
using \hat{G}_{m-1}

r_{nm} = residual of f_{m-1}
using \hat{G}_{m-1}

$$= (r_{nm} - \beta b(x; \sigma))^2$$

So SAM basically fits sequentially to residuals of prev. fit.

What does SAM have to do with boosting?

Ada Boost \approx SAM w/ Exponential Loss

$$L(y, h) = e^{-yh}$$

If I do SAM w/ exp loss then at each step I solve

$$\hat{\beta}_m, \hat{f}_m = \underset{\beta, f_m}{\operatorname{argmin}} \left(\sum_{n=1}^N L(\underbrace{y_n}_y, \underbrace{\hat{G}_{m-1}(x) + \beta f_m(x)}_h) \right)$$

using $L = \text{exp loss}$

$$\leftarrow \hat{G}_{m-1}, [\hat{G}_{m-1}(x) + \beta f_m(x)]$$

$$\sum_n \exp(-y_n [\hat{G}_{m-1}(x) + \beta f_m(x)])$$

$$= \underbrace{\sum_n \exp(-y_n \hat{G}_{m-1}(x))}_{\text{no } f_m, \beta} \exp(-y_n \beta f_m(x))$$

I can regard as a weight
 w_{nm}

$$= \sum_n w_{nm} L(y_n, \beta f_m(x))$$

↑ overall = weighted exp. loss

weights depend on accuracy of G_{m-1} .

For cts y , can do this where we fit each new additive part to residuals of prev. iteration.

Practically:

Tuning params:

— loss function *

- loss function (*)
 - number of trees M
 - Shrinkage factor

$$\hat{G}_m(x) = \hat{G}_{m-1}(x) + \nu \alpha_m \hat{f}_m(x)$$

learning rate $\nu \in [0, 1]$
 - number of splits in my trees \hat{f}_m
-

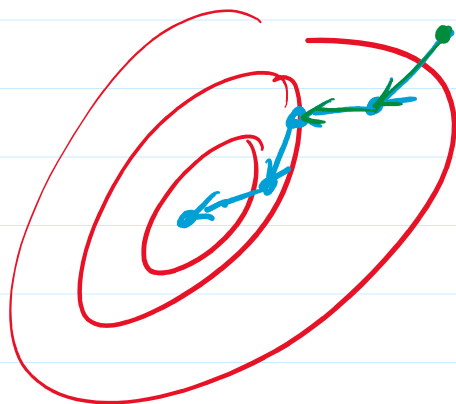
Can generalize this to any differentiable loss function L

Called: Gradient Boosting

Analogy:

Gradient Descent:

Want to optimize function J



For $m = 1, \dots, M$

$$x_m = x_{m-1} - \alpha_m \nabla J|_{x_{m-1}}$$

Boosting:

For $m = 1, \dots, M$

$$\hat{G}_m(x) = \hat{G}_{m-1}(x) + \alpha_m \hat{f}_m(x)$$

Looks like gradient descent if
 $\hat{f} \approx -\text{grad. of loss fn.}$

Neural Networks

NNs are predictive functions \hat{f} that are built via composition of simpler functions: \hat{f}_e

$$\hat{f}(x) = \hat{f}_L(\hat{f}_{L-1}(\hat{f}_{L-2}(\dots \hat{f}_2(\hat{f}_1(x))))).$$

Here L is called the number of layers

Here L is called the number of layers
 \hat{f}_l is called a layer.

The simplest NNs are called feed forward networks b/c the calc. from layer l is fed forward into layer $l+1$.

For FNNs the \hat{f}_l are comprised of two very simple operations:

(1) a linear transformation

$$W^{(l)}h + b^{(l)}$$

weights \nearrow $W^{(l)}$
 \nwarrow generic input to layer h
 $b^{(l)}$ \nwarrow biases

(2) a non-linear transf $g^{(l)}$ \nwarrow activation function
[typically applied element-wise]

$$\hat{n} = (n^{(l)}, \dots, n^{(l)}) \in \mathbb{R}^M$$

$$\text{So } \hat{f}_l(h) = g^{(l)}(w^{(l)}h + b^{(l)})$$

if $h \in \mathbb{R}^K$ then

$$w^{(l)} \in \mathbb{R}^{K \times M}, \quad b^{(l)} \in \mathbb{R}^M$$

and we apply $g^{(l)}$ element-wise

$$\underline{L=1}$$

$$\hat{f}(x) = \hat{f}_1(x) = g^{(1)}(w^{(1)}x + b^{(1)})$$

$$\underline{L=2}$$

$$\hat{f}(x) = \hat{f}_2(\hat{f}_1(x))$$

$$= \hat{f}_2(g^{(1)}(w^{(1)}x + b^{(1)}))$$

$$= g^{(2)}(w^{(2)}(\underbrace{g^{(1)}(w^{(1)}x + b^{(1)})}_{1^{(1)} \quad 0 \dots}) + b^{(2)})$$

$$h^{(1)} = \hat{f}_1(x)$$

$$= g^{(2)}(w^{(2)}h^{(1)} + b^{(2)})$$

Why called NN?

They can be described as a network
(directed acyclic graph - DAG)

