

① Numerical optimization in ~~Function~~ Space.

Assume we are trying to fit a model $f(x) = f(x; \theta)$ parameterized by θ

Risk of the model may be written as: $R(\theta) = E [L(X, f(x; \theta))]$

Given that $R(\theta)$ is $E [L(X, f(x; \theta))]$

at iteration 'm' the estimate of θ is updated from θ^{m-1} to θ^m according to

$$\theta^m = \theta^{(m-1)} + \theta_m \rightarrow \text{Step taken at iteration 'm'}$$

Resulting estimate of ' θ ' after 'M' iterations can be written as a sum:

$$\theta = \theta^{(M)} = \sum_{m=0}^M \theta_m$$

[Gradient Descent:]

Before the update at iteration 'm', the current estimate of ' θ ' is given by $\theta^{(m-1)}$

At this current estimate, the direction of steepest descent of risk is given by

$$-g_m = - \nabla_{\theta} R(\theta) \Big|_{\theta = \theta^{(m-1)}}$$

→ optimal step length, may be

determined by 'Line Search' \Rightarrow Step at iteration 'm' is

$$p_m = \underset{p}{\operatorname{argmin}} R(\theta^{(m-1)} - p g_m) \quad \text{thus} \quad \theta_m = -p_m g_m$$

Newton's Method.

TB-2

Here we estimate the steplength and direction at the same time

Note: We wish to solve $\nabla_{\theta_m} R(\theta^{(m-1)} + \theta_m) = 0$
for optimal θ_m .

Second order Taylor's Expansion: for $R(\theta^{(m-1)} + \theta_m)$

\therefore we have.

$$R(\theta^{(m-1)} + \theta_m) \approx R(\theta^{(m-1)}) + g_m^T \theta_m + \frac{1}{2} \theta_m^T H_m \theta_m.$$

where H_m is the Hessian

matrix at the current estimate.

$$\therefore H_m = \nabla_{\theta}^2 R(\theta) \Big|_{\theta = \theta^{(m-1)}}$$

Thus we get $\nabla_{\theta_m} R(\theta^{(m-1)} + \theta_m) \approx g_m + H_m \theta_m = 0$

$$\Rightarrow \theta_m = -H_m^{-1} g_m.$$

Note: Newton's method is a second order method.
[Gradient descent is a first order Method.]

Numerical optimization in 'Function Space'

We wish to minimize the risk;

$$R(f) = E[L(X, f(X))]$$

Similar to the procedure for parameter optimization; we have here the following. Update at iteration m :

$$f_m^m(x) = f^{(m-1)}(x) + f_m(x)$$

Resulting estimate of 'f' after 'M' iterations

can be written as a sum.

$$\hat{f}(x) = \hat{f}^M(x) = \sum_{m=0}^M \hat{f}_m(x)$$

gradient descent:

$$-g_m(x) = - \left[\frac{\partial R(\hat{f}(x))}{\partial \hat{f}(x)} \right] \hat{f}(x) = \hat{f}^{(m-1)}(x).$$

$$= - \left[\frac{\partial E[L(Y, \hat{f}(x)) | X=x]}{\partial \hat{f}(x)} \right] \hat{f}(x) = \hat{f}^{(m-1)}(x).$$

$$= - E \left[\frac{\partial L(Y, \hat{f}(x))}{\partial \hat{f}(x)} \bigg| X=x \right] \hat{f}(x) = \hat{f}^{(m-1)}(x).$$

The step length P_m to take in steepest descent direction can be determined using line search.

$$P_m = \operatorname{argmin} E [L(Y, \hat{f}^{(m-1)}(x) - P_m g_m(x))]$$

the 'step' at each iter 'm'

is

$$\left[\hat{f}_m(x) = -P_m g_m(x) \right] \leftarrow \text{Performing iterative updates yields gradient descent algorithm in function space.}$$

Newton's Method,

We are trying to solve: $\frac{\partial}{\partial \hat{f}_m(x)} \left[E [L(Y, \hat{f}^{(m-1)}(x) + \hat{f}_m(x) | X=x)] \right] = 0$

→ We wish to find optimal step 'f_m'

Taylor's expansion.

We have:

$$E \left[L(Y, f^{(m-1)}(x) + f_m(x) | X=x) \right] \approx E \left[L(Y, f^{(m-1)}(x) | X=x) \right] + g_m f_m(x) + \frac{1}{2} h_m(x) f_m(x)^2$$

h_m is the Hessian at the current estimate.

Thus we have:

$$h_m(x) = \left[\frac{\partial^2 R(f(x))}{\partial f(x)^2} \right]_{f(x) = f^{(m-1)}(x)}$$

$$= \left[\frac{\partial^2 E[L(Y, f(x) | X=x)]}{\partial f(x)^2} \right]_{f(x) = f^{(m-1)}(x)}$$

Thus we have:

$$\frac{\partial}{\partial f_m(x)} \cdot E \left[L(Y, f^{(m-1)}(x) + f_m(x) | X=x) \right] \approx g_m + h_m(x) f_m(x) = 0$$

\therefore the sdn. is

$$\frac{f_m(x)}{h_m(x)} = \frac{-g_m(x)}{h_m(x)}$$

\rightarrow 'Newton' step in function space.

Boosting Algorithms

[TB-5]

Boosting fits ensemble models of the kind

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

→ Re writing as adaptive basis functions we have.

$$f(x) = \theta_0 + \sum_{m=1}^M \theta_m \phi_m(x) \quad \text{where } f_0(x) = \theta_0$$

$$\text{and } f_m(x) = \theta_m \phi_m(x).$$

for $m = 1, \dots, M$.

Most Boosting algorithms try to solve:

at each iteration:

$$\{\hat{\theta}_m, \hat{\phi}_m\} = \arg \min_{\{\theta_m, \phi_m\}} \sum_{i=1}^n L(y_i, f^{(m-1)}(x_i) + \theta_m \phi_m(x_i))$$

Either exactly OR approximately.

GRADIENT Boosting : Based on Gradient Descent on 'function space'.

We have the empirical version of the -ve gradient given as.

$$-\hat{g}_m(x_i) = - \left[\frac{\partial \hat{R}(f(x_i))}{\partial f(x_i)} \right]_{f(x) = \hat{f}^{(m-1)}(x)}$$

$$= - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x) = \hat{f}^{(m-1)}(x)}$$

→ only defined at data points. $\{x_i\}_{i=1}^n$.

Thus to generalize to other points in \mathcal{X} , prevent overfitting. We need to learn on approximate - no gradient using a restricted set of possible functions.

- We thus constrain set of possible solutions to a set of basis functions ϕ .
- At iteration 'm' the basis function $\phi_m \in \phi$ is learnt from data.
- The basis function we seek; should produce output $\{\phi_m(x_i)\}_{i=1}^n \rightarrow$ which is most highly correlated with -ve gradient $\{-g_m(x_i)\}_{i=1}^n$.

→ This is obtained by: $\hat{\phi}_m = \arg \min_{\phi \in \phi, B} \sum_{i=1}^n [(-g_m(x_i)) - \beta_m \phi(x_i)]^2$

step length: $\hat{\phi}_m = \arg \min_p \sum_{i=1}^n L(y_i, f^{m-1}(x_i) + p \phi_r)$

[Step length] → Additional Note.

Friedman (2001) → Introduced Shrinkage

$$\hat{f}_m(x) = \eta \hat{p}_m \hat{\phi}_m$$

$$0 < \eta \leq 1$$

where $0 < \eta \leq 1$ is the learning rate.

The resulting model may be written as $\hat{f}(x) = \hat{f}^M(x) = \sum_{m=0}^M \hat{f}_m(x)$

This can be seen as an adaptive basis function model.
where:
 $\hat{f}_m(x) = \hat{\theta}_m \hat{\phi}_m$
and $\hat{\theta}_m = \eta \hat{p}_m$.

For $m = 1$ to M .

Gradient Boosting:

Input:

- Data 'D'
- Loss function L
- Base Learner L_ϕ

→ Nos of iterations M

→ Learning rate γ

Initialize:

$$1] \hat{f}^{(0)}(x) = \hat{f}_0(x) = \hat{\theta}_0 = \arg \min_{\theta} \sum_{i=1}^n L(y_i, \theta)$$

2] for $m = 1 \dots M$ do

→ Compute gradient

$$\hat{g}_m(x_i) = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]$$

$$f(x) = \hat{f}^{(m-1)}(x)$$

→ Compute adaptive basis function.

$$\hat{\phi}_m = \arg \min_{\phi} \sum_{i=1}^n \left| \hat{g}_m(x_i) - \beta \phi(x_i) \right|$$

To make it more generic we do not use gradient directly as gradient is 'develop' only on certain data points.

→ Compute step length.

$$\hat{\rho}_m = \arg \min_{\rho} \sum_{i=1}^n L(y_i, \hat{f}^{(m-1)}(x_i) + \rho \hat{\phi}_m(x_i))$$

→ Compute 'step'

$$\hat{f}_m(x) = \gamma \hat{\rho}_m \hat{\phi}_m(x)$$

→ update

$$\hat{f}^m(x) = \hat{f}^{(m-1)}(x) + \hat{f}_m(x)$$

end

↓ output $\hat{f}(x) = \hat{f}^M(x) = \sum_{m=0}^M \hat{f}_m(x)$

Newton's Method for
fn. space: Recap.

at step 'm':] we are trying to solve:

$$\frac{\partial}{\partial f_m(x)} E [L(Y, f^{m-1}(x) + f_m(x)) | X=x] = 0$$

Second order Taylor's expansion,

$$E [L(Y, f^{m-1}(x) + f_m(x)) | X=x] \approx E [L(Y, f^{m-1}(x)) | X=x] + g_m f_m(x) + \frac{1}{2} h_m(x) f_m^2(x)$$

where $h_m(x)$ is Hessian at current estimate:

$$\text{Thus } h_m(x) = E \left[\frac{\partial^2 L(Y, f(x))}{\partial f(x)^2} | X=x \right] f(x) = f^{(m-1)}(x)$$

$$\text{Thus } \frac{\partial}{\partial f_m(x)} E [L(Y, f^{m-1}(x) + f_m(x)) | X=x] \approx g_m(x) + h_m(x) f_m(x)$$

$$\text{Thus Newton step is } f_m(x) = \frac{-g_m(x)}{h_m(x)} = 0$$

Newton Boosting

] Risk is Now Empirical Risk

$$\text{Empirical Hessian} \rightarrow \hat{h}_m(x_i) = \left[\frac{\partial^2 \hat{R}(f(x_i))}{\partial f(x_i)^2} \right] f(x) = \hat{f}^{(m-1)}(x)$$

Newton's step]

$$\hat{\phi}_m = \arg \min_{\phi \in \phi} \sum_{i=1}^n \left[\hat{g}_m(x_i) \phi(x_i) + \frac{1}{2} \hat{h}_m(x_i) \phi(x_i)^2 \right]$$

we can re-write: \rightarrow Weighted Least Square problem.

$$\hat{\phi}_m = \arg \min_{\phi \in \phi} \sum_{i=1}^n \frac{1}{2} \hat{h}_m(x_i) \left[\frac{-\hat{g}_m(x_i)}{\hat{h}_m(x_i)} - \phi(x_i) \right]^2$$

Algorithm For Newton Boosting.

TB-9

Algorithm:

[input:] Data set D
 Loss function L .
 A base learner ϕ .
 Nos of iter. M .
 Learning rate γ .

Step 1] init

$$\hat{f}^{(0)}(x) = \hat{f}_0(x) = \hat{\theta}_0 = \arg \min_{\theta} \sum_{i=1}^n L(y_i, \theta)$$

Step 2] - for $m = 1, 2, \dots, M$

Step 3] do {

$$\hat{g}_m(x_i) = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right] f(x) = \hat{f}^{(m-1)}(x).$$

Step 4]

$$\hat{h}_m(x_i) = \left[\frac{\partial^2 L(y_i, f(x_i))}{\partial f(x_i)^2} \right] f(x) = \hat{f}^{(m-1)}(x).$$

Step 5]

$$\hat{\phi}_m = \arg \min_{\phi \in \phi} \sum_{i=1}^n \frac{1}{2} \hat{h}_m(x_i) \left[\left(\frac{-\hat{g}_m(x_i)}{\hat{h}_m(x_i)} \right) - \phi(x_i) \right]$$

Step 6]

$$\hat{f}_m(x) = \sum \hat{\phi}_m(x)$$

Step 7]

$$\hat{f}^{(m)}(x) = \hat{f}^{(m-1)}(x) + \hat{f}_m(x).$$

8 end]

$$\text{output } \hat{f}(x) = \hat{f}^M(x) = \sum_{m=0}^M \hat{f}_m(x).$$