

Gradient boosting:

- I wish to construct a function approximating some data (x_i, y_i) in some sense.

- So must minimize

$$\int \text{Loss}(f(x), y(x)) P(x) dx$$

(which I can't evaluate - must substitute a sum & eventually.)

But look at the problem it's a variational problem.

I have

$$\min_f J(f(x))$$

Attack by

$$f^{(i)}(x) = \sum_{j=1}^i \alpha_j \phi_j(x, \theta_j)$$

these are parameters

where these are model functions of some sort - eg trees, etc

Assume I have

$f^{(i)}$ - how to get next?

Take functional gradient of $J|_{f^{(i)}}$

- this is a derivative, wrt \underline{f}
- linear function of a small change δf .

$(\nabla_f F)(\phi)$ is best approx
to $F(f + \phi)$

- $(\nabla_f F)(\phi)$ is linear, so it must have the form

$$\int H \cdot \phi \, dx$$

- Notice that if F depends on f , but NOT derivatives of f ,

$$F(f + \phi) = \int \text{loss}(f + \phi, y) P(x) \, dx$$

$$\approx \int \text{loss}(f, y) P(x) \, dx +$$

$$\int \frac{\partial \text{loss}(f, y)}{\partial u} \cdot \phi \cdot P(x) \, dx$$

(A)

This means that the functional gradient is

$$\langle \delta J, \phi \rangle = \int \frac{\partial \text{loss}(f, y)}{\partial u_1} \cdot \phi P(x) dx$$

↑
first argument

(for the case where loss depends only on f not f' , etc)

Now we can do functional gradient descent.

Notice similarity / analogy with gradient

directional derivative in dir a

$$= \nabla f \cdot a$$

dir derivative in dir ϕ

$$= \int \frac{\partial J}{\partial u_1} P(x) \phi dx$$

Gradient Descent:

- assume we have $f^{(n)}$
- functional gradient in dir φ is

$$\int \frac{\delta L}{\delta u_i} \Big|_{u_i = f^{(n)}} \cdot \varphi \cdot P(x) dx.$$

- we cannot represent this accurately, its at least as complex as a function

- Option choose an update from a finite dimensional space of f 's

$$\Delta f^{(n)} \in S$$

↑ this could be trees, etc, etc.

Now, find the element of S that best approximates the gradient

Want to minimize

$$\int \left\| \frac{\partial \mathcal{L}}{\partial u_i} \Big|_{u_i = f^{(u)}} - \Delta f \right\|^2 P(x) dx.$$

By choice of function Δf in S

Notice that, if S consists of functions with bounded norm (some trees, etc).

this is equivalent to

~~max~~
$$\int \frac{\partial \mathcal{L}}{\partial u_i} \Big|_{u_i = f^{(u)}} \cdot \Delta f \cdot P(x) dx$$

analogy: $\frac{\partial \mathcal{L}}{\partial u_i} \Big|_{u_i = f^{(u)}} \equiv \nabla f$

$$\int (\alpha, \beta) P(x) dx \equiv \text{dot product}$$

So we are looking for the Δf that has largest dot product w/ "gradient"

- for many types of function model, this is fairly easy.

eg. $S =$ Decision stumps

$S =$ Decision trees

- Now once we have $\Delta f^{(n)}$, we can choose α_n by line search

$$f^{(n+1)} = f^{(n)} + \alpha_n \Delta f^{(n)}$$

line search

~~$$g(\alpha_n) = \int \text{loss}(f^{(n)}, y)$$~~

Line search:

$$g(\alpha_n) = \mathbb{E}(f \stackrel{(n)}{\leftarrow} \alpha_n \Delta f \stackrel{(n)}{\rightarrow})$$

$$= \int \text{loss}(f \stackrel{(n)}{\leftarrow} \alpha_n \Delta f \stackrel{(n)}{\rightarrow}, y(x)) P(x) dx$$

$$\approx \sum_{\text{ie examples}} \text{loss}(f(x_i) \stackrel{(n)}{\leftarrow} \alpha_n \Delta f(x_i) \stackrel{(n)}{\rightarrow}, y_i)$$

Notice that we can iterate ad nauseam
~~particularly if the~~

- the cost is applied to an ∞ -dim obj.
- we are taking fd steps.

Example

- Classification with stumps.

- $f(x) = \sum_{i \in \text{stumps}} \alpha_i \{ \text{value of leaf at } x \}$

where ~~stumps~~ leaves have values $-1, 1$

- $\text{loss} = \max(0, 1 - y_{\text{known}} \cdot f_{\text{pred}})$

Procedure:

- have $f^{(n)}$
- $\frac{\partial L}{\partial \alpha_i} \Big|_{f^{(n)}} = \begin{cases} 0 & \text{if } y_k \cdot f_{\text{pred}} > 1 \\ -y_{\text{known}} & \text{otherwise} \end{cases}$

\Rightarrow new stump tries to have

items too close to margin on

either side dep. sign y_{known}

and ignores ~~correct~~ answers to the far side of margin.

NUMEROUS other examples in
Friedman paper, q.v.