# Gradient boosting

- can be used for classification or regression

- iterative method (not easily parallelizable)

- ensemble model
- need to choose an appropriate loss function given problem at hand

- need to choose a set of weak learners to choose from at each iteration, $H$
(often trees/stumps) can also be splines or anything else.
- the algo learns a function of the form:

$$ F(x) \in \left\{ \sum_{m=1}^{M} \nu_m h(x; a_m) \mid \nu_m \in R, h(\cdot) \in H \right\} $$

## Idea of gradient boosting

- instead of fitting a parameterized model by employing GD in parameter space to minimize a cost function, run gradient decent in a function space of dims $N \leftarrow \#$ data points $(R^N)$, and at each iteration approximate the gradient by choosing a parameterized weak learner from $H$.

· goal is to find:

$$F^* = \underset{F \in \mathcal{F}}{\text{argmin}} \left\{ \sum_{i=1}^{N} \ell(y^{(i)}, F(x^{(i)})) \right\} \qquad (\text{i.e. minimize empirical risk})$$

let $\vec{F} = [F(x^{(1)}), F(x^{(2)}), \ldots, F(x^{(N)})]^T$

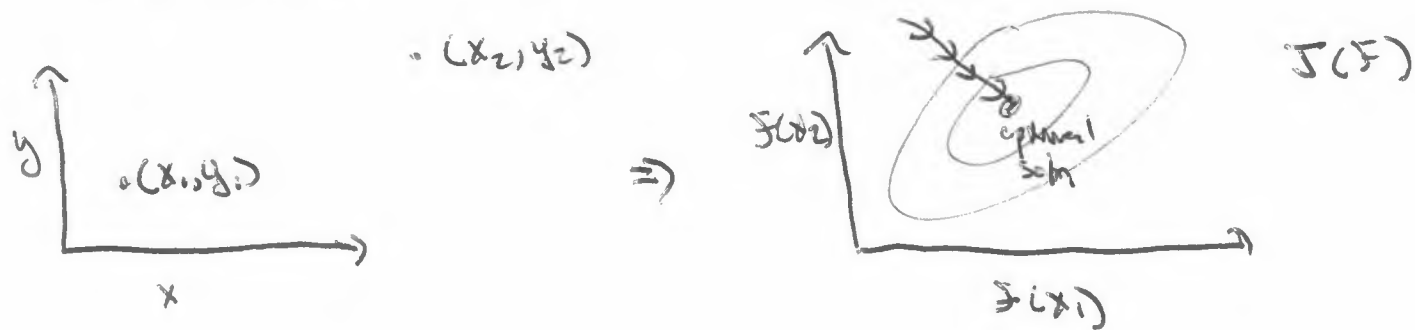$\Rightarrow$ $J(\vec{F}) = \sum_{i=1}^{N} \ell(y^{(i)}, \vec{F}_i)$

$\Rightarrow$ Thus we run GD on the space formed by $F_1, F_2, \ldots, F_N$ to find $\vec{F}^*$ the vector that minimizes the cost

$$\vec{F}^{(m)} = \vec{F}^{(m-1)} - \rho_m \vec{g}_m$$

where $[\vec{g}_m]_j = \frac{\partial}{\partial F_j} \sum_i \ell(y^{(i)}, \vec{F}_i^{(m)}) = \frac{\partial \ell(y^{(i)}, \vec{F}_j^{(m-1)})}{\partial F_j}$

and $\rho_m = \underset{\rho \in \mathbb{R}}{\text{argmin}} \left\{ \sum_i \ell(y^{(i)}, \vec{F}^{m-1} - \rho \vec{g}_m) \right\}$

in pictures, if we have 2 data points in the training set:



- Note that we wish to predict $y$ at more points than just $x_1, x_2$
- to extrapolate to other values of $x$, we guess a parameterized function for the gradient at each boosting iteration. Let $\vec{\phi}_{\vec{a}} = [\phi(x_1; \vec{a}), \phi(x_2; \vec{a}) \dots \phi(x_N; \vec{a})]^T$

then the algo is:

① $F^0(x) = 0$

② For $m = 1$ to $M$:

    - compute $\vec{g}_m$ as given by formula on previous page

    - approximate $\vec{g}_m$ w/ a parameterized weak learner by finding the vector closest in $\mathbb{R}^N$ (in $L^2$ sense):

$$\phi(x, a_m) = \arg\min_a \left\{ \sum_{i=1}^{N} \left( -\vec{g}_m^{(i)} - \vec{\phi}_{\vec{a}}^{(i)} \right)^2 \right\}$$

$\longrightarrow$

$$- \rho_m = \underset{\rho}{\text{argmin}} \left\{ \sum_i \ell(y^{(i)}, \vec{F}^{m-1} - \rho \, \phi(x^{(i)}; \vec{a}_m)) \right\}$$

$$- \vec{F}^m_i = \vec{F}^{m-1} + \rho_m \, \phi(x; \vec{a}_m)$$

output: $F^M$

- For regression: $\hat{y}(x) = \sum_{m=1}^{M} \rho_m \, \phi(x; \vec{a}_m)$     (a weighted sum)

- For classification: $\hat{y} = \text{sign}\left( \sum_{m=1}^{M} \rho_m \, \phi(x; \vec{a}_m) \right)$

(loss function typically defined in such a way that it is easy to compute class probabilities from $\sum \rho_m \, \phi(x; \vec{a}_m)$)

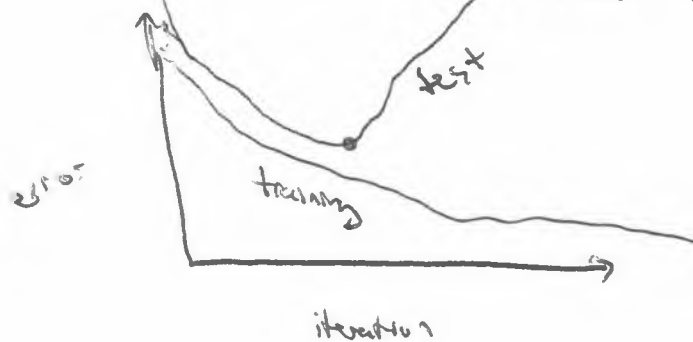- thus we traverse "$F$" space like this

m pictures:

- ways to regularize/improve test set performance

  - shrinkage  let $\nu \in [0,1]$
                 ⤷ step size

  - by making step size smaller test error seems to improve (will need more
                                                           iterations to converge though)

- hold out some of the data on each iteration (fit on less training data, so less
                                                     propensity to overfit)

  - has the advantage that we can
    test error at each iteration on an oob sample to build up a learning
    curve to throw when to stop



- early stopping                                    ⤷

* these regularization methods seen to work well
  from empirical studies.

Pros

- Very good performance (typically)
- lots of flexibility in choosing loss function + weak learner to deal w/ problem at hand
- can learn complicated decision functions
- non-parametric (don't need to have much knowledge about how function should look)

Cons

- can't parallelize (XGBoost can llize trees, but not forest)
- can be difficult to tune
- can overfit if not careful
- training can take a while