

# MATH 342W Lec 23

Bagging: meta-algorithm

$$f_{\text{bag}} := \frac{f_1 + f_2 + \dots + f_M}{M} \text{ for regression}$$

$$f_{\text{bag}} := \text{Mode}[f_1, f_2, \dots, f_M] \text{ for classification}$$

The model  $f_i$  is called the "base-learner"

which could be the result of any algorithm.

Which base learners do well? Those that have low bias and high variance and are as independent as possible

Boasting: <sup>rather</sup> meta-algorithm  $\rightarrow$  ("Strong learners")

$$f_{\text{boost}} := f_1 + f_2 + \dots + f_M$$

where  $f_2$  is fit on where  $\hat{y}_1 = f_1(x)$  doesn't fit,

$f_3 \dots \dots \dots \hat{y}_2 = f_1(x) + f_2(x)$  doesn't fit

$\vdots$   
 $f_M \dots \dots \dots \hat{y}_{M-1} = f_1(x) + f_2(x) + \dots + f_{M-1}(x)$  "etc"

which is "iterative self-correction"

There are many ways to implement boosting. We will study "gradient boosting". What does "gradient" mean.

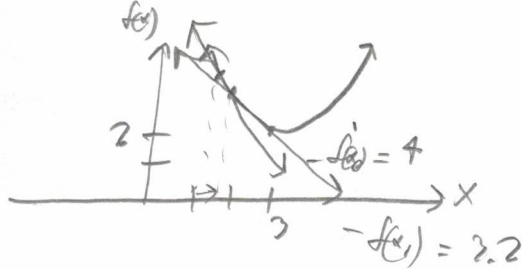
Go back to calculus...

Consider  $f(x) = 2 + (x-3)^2$

we want to find

$x_0 := \operatorname{argmin}_{x \in \mathbb{R}} \{f(x)\}$  or at

least a local minimum



we first solve for derivative  $f'(x) = 2(x-3)$

Pretend that you cannot solve for  $x$  if  $f'(x) = 0$ .

so... start at initial value  $x_0 \stackrel{\text{e.g.}}{=} 1$

Here  $-f'(x_0) = -2(1-3) = 4$

iterate to get closer which is defined as

If a "step" multiple fraction of the negative derivative, we get closer to the local minimum. Let  $\eta = 0.1$ , the step size.

$x_1 = x_0 + \eta(-f'(x_0)) = 1 + 0.1 \cdot 4 = 1.4$ ,  $-f'(x_1) = -2(1.4-3) = 3.2$

we iterate again

$\Rightarrow x_2 = x_1 + \eta(-f'(x_1)) = 1.4 + 0.1 \cdot 3.2 = 1.72$

This is called "gradient descent"

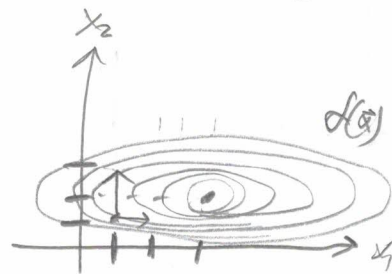
In  $n$  dimensions e.g.  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$f(\vec{x}) = 2 + (x_1-3)^2 + (x_2-2)^2$

$\vec{x}_0 := \operatorname{argmin}\{f(\vec{x})\}$  will be a 2-dim vector

we first solve for  $\vec{\nabla} f(\vec{x}) := \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2(x_1-3) \\ 2(x_2-2) \end{bmatrix}$

we start at  $\vec{x}_0 \stackrel{\text{e.g.}}{=} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ ,  $-\vec{\nabla} f(\vec{x}_0) = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$



Now we step  $\vec{x}_1 = \vec{x}_0 + \eta(-\vec{\nabla} f(\vec{x}_0)) = \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0.1 \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 1.3 \\ 1.4 \end{bmatrix}$  and repeat

Let's transfer this to strong  $\mathcal{D}$  learning.

Gradient Descent  
 $f(\vec{x}) \rightarrow$  <sup>any  $\mathcal{D}$</sup>  loss function on the current fit  $\vec{y}$ ,  $L(\vec{y}, \vec{y})$  which must be differentiable in all  $n$  dimensions  
 $\vec{x}_0 \rightarrow$  starting prediction  $\vec{y}_0$

$$\vec{\nabla} f(\vec{x}) \rightarrow \vec{\nabla} L(\vec{y}, \vec{y}) = \begin{bmatrix} \frac{\partial L}{\partial y_1} \\ \vdots \\ \frac{\partial L}{\partial y_n} \end{bmatrix}$$

$$\vec{y}_1 = \vec{y}_0 - \eta \vec{\nabla} L(\vec{y}, \vec{y}_0)$$

$$\vec{y}_2 = \vec{y}_1 - \eta \vec{\nabla} L(\vec{y}, \vec{y}_1)$$

$\vdots$

$\vec{y}_m \approx \vec{y}_*$ , the predictions that minimize the objective function.

But the goal of supervised learning is not to find best  $\vec{y}$  for  $\mathcal{D}$ , it's to build a generalizable model  $g$ .

So we instead want to use gradient descent in the space of functions  $g: \mathbb{R}^p \rightarrow \mathbb{R}$ . How can we amend the above to do this?

This is why it took until 2001 to make this as it's not so easy to "port" this idea over.

The trick is to fit a model to the gradient itself.  
This is strange is we only fit models to data. So

let's translate:

Fitting 0 with  
Gradient Descent

Gradient boosting for supervised learning using base learner  $A$

$$\hat{y}_0 \longrightarrow f_0, \text{ the default model } \hat{y}_0(x)$$

$$\hat{y}_1 = \hat{y}_0 - \eta \vec{\nabla} L(\hat{y}, y) \longrightarrow G_1 = f_0 + \eta A(\langle X, -\vec{\nabla} L(\hat{y}, \hat{y}_0) \rangle, \mathcal{H})$$

Let  $G_t := f_0 + f_1 + \dots + f_t$   
 $\uparrow$   
 the partial sums

instead of fitting a model to  
the data, you're fitting  
a model to the negative  
gradient

$f_1$  is a step in the direction  
in function space towards  
the "best" model

$$\hat{y}_2 = \hat{y}_1 - \vec{\nabla} L(\hat{y}, \hat{y}_1) \longrightarrow G_2 = G_1 + \eta A(\langle X, -\vec{\nabla} L(\hat{y}, \hat{y}_1) \rangle, \mathcal{H})$$

$$\vdots$$

$$\hat{y}_m = \hat{y}_{m-1} - \vec{\nabla} L(\hat{y}, \hat{y}_{m-1}) \longrightarrow G_m = G_{m-1} + \eta A(\langle X, -\vec{\nabla} L(\hat{y}, \hat{y}_{m-1}) \rangle, \mathcal{H})$$

Which algorithms work well in gradient boosting? Weak learners (high bias, low variance) that can still "span" the space of non-linear and interactive models (interactions among the  $p$  features).

$\Rightarrow$  Trees with low depth / high  $N_0$  node size is  
usually the default choice

We will now apply gradient descent to both regression and prob. estimation for binary outcome.

Regression: let the objective function to be minimized be:

$$L(\vec{y}, \vec{\hat{y}}) := \sum_{i=1}^n (y_i - \hat{y}_i)^2 = SSE. \text{ Needs to be ... just must be differentiable}$$

Prob est: recall the likelihood function

$$\prod_{i=1}^n \hat{p}_i^{y_i} (1 - \hat{p}_i)^{1-y_i}$$

It will be easier to use log odds:

$$\text{let } \hat{y}_i := \ln\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) \Leftrightarrow \hat{p}_i = \frac{e^{\hat{y}_i}}{1+e^{\hat{y}_i}} \quad \begin{matrix} \text{1:1} \\ \text{function} \end{matrix}$$

The likelihood function becomes

$$\prod_{i=1}^n \left(\frac{e^{\hat{y}_i}}{1+e^{\hat{y}_i}}\right)^{y_i} \left(\frac{1}{1+e^{\hat{y}_i}}\right)^{1-y_i} = \prod_{i=1}^n \frac{e^{y_i \hat{y}_i}}{1+e^{\hat{y}_i}}$$

Since maximizing likelihood is the same as max. log likelihood, we can maximize

$$\ln(\cdot) = \sum_{i=1}^n \ln\left(\frac{e^{y_i \hat{y}_i}}{1+e^{\hat{y}_i}}\right) = \sum_{i=1}^n y_i \hat{y}_i - \ln(1+e^{\hat{y}_i})$$

Since we want to minimize function in gradient descent, we now have our objective function:

$$L(\vec{y}, \vec{\hat{y}}) = \sum_{i=1}^n -y_i \hat{y}_i + \ln(1+e^{\hat{y}_i})$$

$$-\vec{\nabla} L(\vec{y}, \vec{\hat{y}}) = \begin{bmatrix} -y_1 + \frac{e^{\hat{y}_1}}{1+e^{\hat{y}_1}} \\ -y_2 + \frac{e^{\hat{y}_2}}{1+e^{\hat{y}_2}} \\ \vdots \\ -y_n + \frac{e^{\hat{y}_n}}{1+e^{\hat{y}_n}} \end{bmatrix}$$

$$= -\vec{y} + \frac{e^{\vec{\hat{y}}}}{1+e^{\vec{\hat{y}}}} \xrightarrow{\text{intuition}} -\vec{y} + \vec{\hat{p}} \xrightarrow{\text{we have}} \langle X, \vec{y} - \vec{\hat{p}} \rangle$$

$$-\vec{\nabla} L(\vec{y}, \vec{\hat{y}}) = \begin{bmatrix} -2(y_1 - \hat{y}_1) \\ -2(y_2 - \hat{y}_2) \\ \vdots \\ -2(y_n - \hat{y}_n) \end{bmatrix}$$

$$= -2 \vec{e}$$

$\Rightarrow$  essentially you fit a shallow regression, then take

$$\langle X, \vec{e} \rangle$$

modulo a constant...

very intuitive

modulo a constant, intuitive!!

↑ end of course

↓ extra

Possible improvement; eliminate hyperparameter  $\eta$ .

$$\text{let } G_t = G_{t-1} + \eta^* A(\langle x, -\vec{\nabla}_{\vec{y}} \tilde{f}_{t-1} \rangle, \mathcal{A})$$

$$\text{where } \eta^* = \underset{\eta \in (0, \infty)}{\operatorname{argmin}} \left\{ L(\vec{y}, G_{t-1}(x) + \eta \tilde{f}_t(x)) \right\}$$

6