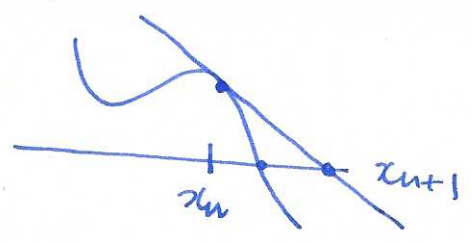


XG Boost

• gradient descent: 1-d version Newton's method:  
looking for zeros.

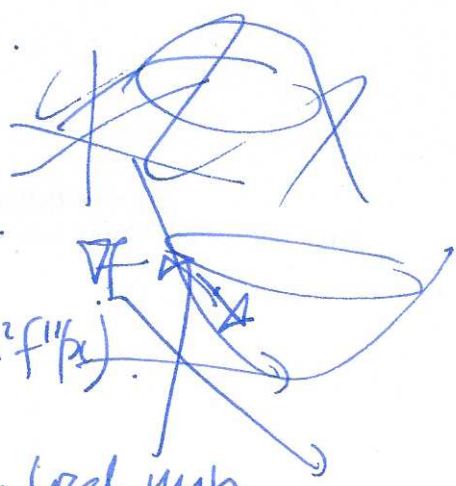


high dimension version:

$$x_{n+1} = x_n - \frac{\nabla f}{f''}$$

• linear approx:  $f(x+t) \approx f(x) + tf'(x)$

• quadratic approx:  $f(x+t) \approx f(x) + tf'(x) + \frac{1}{2}t^2 f''(x)$

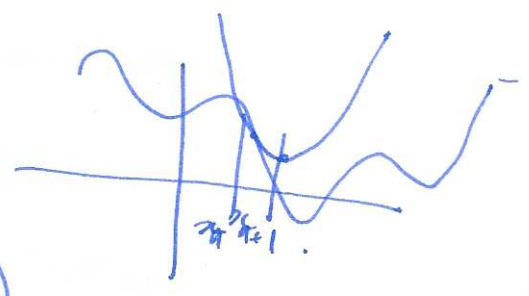


note: if  $f''(x) > 0$  then  $f$  has a local min and you can find it!

quadratic expression  
 $q(x) = f(x) + tf'(x) + \frac{1}{2}t^2 f''(x)$

$$q'(x) = f'(x) + tf''(x) = 0 \leftarrow \text{solve for critical pt}$$

$$t = - \frac{f'(x)}{f''(x)}$$



so set  $x_{n+1} = x_n + t = x_n - \frac{f'(x_n)}{f''(x_n)}$

note if original function actually quadratic this gives the correct answer in 1 step!

high dimensional version:  $f'(x) \Leftrightarrow \nabla f$  - vector of partial derivatives

$$f''(x) = \nabla^2 f(x) = H_f(x) = \text{matrix of 2nd partial derivatives} \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_{i-1} \partial x_i} & \dots & \frac{\partial^2 f}{\partial x_i^2} \end{bmatrix}$$

$$x_{n+1} = x_n - \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_{i-1} \partial x_i} & \dots & \frac{\partial^2 f}{\partial x_i^2} \end{bmatrix}^{-1} = (H_f(x))^{-1} \cdot \nabla f$$


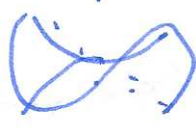
problem . need to compute inverse of matrix, Gaussian elimination methods are  $O(n^3)$ .

• if  $H^{-1}$  very small det,  $H^{-1}$  may be very large, i.e. local min may be far away, so common modification is to fix  $\epsilon > 0$  and  $0 < \epsilon < 1$

we:  $x_{n+1} = x_n - \epsilon H_{f(x_n)}^{-1} \nabla f(x_n)$ .

Fact: if  $f$  is strongly convex and  $H$  is Lipschitz and  $x_0$  is close enough to  $x_{min}$ , get quadratic convergence.

note instead of computing  $H_f^{-1}(x_n)$  can solve  $[H_f(x_n)] h = -f'(x_n)$   
(then  $h = H_f^{-1}(x_n) \nabla f$  if  $H_f^{-1}$  invertible).  
 $n \times n$  matrix       $n$ -vector

(dis) advantages: still  $O(n^3)$ , will still get (family) of solutions if  $\det(H) = 0$   
can use special methods assuming  $H$  is positive definite - if special methods don't work and  $H$  not true definite have problems anyway...  
  $f$  quad true det  $>$   
  $f$  quad saddle - critical point not min!

• problem: if  $\det(H)$  close to zero then method is probably unstable.  
- usual convergence problems for Newton's method / Gradient descent.

initial data:  $X = \{(x_i, y_i)\}_{i=1}^n$   $x_i \in \mathbb{R}^d$

predicting quantitative data

eg  $L(\theta) = \frac{1}{n} \sum_i (y_i - \hat{f}(x_i))^2$   
(e.g. mean sum of squares  
mean sum of the differences)

loss function  $L(y, \hat{f}(x))$

Logistic loss

$L(\theta) = \sum_i (y_i \ln(1 + e^{-\hat{y}_i}) + (1 - y_i) \ln(1 + e^{\hat{y}_i}))$   
↑ predicting probabilistic qualitative data.

regularization function / shrinkage function / pruning function  $\Omega(\theta)$

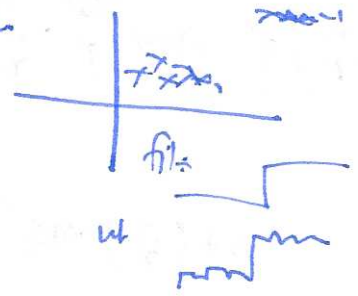
actually minimize

$obj(\theta) = L(\theta) + \Omega(\theta)$

↑ error  
minimization  
term

↑ model  
complexity term.  
(e.g. # of nodes)

example

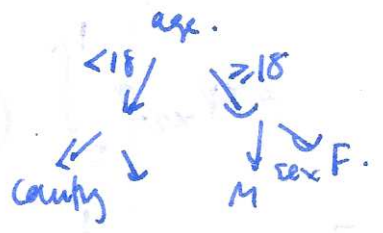


recall decision trees.

e.g. use age, sex, country to predict height.

24, M, dutch.  
71, M, Japanese.

tree



many trees: (eg. random forest:  $\{T_1, T_2, \dots\}$ ).  
take average...

- need to choose trees.
- need to combine trees.

$\hat{f}(x) = \sum_T f_T(x)$  where  $f_T(x) =$   
 $\Omega(\theta)$  could be # of trees

• # of trees + # leaves

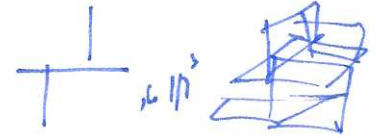
problem: each tree has a certain number of parameters (where to choose each branch, leaf value)

- usually can't do just solve to find optimal parameters
- so need to use gradient descent / Newton Raphson / Hessian method...

simplification: just add one tree at a time!

a single tree  $t$  is:  $f_t(x) = w_{q(x)}$ ,  $w \in \mathbb{R}^T$

$w$  vector of values on each leaf  $q: \mathbb{R}^d \rightarrow \{1, 2, \dots, T\}$

$q$  function sends data points /  $\mathbb{R}^d \rightarrow$  leaves, e.g. in  $\mathbb{R}^2$  

$w(f_t) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$       $\gamma, \lambda$  tuning parameters we need to pick, 0.3?

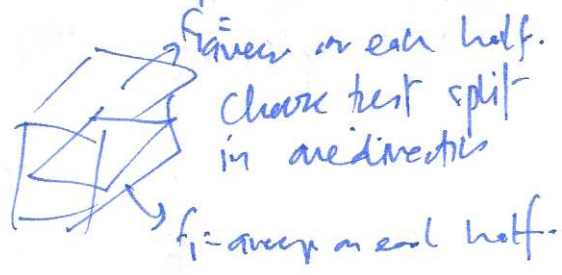
$T = \# \text{leaves.}$

so as you add leaves, if improvement  $< \gamma$  then better off not adding that leaf.

build first tree:



$f_0 = \text{const func}$   
 $= \text{avg}$



$L(f_1) < L(f_0)$ .  
↑  
but only add if at least  $\gamma$  better!

add second tree: we have the function from the first tree  $f_1: \mathbb{R}^d \rightarrow \mathbb{R}$ .

replace data by error in prediction / residuals:  $\{(x_i, y_i)\}_{i=1}^n$   
↓  
 $\{(x_i, y_i - f_1(y_i))\}_{i=1}^n$

now use  $f_1 + \alpha f_2$ .  
↑  
learning parameter (0.3?)

keep going...