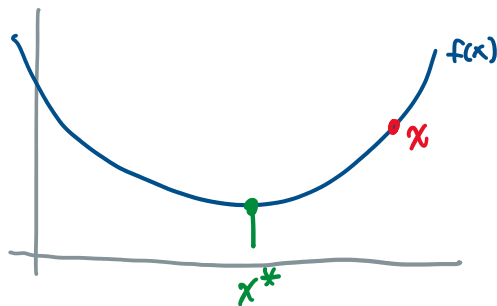


Gradient Descent

Thursday, January 23, 2020 13:26

Let's look at 1-D case: $f(x)$



$$\begin{aligned} \text{if } x > x^* &\Leftrightarrow f'(x) > 0 \\ \text{if } x < x^* &\Leftrightarrow f'(x) < 0 \\ \text{if } x = x^* &\Leftrightarrow f'(x) = 0 \end{aligned}$$

Algorithm:

① Initialize x to some value

② if $f'(x)$ is close to 0, then done ✓

else if $f'(x) > 0$, then $x_{i+1} = x_i - \eta$

else if $f'(x) < 0$, then $x_{i+1} = x_i + \eta$

step size

if step size too small: algorithm too slow

too large: oscillation, no convergence.

typically: step size \downarrow over iterations.

for $d > 1$:

$$f(\underline{x}_1) = f(\underline{x}_0 + \eta \cdot \underline{v})$$

some direction to minimize $f(\underline{x})$

$$= f(\underline{x}_0) + \eta \underline{v}^T \nabla f(\underline{x}_0) + \dots$$

ignore

$$\approx f(\underline{x}_0) + \eta \underline{v}^T \nabla f(\underline{x}_0)$$

learning rate

$$\Rightarrow \underline{v}^* = \underset{\|\underline{v}\|=1}{\operatorname{argmin}} \left(f(\underline{x}_0) + \eta \underline{v}^T \nabla f(\underline{x}_0) \right)$$

not minimizable

minimize this

$$\underline{v}^* = \frac{-\nabla f(\underline{x}_0)}{\|\nabla f(\underline{x}_0)\|}$$

Updating Weights:

adaptive "learning rate"



Req. 1:

$$w(t+1) = w(t) + \eta \cdot \underline{v}_t$$

$$\underline{w}(t+1) = \underline{w}(t) + \eta_t \cdot \underline{v}_t$$

$\eta_t = \eta \cdot \underbrace{\|\nabla E_{in}\|}_{\nabla f(x_t)}$

$\eta_t \underline{v}_t = -\eta \nabla f(x_t)$

Recall in perceptron we initialize to \emptyset

$$\underline{w}(t=0) = \text{normal}(0, I)$$

▷ could be # of iterations

———— change of error : $|\Delta E_n| < \delta,$

① Initialize weights at $t=0$ w/ $\underline{w}(t=0)$

③ Compute $\nabla E_{in}(\underline{w}(t))$

⑤ Update weights

$$\underline{w}(t+1) = \underline{w}(t) + \eta \cdot \underline{V}(t)$$

ELEC 400M sidan 2

Gradient Descent in Logistic Regression

Recall logistic regression cost function:

$$E_{in} = \frac{1}{N} \cdot \sum_{n=1}^N \underbrace{\log(1 + e^{-y_n \cdot \underline{w}^T \underline{x}_n})}_{e_n}$$

$$\begin{aligned} \nabla E_{in} &= \frac{\partial E_{in}}{\partial \underline{w}} = \frac{1}{N} \sum_{n=1}^N \left(\frac{\partial E_{in}}{\partial \underline{w}} \right) \\ &= \frac{1}{N} \sum_{n=1}^N \frac{1}{\log(1 + e^{-y_n \underline{w}^T \underline{x}_n})} \cdot (e^{-y_n \underline{w}^T \underline{x}_n}) \cdot (-y_n \underline{x}_n) \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\partial y_n \underline{x}_n}{1 + e^{y_n \underline{w}^T \underline{x}_n}} \end{aligned}$$

EX. linear regression solve vs. GD.

$$f(\underline{x}) \equiv E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N (\underline{w}^T \underline{x}_n - y_n)^2$$

Recall closed form sol: $\underline{x}^T \underline{x} \underline{w}^* = \underline{x}^T \underline{y}$

$$\underline{w}^* = \underline{x}^T \underline{y} (\underline{x}^T \underline{x})^{-1}$$

computation complexity: $\mathcal{O}(N \cdot d^2 + d^3)$ operations

alternatively we could use gradient descent.

$$\underline{w}(t+1) = \underline{w}(t) - \eta \cdot \frac{2}{N} \sum_{n=1}^N (\underline{w}^T(t) \underline{x}_n - y_n) \cdot \underline{x}_n$$

$$\underline{w}(t+1) = \underline{w}(t) - \eta \cdot \underbrace{\frac{2}{N} \sum_{n=1}^N (\underline{w}^T(t) \underline{x}_n - y_n) \cdot \underline{x}_n}_{\nabla}$$

Complexity: $O(N \times d \times \frac{I}{\epsilon})$ operations
iterations

\Rightarrow gradient descent is faster if d is large
 \Rightarrow gradient descent is not good for large N
 (gradient is expensive for large N)