#### Optimization

Friday, 30 April 2021 07:42

## Gradient Descent

- In general used to find a minimum of any function

- In ML mostly used to find a minimum of the error function (typically MSE) dependent on model parameters

General: 
$$\vec{\chi} = (\chi_{\eta_1, \dots, \chi_n})$$
  $f(\chi_{\eta_1, \dots, \chi_n})$ 

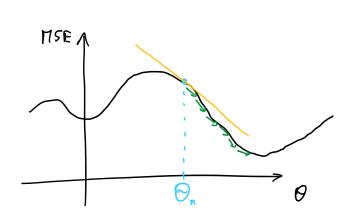
### Example

$$\begin{array}{cccc}
\gamma \pi i n & \sum_{(x_i, y) \in \Omega} (y - (\Theta_i x + \Theta_i))^2 \\
(\theta_i, \theta_i) & (x_i, y) \in \Omega
\end{array}$$

## Idea of 60

1. Start eith any B 2. Iteratively move  $\vec{\Theta}$ 

in the direction opposite to the derivative



The slope of the tangent line is equal to the derivative of the function eith respect to 
$$\Theta$$

# Derivative calculation for 175E

For simplicity we assume a linear model  $f(x|Q_0,Q_1) = f(x,\theta_0,\theta_1) = \theta_1 x + \theta_0$   $f(x|Q_0,Q_1) = f(x,\theta_0,\theta_1) = \theta_1 x + \theta_0$   $f(x|Q_0,Q_1) = f(x,y) \in 0 \quad (y - (\theta_1 x + \theta_0))^2$ 

Since (f+g)'(x) = f'(x) + g'(x) and (af(x))' = af'(x) we can make the derivative calculations for a single element in the sum (single datapoint) and then average.

$$\frac{\partial}{\partial \theta_{o}} \left( y - (\theta_{1} \times + \theta_{o}) \right)^{2} = 2 \left( y - (\theta_{1} \times + \theta_{o}) \right) \cdot \frac{\partial}{\partial \theta_{o}} \left( y - (\theta_{1} \times + \theta_{o}) \right)$$

$$= 2 \left( y - (\theta_{1} \times + \theta_{o}) \right) (-1)$$

$$= -2 \left( y - (\theta_{1} \times + \theta_{o}) \right)$$

$$= -2 \left( y - (\theta_{1} \times + \theta_{o}) \right)$$

$$\frac{\partial}{\partial \theta_{1}} \left( y - (\theta_{1} \times + \theta_{0}) \right)^{2} = 2 \left( y - (\theta_{1} \times + \theta_{0}) \right) \cdot \frac{\partial}{\partial \theta_{1}} \left( y - (\theta_{1} \times + \theta_{0}) \right)$$

$$= 2 \left( y - (\theta_{1} \times + \theta_{0}) \right) (-x)$$

$$= -2 \times \left( y - (\theta_{1} \times + \theta_{0}) \right)$$

$$= -2 \times \left( y - (\theta_{1} \times + \theta_{0}) \right)$$

Updating rule

$$\vec{\partial}^{n} = \vec{\partial}^{n-1} - \lambda \frac{\partial}{\partial \theta} MSE(\vec{\theta})$$

$$\Theta_{0}^{n} = \Theta_{0}^{n-1} - \lambda \frac{\partial}{\partial \Theta_{0}} MSE(\Theta_{0}^{n-1}\Theta_{1}^{n-1})$$

$$= \Theta_{0}^{n-1} - \lambda \frac{1}{|\Theta|} \sum_{(x,y)\in \Theta} (-2)(y - (\Theta_{1}x + \Theta_{0}))$$

$$\Theta_{1}^{n} = \Theta_{1}^{n-1} - \lambda \frac{\Im}{\Im \Theta_{1}} MSE(\Theta_{0}^{n-1}\Theta_{1}^{n-1})$$

$$= \Theta_{1}^{n-1} - \lambda \frac{1}{101} \sum_{(x,y) \in \Theta} (-2x)(y - (\Theta_{1}x + \Theta_{0}))$$

# Stochastic Gradient Wescent (SGO)

- 1. Start with any of
- 2. Iteratively:
  - a) take a datapoint  $(\vec{x}_{i}y) \in 0$
  - b) calculate the derivative of MSE on this single datapoint with respect to  $\vec{\theta}$
  - c) shift 0 in the direction opposite to the derivative

Croblem: SGO can be unstable and diverge

Botch Gradient Descent

Remark: most of the time (
when people say SGO
they mean Batch GO

- 1. Start with any o
- 2. Iteratively:
  - a) take a elatapoint  $(\tilde{X}_{i}y) \in \Omega$
  - b) calculate the derivative of MSE on this single datapoint

- on this single datapoint with respect to  $\vec{\theta}$ c) shift  $\theta$  in the direction opposite to the derivative
- There are many other variants of 560 used in practice:
  - SGD with momentum
  - Rmsgroß
  - NAG
  - Adam (the most popular)
  - Ada Grad
  - A de Delta