

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 formulation :  $\vec{x} = \min_{(x_1, \dots, x_n)} f(x_1, \dots, x_n)$

ML formulation (special case of the above) :  $\vec{\theta} = \min_{(\theta_1, \dots, \theta_n)} \sum_{(\vec{x}, y) \in \mathcal{D}} (y - f(\vec{x} | \vec{\theta}))^2$   
 ↑  
 training data

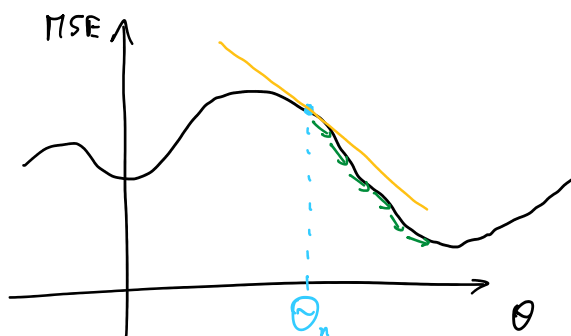
Example

$$\min_{(\theta_0, \theta_1)} \sum_{(x, y) \in \mathcal{D}} (y - (\theta_1 x + \theta_0))^2$$

$$f(x | \theta) = \theta_1 x + \theta_0$$

Idea of GD

1. Start with any  $\vec{\theta}$
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 with respect to  $\theta$

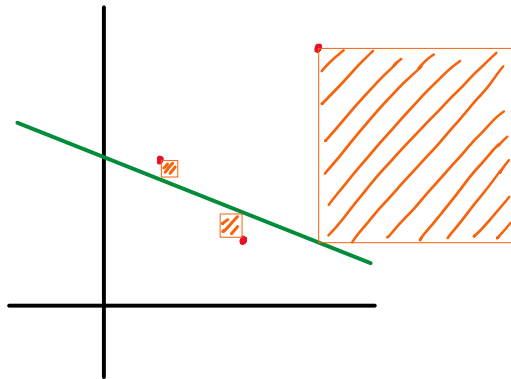
$$\frac{\partial \text{MSE}}{\partial \theta}$$

## Derivative calculation for MSE

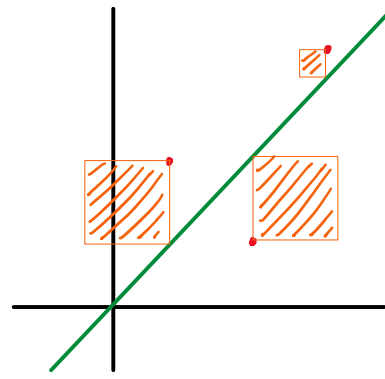
For simplicity we assume a linear model

$$f(x|\theta_0, \theta_1) = f(x, \theta_0, \theta_1) = \theta_1 x + \theta_0$$

$$MSE = \frac{1}{|D|} \sum_{(x,y) \in D} (y - (\theta_1 x + \theta_0))^2$$



Error = [small square] + [small square] + [large square]



Error = [medium square] + [large square] + [small square]

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

$$\begin{aligned} \frac{\partial}{\partial \theta_0} (y - (\theta_1 x + \theta_0))^2 &= 2(y - (\theta_1 x + \theta_0)) \cdot \frac{\partial}{\partial \theta_0} (y - (\theta_1 x + \theta_0)) \\ &= 2(y - (\theta_1 x + \theta_0)) (-1) \\ &= -2(y - (\theta_1 x + \theta_0)) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \theta_1} (y - (\theta_1 x + \theta_0))^2 &= 2(y - (\theta_1 x + \theta_0)) \cdot \frac{\partial}{\partial \theta_1} (y - (\theta_1 x + \theta_0)) \\ &= 2(y - (\theta_1 x + \theta_0)) (-x) \end{aligned}$$

$$= -2x(y - (\theta_1 x + \theta_0))$$

Updating rule

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

$$\theta_0^n = \theta_0^{n-1} - \alpha \frac{\partial}{\partial \theta_0} \text{MSE}(\theta_0^{n-1}, \theta_1^{n-1})$$

$$= \theta_0^{n-1} - \alpha \frac{1}{|D|} \sum_{(x,y) \in D} (-2)(y - (\theta_1 x + \theta_0))$$

$$\theta_1^n = \theta_1^{n-1} - \alpha \frac{\partial}{\partial \theta_1} \text{MSE}(\theta_0^{n-1}, \theta_1^{n-1})$$

$$= \theta_1^{n-1} - \alpha \frac{1}{|D|} \sum_{(x,y) \in D} (-2x)(y - (\theta_1 x + \theta_0))$$

## Stochastic Gradient Descent (SGD)

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

Problem : SGD can be unstable and diverge

## Mini-batch Gradient Descent

1. Start with any  $\vec{\theta}$

Remark : most of the time when people say SGD they mean Mini-batch GD

2. Iteratively :

- take a mini-batch  $\{(\vec{x}, y)\} \subset D$
- calculate the derivative of MSE on this mini batch with respect to  $\vec{\theta}$
- shift  $\theta$  in the direction opposite to the derivative

There are many other variants of SGD used in practice :

- SGD with momentum
- Rmsprop
- NAG
- Adam (the most popular)
- Ada Grad
- Ada Delta

### Stochastic Gradient Descent for matrix factorization

$$\min_{\{p_u, q_i\}} \sum_{(u,i) \in K} \underbrace{(r_{ui} - q_i^T p_u)^2}_{\text{error}(u,i)} + \underbrace{\lambda (\|q_i\|^2 + \|p_u\|^2)}_{L_2 \text{ regularization term}}$$

$\underbrace{\hspace{15em}}_{\text{error}}$

For simplicity consider embedding  $\text{dim} = 2$

$$p_u = (p_{u1}, p_{u2})$$

$$q_i = (q_{i1}, q_{i2})$$

The stochastic gradient descent step looks as follows

$$p_{u1}^{(n)} = p_{u1}^{(n-1)} - \alpha \frac{\partial}{\partial p_{u1}} \text{error}(u,i)$$

$$\rho_{u2}^{(n)} = \rho_{u2}^{(n-1)} - \alpha \frac{\partial}{\partial \rho_{u2}} \text{error}(u, i)$$

$$q_{i1}^{(n)} = q_{i1}^{(n-1)} - \alpha \frac{\partial}{\partial q_{i1}} \text{error}(u, i)$$

$$q_{i2}^{(n)} = q_{i2}^{(n-1)} - \alpha \frac{\partial}{\partial q_{i2}} \text{error}(u, i)$$

Substituting the formula for the error we get

$$\begin{aligned} \frac{\partial}{\partial \rho_{u1}} \text{error} &= \frac{\partial}{\partial \rho_{u1}} \left( (r_{ui} - (q_{i1} \rho_{u1} + q_{i2} \rho_{u2}))^2 + \lambda (q_{i1}^2 + q_{i2}^2 + \rho_{u1}^2 + \rho_{u2}^2) \right) \\ &= 2(r_{ui} - (q_{i1} \rho_{u1} + q_{i2} \rho_{u2}))(-q_{i1}) + 2\lambda \rho_{u1} \end{aligned}$$

Denote

$$e_{ui} = r_{ui} - (q_{i1} \rho_{u1} + q_{i2} \rho_{u2})$$

Then

$$\frac{\partial}{\partial \rho_{u1}} \text{error} = -2e_{ui} q_{i1} + 2\lambda \rho_{u1}$$

Analogously for  $\rho_{u2}$ ,  $q_{i1}$ ,  $q_{i2}$

$$\frac{\partial}{\partial \rho_{u2}} \text{error} = -2e_{ui} q_{i2} + 2\lambda \rho_{u2}$$

$$\frac{\partial}{\partial q_{u1}} \text{error} = -2e_{ui} \rho_{i1} + 2\lambda q_{u1}$$

$$\frac{\partial}{\partial q_{u2}} \text{error} = -2e_{ui} \rho_{i2} + 2\lambda q_{u2}$$

Therefore the updating formula for SGD looks as follows:

$$(\rho_{u1}^{(n)}, \rho_{u2}^{(n)}) = (\rho_{u1}^{(n-1)}, \rho_{u2}^{(n-1)}) - \alpha \left( \frac{\partial}{\partial \rho_{u1}} \text{error}, \frac{\partial}{\partial \rho_{u2}} \text{error} \right)$$

$$(\rho_{u1}^{(n)}, \rho_{u2}^{(n)}) = (\rho_{u1}^{(n-1)}, \rho_{u2}^{(n-1)}) - \mathcal{L} \left( \frac{\partial}{\partial \rho_{u1}} \text{error}, \frac{\partial}{\partial \rho_{u2}} \text{error} \right)$$

$$= (\rho_{u1}^{(n-1)}, \rho_{u2}^{(n-1)}) - \mathcal{L} \left( -2 e_{ui}^{(n-1)} q_{i1}^{(n-1)} + 2\lambda \rho_{u1}^{(n-1)}, -2 e_{ui}^{(n-1)} q_{i2}^{(n-1)} + 2\lambda \rho_{u2}^{(n-1)} \right)$$

$$= (\rho_{u1}^{(n-1)}, \rho_{u2}^{(n-1)}) + 2\mathcal{L} \left( e_{ui}^{(n-1)} q_{i1}^{(n-1)} - \lambda \rho_{u1}^{(n-1)}, e_{ui}^{(n-1)} q_{i2}^{(n-1)} - \lambda \rho_{u2}^{(n-1)} \right)$$

Analogously for  $q_{i1}, q_{i2}$

Reverting the above equations in vector form gives the following formulation (renaming  $\mathcal{L}$  to be  $2\mathcal{L}$ )

$$\rho_u^{(n)} = \rho_u^{(n-1)} + \mathcal{L} (e_{ui}^{(n-1)} q_i^{(n-1)} - \lambda \rho_u^{(n-1)})$$

$$q_i^{(n)} = q_i^{(n-1)} + \mathcal{L} (e_{ui}^{(n-1)} \rho_u^{(n-1)} - \lambda q_i^{(n-1)})$$