Gradient Descent

Nipun Batra July 26, 2025

IIT Gandhinagar

Revision

Contour Plot And Gradients

$$z = f(x, y) = x^{2} + y^{2}$$

Surface Plot
$$z = 40$$

$$z = 40$$

$$z = 0$$

$$y = 5$$

$$y = 5$$

Surface Plot
$$z = 0$$

$$z = 0$$

$$z = 0$$

Surface Plot
$$z = 0$$

$$z = 0$$

$$z = 0$$

Contour Plot
$$z = 0$$

$$z = 0$$

Surface Plot
$$z = 0$$

Surface Plot
$$z = 0$$

Surface Plot
Surface

Contour Plot And Gradients



Gradient denotes the direction of steepest ascent or the direction in which there is a maximum increase in f(x,y)

Contour Plot And Gradients



Gradient denotes the direction of steepest ascent or the direction in which there is a maximum increase in f(x,y)

$$\nabla f(x,y) = \begin{bmatrix} \frac{\partial f(x,y)}{\partial x} \\ \frac{\partial f(x,y)}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

• We often want to minimize/maximize a function

- We often want to minimize/maximize a function
- We wanted to minimize the cost function:

$$f(\theta) = (\mathbf{y} - \mathbf{X}\theta)^{\mathsf{T}} (\mathbf{y} - \mathbf{X}\theta)$$
(1)

- We often want to minimize/maximize a function
- We wanted to minimize the cost function:

$$f(\theta) = (\mathbf{y} - \mathbf{X}\theta)^{T} (\mathbf{y} - \mathbf{X}\theta)$$
(1)

• Note, here heta is the parameter vector

• In general, we have following components:

- In general, we have following components:
- Maximize or Minimize a function subject to some constraints

- In general, we have following components:
- Maximize or Minimize a function subject to some constraints
- Today, we will focus on unconstrained optimization (no constraints)

- In general, we have following components:
- Maximize or Minimize a function subject to some constraints
- Today, we will focus on unconstrained optimization (no constraints)
- We will focus on minimization

- In general, we have following components:
- Maximize or Minimize a function subject to some constraints
- Today, we will focus on unconstrained optimization (no constraints)
- We will focus on minimization
- Goal:

$$\theta^* = \underset{\theta}{\arg\min}f(\theta)$$
 (2)

• Gradient descent is an optimization algorithm

- Gradient descent is an optimization algorithm
- It is used to find the minimum of a function in unconstrained settings

- Gradient descent is an optimization algorithm
- It is used to find the minimum of a function in unconstrained settings
- It is an iterative algorithm

- Gradient descent is an optimization algorithm
- It is used to find the minimum of a function in unconstrained settings
- It is an iterative algorithm
- It is a first order optimization algorithm

- Gradient descent is an optimization algorithm
- It is used to find the minimum of a function in unconstrained settings
- It is an iterative algorithm
- It is a first order optimization algorithm
- It is a local search algorithm/greedy

1. Initialize heta to some random value

- 1. Initialize heta to some random value
- 2. Compute the gradient of the cost function at θ , $\nabla f(\theta)$

- 1. Initialize heta to some random value
- 2. Compute the gradient of the cost function at θ , $\nabla f(\theta)$
- 3. For Iteration i (i = 1, 2, ...) or until convergence:

- 1. Initialize heta to some random value
- 2. Compute the gradient of the cost function at θ , $\nabla f(\theta)$
- 3. For Iteration i (i = 1, 2, ...) or until convergence:

•
$$\boldsymbol{\theta}_i \leftarrow \boldsymbol{\theta}_{i-1} - \alpha \nabla f(\boldsymbol{\theta}_{i-1})$$

• Taylor's series is a way to approximate a function f(x) around a point x_0 using a polynomial

- Taylor's series is a way to approximate a function f(x) around a point x_0 using a polynomial
- The polynomial is given by

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots \quad (3)$$

- Taylor's series is a way to approximate a function f(x) around a point x₀ using a polynomial
- The polynomial is given by

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots \quad (3)$$

• The vector form of the above equation is given by:

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^T \nabla^2 f(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0) + \dots$$
(4)

- Taylor's series is a way to approximate a function f(x) around a point x₀ using a polynomial
- The polynomial is given by

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots \quad (3)$$

• The vector form of the above equation is given by:

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^T \nabla^2 f(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0) + \dots$$
(4)

where ∇² f(x₀) is the Hessian matrix and ∇f(x₀) is the gradient vector

• Let us consider $f(x) = \cos(x)$ and $x_0 = 0$

- Let us consider $f(x) = \cos(x)$ and $x_0 = 0$
- Then, we have:

- Let us consider $f(x) = \cos(x)$ and $x_0 = 0$
- Then, we have:
- $f(x_0) = \cos(0) = 1$

- Let us consider $f(x) = \cos(x)$ and $x_0 = 0$
- Then, we have:
- $f(x_0) = \cos(0) = 1$
- $f'(x_0) = -\sin(0) = 0$

- Let us consider $f(x) = \cos(x)$ and $x_0 = 0$
- Then, we have:
- $f(x_0) = \cos(0) = 1$
- $f'(x_0) = -\sin(0) = 0$
- $f''(x_0) = -\cos(0) = -1$

- Let us consider $f(x) = \cos(x)$ and $x_0 = 0$
- Then, we have:
- $f(x_0) = \cos(0) = 1$
- $f'(x_0) = -\sin(0) = 0$
- $f''(x_0) = -\cos(0) = -1$
- We can write the second order Taylor's series as:

- Let us consider $f(x) = \cos(x)$ and $x_0 = 0$
- Then, we have:

•
$$f(x_0) = \cos(0) = 1$$

•
$$f'(x_0) = -\sin(0) = 0$$

•
$$f''(x_0) = -\cos(0) = -1$$

• We can write the second order Taylor's series as:

•
$$f(x) = 1 + 0(x - 0) + \frac{-1}{2!}(x - 0)^2 = 1 - \frac{x^2}{2}$$

• Let us consider another example: $f(x) = x^2 + 2$ and $x_0 = 2$
- Let us consider another example: $f(x) = x^2 + 2$ and $x_0 = 2$
- Question: How does the first order Taylor's series approximation look like?

- Let us consider another example: $f(x) = x^2 + 2$ and $x_0 = 2$
- Question: How does the first order Taylor's series approximation look like?
- First order Taylor's series approximation is given by:

- Let us consider another example: $f(x) = x^2 + 2$ and $x_0 = 2$
- Question: How does the first order Taylor's series approximation look like?
- First order Taylor's series approximation is given by:
- $f(x) = f(x_0) + f'(x_0)(x x_0) = 6 + 4(x 2) = 4x 2$

• We have:

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots \quad (5)$$

• We have:

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots \quad (5)$$

• Let us consider $x = x_0 + \Delta x$ where Δx is a small quantity

• We have:

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots \quad (5)$$

- Let us consider $x = x_0 + \Delta x$ where Δx is a small quantity
- Then, we have:

$$f(x_0 + \Delta x) = f(x_0) + \frac{f'(x_0)}{1!} \Delta x + \frac{f''(x_0)}{2!} \Delta x^2 + \dots \quad (6)$$

• We have:

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots \quad (5)$$

- Let us consider $x = x_0 + \Delta x$ where Δx is a small quantity
- Then, we have:

$$f(x_0 + \Delta x) = f(x_0) + \frac{f'(x_0)}{1!} \Delta x + \frac{f''(x_0)}{2!} \Delta x^2 + \dots$$
 (6)

 Let us assume Δx is small enough such that Δx² and higher order terms can be ignored

• We have:

$$f(x) = f(x_0) + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots \quad (5)$$

- Let us consider $x = x_0 + \Delta x$ where Δx is a small quantity
- Then, we have:

$$f(x_0 + \Delta x) = f(x_0) + \frac{f'(x_0)}{1!} \Delta x + \frac{f''(x_0)}{2!} \Delta x^2 + \dots$$
 (6)

- Let us assume Δx is small enough such that Δx² and higher order terms can be ignored
- Then, we have: $f(x_0 + \Delta x) \approx f(x_0) + \frac{f'(x_0)}{1!} \Delta x$

• Then, we have: $f(x_0 + \Delta x) \approx f(x_0) + \frac{f'(x_0)}{1!} \Delta x$

- Then, we have: $f(x_0 + \Delta x) \approx f(x_0) + \frac{f'(x_0)}{1!} \Delta x$
- Or, in vector form: $f(\mathbf{x}_0 + \Delta \mathbf{x}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$

- Then, we have: $f(x_0 + \Delta x) \approx f(x_0) + \frac{f'(x_0)}{1!} \Delta x$
- Or, in vector form: $f(\mathbf{x}_0 + \Delta \mathbf{x}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$
- Goal: Find $\Delta \mathbf{x}$ such that $f(\mathbf{x}_0 + \Delta \mathbf{x})$ is minimized

- Then, we have: $f(x_0 + \Delta x) \approx f(x_0) + \frac{f'(x_0)}{1!} \Delta x$
- Or, in vector form: $f(\mathbf{x}_0 + \Delta \mathbf{x}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$
- Goal: Find $\Delta \mathbf{x}$ such that $f(\mathbf{x}_0 + \Delta \mathbf{x})$ is minimized
- This is equivalent to minimizing $f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$

- Then, we have: $f(x_0 + \Delta x) \approx f(x_0) + \frac{f'(x_0)}{1!} \Delta x$
- Or, in vector form: $f(\mathbf{x}_0 + \Delta \mathbf{x}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$
- Goal: Find $\Delta \mathbf{x}$ such that $f(\mathbf{x}_0 + \Delta \mathbf{x})$ is minimized
- This is equivalent to minimizing $f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$
- This happens when vectors $abla f(\mathbf{x}_0)$ and $\Delta \mathbf{x}$ are at phase angle of 180°

- Then, we have: $f(x_0 + \Delta x) \approx f(x_0) + \frac{f'(x_0)}{1!} \Delta x$
- Or, in vector form: $f(\mathbf{x}_0 + \Delta \mathbf{x}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$
- Goal: Find $\Delta \mathbf{x}$ such that $f(\mathbf{x}_0 + \Delta \mathbf{x})$ is minimized
- This is equivalent to minimizing $f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$
- This happens when vectors ∇f(x₀) and Δx are at phase angle of 180°
- This happens when $\Delta \mathbf{x} = -\alpha \nabla f(\mathbf{x}_0)$ where α is a scalar

- Then, we have: $f(x_0 + \Delta x) \approx f(x_0) + \frac{f'(x_0)}{1!} \Delta x$
- Or, in vector form: $f(\mathbf{x}_0 + \Delta \mathbf{x}) \approx f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$
- Goal: Find $\Delta \mathbf{x}$ such that $f(\mathbf{x}_0 + \Delta \mathbf{x})$ is minimized
- This is equivalent to minimizing $f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \Delta \mathbf{x}$
- This happens when vectors ∇f(x₀) and Δx are at phase angle of 180°
- This happens when $\Delta \mathbf{x} = -\alpha \nabla f(\mathbf{x}_0)$ where α is a scalar
- This is the gradient descent algorithm: $\mathbf{x}_1 = \mathbf{x}_0 \alpha \nabla f(\mathbf{x}_0)$







High learning rate $\alpha = 0.8$: Converges quickly, but might overshoot





Appropriate learning rate $\alpha=0.1$



Gradient Descent for linear regression

• Loss function is usually a function defined on a data point, prediction and label, and measures the penalty.

- Loss function is usually a function defined on a data point, prediction and label, and measures the penalty.
- square loss $I(f(x_i; \theta), y_i) = (f(x_i; \theta) y_i)^2$, used in linear regression

- Loss function is usually a function defined on a data point, prediction and label, and measures the penalty.
- square loss $I(f(x_i; \theta), y_i) = (f(x_i; \theta) y_i)^2$, used in linear regression
- **Cost function** is usually more general. It might be a sum of loss functions over your training set plus some model complexity penalty (regularization). For example:

- Loss function is usually a function defined on a data point, prediction and label, and measures the penalty.
- square loss $I(f(x_i; \theta), y_i) = (f(x_i; \theta) y_i)^2$, used in linear regression
- **Cost function** is usually more general. It might be a sum of loss functions over your training set plus some model complexity penalty (regularization). For example:
- Mean Squared Error $MSE(\theta) = \frac{1}{n} \sum_{i=1}^{n} (f(x_i; \theta) y_i)^2$

- Loss function is usually a function defined on a data point, prediction and label, and measures the penalty.
- square loss $I(f(x_i; \theta), y_i) = (f(x_i; \theta) y_i)^2$, used in linear regression
- **Cost function** is usually more general. It might be a sum of loss functions over your training set plus some model complexity penalty (regularization). For example:
- Mean Squared Error $MSE(\theta) = \frac{1}{n} \sum_{i=1}^{n} (f(x_i; \theta) y_i)^2$
- **Objective function** is the most general term for any function that you optimize during training.

Learn $y = \theta_0 + \theta_1 x$ on following dataset, using gradient descent where initially $(\theta_0, \theta_1) = (4, 0)$ and step-size, $\alpha = 0.1$, for 2 iterations.

x	у
1	1
2	2
3	3

Our predictor, $\hat{y} = \theta_0 + \theta_1 x$

Error for *i*th datapoint,
$$\epsilon_i = y_i - \hat{y}_i$$

 $\epsilon_1 = 1 - \theta_0 - \theta_1$
 $\epsilon_2 = 2 - \theta_0 - 2\theta_1$
 $\epsilon_3 = 3 - \theta_0 - 3\theta_1$

$$\mathsf{MSE} = \frac{\epsilon_1^2 + \epsilon_2^2 + \epsilon_3^2}{3} = \frac{14 + 3\theta_0^2 + 14\theta_1^2 - 12\theta_0 - 28\theta_1 + 12\theta_0\theta_1}{3}$$

$\sum \epsilon_i^2$ increases as the number of examples increase So, we use MSE

$$MSE = \frac{1}{n} \sum \epsilon_i^2$$

Here n denotes the number of samples

Gradient Descent : Example

$$\frac{\partial \mathsf{MSE}}{\partial \theta_0} = \frac{2\sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)(-1)}{n} = \frac{2\sum_{i=1}^n \epsilon_i(-1)}{n}$$

$$\frac{\partial \mathsf{MSE}}{\partial \theta_1} = \frac{2\sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)(-x_i)}{n} = \frac{2\sum_{i=1}^n \epsilon_i(-x_i)}{n}$$

Iteration 1

$$\theta_{0} = \theta_{0} - \alpha \frac{\partial MSE}{\partial \theta_{0}}$$
$$\theta_{1} = \theta_{1} - \alpha \frac{\partial MSE}{\partial \theta_{1}}$$

Iteration 1

$$\theta_0 = \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0}$$

$$\theta_0 = 4 - 0.2 \frac{((1 - (4 + 0))(-1) + (2 - (4 + 0))(-1) + (3 - (4 + 0))(-1))}{3}$$

$$\theta_0 = 3.6$$

$$\theta_1 = \theta_1 - \alpha \frac{\partial \textit{MSE}}{\partial \theta_1}$$

Gradient Descent : Example

Iteration 1

 $\theta_1 = -0.67$

$$\begin{aligned} \theta_0 &= \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0} \\ \theta_0 &= 4 - 0.2 \frac{((1 - (4 + 0))(-1) + (2 - (4 + 0))(-1) + (3 - (4 + 0))(-1))}{3} \\ \theta_0 &= 3.6 \\ \theta_1 &= \theta_1 - \alpha \frac{\partial MSE}{\partial \theta_1} \\ \theta_1 &= 0 - 0.2 \frac{((1 - (4 + 0))(-1) + (2 - (4 + 0))(-2) + (3 - (4 + 0))(-3))}{3} \end{aligned}$$

Iteration 2

$$\theta_0 = \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0}$$

$$\theta_1 = \theta_1 - \alpha \frac{\partial MSE}{\partial \theta_1}$$

Iteration 2

$$\theta_0 = \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0}$$

$$\theta_0 = \\ 3.6 - 0.2 \frac{((1 - (3.6 - 0.67))(-1) + (2 - (3.6 - 0.67 \times 2))(-1) + (3 - (3.6 - 0.67 \times 3))(-1))}{3}$$

 $\theta_0 = 3.54$

 $\theta_1 = \theta_1 - \alpha \frac{\partial \textit{MSE}}{\partial \theta_1}$

Gradient Descent : Example

Iteration 2

$$\theta_0 = \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0}$$

$$\theta_0 = 3.6 - 0.2 \frac{((1 - (3.6 - 0.67))(-1) + (2 - (3.6 - 0.67 \times 2))(-1) + (3 - (3.6 - 0.67 \times 3))(-1))}{3}$$

$$\theta_0 = 3.54$$

$$\theta_1 = \theta_1 - \alpha \frac{\partial MSE}{\partial \theta_1}$$

$$\theta_0 = 3.6 - 0.2 \frac{((1 - (3.6 - 0.67))(-1) + (2 - (3.6 - 0.67 \times 2))(-2) + (3 - (3.6 - 0.67 \times 3))(-3))}{3}$$

 $\theta_{0} = -0.55$

Gradient Descent : Example (Iteraion 0)


Gradient Descent : Example (Iteraion 2)



Gradient Descent : Example (Iteraion 4)



Gradient Descent : Example (Iteraion 6)



Gradient Descent : Example (Iteraion 8)



Gradient Descent : Example (Iteraion 10)



Gradient Descent : Example (Iteraion 12)



Gradient Descent : Example (Iteraion 14)



Gradient Descent : Example (Iteraion 16)



Gradient Descent : Example (Iteraion 18)



Gradient Descent : Example (Iteraion 20)



Gradient Descent : Example (Iteraion 22)



Gradient Descent : Example (Iteraion 24)



Gradient Descent : Example (Iteraion 26)



Gradient Descent : Example (Iteraion 28)



Gradient Descent : Example (Iteraion 30)



Gradient Descent : Example (Iteraion 32)



Gradient Descent : Example (Iteraion 34)



Gradient Descent : Example (Iteraion 36)



Gradient Descent : Example (Iteraion 38)



Gradient Descent : Example (Iteraion 40)



Iteration vs Epochs for gradient descent

• Iteration: Each time you update the parameters of the model

- Iteration: Each time you update the parameters of the model
- Epoch: Each time you have seen all the set of examples

• Dataset:
$$\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$$
 of size n

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Predict $\hat{\mathbf{y}} = \text{pred}(\mathbf{X}, \boldsymbol{\theta})$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Predict $\hat{\mathbf{y}} = \text{pred}(\mathbf{X}, \boldsymbol{\theta})$
 - Compute loss: $J(\theta) = loss(\mathbf{y}, \hat{\mathbf{y}})$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Predict $\hat{\mathbf{y}} = \operatorname{pred}(\mathbf{X}, \boldsymbol{\theta})$
 - Compute loss: $J(\theta) = loss(\mathbf{y}, \hat{\mathbf{y}})$
 - Compute gradient: $\nabla J(\theta) = \operatorname{grad}(J)(\theta)$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Predict $\hat{\mathbf{y}} = \operatorname{pred}(\mathbf{X}, \boldsymbol{\theta})$
 - Compute loss: $J(\theta) = loss(\mathbf{y}, \hat{\mathbf{y}})$
 - Compute gradient: $\nabla J(\theta) = \operatorname{grad}(J)(\theta)$
 - Update: $\theta = \theta \alpha \nabla J(\theta)$

• Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - $\bullet \ \ \mathsf{Shuffle} \ \mathcal{D}$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Shuffle ${\cal D}$
 - For *i* in [1, *n*]

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - $\bullet \ \ \mathsf{Shuffle} \ \mathcal{D}$
 - For *i* in [1, *n*]
 - Predict $\hat{\mathbf{y}}_i = \text{pred}(\mathbf{x}_i, \boldsymbol{\theta})$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - $\bullet \ \ \mathsf{Shuffle} \ \mathcal{D}$
 - For *i* in [1, *n*]
 - Predict $\hat{\mathbf{y}}_i = \text{pred}(\mathbf{x}_i, \boldsymbol{\theta})$
 - Compute loss: $J(\boldsymbol{\theta}) = loss(y_i, \hat{\mathbf{y}}_i)$
Stochastic Gradient Descent (SGD)

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - $\bullet \ \ \mathsf{Shuffle} \ \mathcal{D}$
 - For *i* in [1, *n*]
 - Predict $\hat{\mathbf{y}}_i = \text{pred}(\mathbf{x}_i, \boldsymbol{\theta})$
 - Compute loss: $J(\boldsymbol{\theta}) = loss(y_i, \hat{\mathbf{y}}_i)$
 - Compute gradient: $\nabla J(\theta) = \operatorname{grad}(J)(\theta)$

Stochastic Gradient Descent (SGD)

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - $\bullet \ \ \mathsf{Shuffle} \ \mathcal{D}$
 - For *i* in [1, *n*]
 - Predict $\hat{\mathbf{y}}_i = \text{pred}(\mathbf{x}_i, \boldsymbol{\theta})$
 - Compute loss: $J(\boldsymbol{\theta}) = loss(y_i, \hat{\mathbf{y}}_i)$
 - Compute gradient: $\nabla J(\theta) = \operatorname{grad}(J)(\theta)$
 - Update: $\theta = \theta \alpha \nabla J(\theta)$

• Dataset:
$$\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$$
 of size *n*

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - $\bullet \ \ Shuffle \ \mathcal{D}$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - $\bullet \ \ Shuffle \ \mathcal{D}$
 - Batches = make_batches(D, B)

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Shuffle \mathcal{D}
 - Batches = make_batches(D, B)
 - For b in Batches

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Shuffle \mathcal{D}
 - Batches = make_batches(D, B)
 - For b in Batches
 - $\mathbf{X}_b, \mathbf{y}_b = b$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Shuffle \mathcal{D}
 - Batches = make_batches(D, B)
 - For b in Batches
 - $\mathbf{X}_b, \mathbf{y}_b = b$
 - Predict $\hat{\mathbf{y}}_b = \operatorname{pred}(\mathbf{X}_b, \boldsymbol{\theta})$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Shuffle \mathcal{D}
 - Batches = make_batches(D, B)
 - For b in Batches
 - $\mathbf{X}_b, \mathbf{y}_b = b$
 - Predict $\hat{\mathbf{y}}_b = \text{pred}(\mathbf{X}_b, \boldsymbol{\theta})$
 - Compute loss: $J(\theta) = loss(\mathbf{y}_b, \hat{\mathbf{y}}_b)$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Shuffle \mathcal{D}
 - Batches = make_batches(D, B)
 - For b in Batches
 - $\mathbf{X}_b, \mathbf{y}_b = b$
 - Predict $\hat{\mathbf{y}}_b = \operatorname{pred}(\mathbf{X}_b, \boldsymbol{\theta})$
 - Compute loss: $J(\theta) = loss(\mathbf{y}_b, \hat{\mathbf{y}}_b)$
 - Compute gradient: $\nabla J(\theta) = \operatorname{grad}(J)(\theta)$

- Dataset: $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$ of size n
- Initialize θ
- For epoch e in [1, E]
 - Shuffle \mathcal{D}
 - Batches = make_batches(D, B)
 - For b in Batches
 - $\mathbf{X}_b, \mathbf{y}_b = b$
 - Predict $\hat{\mathbf{y}}_b = \text{pred}(\mathbf{X}_b, \boldsymbol{\theta})$
 - Compute loss: $J(\theta) = loss(\mathbf{y}_b, \hat{\mathbf{y}}_b)$
 - Compute gradient: $\nabla J(\theta) = \operatorname{grad}(J)(\theta)$
 - Update: $\theta = \theta \alpha \nabla J(\theta)$

Vanilla Gradient Descent

• in Vanilla (Batch) gradient descent: We update params after going through all the data

Vanilla Gradient Descent

- in Vanilla (Batch) gradient descent: We update params after going through all the data
- Smooth curve for Iteration vs Cost

Vanilla Gradient Descent

- in Vanilla (Batch) gradient descent: We update params after going through all the data
- Smooth curve for Iteration vs Cost
- For a single update, it needs to compute the gradient over all the samples. Hence takes more time

Vanilla Gradient Descent

- in Vanilla (Batch) gradient descent: We update params after going through all the data
- Smooth curve for Iteration vs Cost
- For a single update, it needs to compute the gradient over all the samples. Hence takes more time

Vanilla Gradient Descent

- in Vanilla (Batch) gradient descent: We update params after going through all the data
- Smooth curve for Iteration vs Cost
- For a single update, it needs to compute the gradient over all the samples. Hence takes more time

Stochastic Gradient Descent

• In SGD, we update parameters after seeing each each point

Vanilla Gradient Descent

- in Vanilla (Batch) gradient descent: We update params after going through all the data
- Smooth curve for Iteration vs Cost
- For a single update, it needs to compute the gradient over all the samples. Hence takes more time

Stochastic Gradient Descent

- In SGD, we update parameters after seeing each each point
- Noisier curve for iteration vs cost

Vanilla Gradient Descent

- in Vanilla (Batch) gradient descent: We update params after going through all the data
- Smooth curve for Iteration vs Cost
- For a single update, it needs to compute the gradient over all the samples. Hence takes more time

Stochastic Gradient Descent

- In SGD, we update parameters after seeing each each point
- Noisier curve for iteration vs cost
- For a single update, it computes the gradient over one example. Hence lesser time

Learn $y = \theta_0 + \theta_1 x$ on following dataset, using SGD where initially $(\theta_0, \theta_1) = (4, 0)$ and step-size, $\alpha = 0.1$, for 1 epoch (3 iterations).

x	у
2	2
3	3
1	1

Stochastic Gradient Descent : Example

Our predictor, $\hat{y} = \theta_0 + \theta_1 x$

Error for *i*th datapoint,
$$e_i = y_i - \hat{y}_i$$

 $\epsilon_1 = 2 - \theta_0 - 2\theta_1$
 $\epsilon_2 = 3 - \theta_0 - 3\theta_1$
 $\epsilon_3 = 1 - \theta_0 - \theta_1$

While using SGD, we compute the MSE using only 1 datapoint per iteration.

So MSE is ϵ_1^2 for iteration 1 and ϵ_2^2 for iteration 2.

Stochastic Gradient Descent : Example

Contour plot of the cost functions for the three datapoints



For Iteration *i*

.

$$\frac{\partial MSE}{\partial \theta_0} = 2\left(y_i - \theta_0 - \theta_1 x_i\right)(-1) = 2\epsilon_i (-1)$$

$$\frac{\partial MSE}{\partial \theta_1} = 2\left(y_i - \theta_0 - \theta_1 x_i\right)\left(-x_i\right) = 2\epsilon_i\left(-x_i\right)$$

Iteration 1

$$\theta_{0} = \theta_{0} - \alpha \frac{\partial MSE}{\partial \theta_{0}}$$
$$\theta_{1} = \theta_{1} - \alpha \frac{\partial MSE}{\partial \theta_{1}}$$

Iteration 1

$$\theta_0 = \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0}$$
$$\theta_0 = 4 - 0.1 \times 2 \times (2 - (4 + 0)) (-1)$$
$$\theta_0 = 3.6$$

)

$$\theta_1 = \theta_1 - \alpha \frac{\partial MSE}{\partial \theta_1}$$

Stochastic Gradient Descent : Example

Iteration 1

$$\begin{aligned} \theta_0 &= \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0} \\ \theta_0 &= 4 - 0.1 \times 2 \times (2 - (4 + 0)) (-1) \\ \theta_0 &= 3.6 \\ \theta_1 &= \theta_1 - \alpha \frac{\partial MSE}{\partial \theta_1} \\ \theta_1 &= 0 - 0.1 \times 2 \times (2 - (4 + 0)) (-2) \\ \theta_1 &= -0.8 \end{aligned}$$

Iteration 2

$$\theta_{0} = \theta_{0} - \alpha \frac{\partial MSE}{\partial \theta_{0}}$$
$$\theta_{1} = \theta_{1} - \alpha \frac{\partial MSE}{\partial \theta_{1}}$$

Iteration 2

$$\theta_0 = \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0}$$

$$\theta_0 = 3.6 - 0.1 \times 2 \times (3 - (3.6 - 0.8 \times 3))(-1)$$

 $\theta_0 = 3.96$

$$\theta_1 = \theta_1 - \alpha \frac{\partial MSE}{\partial \theta_1}$$

Stochastic Gradient Descent : Example

Iteration 2

$$\begin{aligned} \theta_{0} &= \theta_{0} - \alpha \frac{\partial MSE}{\partial \theta_{0}} \\ \theta_{0} &= 3.6 - 0.1 \times 2 \times (3 - (3.6 - 0.8 \times 3)) (-1) \\ \theta_{0} &= 3.96 \\ \theta_{1} &= \theta_{1} - \alpha \frac{\partial MSE}{\partial \theta_{1}} \\ \theta_{0} &= -0.8 - 0.1 \times 2 \times (3 - (3.6 - 0.8 \times 3)) (-3) \\ \theta_{1} &= 0.28 \end{aligned}$$

Iteration 3

$$\theta_{0} = \theta_{0} - \alpha \frac{\partial MSE}{\partial \theta_{0}}$$
$$\theta_{1} = \theta_{1} - \alpha \frac{\partial MSE}{\partial \theta_{1}}$$

Iteration 3

$$\theta_0 = \theta_0 - \alpha \frac{\partial MSE}{\partial \theta_0}$$

$$\theta_0 = 3.96 - 0.1 \times 2 \times (1 - (3.96 + 0.28 \times 1))(-1)$$

 $\theta_0 = 3.312$

$$\theta_1 = \theta_1 - \alpha \frac{\partial MSE}{\partial \theta_1}$$

Stochastic Gradient Descent : Example

Iteration 3

$$\begin{aligned} \theta_{0} &= \theta_{0} - \alpha \frac{\partial MSE}{\partial \theta_{0}} \\ \theta_{0} &= 3.96 - 0.1 \times 2 \times (1 - (3.96 + 0.28 \times 1)) (-1) \\ \theta_{0} &= 3.312 \\ \theta_{1} &= \theta_{1} - \alpha \frac{\partial MSE}{\partial \theta_{1}} \\ \theta_{0} &= 0.28 - 0.1 \times 2 \times (1 - (3.96 + 0.28 \times 1)) (-1) \\ \theta_{1} &= -0.368 \end{aligned}$$

)

)

Stochastic gradient is an unbiased estimator of the true gradient

True Gradient

Based on Estimation Theory and Machine Learning by Florian Hartmann

• Let us say we have a dataset \mathcal{D} containing input output pairs $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

True Gradient

Based on Estimation Theory and Machine Learning by Florian Hartmann

- Let us say we have a dataset \mathcal{D} containing input output pairs $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$
- We can define overall loss as:

$$L(\theta) = \frac{1}{N} \sum_{i=1}^{N} loss(f(x_i, \theta), y_i)$$
True Gradient

Based on Estimation Theory and Machine Learning by Florian Hartmann

- Let us say we have a dataset \mathcal{D} containing input output pairs $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$
- We can define overall loss as:

$$L(\theta) = \frac{1}{N} \sum_{i=1}^{N} loss(f(x_i, \theta), y_i)$$

• loss can be any loss function such as squared loss, cross-entropy loss etc.

$$loss(f(x_i, \theta), y_i) = (f(x_i, \theta) - y_i)^2$$

• The true gradient of the loss function is given by:

$$\nabla L = \nabla \frac{1}{n} \sum_{i=1}^{n} loss(f(x_i), y_i)$$
$$= \frac{1}{n} \sum_{i=1}^{n} \nabla loss(f(x_i), y_i)$$

• The true gradient of the loss function is given by:

$$\nabla L = \nabla \frac{1}{n} \sum_{i=1}^{n} loss(f(x_i), y_i)$$
$$= \frac{1}{n} \sum_{i=1}^{n} \nabla loss(f(x_i), y_i)$$

• The above is a consequence of linearity of the gradient operator.

• In practice, we do not have access to the true gradient

- In practice, we do not have access to the true gradient
- We can only estimate the true gradient using a subset of the data

- In practice, we do not have access to the true gradient
- We can only estimate the true gradient using a subset of the data
- For SGD, we use a single example to estimate the true gradient, for mini-batch gradient descent, we use a mini-batch of examples to estimate the true gradient

- In practice, we do not have access to the true gradient
- We can only estimate the true gradient using a subset of the data
- For SGD, we use a single example to estimate the true gradient, for mini-batch gradient descent, we use a mini-batch of examples to estimate the true gradient
- Let us say we have a sample: (x, y)

- In practice, we do not have access to the true gradient
- We can only estimate the true gradient using a subset of the data
- For SGD, we use a single example to estimate the true gradient, for mini-batch gradient descent, we use a mini-batch of examples to estimate the true gradient
- Let us say we have a sample: (x, y)
- The estimated gradient is given by:

$$\nabla \tilde{L} = \nabla \operatorname{loss}(f(x), y)$$

Bias of the estimator

 One measure for the quality of an estimator X
 is its bias or how far off its estimate is on average from the true value X :

$$\mathsf{bias}(X) = \mathbb{E}[\tilde{X}] - X$$

Bias of the estimator

 One measure for the quality of an estimator X is its bias or how far off its estimate is on average from the true value X :

$$\mathsf{bias}(X) = \mathbb{E}[\tilde{X}] - X$$

• Using the rules of expectation, we can show that the expected value of the estimated gradient is the true gradient:

$$\mathbb{E}[\nabla \tilde{L}] = \sum_{i=1}^{n} \frac{1}{n} \nabla \log (f(x_i), y_i)$$
$$= \frac{1}{n} \nabla \sum_{i=1}^{n} \log (f(x_i), y_i)$$
$$= \nabla L$$

Bias of the estimator

 One measure for the quality of an estimator X is its bias or how far off its estimate is on average from the true value X :

$$\mathsf{bias}(X) = \mathbb{E}[\tilde{X}] - X$$

• Using the rules of expectation, we can show that the expected value of the estimated gradient is the true gradient:

$$\mathbb{E}[\nabla \tilde{L}] = \sum_{i=1}^{n} \frac{1}{n} \nabla \log (f(x_i), y_i)$$
$$= \frac{1}{n} \nabla \sum_{i=1}^{n} \log (f(x_i), y_i)$$
$$= \nabla L$$

• Thus, the estimated gradient is an unbiased estimator of the true gradient

Time Complexity: Gradient Descent vs Normal Equation for Linear Regression

• Consider $\mathbf{X} \in \mathbb{R}^{n \times d}$

- Consider $\mathbf{X} \in \mathbb{R}^{n \times d}$
- *n* examples and *d* dimensions

- Consider $\mathbf{X} \in \mathbb{R}^{n \times d}$
- *n* examples and *d* dimensions
- What is the time complexity of solving the normal equation $\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$?

• **X** has dimensions $n \times d$, **X**^T has dimensions $d \times n$

- **X** has dimensions $n \times d$, **X**^T has dimensions $d \times n$
- X^TX is a matrix product of matrices of size: d × n and n × d, which is O(d²n)

- **X** has dimensions $n \times d$, **X**^T has dimensions $d \times n$
- X^TX is a matrix product of matrices of size: d × n and n × d, which is O(d²n)
- Inversion of $\mathbf{X}^T \mathbf{X}$ is an inversion of a $d \times d$ matrix, which is $\mathcal{O}(d^3)$

- **X** has dimensions $n \times d$, **X**^T has dimensions $d \times n$
- X^TX is a matrix product of matrices of size: d × n and n × d, which is O(d²n)
- Inversion of $\mathbf{X}^T \mathbf{X}$ is an inversion of a $d \times d$ matrix, which is $\mathcal{O}(d^3)$
- **X**^T**y** is a matrix vector product of size $d \times n$ and $n \times 1$, which is $\mathcal{O}(dn)$

- **X** has dimensions $n \times d$, **X**^T has dimensions $d \times n$
- X^TX is a matrix product of matrices of size: d × n and n × d, which is O(d²n)
- Inversion of $\mathbf{X}^T \mathbf{X}$ is an inversion of a $d \times d$ matrix, which is $\mathcal{O}(d^3)$
- **X**^T**y** is a matrix vector product of size $d \times n$ and $n \times 1$, which is $\mathcal{O}(dn)$
- (X^TX)⁻¹X^Ty is a matrix product of a d × d matrix and d × 1 matrix, which is O(d²)

- **X** has dimensions $n \times d$, **X**^T has dimensions $d \times n$
- X^TX is a matrix product of matrices of size: d × n and n × d, which is O(d²n)
- Inversion of $\mathbf{X}^T \mathbf{X}$ is an inversion of a $d \times d$ matrix, which is $\mathcal{O}(d^3)$
- **X**^T**y** is a matrix vector product of size $d \times n$ and $n \times 1$, which is $\mathcal{O}(dn)$
- (X^TX)⁻¹X^Ty is a matrix product of a d × d matrix and d × 1 matrix, which is O(d²)
- Overall complexity: $\mathcal{O}(d^2n) + \mathcal{O}(d^3) + \mathcal{O}(dn) + \mathcal{O}(d^2) = \mathcal{O}(d^2n) + \mathcal{O}(d^3)$

- **X** has dimensions $n \times d$, **X**^T has dimensions $d \times n$
- X^TX is a matrix product of matrices of size: d × n and n × d, which is O(d²n)
- Inversion of $\mathbf{X}^T \mathbf{X}$ is an inversion of a $d \times d$ matrix, which is $\mathcal{O}(d^3)$
- **X**^T**y** is a matrix vector product of size $d \times n$ and $n \times 1$, which is $\mathcal{O}(dn)$
- (X^TX)⁻¹X^Ty is a matrix product of a d × d matrix and d × 1 matrix, which is O(d²)
- Overall complexity: $\mathcal{O}(d^2n) + \mathcal{O}(d^3) + \mathcal{O}(dn) + \mathcal{O}(d^2) = \mathcal{O}(d^2n) + \mathcal{O}(d^3)$
- Scales cubic in the number of columns/features of ${\bf X}$

Start with random values of θ_0 and θ_1 Till convergence

•
$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (\sum \epsilon_i^2)$$

Start with random values of θ_0 and θ_1 Till convergence

- $\theta_0 = \theta_0 \alpha \frac{\partial}{\partial \theta_0} (\sum \epsilon_i^2)$
- $\theta_1 = \theta_1 \alpha \frac{\partial}{\partial \theta_1} (\sum \epsilon_i^2)$

Start with random values of θ_0 and θ_1 Till convergence

•
$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (\sum \epsilon_i^2)$$

- $\theta_1 = \theta_1 \alpha \frac{\partial}{\partial \theta_1} (\sum \epsilon_i^2)$
- Question: Can you write the above for *d* dimensional data in vectorised form?

Start with random values of θ_0 and θ_1 Till convergence

•
$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (\sum \epsilon_i^2)$$

•
$$\theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} (\sum \epsilon_i^2)$$

• Question: Can you write the above for *d* dimensional data in vectorised form?

•
$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})$$

 $\theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})$
 \vdots
 $\theta_d = \theta_d - \alpha \frac{\partial}{\partial \theta_d} (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})$

Start with random values of θ_0 and θ_1 Till convergence

•
$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (\sum \epsilon_i^2)$$

•
$$\theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} (\sum \epsilon_i^2)$$

• Question: Can you write the above for *d* dimensional data in vectorised form?

•
$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})$$

 $\theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})$
:
 $\theta_d = \theta_d - \alpha \frac{\partial}{\partial \theta_d} (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})$
• $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \frac{\partial}{\partial \theta} (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})$

$$\begin{split} & \frac{\partial}{\partial \theta} (\mathbf{y} - \mathbf{X} \theta)^\top (\mathbf{y} - \mathbf{X} \theta) \\ &= \frac{\partial}{\partial \theta} \left(\mathbf{y}^\top - \theta^\top \mathbf{X}^\top \right) (\mathbf{y} - \mathbf{X} \theta) \\ &= \frac{\partial}{\partial \theta} \left(\mathbf{y}^\top \mathbf{y} - \theta^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \theta + \theta^\top \mathbf{X}^\top \mathbf{X} \theta \right) \\ &= -2 \mathbf{X}^\top \mathbf{y} + 2 \mathbf{X}^\top \mathbf{X} \theta \\ &= 2 \mathbf{X}^\top (\mathbf{X} \theta - \mathbf{y}) \end{split}$$

We can write the vectorised update equation as follows, for each iteration $% \left({{{\mathbf{r}}_{i}}} \right)$

 $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

We can write the vectorised update equation as follows, for each iteration

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$

For t iterations, what is the computational complexity of our gradient descent solution?

We can write the vectorised update equation as follows, for each iteration

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$

For t iterations, what is the computational complexity of our gradient descent solution?

Hint, rewrite the above as: $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} + \alpha \mathbf{X}^{\top} \mathbf{y}$

We can write the vectorised update equation as follows, for each iteration

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$

For t iterations, what is the computational complexity of our gradient descent solution?

Hint, rewrite the above as: $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} + \alpha \mathbf{X}^{\top} \mathbf{y}$

Complexity of computing $\mathbf{X}^{\top}\mathbf{y}$ is $\mathcal{O}(dn)$

We can write the vectorised update equation as follows, for each iteration

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$

For *t* iterations, what is the computational complexity of our gradient descent solution?

Hint, rewrite the above as: $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} + \alpha \mathbf{X}^{\top} \mathbf{y}$

Complexity of computing $\mathbf{X}^{\top}\mathbf{y}$ is $\mathcal{O}(dn)$

Complexity of computing $\alpha \mathbf{X}^{\top} \mathbf{y}$ once we have $\mathbf{X}^{\top} \mathbf{y}$ is $\mathcal{O}(d)$ since $\mathbf{X}^{\top} \mathbf{y}$ has d entries

We can write the vectorised update equation as follows, for each iteration

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$

For *t* iterations, what is the computational complexity of our gradient descent solution?

Hint, rewrite the above as: $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} + \alpha \mathbf{X}^{\top} \mathbf{y}$

Complexity of computing $\mathbf{X}^{\top}\mathbf{y}$ is $\mathcal{O}(dn)$

Complexity of computing $\alpha \mathbf{X}^{\top} \mathbf{y}$ once we have $\mathbf{X}^{\top} \mathbf{y}$ is $\mathcal{O}(d)$ since $\mathbf{X}^{\top} \mathbf{y}$ has d entries

Complexity of computing $\mathbf{X}^{\top}\mathbf{X}$ is $\mathcal{O}(d^2n)$ and then multiplying with α is $\mathcal{O}(d^2)$

We can write the vectorised update equation as follows, for each iteration

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$$

For *t* iterations, what is the computational complexity of our gradient descent solution?

Hint, rewrite the above as: $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} + \alpha \mathbf{X}^{\top} \mathbf{y}$

Complexity of computing $\mathbf{X}^{\top}\mathbf{y}$ is $\mathcal{O}(dn)$

Complexity of computing $\alpha \mathbf{X}^{\top} \mathbf{y}$ once we have $\mathbf{X}^{\top} \mathbf{y}$ is $\mathcal{O}(d)$ since $\mathbf{X}^{\top} \mathbf{y}$ has d entries

Complexity of computing $\mathbf{X}^{\top}\mathbf{X}$ is $\mathcal{O}(d^2n)$ and then multiplying with α is $\mathcal{O}(d^2)$

All of the above need only be calculated once!
The remaining subtraction/addition can be done in $\mathcal{O}(d)$ for each iteration.

The remaining subtraction/addition can be done in $\mathcal{O}(d)$ for each iteration.

What is overall computational complexity?

The remaining subtraction/addition can be done in $\mathcal{O}(d)$ for each iteration.

What is overall computational complexity?

 $\mathcal{O}(td^2) + \mathcal{O}(d^2n) = \mathcal{O}((t+n)d^2)$

If we do not rewrite the expression $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

For each iteration, we have:

• Computing $\mathbf{X}\boldsymbol{\theta}$ is $\mathcal{O}(nd)$

If we do not rewrite the expression $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

- Computing $\mathbf{X}\boldsymbol{\theta}$ is $\mathcal{O}(nd)$
- Computing $\mathbf{X}\boldsymbol{\theta} \mathbf{y}$ is $\mathcal{O}(n)$

If we do not rewrite the expression $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

- Computing $\mathbf{X}\boldsymbol{\theta}$ is $\mathcal{O}(nd)$
- Computing $\mathbf{X}\boldsymbol{\theta} \mathbf{y}$ is $\mathcal{O}(n)$
- Computing $\alpha \mathbf{X}^{\top}$ is $\mathcal{O}(nd)$

If we do not rewrite the expression $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

- Computing $\mathbf{X}\boldsymbol{\theta}$ is $\mathcal{O}(nd)$
- Computing $\mathbf{X}\boldsymbol{\theta} \mathbf{y}$ is $\mathcal{O}(n)$
- Computing $\alpha \mathbf{X}^{\top}$ is $\mathcal{O}(nd)$
- Computing $\alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(nd)$

If we do not rewrite the expression $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

- Computing $\mathbf{X}\boldsymbol{\theta}$ is $\mathcal{O}(nd)$
- Computing $\mathbf{X}\boldsymbol{\theta} \mathbf{y}$ is $\mathcal{O}(n)$
- Computing $\alpha \mathbf{X}^{\top}$ is $\mathcal{O}(nd)$
- Computing $\alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(nd)$
- Computing $\boldsymbol{\theta} = \boldsymbol{\theta} \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(n)$

If we do not rewrite the expression $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

- Computing $\mathbf{X}\boldsymbol{\theta}$ is $\mathcal{O}(nd)$
- Computing $\mathbf{X}\boldsymbol{\theta} \mathbf{y}$ is $\mathcal{O}(n)$
- Computing $\alpha \mathbf{X}^{\top}$ is $\mathcal{O}(nd)$
- Computing $\alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(nd)$
- Computing $\boldsymbol{\theta} = \boldsymbol{\theta} \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(n)$

If we do not rewrite the expression $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

For each iteration, we have:

- Computing $\mathbf{X}\boldsymbol{\theta}$ is $\mathcal{O}(nd)$
- Computing $\mathbf{X}\boldsymbol{\theta} \mathbf{y}$ is $\mathcal{O}(n)$
- Computing $\alpha \mathbf{X}^{\top}$ is $\mathcal{O}(nd)$
- Computing $\alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(nd)$
- Computing $\boldsymbol{\theta} = \boldsymbol{\theta} \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(n)$

What is overall computational complexity?

If we do not rewrite the expression $\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} - \mathbf{y})$

For each iteration, we have:

- Computing $\mathbf{X}\boldsymbol{\theta}$ is $\mathcal{O}(nd)$
- Computing $\mathbf{X}\boldsymbol{\theta} \mathbf{y}$ is $\mathcal{O}(n)$
- Computing $\alpha \mathbf{X}^{\top}$ is $\mathcal{O}(nd)$
- Computing $\alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(nd)$
- Computing $\boldsymbol{\theta} = \boldsymbol{\theta} \alpha \mathbf{X}^{\top} (\mathbf{X} \boldsymbol{\theta} \mathbf{y})$ is $\mathcal{O}(n)$

What is overall computational complexity?

 $\mathcal{O}(ndt)$