# Variants of Gradient Descent

# What we will be covering in this module?

- Variants of Gradient Descent

# What we will be covering in this module?

- Variants of Gradient Descent

- Challenges / Problems with Gradient Descent

Analytics Vidhya
Learn everything about analytics

# What we will be covering in this module?

- Variants of Gradient Descent

- Challenges / Problems with Gradient Descent

- Different types of optimizers

# What we will be covering in this module?

- Variants of Gradient Descent

- Challenges / Problems with Gradient Descent

- Different types of optimizers

- Implementing optimizers from scratch

# Gradient Descent Recap

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$
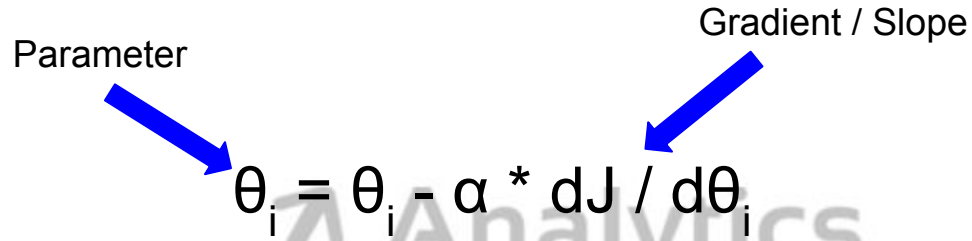
# Gradient Descent Recap

Parameter

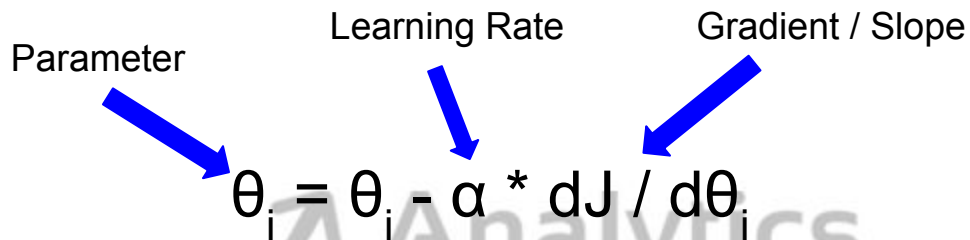$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Analytics Vidhya
Learn everything about analytics

# Gradient Descent Recap
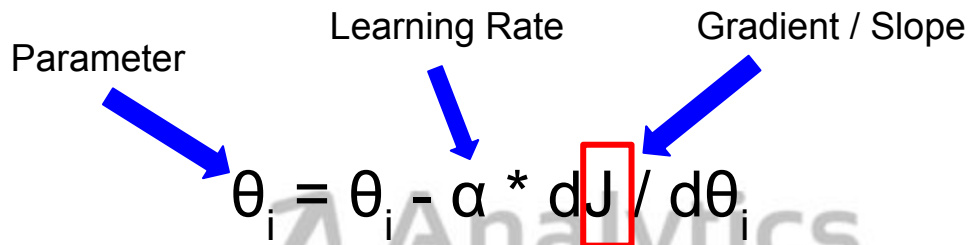
Parameter

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

# Gradient Descent Recap

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Analytics Vidhya
Learn everything about analytics

# Variants of Gradient Descent

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

**Analytics Vidhya**
Learn everything about analytics

# Variants of Gradient Descent

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Entire Training Set (m)

Batch Gradient Descent

# Variants of Gradient Descent

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

| | |
|---|---|
| Entire Training Set (m) | Batch Gradient Descent |
| Single Observation (1) | Stochastic Gradient Descent (SGD) |

# Variants of Gradient Descent

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

| | | |
|---|---|---|
| 1 | 2 | 1 |
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Observations

# Stochastic Gradient Descent (SGD)

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

Observations

| 1 | 2 | 1 |
|---|---|---|
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

# Stochastic Gradient Descent (SGD)

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

Observations

| 1 | 2 | 1 |
|---|---|---|
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Neural Network

Analytics Vidhya
Learn everything about analytics

# Stochastic Gradient Descent (SGD)

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

| | | |
|---|---|---|
| 1 | 2 | 1 |
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Observations

Neural Network → Error

Analytics Vidhya
Learn everything about analytics

# Stochastic Gradient Descent (SGD)

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

| 1 | 2 | 1 |
|---|---|---|
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Observations

Neural Network → Error → Update Parameters

# Stochastic Gradient Descent (SGD)

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

Observations

| 1 | 2 | 1 |
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Neural Network → Error → Update Parameters

Analytics Vidhya
Learn everything about analytics

# Stochastic Gradient Descent (SGD)

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

Observations =5

| 1 | 2 | 1 |
|---|---|---|
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Neural Network → Error → Update Parameters

# Variants of Gradient Descent

Parameter

Learning Rate

Gradient / Slope

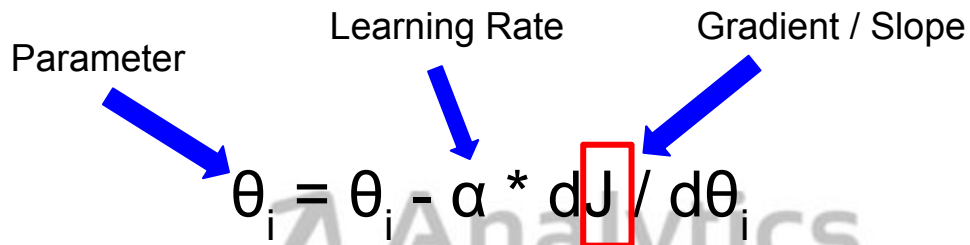$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

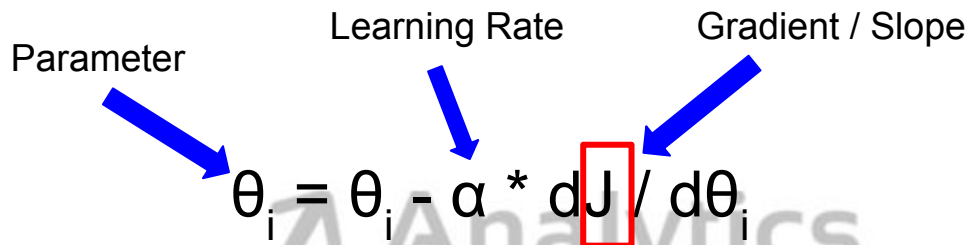Entire Training Set (m)

Batch Gradient Descent

Single Observation (1)

Stochastic Gradient Descent (SGD)

# Variants of Gradient Descent

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

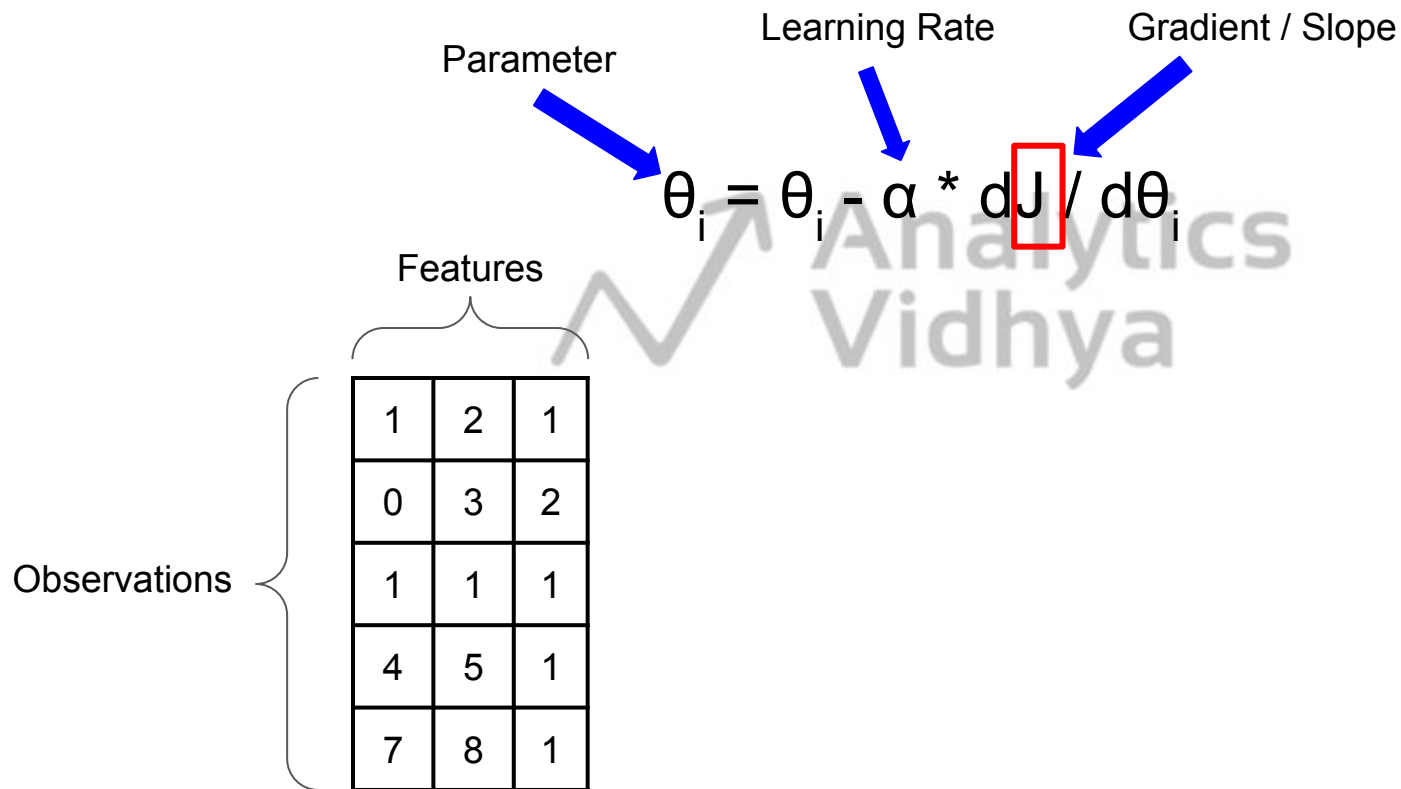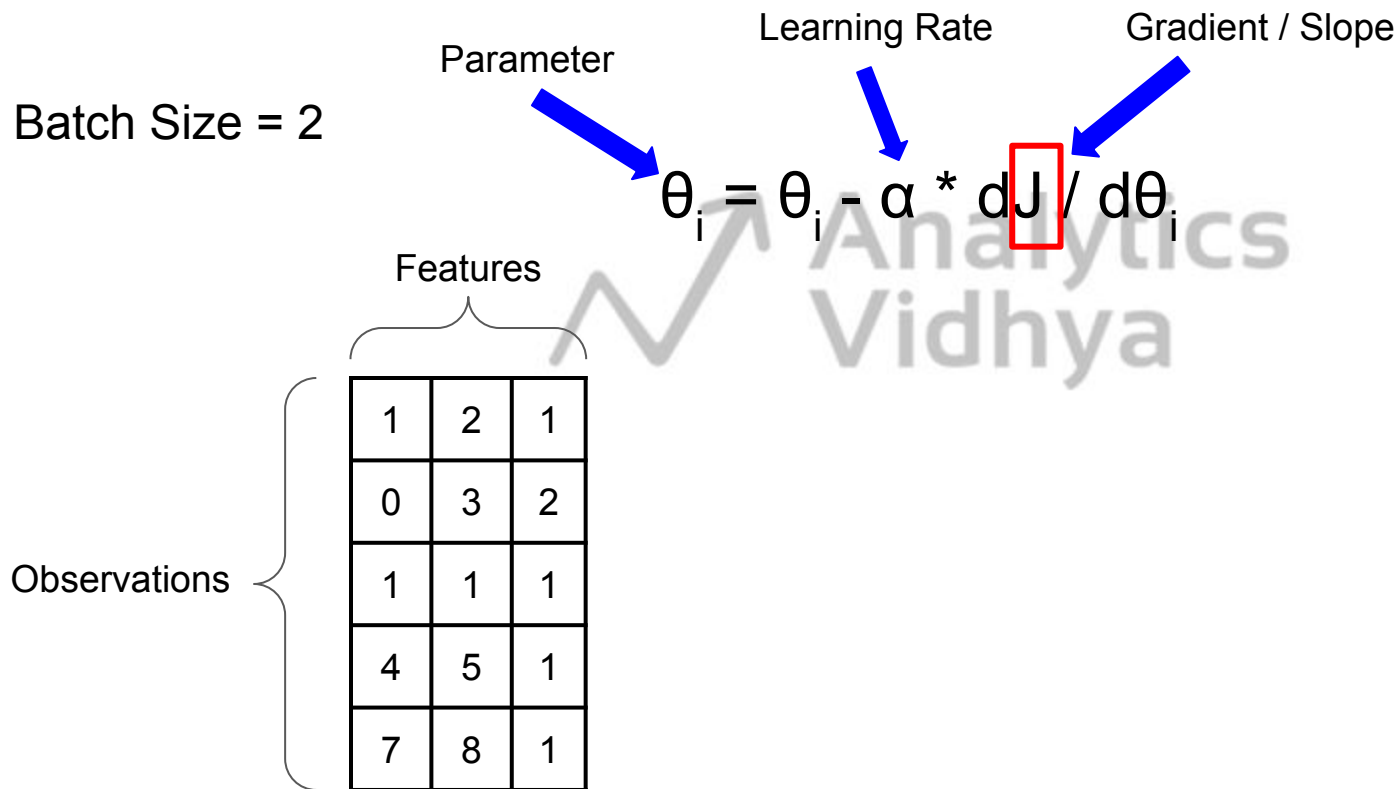| Entire Training Set (m) | Batch Gradient Descent |
|---|---|
| Single Observation (1) | Stochastic Gradient Descent (SGD) |
| $1 < x < m$ | Mini-Batch Gradient Descent |

# Mini-Batch Gradient Descent

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

| 1 | 2 | 1 |
|---|---|---|
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Observations

Analytics Vidhya
Learn everything about analytics

# Mini-Batch Gradient Descent

Batch Size = 2

Parameter     Learning Rate     Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

| | | |
|---|---|---|
| 1 | 2 | 1 |
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Observations

# Mini-Batch Gradient Descent

Batch Size = 2

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

Observations

| 1 | 2 | 1 |
|---|---|---|
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

# Mini-Batch Gradient Descent

Batch Size = 2

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

Observations

| 1 | 2 | 1 |
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Neural Network → Error

Analytics Vidhya
Learn everything about analytics

# Mini-Batch Gradient Descent

Batch Size = 2

Parameter

Learning Rate

Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

Features

Observations

| | | |
|---|---|---|
| 1 | 2 | 1 |
| 0 | 3 | 2 |
| 1 | 1 | 1 |
| 4 | 5 | 1 |
| 7 | 8 | 1 |

Neural Network → Error → Update Parameters
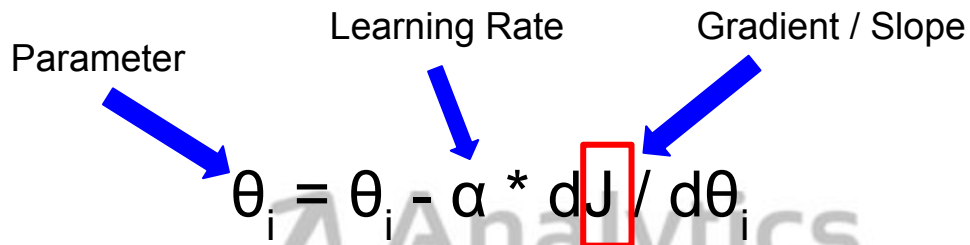
# Variants of Gradient Descent

Parameter  Learning Rate  Gradient / Slope

$$\theta_i = \theta_i - \alpha * dJ / d\theta_i$$

| | |
|---|---|
| Entire Training Set (m) | Batch Gradient Descent |
| Single Observation (1) | Stochastic Gradient Descent (SGD) |
| 1 < x < m | Mini-Batch Gradient Descent |

**Analytics Vidhya**
Learn everything about analytics

# Comparison: Variants of Gradient Descent

**Batch Gradient Descent**

**Stochastic Gradient Descent (SGD)**

**Mini-Batch Gradient Descent**

# Comparison: Number of observations used for updation

| Batch Gradient Descent | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent |
|---|---|---|

- Entire dataset for updation

# Comparison: Number of observations used for updation

| Batch Gradient Descent | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent |
|---|---|---|

- Entire dataset for updation

- Single observation for updation

# Comparison: Number of observations used for updation

| Batch Gradient Descent | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent |
|---|---|---|
| ● Entire dataset for updation | ● Single observation for updation | ● Subset of data for updation |

# Comparison: Cost Function

| Batch Gradient Descent | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent |

- Entire dataset for updation

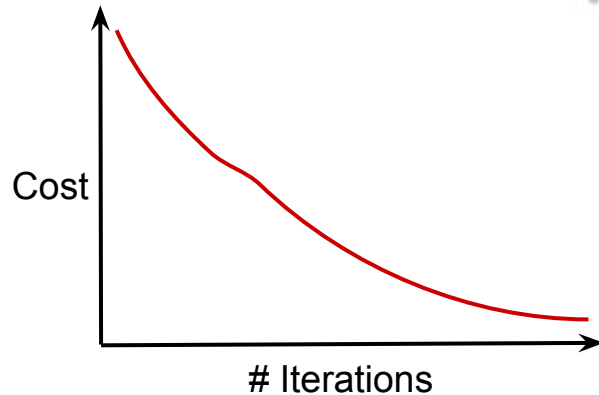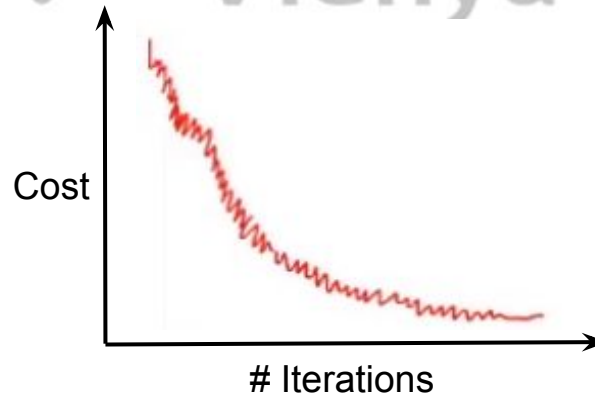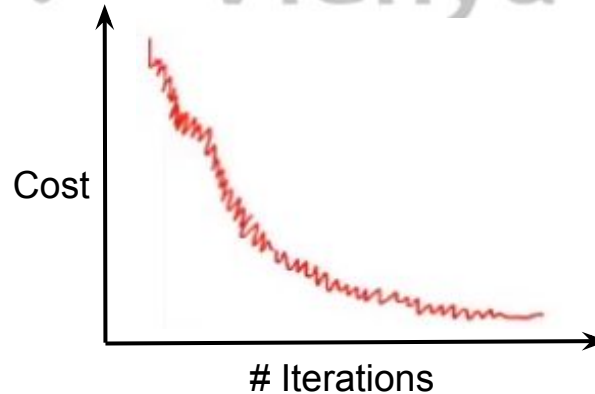- Cost function reduces smoothly

- Single observation for updation

- Subset of data for updation



Cost

\# Iterations

# Comparison: Cost Function

| Batch Gradient Descent | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent |
|---|---|---|
| ● Entire dataset for updation | ● Single observation for updation | ● Subset of data for updation |
| ● Cost function reduces smoothly | ● Lot of variations in cost function | ● Smoother cost function as compared to SGD |

Cost

# Iterations

Cost

# Iterations

Cost

# Iterations

# Comparison: Computation cost and time

| Batch Gradient Descent | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent |
|---|---|---|
| ● Entire dataset for updation | ● Single observation for updation | ● Subset of data for updation |
| ● Cost function reduces smoothly | ● Lot of variations in cost function | ● Smoother cost function as compared to SGD |
| ● Computation cost is very high | | |

# Comparison: Computation cost and time

| Batch Gradient Descent | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent |
|---|---|---|
| ● Entire dataset for updation | ● Single observation for updation | ● Subset of data for updation |
| ● Cost function reduces smoothly | ● Lot of variations in cost function | ● Smoother cost function as compared to SGD |
| ● Computation cost is very high | ● Computation time is more | |


Analytics Vidhya
Learn everything about analytics

# Comparison: Computation cost and time

| Batch Gradient Descent | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent |
|---|---|---|
| • Entire dataset for updation | • Single observation for updation | • Subset of data for updation |
| • Cost function reduces smoothly | • Lot of variations in cost function | • Smoother cost function as compared to SGD |
| • Computation cost is very high | • Computation time is more | • Computation time is lesser than SGD |
| | | • Computation cost is lesser than Batch Gradient Descent |

Thank You