

# Introduction to Machine Learning

## Stochastic Gradient Descent

Andres Mendez-Vazquez

January 28, 2023

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

# Outline

## 1. Introduction

- **Review Gradient Descent**
  - The Problems of Gradient Descent with Large Data Sets
  - Convergence of gradient descent with fixed step size
  - Convergence Rate
    - Convex Functions
    - Back to the Main Problem
  - Accelerating the Gradient Descent
  - Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

# Gradient Descent [1, 2]

The basic procedure is as follow

- 1 Start with a random weight vector  $w_0$ .

# Gradient Descent [1, 2]

The basic procedure is as follow

- 1 Start with a random weight vector  $w_0$ .
- 2 Compute the gradient vector  $\nabla J(w_0)$ .

# Gradient Descent [1, 2]

The basic procedure is as follow

- 1 Start with a random weight vector  $w_0$ .
- 2 Compute the gradient vector  $\nabla J(w_0)$ .
- 3 Obtain value  $w_1$  by moving from  $w_0$  in the direction of the steepest descent:

# Gradient Descent [1, 2]

The basic procedure is as follow

- 1 Start with a random weight vector  $\mathbf{w}_0$ .
- 2 Compute the gradient vector  $\nabla J(\mathbf{w}_0)$ .
- 3 Obtain value  $\mathbf{w}_1$  by moving from  $\mathbf{w}_0$  in the direction of the steepest descent:

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta_n \nabla J(\mathbf{w}_n) \quad (1)$$

# Gradient Descent [1, 2]

The basic procedure is as follow

- 1 Start with a random weight vector  $\mathbf{w}_0$ .
- 2 Compute the gradient vector  $\nabla J(\mathbf{w}_0)$ .
- 3 Obtain value  $\mathbf{w}_1$  by moving from  $\mathbf{w}_0$  in the direction of the steepest descent:

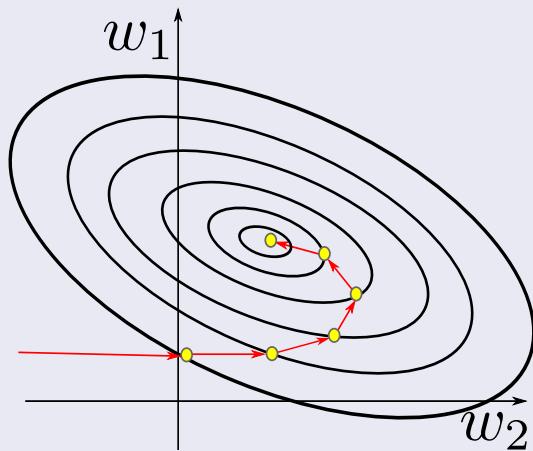
$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta_n \nabla J(\mathbf{w}_n) \quad (1)$$

$\eta_n$  is a positive scale factor or learning rate!!!



# Geometrically

We have the following



$$w(k+1) = w(k) - \eta(k) \nabla J(w(k))$$

# Outline

## 1. Introduction

- Review Gradient Descent
- **The Problems of Gradient Descent with Large Data Sets**
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

# Although

It is possible to prove

- That the gradient direction gives the greatest increase direction!!!

# Although

It is possible to prove

- That the gradient direction gives the greatest increase direction!!!

We have a problem in cost functions like in Deep Neural Networks

$$J(\mathbf{w}) = \sum_{i=1}^N (y_i - f(\mathbf{w}, \mathbf{x}_i))^2$$

- Where, we have that  $f(\mathbf{w}, \mathbf{x}_i) = f_1 \circ f_2 \circ f_3 \circ \dots \circ f_T(\mathbf{w}, \mathbf{x}_i)$

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- **Convergence of gradient descent with fixed step size**
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

Do you remember the problem of the  $\eta$  step size?

### Gradient Descent with fixed step size

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla J(\mathbf{w}_n)$$

Do you remember the problem of the  $\eta$  step size?

### Gradient Descent with fixed step size

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla J(\mathbf{w}_n)$$

### Why to worry about this?

- Because, we want to know how fast Gradient Descent will find the answer...

We have

### Lipschitz Continuous [3]

- Lipschitz continuity, named after Rudolf Lipschitz, is a strong form of uniform continuity for functions.



We have

### Lipschitz Continuous [3]

- Lipschitz continuity, named after Rudolf Lipschitz, is a strong form of uniform continuity for functions.

### Uniform continuity

- The function  $f : A \rightarrow \mathbb{R}$  is said to be uniformly continuous on  $A$  iff for every  $\epsilon > 0$ ,  $\exists \delta > 0$  such that  $|x - y| < \delta$  implies  $|f(x) - f(y)| < \epsilon$ .

# Lipschitz Continuous

## Definition

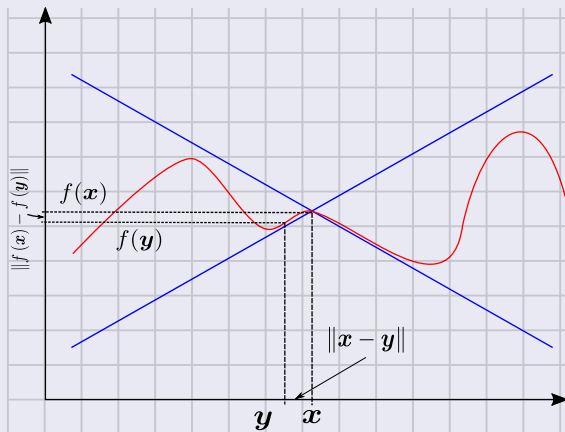
- A function  $f : S \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  satisfies the Lipschitz Continuous at  $x \in S$ , if there is a such constant  $L > 0$  such that

$$\|f(x) - f(y)\| \leq L \|x - y\|$$

for all  $y \in S$  sufficiently near to  $x$ . **Lipschitz continuity can be seen as a refinement of continuity.**

# Example when you see $L$ as the slope

Here the function  $f : \mathbb{R} \rightarrow \mathbb{R}$



## An interesting property of such setup

The derivative of the function cannot exceed  $L$  (Example,  $f : \mathbb{R} \rightarrow \mathbb{R}$ )

$$f'(x) = \lim_{\delta \rightarrow \infty} \frac{f(x + \delta) - f(x)}{\delta}$$

## An interesting property of such setup

The derivative of the function cannot exceed  $L$  (Example,  $f : \mathbb{R} \rightarrow \mathbb{R}$ )

$$f'(x) = \lim_{\delta \rightarrow \infty} \frac{f(x + \delta) - f(x)}{\delta}$$

Then, we have that

$$f'(x) = \lim_{\delta \rightarrow \infty} \frac{f(x) - f(y)}{x - y} \leq \lim_{\delta \rightarrow \infty} \frac{|f(x) - f(y)|}{|x - y|} \leq L$$

Therefore

Lipschitz Continuity implies

$$|f'(x)| < L$$

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- **Convergence Rate**
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

# Convergence idea

## Definition (Big $O$ - Upper Bound) [4]

For a given function  $g(n)$ :

$$O(g(n)) = \{f(n) \mid \text{There exists } c > 0 \text{ and } n_0 > 0 \\ \text{s.t. } 0 \leq f(n) \leq cg(n) \forall n \geq n_0\}$$



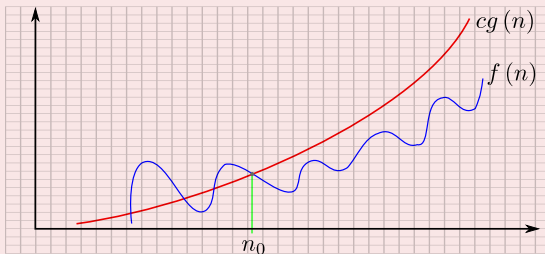
# Convergence idea

## Definition (Big $O$ - Upper Bound) [4]

For a given function  $g(n)$ :

$$O(g(n)) = \{f(n) \mid \text{There exists } c > 0 \text{ and } n_0 > 0 \\ \text{s.t. } 0 \leq f(n) \leq cg(n) \forall n \geq n_0\}$$

## Example



# What are the implications?

## Definition [3]

- Suppose that the sequence  $\{x_k\}$  converges to the number  $L$ :

# What are the implications?

## Definition [3]

- Suppose that the sequence  $\{x_k\}$  converges to the number  $L$ :

We say that this sequence converges linearly to  $L$ , if there exists a number  $\frac{1}{n} \in (0, 1)$  such that

$$\lim_{k \rightarrow \infty} \frac{|x_{n+1} - L|}{|x_n - L|} = \frac{1}{n}$$

# What are the implications?

## Definition [3]

- Suppose that the sequence  $\{x_k\}$  converges to the number  $L$ :

We say that this sequence converges linearly to  $L$ , if there exists a number  $\frac{1}{n} \in (0, 1)$  such that

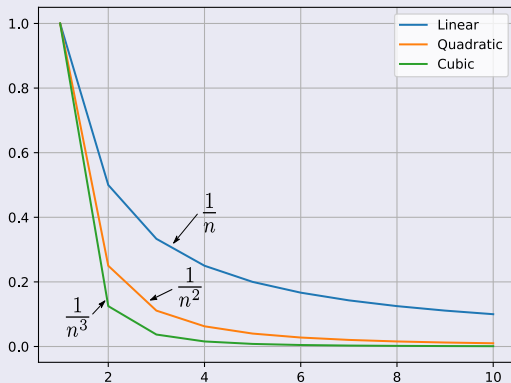
$$\lim_{k \rightarrow \infty} \frac{|x_{n+1} - L|}{|x_n - L|} = \frac{1}{n}$$

Thus, Gradient Descent has a linear convergence speed

- If you do a comparison with quadratic convergence...

# Example

As you can see the quadratic is faster than linear in convergence



# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- **Convergence Rate**
  - **Convex Functions**
    - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

# Why the importance of Convex Functions?

There is an interest on the rates of convergence for many optimization algorithms

- And they are affected by the different cost function that can be used:
  - ▶ Lipschitz-continuity, convexity, strong convexity, and smoothness

# Why the importance of Convex Functions?

There is an interest on the rates of convergence for many optimization algorithms

- And they are affected by the different cost function that can be used:
  - ▶ Lipschitz-continuity, convexity, strong convexity, and smoothness

There are different rates of convergence for the Gradient Descent

- For example when a function is strongly convex

$$\nabla^2 f(x) \succeq \alpha I \iff \nabla^2 f(x) - \alpha I \succeq 0 \text{ (Matrix greater of equal)}$$



# Why the importance of Convex Functions?

There is an interest on the rates of convergence for many optimization algorithms

- And they are affected by the different cost function that can be used:
  - ▶ Lipschitz-continuity, convexity, strong convexity, and smoothness

There are different rates of convergence for the Gradient Descent

- For example when a function is strongly convex

$$\nabla^2 f(x) \succeq \alpha I \iff \nabla^2 f(x) - \alpha I \succeq 0 \text{ (Matrix greater of equal)}$$

This means that

- The curvature of  $f(x)$  is not very close to zero, making possible to accelerate the convergence

# Convex Sets

## Definition

- For a convex set  $X$ , for any two points  $x$  and  $y$  such that  $x, y \in X$ , the line between them lies within the set

$$z = \lambda x + (1 - \lambda) y, \forall \lambda \in (0, 1) \text{ then } z \in X$$

- ▶ The sum  $\lambda x + (1 - \lambda) y$  is termed as convex linear combination.

# Convex Functions

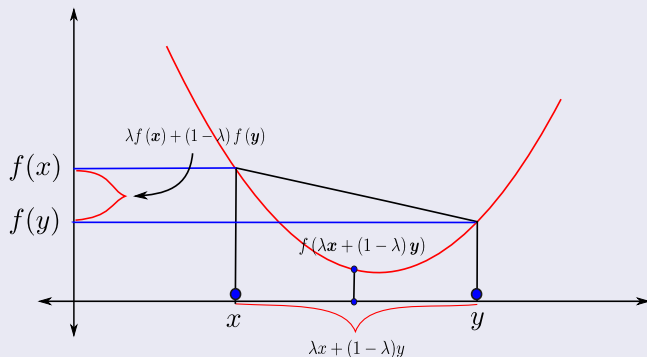
## Definition

- A function  $f(\mathbf{x})$  is convex if the following holds:
  - 1 The Domain of  $f$  is convex
  - 2  $\forall \mathbf{x}, \mathbf{y}$  in the Domain of  $f$  and  $\lambda \in (0, 1)$

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y})$$

# Graphically

This can further expanded to functions



# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- **Convergence Rate**
  - Convex Functions
  - **Back to the Main Problem**
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

# Convergence of gradient descent with **fixed step size**

## Theorem

- Suppose the function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex and differentiable, and we have that  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|$  (Lipschitz Continuous Gradient) for any  $\mathbf{x}, \mathbf{y}$  and  $L > 0$ .

# Convergence of gradient descent with **fixed step size**

## Theorem

- Suppose the function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex and differentiable, and we have that  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|$  (Lipschitz Continuous Gradient) for any  $\mathbf{x}, \mathbf{y}$  and  $L > 0$ .

## We have that

- Then, if we run the **gradient descent** for  $k$  iterations with a fixed step size  $\eta \leq \frac{1}{L}$ , it will yield a solution  $f_n$  which satisfies

$$f(x_n) - f(x^*) \leq \frac{\|x_{(0)} - x^*\|_2^2}{2\eta n}$$

where  $f(x^*)$  is the optimal value.

# Proof

$f(\mathbf{x})$  is Lipschitz continuous with constant  $L$  implies  
( $\|\mathbf{y} - \mathbf{x}\|^2 = \|\mathbf{y} - \mathbf{x}\|_2^2$ )

$\nabla^2 f(\mathbf{x}) - LI$  as semi-definite matrix



# Proof

$f(\mathbf{x})$  is Lipschitz continuous with constant  $L$  implies  
( $\|\mathbf{y} - \mathbf{x}\|^2 = \|\mathbf{y} - \mathbf{x}\|_2^2$ )

$\nabla^2 f(\mathbf{x}) - LI$  as semi-definite matrix

We have the following inequality

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2} \nabla^2 f(\mathbf{x}) \|\mathbf{y} - \mathbf{x}\|^2 \\ &\leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2} L \|\mathbf{y} - \mathbf{x}\|^2 \end{aligned}$$

# Proof

Now, if we apply the Gradient update  $\mathbf{y} = \mathbf{x}^+ = \mathbf{x} - \eta \nabla f(\mathbf{x})$

$$\begin{aligned} f(\mathbf{x}^+) &\leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{x}^+ - \mathbf{x}) + \frac{1}{2} L \|\mathbf{x}^+ - \mathbf{x}\|^2 \\ &= f(\mathbf{x}) - \eta \|\nabla f(\mathbf{x})\|^2 + \frac{1}{2} L \eta^2 \|\nabla f(\mathbf{x})\|^2 \\ &= f(\mathbf{x}) - \left(1 - \frac{1}{2} L \eta\right) \eta \|\nabla f(\mathbf{x})\|^2 \end{aligned}$$

# Proof

Now, if we apply the Gradient update  $\mathbf{y} = \mathbf{x}^+ = \mathbf{x} - \eta \nabla f(\mathbf{x})$

$$\begin{aligned} f(\mathbf{x}^+) &\leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{x}^+ - \mathbf{x}) + \frac{1}{2} L \|\mathbf{x}^+ - \mathbf{x}\|^2 \\ &= f(\mathbf{x}) - \eta \|\nabla f(\mathbf{x})\|^2 + \frac{1}{2} L \eta^2 \|\nabla f(\mathbf{x})\|^2 \\ &= f(\mathbf{x}) - \left(1 - \frac{1}{2} L \eta\right) \eta \|\nabla f(\mathbf{x})\|^2 \end{aligned}$$

Using  $\eta \leq \frac{1}{L}$

$$-\left(1 - \frac{1}{2} L \eta\right) \leq -\frac{1}{2}$$

Therefore

We have that

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) - \frac{1}{2}\eta \|\nabla f(\mathbf{x})\|^2 \quad (2)$$

Therefore

We have that

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) - \frac{1}{2}\eta \|\nabla f(\mathbf{x})\|^2 \quad (2)$$

Implying that

- This inequality implies that the objective function value strictly decreases until it reaches the optimal value

Therefore

We have that

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) - \frac{1}{2}\eta \|\nabla f(\mathbf{x})\|^2 \quad (2)$$

Implying that

- This inequality implies that the objective function value strictly decreases until it reaches the optimal value

This only holds when  $\eta$  is small enough

- This explains why we observe in practice that gradient descent diverges when the step size is too large.

Since  $f$  is convex

We can write

$$\begin{aligned}f(\mathbf{x}^*) &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{x}^* - \mathbf{x}) \\f(\mathbf{x}) &\leq f(\mathbf{x}^*) + \nabla f(\mathbf{x})^T (\mathbf{x} - \mathbf{x}^*)\end{aligned}$$

Since  $f$  is convex

We can write

$$\begin{aligned}f(\mathbf{x}^*) &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{x}^* - \mathbf{x}) \\f(\mathbf{x}) &\leq f(\mathbf{x}^*) + \nabla f(\mathbf{x})^T (\mathbf{x} - \mathbf{x}^*)\end{aligned}$$

This comes from the “First order condition for convexity”

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x})$$



Then

Plugging this in to (Equation 2)

$$f(\mathbf{x}^+) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x})^T (\mathbf{x} - \mathbf{x}^*) - \frac{1}{2} \eta \|\nabla f(\mathbf{x})\|^2$$

Then

Plugging this in to (Equation 2)

$$f(\mathbf{x}^+) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x})^T (\mathbf{x} - \mathbf{x}^*) - \frac{1}{2}\eta \|\nabla f(\mathbf{x})\|^2$$

Therefore

$$f(\mathbf{x}^+) - f(\mathbf{x}^*) \leq \frac{1}{2\eta} \left[ \|\mathbf{x} - \mathbf{x}^*\|^2 - \|\mathbf{x} - \eta \nabla f(\mathbf{x}) - \mathbf{x}^*\|^2 \right]$$

Then

Plugging this in to (Equation 2)

$$f(\mathbf{x}^+) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x})^T (\mathbf{x} - \mathbf{x}^*) - \frac{1}{2}\eta \|\nabla f(\mathbf{x})\|^2$$

Therefore

$$f(\mathbf{x}^+) - f(\mathbf{x}^*) \leq \frac{1}{2\eta} \left[ \|\mathbf{x} - \mathbf{x}^*\|^2 - \|\mathbf{x} - \eta \nabla f(\mathbf{x}) - \mathbf{x}^*\|^2 \right]$$

Then plugging this  $\mathbf{x}^+ = \mathbf{x} - \eta \nabla f(\mathbf{x})$  into

$$f(\mathbf{x}^+) - f(\mathbf{x}^*) \leq \frac{1}{2\eta} \left[ \|\mathbf{x} - \mathbf{x}^*\|^2 - \|\mathbf{x}^+ - \mathbf{x}^*\|^2 \right]$$

Then

Summing over all iterations and the telescopic sum in the right side

$$\sum_{i=1}^n \left[ f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \leq \frac{1}{2\eta} \left[ \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 \right]$$

Then

Summing over all iterations and the telescopic sum in the right side

$$\sum_{i=1}^n \left[ f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \leq \frac{1}{2\eta} \left[ \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 \right]$$

Finally, using the fact that  $f$  decreasing on every iteration

$$f(\mathbf{x}^{(n)}) - f(\mathbf{x}^*) \leq \frac{1}{n} \sum_{i=1}^n \left[ f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \leq \frac{1}{2\eta n} \left[ \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 \right] = C \times \frac{1}{n}$$

Therefore

It converges with rate

$$O\left(\frac{1}{n}\right)$$

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- **Accelerating the Gradient Descent**
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

# Accelerating the Gradient Descent

It is possible to modify the Batch Gradient Descent

- In order to accelerate it several modifications have been proposed



# Accelerating the Gradient Descent

It is possible to modify the Batch Gradient Descent

- In order to accelerate it several modifications have been proposed

## Possible Methods

- Polyak's Momentum Method or Heavy-Ball Method (1964)
- Nesterov's Proposal (1983)
- Stochastic Gradient Descent (1951)

# Polyak's Momentum Method

## Polyak's Step Size

- He Proposed that the step size could be modified to

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \alpha \nabla f(\mathbf{w}_n) + \mu (\mathbf{w}_n - \mathbf{w}_{n-1}) \text{ with } \mu \in [0, 1], \alpha > 0$$

# Polyak's Momentum Method

## Polyak's Step Size

- He Proposed that the step size could be modified to

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \alpha \nabla f(\mathbf{w}_n) + \mu (\mathbf{w}_n - \mathbf{w}_{n-1}) \text{ with } \mu \in [0, 1], \alpha > 0$$

Basically, the method uses the previous gradient information through the step difference  $(\mathbf{w}_n - \mathbf{w}_{n-1})$

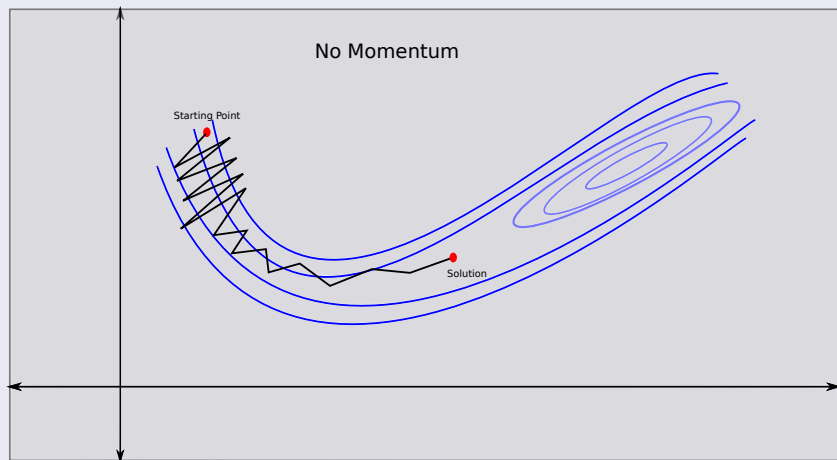
- By the discretization of the second order ODE

$$\ddot{\mathbf{w}} + a\dot{\mathbf{w}} + b\nabla f(\mathbf{w}) = 0$$

- ▶ **which models the motion of a body in a potential field given by  $f$  with friction.**

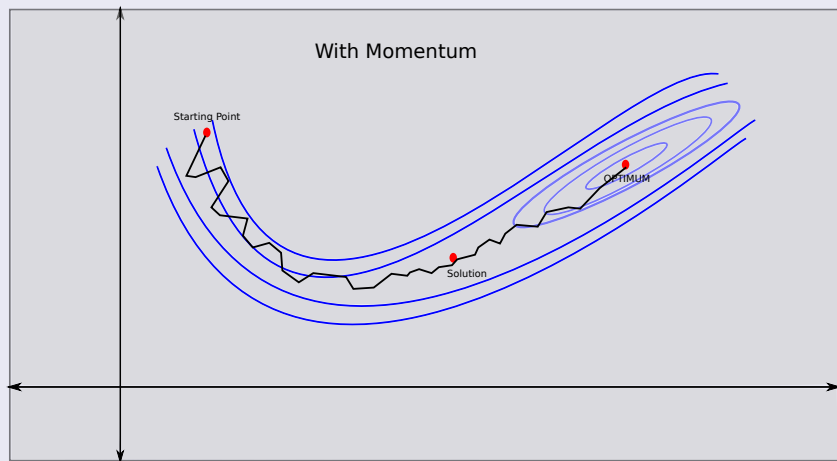
# The Momentum helps to stabilize the GD

## If we do not have Momentum



# Then, with Momentum

## If we have Momentum



# Problem

It has been proved that the method has problems

- L. Lessard, B. Recht, and A. Packard. Analysis and Design of Optimization Algorithms via Integral Quadratic Constraints. ArXiv e-prints, Aug. 2014.

# Problem

It has been proved that the method has problems

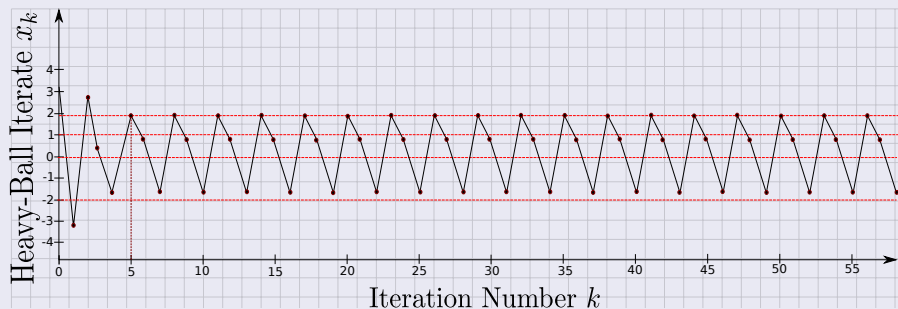
- L. Lessard, B. Recht, and A. Packard. Analysis and Design of Optimization Algorithms via Integral Quadratic Constraints. ArXiv e-prints, Aug. 2014.

Under the function

$$\nabla f(x) = \begin{cases} 25x & \text{if } x < 1 \\ x + 24 & \text{if } 1 \leq x \leq 2 \\ 25x - 24 & \text{if otherwise} \end{cases}$$

In Lessard et al.

We have a non-convergence (Original Lessard et al.) [5]





# Nesterov's Proposal

He proposed a Quasi-Convex Combination

- Instead to use

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \alpha \nabla f(\mathbf{w}_n) + \mu (\mathbf{w}_n - \mathbf{w}_{n-1})$$

# Nesterov's Proposal

He proposed a Quasi-Convex Combination

- Instead to use

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \alpha \nabla f(\mathbf{w}_n) + \mu (\mathbf{w}_n - \mathbf{w}_{n-1})$$

Have an intermediate step to update  $\mathbf{w}_{n+1}$

$$\mathbf{w}_{n+1} = (1 - \gamma_n) \mathbf{y}_{n+1} + \gamma_n \mathbf{y}_n$$

# Nesterov's Proposal

He proposed a Quasi-Convex Combination

- Instead to use

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \alpha \nabla f(\mathbf{w}_n) + \mu (\mathbf{w}_n - \mathbf{w}_{n-1})$$

Have an intermediate step to update  $\mathbf{w}_{n+1}$

$$\mathbf{w}_{n+1} = (1 - \gamma_n) \mathbf{y}_{n+1} + \gamma_n \mathbf{y}_n$$

This allow to weight the actual original gradient change

- with the previous gradient change... making possible to avoid the original problem by Polyak... Which is based in Lyapunov Analysis

## Nesterov's Proposal [6]

### Nesterov's Accelerated Gradient Descent (A Quasi-Convex Modification)

$$\mathbf{y}_{n+1} = \mathbf{w}_n - \frac{1}{\beta} \nabla J(\mathbf{w}_n)$$

$$\mathbf{w}_{n+1} = (1 - \gamma_n) \mathbf{y}_{n+1} + \gamma_n \mathbf{y}_n$$

## Nesterov's Proposal [6]

### Nesterov's Accelerated Gradient Descent (A Quasi-Convex Modification)

$$\mathbf{y}_{n+1} = \mathbf{w}_n - \frac{1}{\beta} \nabla J(\mathbf{w}_n)$$

$$\mathbf{w}_{n+1} = (1 - \gamma_n) \mathbf{y}_{n+1} + \gamma_n \mathbf{y}_n$$

Where, we use the following constants

$$\lambda_0 = 0$$

$$\lambda_n = \frac{1 + \sqrt{1 + 4\lambda_{n-1}^2}}{2}$$

$$\gamma_n = \frac{1 - \lambda_n}{\lambda_{n+1}}$$

# Nesterov's Algorithm

## Nesterov Accelerated Gradient

**Input:** Training Time  $T$ , Learning Rate  $\beta$ , an initialization  $w_0$

- 1  $y_0 \leftarrow w_0$
- 2  $\lambda_0 \leftarrow 0$

# Nesterov's Algorithm

## Nesterov Accelerated Gradient

**Input:** Training Time  $T$ , Learning Rate  $\beta$ , an initialization  $w_0$

- ①  $y_0 \leftarrow w_0$
- ②  $\lambda_0 \leftarrow 0$
- ③ **for**  $t = 0$  **to**  $T - 1$  **do**

# Nesterov's Algorithm

## Nesterov Accelerated Gradient

**Input:** Training Time  $T$ , Learning Rate  $\beta$ , an initialization  $w_0$

- ①  $y_0 \leftarrow w_0$
- ②  $\lambda_0 \leftarrow 0$
- ③ **for**  $t = 0$  **to**  $T - 1$  **do**
- ④      $y_{n+1} = w_n - \frac{1}{\beta} \nabla J(w_n)$



# Nesterov's Algorithm

## Nesterov Accelerated Gradient

**Input:** Training Time  $T$ , Learning Rate  $\beta$ , an initialization  $\mathbf{w}_0$

- ①  $\mathbf{y}_0 \leftarrow \mathbf{w}_0$
- ②  $\lambda_0 \leftarrow 0$
- ③ **for**  $t = 0$  **to**  $T - 1$  **do**
- ④      $\mathbf{y}_{n+1} = \mathbf{w}_n - \frac{1}{\beta} \nabla J(\mathbf{w}_n)$
- ⑤      $\lambda_n = \frac{1 + \sqrt{1 + 4\lambda_{n-1}^2}}{2}$

# Nesterov's Algorithm

## Nesterov Accelerated Gradient

**Input:** Training Time  $T$ , Learning Rate  $\beta$ , an initialization  $\mathbf{w}_0$

- ①  $\mathbf{y}_0 \leftarrow \mathbf{w}_0$
- ②  $\lambda_0 \leftarrow 0$
- ③ **for**  $t = 0$  **to**  $T - 1$  **do**
- ④      $\mathbf{y}_{n+1} = \mathbf{w}_n - \frac{1}{\beta} \nabla J(\mathbf{w}_n)$
- ⑤      $\lambda_n = \frac{1 + \sqrt{1 + 4\lambda_{n-1}^2}}{2}$
- ⑥      $\lambda_{n+1} = \frac{1 + \sqrt{1 + 4\lambda_n^2}}{2}$

# Nesterov's Algorithm

## Nesterov Accelerated Gradient

**Input:** Training Time  $T$ , Learning Rate  $\beta$ , an initialization  $\mathbf{w}_0$

- ①  $\mathbf{y}_0 \leftarrow \mathbf{w}_0$
- ②  $\lambda_0 \leftarrow 0$
- ③ **for**  $t = 0$  **to**  $T - 1$  **do**
- ④      $\mathbf{y}_{n+1} = \mathbf{w}_n - \frac{1}{\beta} \nabla J(\mathbf{w}_n)$
- ⑤      $\lambda_n = \frac{1 + \sqrt{1 + 4\lambda_{n-1}^2}}{2}$
- ⑥      $\lambda_{n+1} = \frac{1 + \sqrt{1 + 4\lambda_n^2}}{2}$
- ⑦      $\gamma_n = \frac{1 - \lambda_n}{\lambda_{n+1}}$

# Nesterov's Algorithm

## Nesterov Accelerated Gradient

**Input:** Training Time  $T$ , Learning Rate  $\beta$ , an initialization  $\mathbf{w}_0$

- 1  $\mathbf{y}_0 \leftarrow \mathbf{w}_0$
- 2  $\lambda_0 \leftarrow 0$
- 3 **for**  $t = 0$  **to**  $T - 1$  **do**
- 4      $\mathbf{y}_{n+1} = \mathbf{w}_n - \frac{1}{\beta} \nabla J(\mathbf{w}_n)$
- 5      $\lambda_n = \frac{1 + \sqrt{1 + 4\lambda_{n-1}^2}}{2}$
- 6      $\lambda_{n+1} = \frac{1 + \sqrt{1 + 4\lambda_n^2}}{2}$
- 7      $\gamma_n = \frac{1 - \lambda_n}{\lambda_{n+1}}$
- 8      $\mathbf{w}_{n+1} = (1 - \gamma_n) \mathbf{y}_{n+1} + \gamma_n \mathbf{y}_n$

## With the following complexity

### Theorem (Nesterov 1983)

- Let  $f$  be a convex and  $\beta$ -smooth function ( $\nabla f$  is  $\beta$ -Lipschitz continuous), then Nesterov's Accelerated Gradient Descent satisfies:

$$f(\mathbf{y}_{n+1}) - f(\mathbf{w}^*) \leq \frac{2\beta \|\mathbf{w}_1 - \mathbf{w}^*\|^2}{n^2}$$

## With the following complexity

### Theorem (Nesterov 1983)

- Let  $f$  be a convex and  $\beta$ -smooth function ( $\nabla f$  is  $\beta$ -Lipschitz continuous), then Nesterov's Accelerated Gradient Descent satisfies:

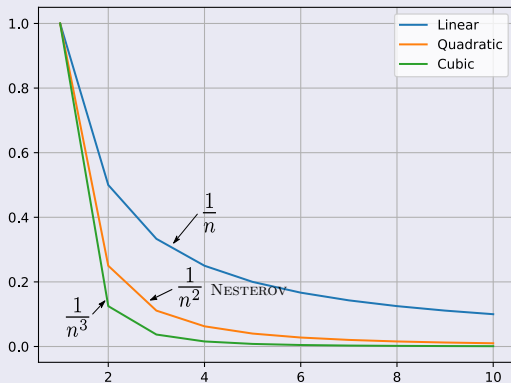
$$f(\mathbf{y}_{n+1}) - f(\mathbf{w}^*) \leq \frac{2\beta \|\mathbf{w}_1 - \mathbf{w}^*\|^2}{n^2}$$

It converges with rate

$$O\left(\frac{1}{n^2}\right)$$

# Example

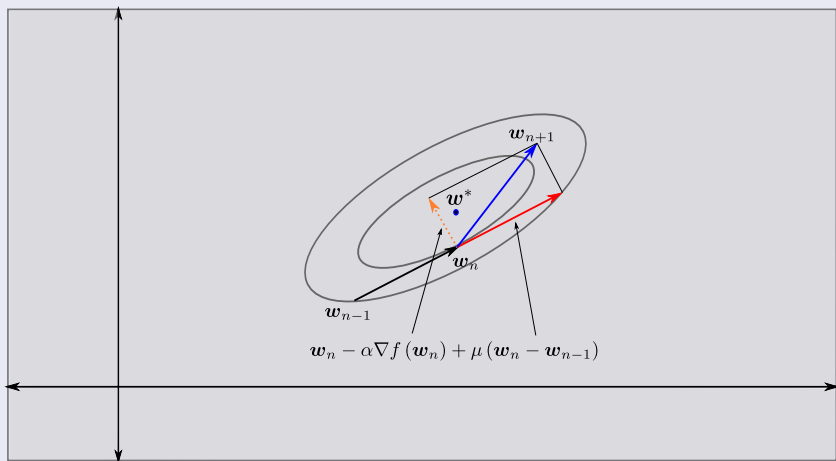
As you can see Nesterov is faster...



## Remark, Polyak vs Nesterov

### We have a remarkable difference

- The gradient descent step (orange arrow) is perpendicular to the level set before applying momentum to  $w_1$  (red arrow) in Polyak's algorithm





## In the case of Nesterov

If we rewrite the equations

$$\begin{aligned}\mathbf{w}_{n+1} &= (1 - \gamma_n) \left[ \mathbf{w}_n - \frac{1}{\beta} \nabla J(\mathbf{w}_n) \right] + \gamma_n \mathbf{y}_n \\ &= \mathbf{w}_n - \gamma_n \mathbf{w}_n - \frac{1}{\beta} \nabla J(\mathbf{w}_n) + \frac{\gamma_n}{\beta} \nabla J(\mathbf{w}_n) + \gamma_n \mathbf{w}_{n-1} - \frac{\gamma_n}{\beta} \nabla J(\mathbf{w}_{n-1}) \\ &= \mathbf{w}_n - \gamma_n (\mathbf{w}_n - \mathbf{w}_{n-1}) - \frac{1}{\beta} [\nabla J(\mathbf{w}_n) + \gamma_n \nabla J(\mathbf{w}_n) - \gamma_n \nabla J(\mathbf{w}_{n-1})] \\ &= \mathbf{w}_n - \gamma_n (\mathbf{w}_n - \mathbf{w}_{n-1}) - \frac{1}{\beta} [\nabla J(\mathbf{w}_n + \gamma_n [\mathbf{w}_n - \mathbf{w}_{n-1}])]\end{aligned}$$



There is a dependence with respect with different properties of  $f$

In this table, we can see upper bounds for the convergences  $D = \|\mathbf{x}_1 - \mathbf{x}^*\|_2$  and  $\lambda$  regularization term [7]

Properties of the Objective Function	Upper Bound for Gradient Descent
convex and $L$ -Lipschitz	$\frac{D_1 L}{\sqrt{n}}$
convex and $\beta$ -smooth	$\frac{\beta D_1^2}{n}$
$\alpha$ -strongly convex and $L$ -Lipschitz	$\frac{L^2}{\alpha n}$
$\alpha$ -strongly convex and $\beta$ -smooth	$\beta D_1^2 \exp\left(-\frac{4n}{\beta/\lambda}\right)$

# A Hierarchy can be established (Black Box Model)

## Based on the following idea

- A black box model assumes that the algorithm does not know the objective function  $f$  being minimized.

# A Hierarchy can be established (Black Box Model)

## Based on the following idea

- A black box model assumes that the algorithm does not know the objective function  $f$  being minimized.

## Not only that

- Information about the objective function can only be accessed by querying an oracle.

# A Hierarchy can be established (Black Box Model)

## Based on the following idea

- A black box model assumes that the algorithm does not know the objective function  $f$  being minimized.

## Not only that

- Information about the objective function can only be accessed by querying an oracle.

## Remarks

- The oracle serves as a bridge between the unknown objective function and the optimizer.

## Furthermore

At any given step, the optimizer queries the oracle with a guess  $x$

- The oracle responds with information about the function around  $x$

## Furthermore

At any given step, the optimizer queries the oracle with a guess  $x$

- The oracle responds with information about the function around  $x$

For Example

- Value of the Cost function, Gradient, Hessian, etc.



Then, we have

## Zeroth Order Methods

- These methods only require the value of function  $f$  at the current guess  $x$ .
  - ▶ The Bisection, Genetic Algorithms, Simulated Annealing and Metropolis-Hastings methods

Then, we have

## Zeroth Order Methods

- These methods only require the value of function  $f$  at the current guess  $x$ .
  - ▶ The Bisection, Genetic Algorithms, Simulated Annealing and Metropolis-Hastings methods

## First Order Methods

- These methods can inquire the value of the function  $f$  and its first derivative.
  - ▶ Gradient descent, Nesterov's and Polyak's

Then, we have

## Zeroth Order Methods

- These methods only require the value of function  $f$  at the current guess  $x$ .
  - ▶ The Bisection, Genetic Algorithms, Simulated Annealing and Metropolis-Hastings methods

## First Order Methods

- These methods can inquire the value of the function  $f$  and its first derivative.
  - ▶ Gradient descent, Nesterov's and Polyak's

## Second Order Methods

- These methods require the value of the function  $f$ , its first derivative (gradient), and its second derivative (Hessian).
  - ▶ Newton's method. Improving the efficiency of these algorithms is an active area of research.

# One of the Last Hierarchy

## Adaptive Methods and Conjugate Gradients

- The methods we mentioned until this point assume that all dimensions of a vector-valued variable have a common set of hyperparameters.

# One of the Last Hierarchy

## Adaptive Methods and Conjugate Gradients

- The methods we mentioned until this point assume that all dimensions of a vector-valued variable have a common set of hyperparameters.

## Adaptive methods relax this assumption

- They allow for every variable to have its own set of hyperparameters.

# One of the Last Hierarchy

## Adaptive Methods and Conjugate Gradients

- The methods we mentioned until this point assume that all dimensions of a vector-valued variable have a common set of hyperparameters.

## Adaptive methods relax this assumption

- They allow for every variable to have its own set of hyperparameters.

## Some popular methods under this paradigm

- AdaGrad and ADAM

# Finally, but not less important

## Lower Bounds

- Lower bounds are useful because they tell us what's the best possible rate of convergence we can have given a category of optimizer.

# Finally, but not less important

## Lower Bounds

- Lower bounds are useful because they tell us what's the best possible rate of convergence we can have given a category of optimizer.

## Something Notable

- Without lower bounds, an unnecessary amount of research energy would be spent in designing better optimizers
  - ▶ Even if convergence rate improvement is impossible within this category of algorithm



# Finally, but not less important

## Lower Bounds

- Lower bounds are useful because they tell us what's the best possible rate of convergence we can have given a category of optimizer.

## Something Notable

- Without lower bounds, an unnecessary amount of research energy would be spent in designing better optimizers
  - ▶ Even if convergence rate improvement is impossible within this category of algorithm

However, if we prove that each procedure has a lower bounded rate of convergence

- We do not know if a specific method reaches this bound.

However

Please, take a look

- Convex Optimization: Algorithms and Complexity by Sébastien Bubeck - Theory Group, Microsoft Research [7]

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- **Even with such Speeds**

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

## In our classic Convex Scenario [2]

Least Square Problem looking to minimize the average of the LSE

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2M} \sum_{m=1}^M \left( \mathbf{w}^T \mathbf{x}_m - y_m \right)^2 = \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2M} \|X\mathbf{w} - Y\|^2$$

Therefore

### Calculating the Gradient

$$\nabla_{\mathbf{w}} f(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M (\mathbf{w}^T \mathbf{x}_m - y_m) \mathbf{x}_m$$

# Observations

It is easy to verify that the complexity per iteration is  $O(dM)$

- With  $M$  is for the sum and  $d$  is for  $\mathbf{w}^T \mathbf{x}_m$ .

# Drawbacks

When the number of samples  $M$  is Large

- Even with a rate of linear convergence, Gradient Descent

# Drawbacks

When the number of samples  $M$  is Large

- Even with a rate of linear convergence, Gradient Descent

Not only that but in the On-line Learning scenario

- The data  $(\mathbf{x}_i, y_i)$  is coming one by one making the gradient not computable.



Therefore

Thus, the need to look for something faster

- Two possibilities:

# Therefore

Thus, the need to look for something faster

- Two possibilities:
  - ▶ Accelerating Gradient Decent Using Stochastic Gradient Descent!!!

# Therefore

Thus, the need to look for something faster

- Two possibilities:
  - ▶ Accelerating Gradient Descent Using Stochastic Gradient Descent!!!
  - ▶ Accelerating Gradient Descent Using The Best of Both World, Min-Batch!!!

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- **Robbins-Monro Theorem**
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

## We have that the Robbins-Monro Theorem[8]

The origins of such techniques are traced back to 1951

- When Robbins and Monro introduced the method of stochastic approximation
  - ▶ DARPA project!!!

## We have that the Robbins-Monro Theorem[8]

The origins of such techniques are traced back to 1951

- When Robbins and Monro introduced the method of stochastic approximation
  - ▶ DARPA project!!!

Setup, given a function  $M(\mathbf{w})$  and a constant  $\alpha$  such that the equation

$$M(\mathbf{w}) = \alpha$$

- It has a unique root  $\mathbf{w} = \mathbf{w}^*$

# Goal

We want to compute the root,  $w$ , of such equation

$$M(w^*) = \alpha$$

# Goal

We want to compute the root,  $w$ , of such equation

$$M(w^*) = \alpha$$

Then, we want to generate values  $w_1, w_2, \dots, w_{n-1}$  thus, we generate  $w_n$  from

- 1  $M(w_1), M(w_2), \dots, M(w_{n-1})$
- 2 and the possible derivatives  $M'(w_1), M'(w_2), \dots, M'(w_{n-1})$



# Goal

We want to compute the root,  $w$ , of such equation

$$M(w^*) = \alpha$$

Then, we want to generate values  $w_1, w_2, \dots, w_{n-1}$  thus, we generate  $w_n$  from

- 1  $M(w_1), M(w_2), \dots, M(w_{n-1})$
- 2 and the possible derivatives  $M'(w_1), M'(w_2), \dots, M'(w_{n-1})$

Thus, we would love that

$$\lim_{n \rightarrow \infty} w_n = w^*$$

Instead, we suppose that for each  $\boldsymbol{w}$  corresponds a Random Variable  $Y = Y(\boldsymbol{w})$

This Random Variable has a distribution function

$$Pr[Y(\boldsymbol{w}) \leq y] = H(y|\boldsymbol{w})$$

Instead, we suppose that for each  $\mathbf{w}$  corresponds a Random Variable  $Y = Y(\mathbf{w})$

This Random Variable has a distribution function

$$Pr[Y(\mathbf{w}) \leq y] = H(y|\mathbf{w})$$

Such that

$$M(\mathbf{w}) = \int_{-\infty}^{\infty} y dH(y|\mathbf{w})$$

# We Postulate

First a bound to the  $M(\mathbf{w})$

$$|M(\mathbf{w})| \leq C < \infty, \quad \int_{-\infty}^{\infty} (y - M(\mathbf{w}))^2 dH(y|\mathbf{w}) \leq \sigma^2 < \infty$$

# IMPORTANT

Neither the exact nature of  $H(y|\mathbf{w})$  nor that of  $M(\mathbf{w})$  is known

- But an important assumption is that

$$M(\mathbf{w}) - \alpha = 0$$

It has only one root

# IMPORTANT

Neither the exact nature of  $H(y|\mathbf{w})$  nor that of  $M(\mathbf{w})$  is known

- But an important assumption is that

$$M(\mathbf{w}) - \alpha = 0$$

It has only one root

Here is we use the  $\alpha$  value to generate the root by assuming

- $M(\mathbf{w}) - \alpha \leq 0$  for  $\mathbf{w} \leq \mathbf{w}^*$  and  $M(\mathbf{w}) - \alpha \geq 0$  for  $\mathbf{w} > \mathbf{w}^*$ .

Now, For a positive  $\delta$

$M(\mathbf{w})$  is strictly increasing if

$$\|\mathbf{w}^* - \mathbf{w}\| < \delta$$

Now, For a positive  $\delta$

$M(\mathbf{w})$  is strictly increasing if

$$\|\mathbf{w}^* - \mathbf{w}\| < \delta$$

And Finally

$$\inf_{\|\mathbf{w}^* - \mathbf{w}\| \geq \delta} |M(\mathbf{w}) - \alpha| > 0$$



Now choose a sequence  $\{\mu_i\}$

Such that

$$\sum_{i=1}^{\infty} \mu_i^2 = A < \infty \text{ and } \sum_{i=1}^{\infty} \mu_i = \infty$$

Now choose a sequence  $\{\mu_i\}$

Such that

$$\sum_{i=1}^{\infty} \mu_i^2 = A < \infty \text{ and } \sum_{i=1}^{\infty} \mu_i = \infty$$

Now, we define a non-stationary Markov Chain  $\{w_n\}$

$$w_{n+1} - w_n = \mu_n (\alpha - y_n)$$

Now choose a sequence  $\{\mu_i\}$

Such that

$$\sum_{i=1}^{\infty} \mu_i^2 = A < \infty \text{ and } \sum_{i=1}^{\infty} \mu_i = \infty$$

Now, we define a non-stationary Markov Chain  $\{\mathbf{w}_n\}$

$$\mathbf{w}_{n+1} - \mathbf{w}_n = \mu_n (\alpha - y_n)$$

Where  $y_n$  is a random variable such that

$$Pr[y_n \leq y | \mathbf{w}_n] = H(y | \mathbf{w}_n)$$

Using the expected value!!!

Here, we define  $b_n$

$$b_n = E [\mathbf{w}_n - \mathbf{w}^*]^2$$

## Using the expected value!!!

Here, we define  $b_n$

$$b_n = E[\mathbf{w}_n - \mathbf{w}^*]^2$$

We want conditions where this variance goes to zero

$$\lim_{n \rightarrow \infty} b_n = 0$$

- No matter what is the initial value  $\mathbf{w}_0$ .

We have then

Based on

$$\mathbf{w}_{n+1} - \mathbf{w}_n = \mu_n (\alpha - y_n)$$

We have then

Based on

$$\mathbf{w}_{n+1} - \mathbf{w}_n = \mu_n (\alpha - y_n)$$

We have then

$$\begin{aligned} b_{n+1} &= E [\mathbf{w}_{n+1} - \mathbf{w}^*]^2 = E [E [\mathbf{w}_{n+1} - \mathbf{w}^*]^2 | \mathbf{w}_n] \\ &= E \left[ \int_{-\infty}^{\infty} [\mathbf{w}_n - \mathbf{w}^* - \mu_n (y - \alpha)]^2 dH(y | \mathbf{w}_n) \right] \\ &= b_n + \mu_n E \left[ \int_{-\infty}^{\infty} (y - \alpha)^2 dH(y | \mathbf{w}_n) \right] - 2\mu_n E [(\mathbf{w}_n - \mathbf{w}^*) (M(\mathbf{w}_n) - \alpha)] \\ &= b_n + \mu_n^2 e_n - 2\mu_n d_n \end{aligned}$$

# With Values

We have

$$d_n = E [(\mathbf{w}_n - \mathbf{w}^*) (M(\mathbf{w}_n) - \alpha)]$$

$$e_n = E \left[ \int_{-\infty}^{\infty} (y - \alpha)^2 dH(y|\mathbf{w}_n) \right]$$



## With Values

We have

$$d_n = E[(\mathbf{w}_n - \mathbf{w}^*)(M(\mathbf{w}_n) - \alpha)]$$

$$e_n = E\left[\int_{-\infty}^{\infty} (y - \alpha)^2 dH(y|\mathbf{w}_n)\right]$$

From  $M(\mathbf{w}) \leq \alpha$  for  $\mathbf{w} \leq \mathbf{w}^*$  and  $M(\mathbf{w}) \geq \alpha$  for  $\mathbf{w} > \mathbf{w}^*$

$$d_n \geq 0$$

Additionally

Now, assuming that exist  $C$  such that

$$Pr [|Y(\mathbf{w})| \leq C] = \int_{-C}^C dH(y|\mathbf{w}) = 1 \quad \forall x$$

## Additionally

Now, assuming that exist  $C$  such that

$$Pr [|Y(\mathbf{w})| \leq C] = \int_{-C}^C dH(y|\mathbf{w}) = 1 \quad \forall x$$

We can prove that

$$0 \leq e_n \leq [C + |\alpha|^2] < \infty$$

Additionally

Now, assuming that exist  $C$  such that

$$Pr [|Y(\mathbf{w})| \leq C] = \int_{-C}^C dH(y|\mathbf{w}) = 1 \quad \forall x$$

We can prove that

$$0 \leq e_n \leq [C + |\alpha|^2] < \infty$$

Now, given

$$\sum_{i=1}^{\infty} \mu_i^2 = A < \infty \text{ and } \sum_{i=1}^{\infty} \mu_i = \infty$$

Therefore  $\sum_{i=1}^{\infty} \mu_i^2 e_i$  converges

Then, summing over  $i$  we obtain

$$b_{n+1} = b_1 + \sum_{i=1}^n \mu_i^2 e_i - 2 \sum_{i=1}^n \mu_i d_i$$

Therefore  $\sum_{i=1}^{\infty} \mu_i^2 e_i$  converges

Then, summing over  $i$  we obtain

$$b_{n+1} = b_1 + \sum_{i=1}^n \mu_i^2 e_i - 2 \sum_{i=1}^n \mu_i d_i$$

Since  $b_{n+1} \geq 0$

$$\sum_{i=1}^n \mu_i d_i \leq \frac{1}{2} \left[ b_1 + \sum_{i=1}^n \mu_i^2 e_i \right] < \infty$$

Then

Hence the positive-term series

$$\sum_{i=1}^{\infty} \mu_i d_i \text{ converges}$$

Then

Hence the positive-term series

$$\sum_{i=1}^{\infty} \mu_i d_i \text{ converges}$$

Then,  $\lim_{n \rightarrow \infty} b_n$  exists and...

$$\lim_{n \rightarrow \infty} b_n = b_1 + \sum_{i=1}^{\infty} \mu_i^2 e_i - 2 \sum_{i=1}^{\infty} \mu_i d_i = b$$



Therefore

If a sequence of  $\{k_i\}$  of non-negative constants such that

$$d_i \geq k_i b_i, \quad \sum_{i=1}^{\infty} \mu_i k_i = \infty$$

Therefore

If a sequence of  $\{k_i\}$  of non-negative constants such that

$$d_i \geq k_i b_i, \quad \sum_{i=1}^{\infty} \mu_i k_i = \infty$$

We want to prove that

$$\sum_{i=1}^{\infty} \mu_i k_i b_i < \infty$$

For this

We know that

$$\sum_{i=1}^{\infty} \mu_i d_i \text{ converges}$$

For this

We know that

$$\sum_{i=1}^{\infty} \mu_i d_i \text{ converges}$$

Therefore

$$k_i b_i \leq d_i \Rightarrow \mu_i k_i b_i \leq \mu_i d_i$$

Then

We have that

$$\sum_{i=1}^{\infty} \mu_i k_i b_i \leq \sum_{i=1}^{\infty} \mu_i d_i < \infty$$

Then

We have that

$$\sum_{i=1}^{\infty} \mu_i k_i b_i \leq \sum_{i=1}^{\infty} \mu_i d_i < \infty$$

Then, we have that

$$\sum_{i=1}^{\infty} \mu_i k_i b_i < \infty, \quad \sum_{i=1}^{\infty} \mu_i k_i = \infty$$

## Finally

For any  $\epsilon > 0$  there must be infinitely values  $i$  such that  $b_i < \epsilon$

- Therefore given that  $\lim_{n \rightarrow \infty} b_n = b$  then  $b = 0$ .

# Robbins and Monro Theorem (Original)

If  $\{\mu_n\}$  is of type  $\frac{1}{n}$

- Given a family of conditional probabilities

$$\{H(y|\mathbf{w}) = \Pr(Y(\mathbf{w}) \leq y|\mathbf{w})\}$$



# Robbins and Monro Theorem (Original)

If  $\{\mu_n\}$  is of type  $\frac{1}{n}$

- Given a family of conditional probabilities

$$\{H(y|\mathbf{w}) = \Pr(Y(\mathbf{w}) \leq y|\mathbf{w})\}$$

We have the following Expected Risk

$$M(\mathbf{w}) = \int_{-\infty}^{\infty} y dH(y|\mathbf{w})$$

Now

If we additionally have that

$$\Pr(|Y(\mathbf{w})| \leq C) = \int_{-C}^C dH(y|\mathbf{w}) = 1 \quad (3)$$

Then under the following constraints

For some  $\delta > 0$

$$M(w) \leq \alpha - \delta \text{ for } w < w^*$$

$$M(w) \geq \alpha + \delta \text{ for } w > w^* \quad (4)$$

Then under the following constraints

For some  $\delta > 0$

$$\begin{aligned}M(w) &\leq \alpha - \delta \text{ for } w < w^* \\ M(w) &\geq \alpha + \delta \text{ for } w > w^*\end{aligned}\tag{4}$$

Or Else

$$\begin{aligned}M(w) &< \alpha \text{ for } w < w^* \\ M(w^*) &= \alpha \\ M(w) &> \alpha \text{ for } w > w^*\end{aligned}\tag{5}$$

## Next

### Furthermore

$M(\boldsymbol{w})$  is strictly increasing if  $|\boldsymbol{w} - \boldsymbol{w}^*| < \delta$  (6)

## Next

### Furthermore

$M(\boldsymbol{w})$  is strictly increasing if  $|\boldsymbol{w} - \boldsymbol{w}^*| < \delta$  (6)

### And

$$\inf_{|\boldsymbol{w} - \boldsymbol{w}^*| \geq \delta} |M(\boldsymbol{w}) - \alpha| > 0 \quad (7)$$

## Next

Furthermore

$$M(\mathbf{w}) \text{ is strictly increasing if } |\mathbf{w} - \mathbf{w}^*| < \delta \quad (6)$$

And

$$\inf_{|\mathbf{w} - \mathbf{w}^*| \geq \delta} |M(\mathbf{w}) - \alpha| > 0 \quad (7)$$

And Let  $\{\mu_i\}$  be a sequence of positive numbers such that

$$\sum_{n=1}^{\infty} \mu_n = \infty \text{ and } \sum_{n=1}^{\infty} \mu_n^2 < \infty \quad (8)$$

Then

Let  $x_1$  an arbitrary number, then under the recursion

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu_n (\alpha - y_n)$$

- Where  $y_n \sim P(y|\mathbf{w}_n)$



Then

Let  $x_1$  an arbitrary number, then under the recursion

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu_n (\alpha - y_n)$$

- Where  $y_n \sim P(y|\mathbf{w}_n)$

## Theorem

- If (3) and (8), either (4) or (5,6,7) hold, then  $\mathbf{w}_n$  converges stochastically to  $\mathbf{w}^*$  given that  $b = 0$ .

# Recap of Robbins-Monro Proposal

Given the following function

$$f(\mathbf{w}) = E[\phi(\mathbf{w}, \eta)], \mathbf{w} \in \mathbb{R}^{d+1}$$

# Recap of Robbins-Monro Proposal

Given the following function

$$f(\mathbf{w}) = E[\phi(\mathbf{w}, \eta)], \mathbf{w} \in \mathbb{R}^{d+1}$$

Given a series of i.i.d. observations  $x_0, x_1, \dots$

- The following iterative procedure (Robbins-Monro Scheme)

$$\mathbf{w}_n = \mathbf{w}_{n-1} - \mu_n \phi(\mathbf{w}_{n-1}, \mathbf{x}_n)$$

# Robbins-Monro Proposal

Starting from an arbitrary initial condition,  $w_0$

- It converges to a root of  $M(w) = \alpha$

# Robbins-Monro Proposal

Starting from an arbitrary initial condition,  $w_0$

- It converges to a root of  $M(w) = \alpha$

Under some general conditions about the step size

$$\sum_{i=0}^{\infty} \mu_i^2 < \infty$$

$$\sum_{i=0}^{\infty} \mu_i \rightarrow \infty$$

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- **Robbins-Monro Scheme for Minimum-Square Error**
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

## Mean-Square Error [2]

### Cost function for MSE

$$J(\mathbf{w}) = E[\mathcal{L}(\mathbf{w}, \mathbf{x}, y)]$$

- Also known as the expected risk or the expected loss.

# Mean-Square Error [2]

## Cost function for MSE

$$J(\mathbf{w}) = E[\mathcal{L}(\mathbf{w}, \mathbf{x}, y)]$$

- Also known as the expected risk or the expected loss.

## Then, our objective is the reduction of the Expected Risk!!!

- Thus, the simple thing to do is to derive the function and make such gradient equal to zero.



Therefore

We can get the Gradient of the Expected Cost Function

$$\nabla J(\mathbf{w}) = E[\nabla \mathcal{L}(\mathbf{w}, \mathbf{x}, y)]$$

- where the expectation is w.r.t. the pair  $(\mathbf{x}, y)$

# Therefore

We can get the Gradient of the Expected Cost Function

$$\nabla J(\mathbf{w}) = E[\nabla \mathcal{L}(\mathbf{w}, \mathbf{x}, y)]$$

- where the expectation is w.r.t. the pair  $(\mathbf{x}, y)$

Therefore, everything depends on the form of the Loss function

$$\mathcal{L}_1(\mathbf{w}, \mathbf{x}, y) = \frac{1}{2} \left\| \mathbf{w}^T \mathbf{x} - y \right\|_2^2 \quad (\text{Least Squared Loss})$$

$$\mathcal{L}_2(\mathbf{w}, \mathbf{x}, y) = \left[ \frac{1}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}} \right]^{1-y} \left[ \frac{\exp\{\mathbf{w}^T \mathbf{x}\}}{1 + \exp\{\mathbf{w}^T \mathbf{x}\}} \right]^y \quad (\text{Logistic Loss})$$

$$\mathcal{L}_3(\mathbf{w}, \mathbf{x}, y) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log(y_{nk}^{(l)}) \quad (\text{Cross-Entropy Loss})$$

Therefore

We simply take  $\alpha = 0$  then

$$\nabla J(\mathbf{w}) = E[\nabla \mathcal{L}(\mathbf{w}, \mathbf{x}, y)] = 0$$

Therefore

We simply take  $\alpha = 0$  then

$$\nabla J(\mathbf{w}) = E[\nabla \mathcal{L}(\mathbf{w}, \mathbf{x}, y)] = 0$$

Then, we apply the Robbins-Monroe Schema to the function

$$f(\mathbf{w}) = \nabla J(\mathbf{w}) = 0$$

Then

Given the sequence of observations  $\{(\mathbf{x}_i, y_i)\}_{i=1,2,\dots}$  and values  $\{\mu_i\}_{i=1,2,\dots}$

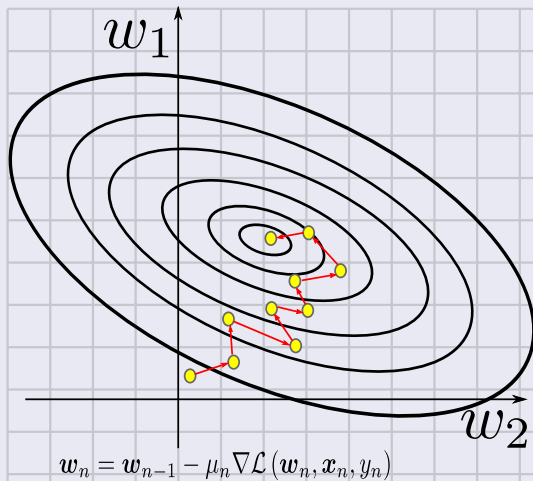
- We have that the iterative procedure becomes:

$$\mathbf{w}_n = \mathbf{w}_{n-1} - \mu_n \nabla \mathcal{L}(\mathbf{w}_n, \mathbf{x}_n, y_n)$$

- ▶ The Well known Vanilla Stochastic Gradient Descent (SGD)

# Geometrically

We have the following



# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- **Convergence**

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

Therefore

However, although the theorem is important

- it is not by itself enough.



# Therefore

However, although the theorem is important

- it is not by itself enough.

One has to know something more concerning

- The rate of convergence of such a scheme.

# Therefore

However, although the theorem is important

- it is not by itself enough.

One has to know something more concerning

- The rate of convergence of such a scheme.

It has been shown that

$$\mu_n = O\left(\frac{1}{n}\right)$$

## Additionally

Assuming that iterations have brought the estimate close to the optimal value

$$E(\mathbf{w}_n) = \mathbf{w}^* + \frac{1}{n}\mathbf{c}$$

## Additionally

Assuming that iterations have brought the estimate close to the optimal value

$$E(\mathbf{w}_n) = \mathbf{w}^* + \frac{1}{n}\mathbf{c}$$

And

$$Cov(\mathbf{w}_n) = \frac{1}{n}V + O\left(\frac{1}{n^2}\right)$$

- Where  $\mathbf{c}$  and  $V$  are constants that depend on the form of the expected risk.

# Meaning

## Therefore

- These formulas indicate that the parameter vector estimate fluctuates around the optimal value.

# Meaning

## Therefore

- These formulas indicate that the parameter vector estimate fluctuates around the optimal value.

## However

- Low complexity requirements makes this algorithmic family to be the one that is selected in a number of practical applications.
  - ▶ Given the problem with Batch Gradient Descent (BGD)

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

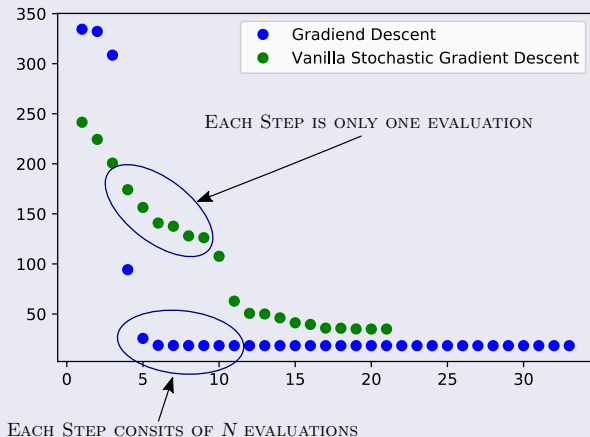
- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- **Example of SGD Vs BGD**
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

## Example of SGD for, $\frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x} - \mathbf{y})^2$

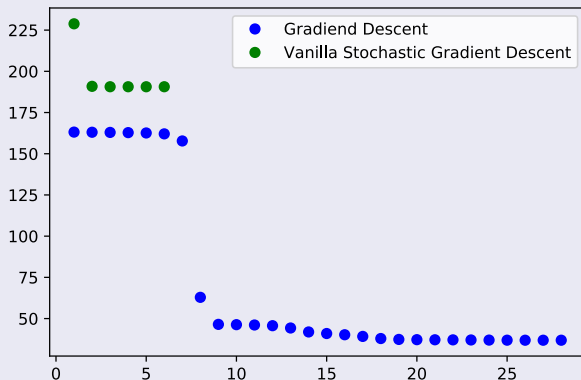
We can see how from the Vanilla SGD improves over the Batch GD with respect to Speed of Evaluation





# Problems

However, we need to improve such Vanilla Stochastic Gradient Descent



# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

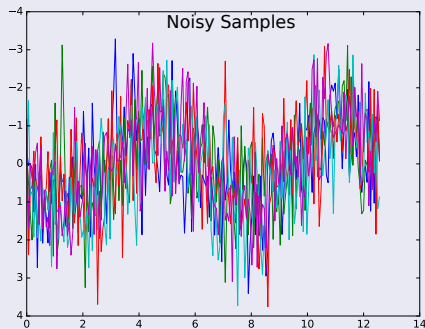
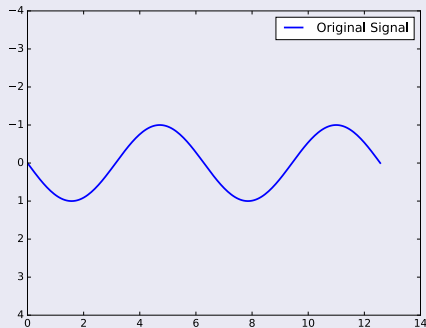
- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- **Using The Expected Value, The Mini-Batch**
- The Least-Mean Squares Adaptive Algorithm
- Conclusions

# Do you Remember?

Imagine the following signal from  $\sin(\theta)$



## What if we know the noise?

Given a series of observed samples  $\{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N\}$  with noise  $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{x}_i = x_i + \epsilon$$

## What if we know the noise?

Given a series of observed samples  $\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_N\}$  with noise  $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \epsilon$$

We can use our knowledge of probability to remove such noise

$$E[\hat{\mathbf{x}}_i] = E[\mathbf{x}_i + \epsilon] = E[\mathbf{x}_i] + E[\epsilon]$$

## What if we know the noise?

Given a series of observed samples  $\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_N\}$  with noise  $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \epsilon$$

We can use our knowledge of probability to remove such noise

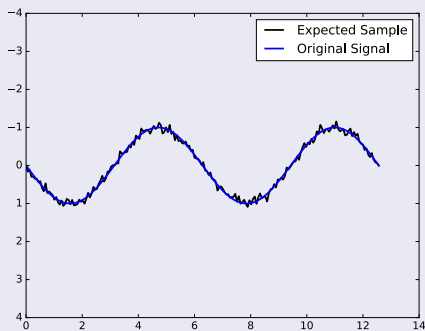
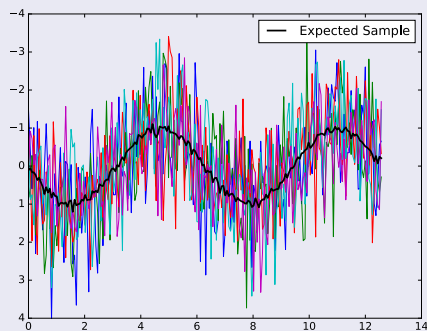
$$E[\hat{\mathbf{x}}_i] = E[\mathbf{x}_i + \epsilon] = E[\mathbf{x}_i] + E[\epsilon]$$

Then, because  $E[\epsilon] = 0$

$$E[\mathbf{x}_i] = E[\hat{\mathbf{x}}_i] \approx \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{x}}_i$$

# In our example

We have a nice result



Thus

Using a similar idea, you could use an average [9]

$$\nabla J(\mathbf{w}_{k-1} | \mathbf{x}_{i:i+m}, y_{i:i+m}) = \dots$$
$$\frac{1}{m} \sum_{i=1}^m \nabla J(\mathbf{w}_{k-1}, \mathbf{x}_i, y_i)$$



Thus

Using a similar idea, you could use an average [9]

$$\nabla J(\mathbf{w}_{k-1} | \mathbf{x}_{i:i+m}, y_{i:i+m}) = \dots$$
$$\frac{1}{m} \sum_{i=1}^m \nabla J(\mathbf{w}_{k-1}, \mathbf{x}_i, y_i)$$

This allows to reduce the variance of the original Stochastic Gradient

- It reduces the variance of the parameter updates, which can lead to more stable convergence.

Thus

Using a similar idea, you could use an average [9]

$$\nabla J(\mathbf{w}_{k-1} | \mathbf{x}_{i:i+m}, y_{i:i+m}) = \dots$$
$$\frac{1}{m} \sum_{i=1}^m \nabla J(\mathbf{w}_{k-1}, \mathbf{x}_i, y_i)$$

This allows to reduce the variance of the original Stochastic Gradient

- It reduces the variance of the parameter updates, which can lead to more stable convergence.
- It can make use of highly optimized matrix optimizations common to state-of-the-art deep learning libraries that make computing the gradient w.r.t. a mini-batch very efficient.

There are other more efficient options

We can update the  $w(k)$

- By Batches per epoch...

There are other more efficient options

We can update the  $w(k)$

- By Batches per epoch...

Therefore

- 1 for  $i$  in batch  $k$

$$w_k = w_{k-1} - \alpha \nabla J(w_{k-1}, x_i, y_i)$$

# Mini-batch gradient descent finally takes the best of both worlds

## Min-Batch( $X$ )

Input:

- Initialize  $w_0$  , Set number of epochs,  $L$ , Set learning rate  $\alpha$

- 1 for  $k = 1$  to  $L$ :
- 2     Randomly pick a mini batch of size  $m$ .
- 3     for  $i = 1$  to  $m$  do:
- 4         Evaluate  $g(k) = \nabla J(w_{k-1}, x_i, y_i)$
- 5          $w_k = w_{k-1} - \alpha g(k)$

Remark, for  $\alpha = \frac{1}{m}$ , the method is equivalent to average sample way

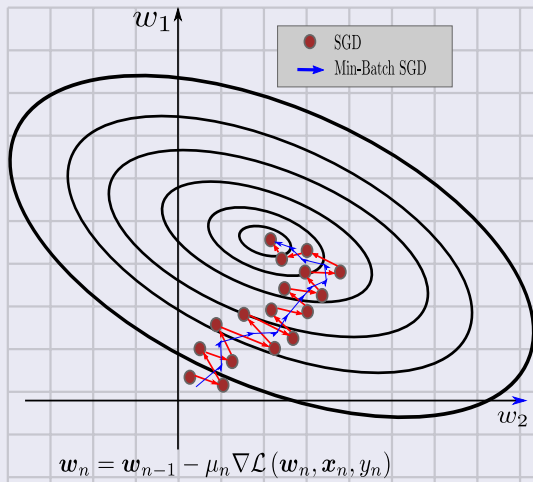
$$\begin{aligned} \mathbf{w}_k &= \mathbf{w}_{k-1} - \alpha \nabla J(\mathbf{w}_{k-1}, \mathbf{x}_i, y_i) - \dots \\ &\quad \alpha \nabla J(\mathbf{w}_{k-1}, \mathbf{x}_{i+1}, y_{i+1}) - \dots \\ &\quad \alpha \nabla J(\mathbf{w}_{k-1}, \mathbf{x}_{i+m}, y_{i+m}) \\ &= \mathbf{w}_{k-1} - \frac{1}{m} \sum_{i=1}^m \nabla J(\mathbf{w}_{k-1}, \mathbf{x}_i, y_i) \end{aligned}$$

## We have the following

- Common mini-batch sizes range between 50 and 256, but can vary for different applications.
- Mini-batch gradient descent is typically the algorithm of choice when training a neural network.

# A Small Intuition

We have smoother version of the Stochastic Gradient Descent





# Drawbacks

## Choosing a proper learning rate can be difficult

- A learning rate that is too small leads to painfully slow convergence,
- Too large can hinder convergence and cause the loss function to fluctuate around the minimum or even to diverge.

# Drawbacks

## Choosing a proper learning rate can be difficult

- A learning rate that is too small leads to painfully slow convergence,
- Too large can hinder convergence and cause the loss function to fluctuate around the minimum or even to diverge.

## Learning Rate Schedules

- To adjust the learning rate during training by e.g. annealing
- These schedules and thresholds, however, have to be defined in advance not on-line

# Drawbacks

## Choosing a proper learning rate can be difficult

- A learning rate that is too small leads to painfully slow convergence,
- Too large can hinder convergence and cause the loss function to fluctuate around the minimum or even to diverge.

## Learning Rate Schedules

- To adjust the learning rate during training by e.g. annealing
- These schedules and thresholds, however, have to be defined in advance not on-line

## Another key challenge of minimizing highly non-convex error functions

- For example, neural networks, it is avoiding getting trapped in their numerous suboptimal local minima.

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- **The Least-Mean Squares Adaptive Algorithm**
- Conclusions

# The MSE Linear Estimation, the Normal Equations

It was proved in slide set 2

- The optimal **Mean-Square Error estimate** of  $y$  given the value  $X = \mathbf{x}$  is

$$E[y|\mathbf{x}] = \hat{y}$$

- ▶ In general, a nonlinear function.

# The MSE Linear Estimation, the Normal Equations

It was proved in slide set 2

- The optimal **Mean-Square Error estimate** of  $y$  given the value  $X = \mathbf{x}$  is

$$E[y|\mathbf{x}] = \hat{y}$$

- ▶ In general, a nonlinear function.

For Linear Estimators, in  $(\mathbf{x}, y) \in \mathbb{R}^d \times \mathbb{R}$  joint distributed random variables of zero mean values

- Our goal is to obtain an estimate of  $\mathbf{w} \in \mathbb{R}^d$  (Our Unknown  $\theta$ ) in the linear estimator model

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

Thus, using MSE as the Cost Equation

### Cost Function

$$J(\mathbf{w}) = E \left[ (y - \hat{y})^2 \right]$$

Thus, using MSE as the Cost Equation

Cost Function

$$J(\mathbf{w}) = E \left[ (y - \hat{y})^2 \right]$$

Thus, we are looking for an estimator that minimize the variance of the error

$$\epsilon = y - \hat{y}$$



Thus, using MSE as the Cost Equation

### Cost Function

$$J(\mathbf{w}) = E \left[ (y - \hat{y})^2 \right]$$

Thus, we are looking for an estimator that minimize the variance of the error

$$\epsilon = y - \hat{y}$$

We want to **Minimize** the cost function  $J(\mathbf{w})$  by finding an optimal  $\mathbf{w}^*$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$$

Then, we can simply use  $\nabla J(\mathbf{w}) = 0$

We have

$$\begin{aligned}\nabla J(\mathbf{w}) &= \nabla E \left[ \left( y - \mathbf{w}^T \mathbf{x} \right)^2 \right] \\ &= \nabla E \left[ \left( y - \mathbf{w}^T \mathbf{x} \right) \left( y - \mathbf{w}^T \mathbf{x} \right) \right] \\ &= \nabla \left\{ E \left[ y^2 \right] - 2\mathbf{w}^T E \left[ \mathbf{x} y \right] + \mathbf{w}^T E \left[ \mathbf{x} \mathbf{x}^T \right] \mathbf{w} \right\} \\ &= -2\mathbf{p} + 2\Sigma_x \mathbf{w} = 0\end{aligned}$$

Where, we have

$$\begin{aligned}\mathbf{p} &= [E[yx_1], E[yx_2], \dots, E[yx_d]] = E[\mathbf{x}y] \\ \Sigma_x &= E[\mathbf{x}\mathbf{x}^T]\end{aligned}$$

This generates what is known as

Then, we get the Normal Equations

$$\Sigma_x \mathbf{w}^* = \mathbf{p}$$

# The Stochastic Gradient Descent

Imagine the follow

- We assume that the covariance matrix and the cross-correlation vector are unknown.

# The Stochastic Gradient Descent

Imagine the follow

- We assume that the covariance matrix and the cross-correlation vector are unknown.

We have that for a single sample

$$\mathcal{L}(\mathbf{w}, y, \mathbf{x}) = \frac{1}{2} \left( \mathbf{w}^T \mathbf{x} - y \right)^2$$

# Therefore

We know

- The solution corresponds to the root of the gradient of the cost function:

$$\Sigma_x \mathbf{w} - \mathbf{p} = E \left[ \mathbf{x} \left( \mathbf{x}^T \mathbf{w} - y \right) \right] = 0$$

# Therefore

We know

- The solution corresponds to the root of the gradient of the cost function:

$$\Sigma_x \mathbf{w} - \mathbf{p} = E \left[ \mathbf{x} \left( \mathbf{x}^T \mathbf{w} - y \right) \right] = 0$$

We have

$$\nabla J(\mathbf{w}) = \Sigma_x \mathbf{w} - \mathbf{p} = E \left[ \mathbf{x} \left( \mathbf{x}^T \mathbf{w} - y \right) \right] = 0$$

# Therefore

## We know

- The solution corresponds to the root of the gradient of the cost function:

$$\Sigma_x \mathbf{w} - \mathbf{p} = E \left[ \mathbf{x} \left( \mathbf{x}^T \mathbf{w} - y \right) \right] = 0$$

## We have

$$\nabla J(\mathbf{w}) = \Sigma_x \mathbf{w} - \mathbf{p} = E \left[ \mathbf{x} \left( \mathbf{x}^T \mathbf{w} - y \right) \right] = 0$$

## Then

$$\mathbf{w}_n = \mathbf{w}_{n-1} + \mu_n \mathbf{x}_n \left( \mathbf{x}_n^T \mathbf{w}_{n-1} - y_n \right)$$



# The Least-Mean Squares Adaptive Algorithm

## The stochastic gradient algorithm for MSE

- It converges to the optimal mean-square error solution provided that  $\mu_n$  satisfies the two convergence conditions.

# The Least-Mean Squares Adaptive Algorithm

## The stochastic gradient algorithm for MSE

- It converges to the optimal mean-square error solution provided that  $\mu_n$  satisfies the two convergence conditions.

## Once the algorithm has converged

- It “locks” at the obtained solution.

# The Least-Mean Squares Adaptive Algorithm

## The stochastic gradient algorithm for MSE

- It converges to the optimal mean-square error solution provided that  $\mu_n$  satisfies the two convergence conditions.

## Once the algorithm has converged

- It “locks” at the obtained solution.

## In a case where the statistics of the involved process changes

- The algorithm cannot track the changes.

Therefore

if such changes occur, the error term

$$e_n = y_n - \mathbf{x}_n^T \mathbf{w}_{n-1}$$

- It will get larger values.

# Therefore

if such changes occur, the error term

$$e_n = y_n - \mathbf{x}_n^T \mathbf{w}_{n-1}$$

- It will get larger values.

# However

- Because  $\mu_n$  is very small, the increased value of the error will not lead to corresponding changes of the estimate at time  $n$ .

## Solution

This can be overcome if one sets the value of  $\mu_n$

- To a preselected fixed value,  $\mu$ .

# Solution

This can be overcome if one sets the value of  $\mu_n$

- To a preselected fixed value,  $\mu$ .

## The celebrated Least-Mean-Squares Algorithms

- Algorithm LMS

- 1  $\mathbf{w}_{-1} = \mathbf{0} \in \mathbb{R}^d$
- 2 Select a value  $\mu$
- 3 **for**  $n = 0, 1, \dots$  **do**
- 4      $e_n = y_n - \mathbf{x}_n^T \mathbf{w}_{n-1}$
- 5      $\mathbf{w}_n = \mathbf{w}_{n-1} + \mu e_n \mathbf{x}_n$

# Complexity

## Something Notable

- The complexity of the algorithm amounts to  $2d$  multiplications/additions (MADs) per time update.



# Complexity

## Something Notable

- The complexity of the algorithm amounts to  $2d$  multiplications/additions (MADs) per time update.

## However

- As the algorithm converges close the solution

# Complexity

## Something Notable

- The complexity of the algorithm amounts to  $2d$  multiplications/additions (MADs) per time update.

## However

- As the algorithm converges close the solution

## Thus

- The error term is expected to take small values making the updates to remain close the solution

# Important

Given that  $\mu$  has a constant value

- The algorithm has now the “agility” to update the estimates
  - ▶ In an attempt to “push” the error to lower values.

# Important

## Given that $\mu$ has a constant value

- The algorithm has now the “agility” to update the estimates
  - ▶ In an attempt to “push” the error to lower values.

## Something Notable

- This small variation of the iterative scheme has important implications.

# Important

## Given that $\mu$ has a constant value

- The algorithm has now the “agility” to update the estimates
  - ▶ In an attempt to “push” the error to lower values.

## Something Notable

- This small variation of the iterative scheme has important implications.

## No More a Robbins-Monro stochastic family

- The resulting algorithm is no more a member of the Robbins-Monro stochastic approximation family.

# Outline

## 1. Introduction

- Review Gradient Descent
- The Problems of Gradient Descent with Large Data Sets
- Convergence of gradient descent with fixed step size
- Convergence Rate
  - Convex Functions
  - Back to the Main Problem
- Accelerating the Gradient Descent
- Even with such Speeds

## 2. Accelerating Gradient Descent

- Robbins-Monro Theorem
- Robbins-Monro Scheme for Minimum-Square Error
- Convergence

## 3. Improving and Measuring Stochastic Gradient Descent

- Example of SGD Vs BGD
- Using The Expected Value, The Mini-Batch
- The Least-Mean Squares Adaptive Algorithm
- **Conclusions**

# Conclusions

## In Machine Learning

- We need to have the best speedups to handle the problem dealing with Big Data...

# Conclusions

## In Machine Learning

- We need to have the best speedups to handle the problem dealing with Big Data...

## As we get more and more algorithms

- It is clear that optimization for Big Data is one of the hottest trends in Machine Learning



# Conclusions

## In Machine Learning

- We need to have the best speedups to handle the problem dealing with Big Data...

## As we get more and more algorithms

- It is clear that optimization for Big Data is one of the hottest trends in Machine Learning

## Take a look to

- Leon Bottou, Frank E. Curtis, Jorge Nocedal: **Optimization Methods for Large-Scale Machine Learning**. SIAM Review 60(2): 223-311 (2018)



C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*.

Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.



S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*.

Academic Press, 1st ed., 2015.



W. Rudin, *Real and Complex Analysis, 3rd Ed.*

New York, NY, USA: McGraw-Hill, Inc., 1987.







T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms, Third Edition*.

The MIT Press, 3rd ed., 2009.



L. Lessard, B. Recht, and A. Packard, “Analysis and design of optimization algorithms via integral quadratic constraints,” *SIAM Journal on Optimization*, vol. 26, no. 1, pp. 57–95, 2016.

-  Y. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Publishing Company, Incorporated, 1 ed., 2014.
-  S. Bubeck, “Convex optimization: Algorithms and complexity,” *arXiv preprint arXiv:1405.4980*, 2014.
-  H. Robbins and S. Monro, “A stochastic approximation method,” *Ann. Math. Statist.*, vol. 22, pp. 400–407, 09 1951.
-  M. Li, T. Zhang, Y. Chen, and A. J. Smola, “Efficient mini-batch training for stochastic optimization,” in *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, (New York, NY, USA), pp. 661–670, ACM, 2014.