



10-301/601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Linear Regression + Optimization for ML

Matt Gormley & Henry Chai

Lecture 8

Sep. 22, 2021

Reminders

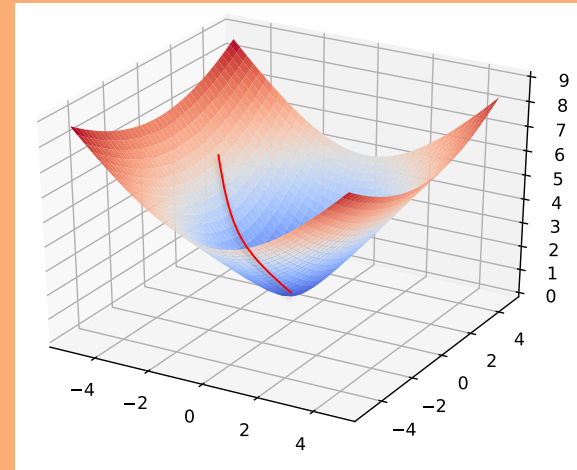
- Homework 3: KNN, Perceptron, Lin.Reg.
 - Out: Mon, Sep. 20
 - Due: Sun, Sep. 26 at 11:59pm
 - **IMPORTANT: you may only use 2 grace days on Homework 3**

OPTIMIZATION METHOD #1: GRADIENT DESCENT

Gradient Descent

Algorithm 1 Gradient Descent

```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$ 
5:   return  $\theta$ 
```



There are many possible ways to detect **convergence**. For example, we could check whether the L2 norm of the gradient is below some small tolerance.

$$\|\nabla_{\theta} J(\theta)\|_2 \leq \epsilon$$

Alternatively we could check that the reduction in the objective function from one iteration to the next is small.

GRADIENT DESCENT FOR LINEAR REGRESSION

Linear Regression as Function Approximation

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$$

where $\mathbf{x} \in \mathbb{R}^M$ and $y \in \mathbb{R}$

1. Assume \mathcal{D} generated as:

$$\begin{aligned}\mathbf{x}^{(i)} &\sim p^*(\cdot) \\ y^{(i)} &= h^*(\mathbf{x}^{(i)})\end{aligned}$$

2. Choose hypothesis space, \mathcal{H} :
all linear functions in M -dimensional space

$$\mathcal{H} = \{h_{\boldsymbol{\theta}} : h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}, \boldsymbol{\theta} \in \mathbb{R}^M\}$$

3. Choose an objective function:
mean squared error (MSE)

$$\begin{aligned}J(\boldsymbol{\theta}) &= \frac{1}{N} \sum_{i=1}^N e_i^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})\right)^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)}\right)^2\end{aligned}$$

4. Solve the unconstrained optimization problem via favorite method:

- gradient descent
- closed form
- stochastic gradient descent
- ...

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$$

5. Test time: given a new \mathbf{x} , make prediction \hat{y}

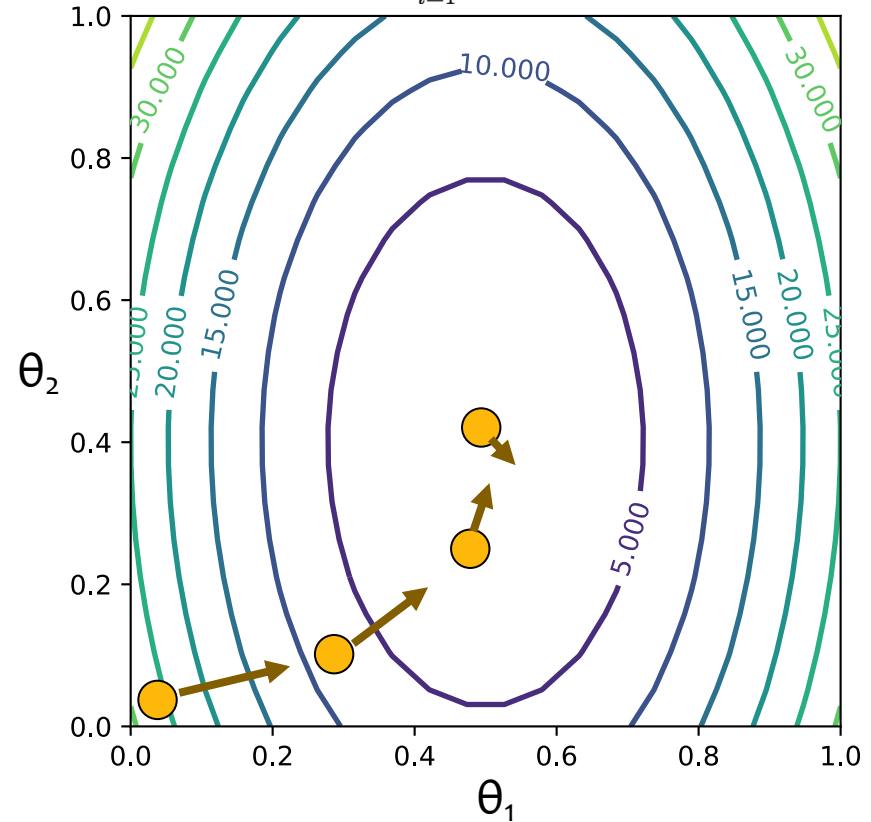
$$\hat{y} = h_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}^T \mathbf{x}$$

Linear Regression by Gradient Desc.

Optimization Method #1: Gradient Descent

1. Pick a random θ
2. Repeat:
 - a. Evaluate gradient $\nabla J(\theta)$
 - b. Step opposite gradient
3. Return θ that gives smallest $J(\theta)$

$$J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$$



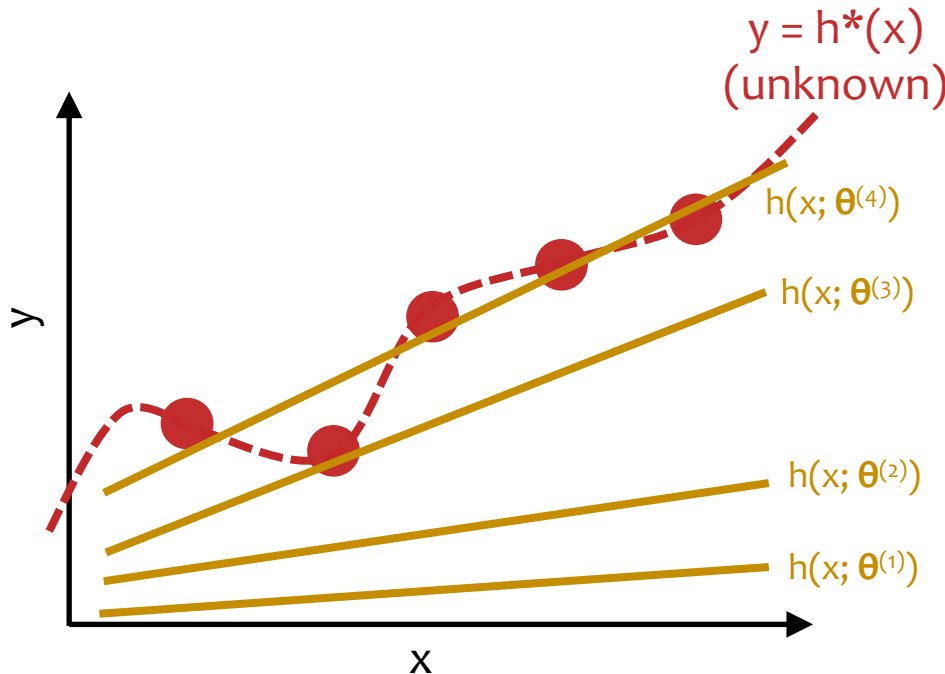
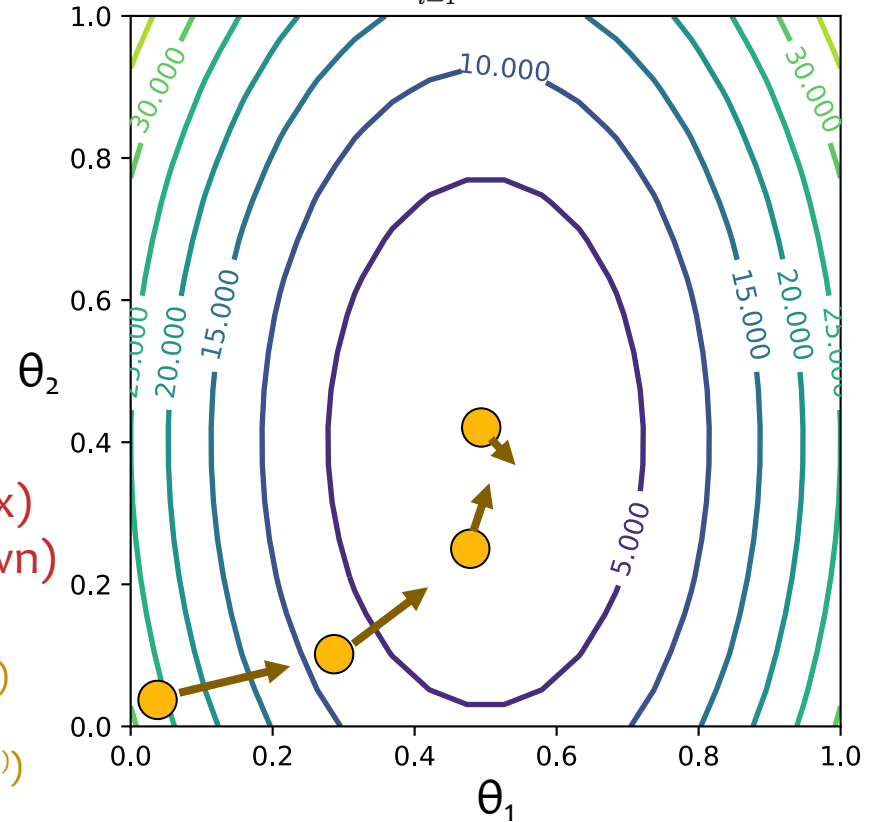
t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

Linear Regression by Gradient Desc.

$$J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$$

Optimization Method #1: Gradient Descent

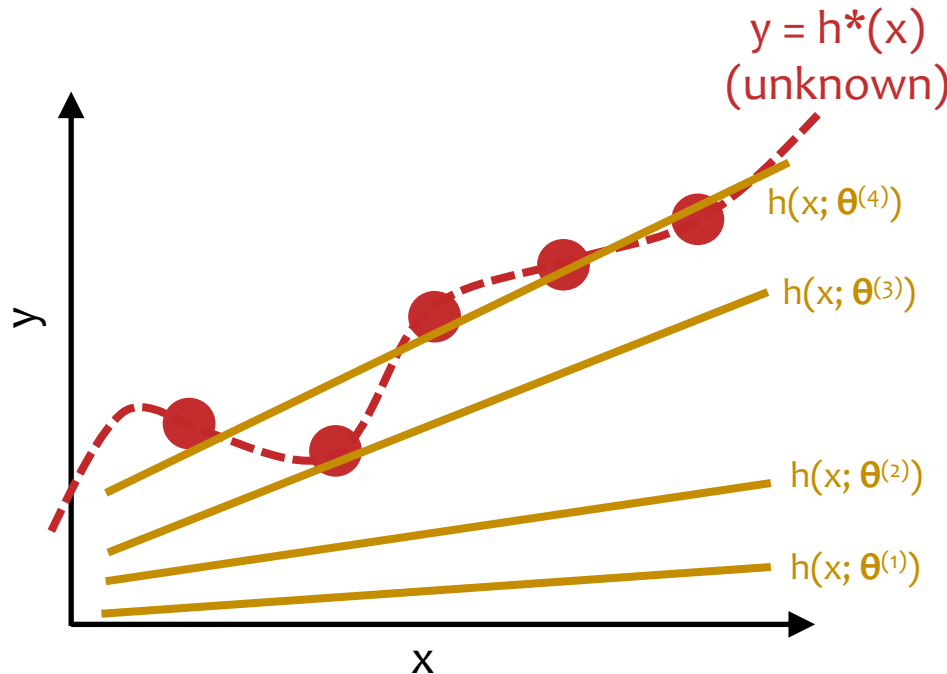
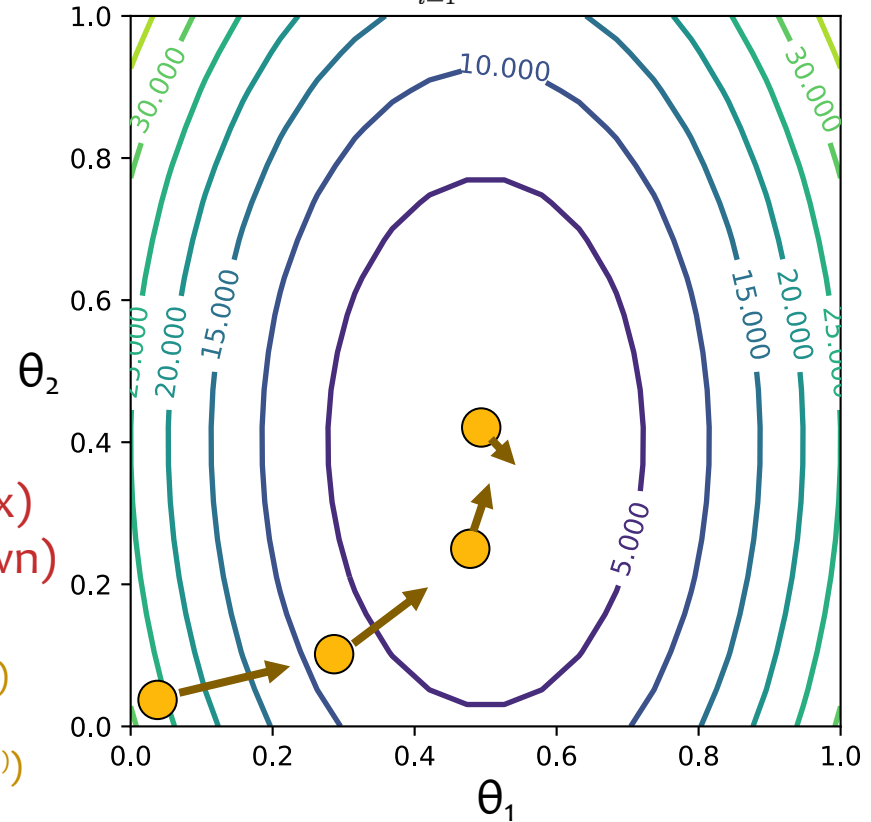
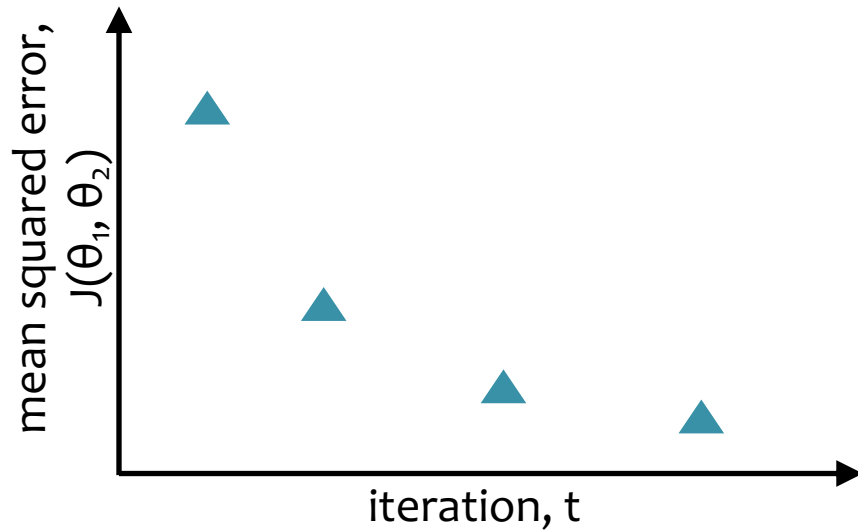
1. Pick a random θ
2. Repeat:
 - a. Evaluate gradient $\nabla J(\theta)$
 - b. Step opposite gradient
3. Return θ that gives smallest $J(\theta)$



t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

Linear Regression by Gradient Desc.

$$J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T x^{(i)})^2$$



t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

Optimization for Linear Regression

Chalkboard

- Computing the gradient for Linear Regression
- Gradient Descent for Linear Regression

GD for Linear Regression

Gradient Descent for Linear Regression repeatedly takes steps opposite the gradient of the objective function

Algorithm 1 GD for Linear Regression

```
1: procedure GDLR( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$                                 ▷ Initialize parameters
3:   while not converged do
4:      $\mathbf{g} \leftarrow \sum_{i=1}^N (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$     ▷ Compute gradient
5:      $\theta \leftarrow \theta - \gamma \mathbf{g}$                                 ▷ Update parameters
6:   return  $\theta$ 
```

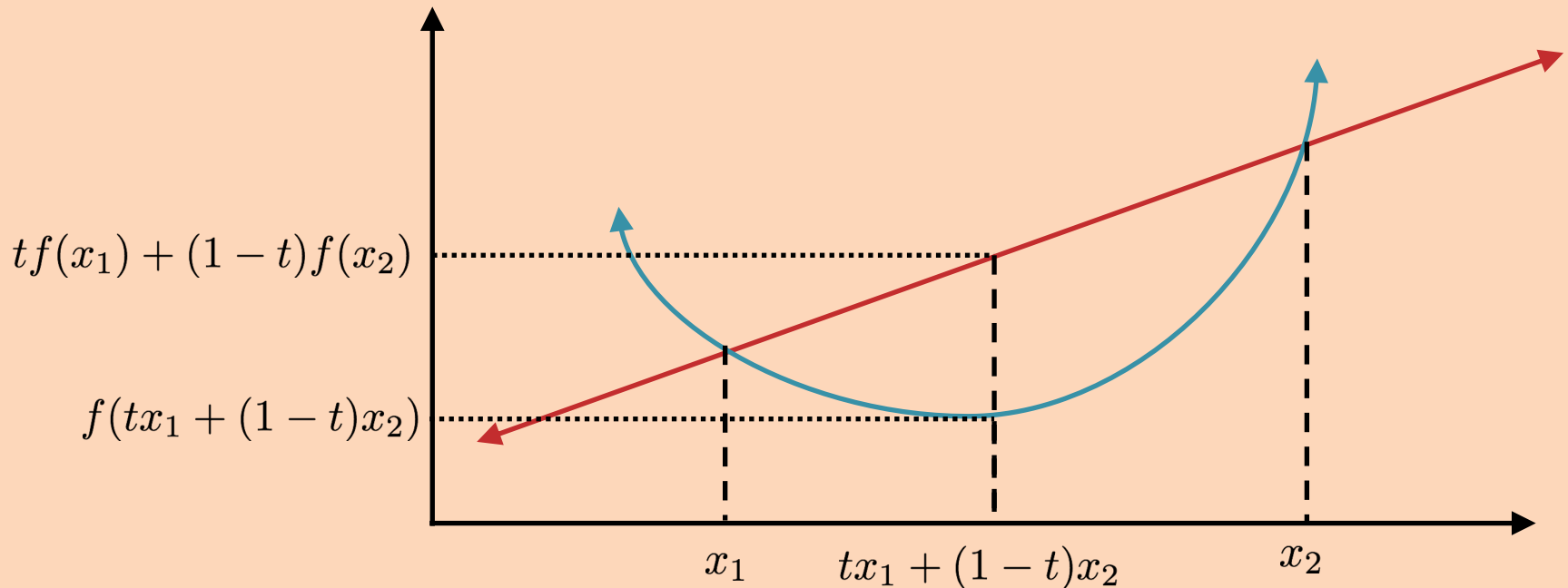
CONVEXITY

Convexity

Function $f : \mathbb{R}^M \rightarrow \mathbb{R}$ is **convex**

if $\forall \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \leq t \leq 1$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$

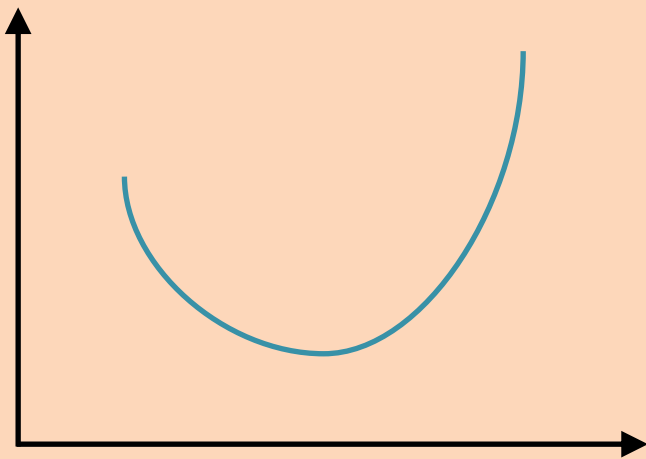


Convexity

Suppose we have a function $f(x) : \mathcal{X} \rightarrow \mathcal{Y}$.

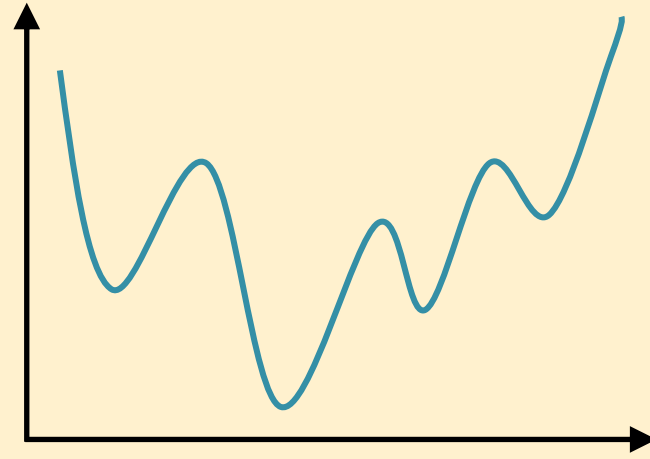
- The value x^* is a **global minimum** of f iff $f(x^*) \leq f(x), \forall x \in \mathcal{X}$.
- The value x^* is a **local minimum** of f iff $\exists \epsilon$ s.t. $f(x^*) \leq f(x), \forall x \in [x^* - \epsilon, x^* + \epsilon]$.

Convex Function



- Each **local minimum** is a **global minimum**

Nonconvex Function

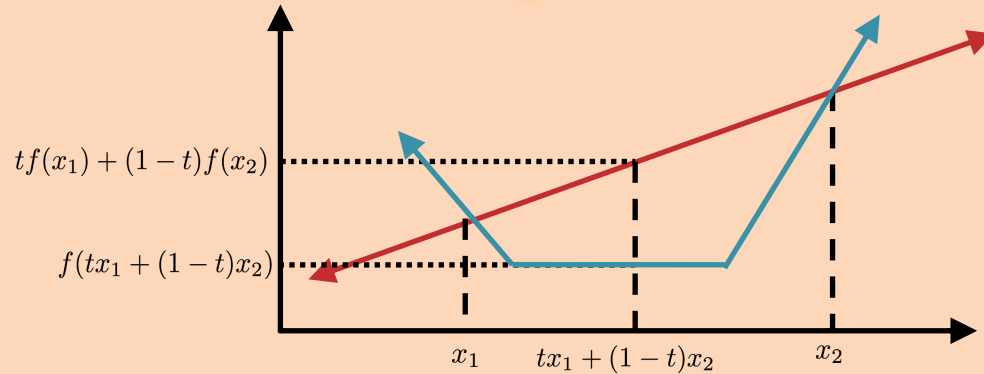


- A **nonconvex** function is **not convex**
- Each **local minimum** is **not necessarily a global minimum**

Convexity

Function $f : \mathbb{R}^M \rightarrow \mathbb{R}$ is **convex**
if $\forall \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \leq t \leq 1$:

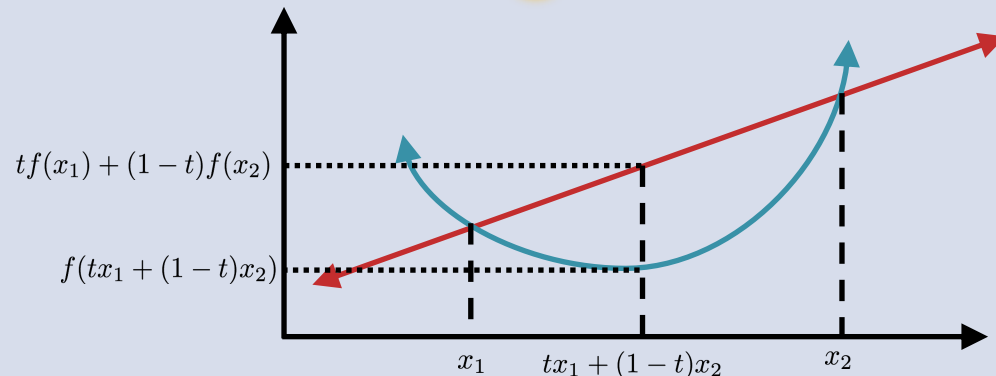
$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



Each **local**
minimum of a
convex function is
also a **global**
minimum.

Function $f : \mathbb{R}^M \rightarrow \mathbb{R}$ is **strictly convex**
if $\forall \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \leq t \leq 1$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) < tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



A **strictly convex**
function has a
unique global
minimum.

CONVEXITY AND LINEAR REGRESSION

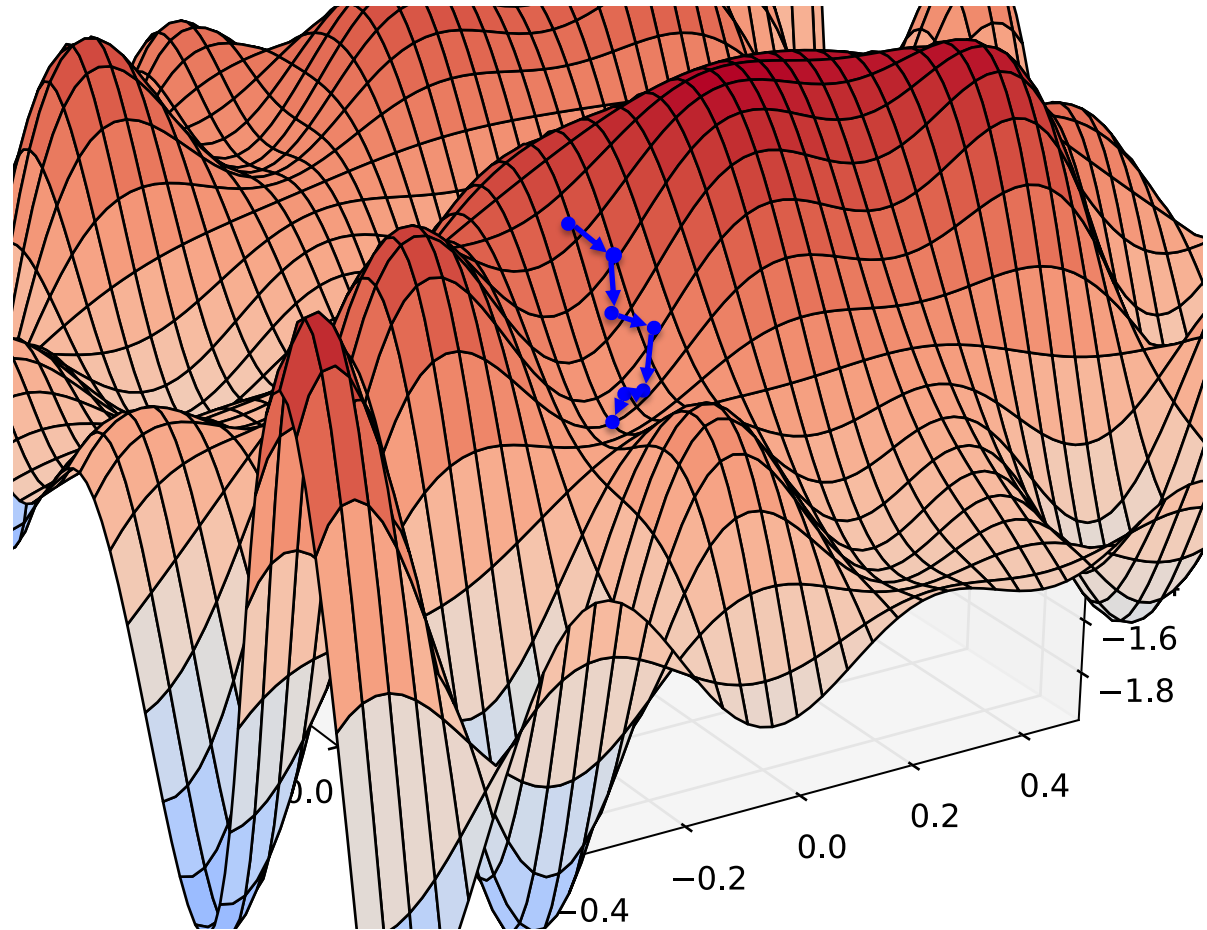
Convexity and Linear Regression

The **Mean Squared Error** function, which we minimize for learning the parameters of Linear Regression, **is convex!**

... but in the general case it is **not strictly convex.**

Gradient Descent & Convexity

- Gradient descent is a **local optimization algorithm**
- If the function is **nonconvex**, it will find a local minimum, not necessarily a global minimum
- If the function is **convex**, it will find a global minimum



Regression Loss Functions

In-Class Exercise:

Which of the following could be used as loss functions for training a linear regression model?

Select all that apply.

A. $\ell(\hat{y}, y) = \|\hat{y} - y\|_2$

B. $\ell(\hat{y}, y) = |\hat{y} - y|$

C. $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

D. $\ell(\hat{y}, y) = \frac{1}{4}(\hat{y} - y)^4$

E. $\ell(\hat{y}, y) = \begin{cases} \frac{1}{2}(\hat{y} - y)^2 & \text{if } |\hat{y} - y| \leq \delta \\ \delta|\hat{y} - y| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}$

F. $\ell(\hat{y}, y) = \log(\cosh(\hat{y} - y))$

Question 1

A

B

C

D

E

F

OPTIMIZATION METHOD #2: CLOSED FORM SOLUTION

Calculus and Optimization

In-Class Exercise

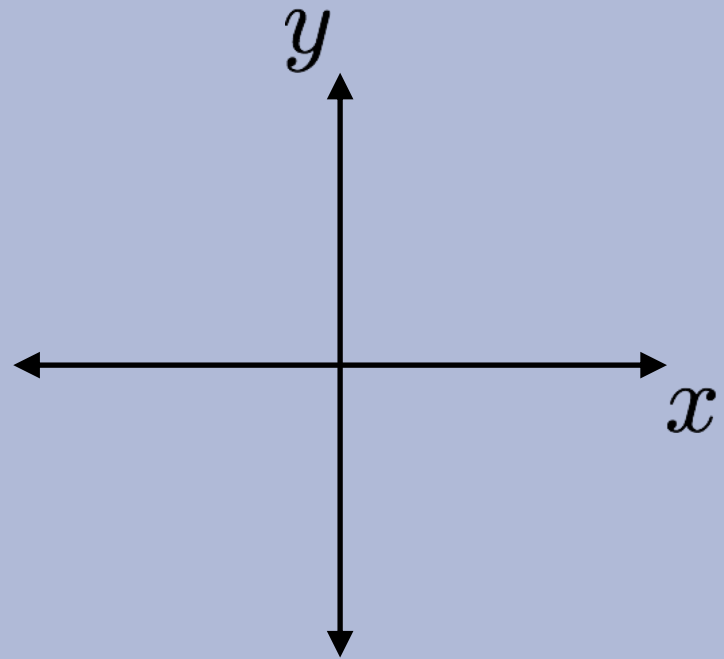
Plot three functions:

1. $f(x) = x^3 - x$

2. $f'(x) = \frac{\partial y}{\partial x}$

3. $f''(x) = \frac{\partial^2 y}{\partial x^2}$

Answer Here:



Optimization: Closed form solutions

Chalkboard

- Zero Derivatives
- Example: 1-D function
- Example: higher dimensions

CLOSED FORM SOLUTION FOR LINEAR REGRESSION

Linear Regression as Function Approximation

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$$

where $\mathbf{x} \in \mathbb{R}^M$ and $y \in \mathbb{R}$

1. Assume \mathcal{D} generated as:

$$\begin{aligned}\mathbf{x}^{(i)} &\sim p^*(\cdot) \\ y^{(i)} &= h^*(\mathbf{x}^{(i)})\end{aligned}$$

2. Choose hypothesis space, \mathcal{H} :
all linear functions in M -dimensional space

$$\mathcal{H} = \{h_{\boldsymbol{\theta}} : h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}, \boldsymbol{\theta} \in \mathbb{R}^M\}$$

3. Choose an objective function:
mean squared error (MSE)

$$\begin{aligned}J(\boldsymbol{\theta}) &= \frac{1}{N} \sum_{i=1}^N e_i^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})\right)^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)}\right)^2\end{aligned}$$

4. Solve the unconstrained optimization problem via favorite method:

- gradient descent
- closed form
- stochastic gradient descent
- ...

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$$

5. Test time: given a new \mathbf{x} , make prediction \hat{y}

$$\hat{y} = h_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}^T \mathbf{x}$$

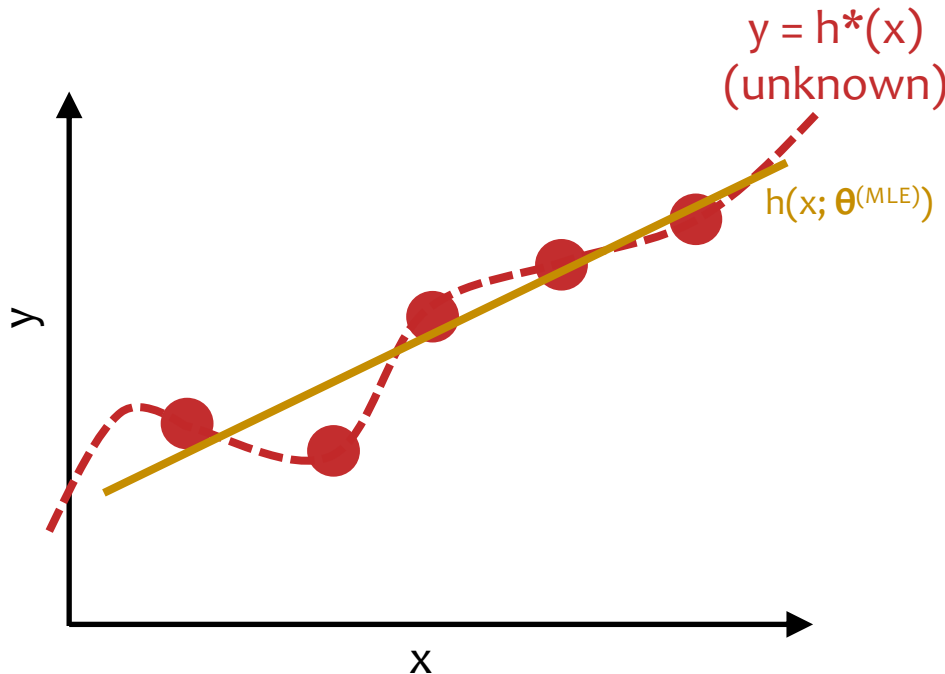
Linear Regression: Closed Form

Optimization Method #2: Closed Form

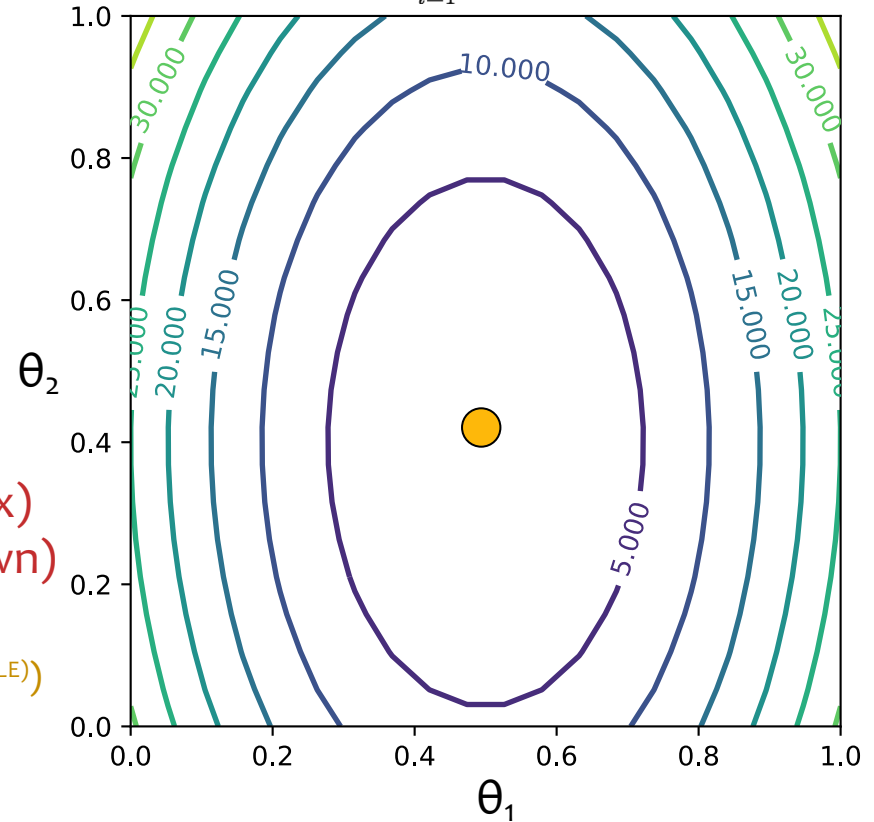
1. Evaluate

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

2. Return θ^{MLE}



$$J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2$$



t	θ_1	θ_2	$J(\theta_1, \theta_2)$
MLE	0.59	0.43	0.2

Optimization for Linear Regression

Chalkboard

- Closed-form (Normal Equations)

COMPUTATIONAL COMPLEXITY

Computational Complexity of OLS

To solve the Ordinary Least Squares problem we compute:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - (\theta^T \mathbf{x}^{(i)}))^2$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{Y})$$

The resulting shape of the matrices:

$$\underbrace{\left(\begin{array}{cc} \mathbf{X}^T & \mathbf{X} \\ \hline M \times N & N \times M \end{array} \right)^{-1}}_{M \times M} \underbrace{\left(\begin{array}{cc} \mathbf{X}^T & \mathbf{Y} \\ \hline M \times N & N \times 1 \end{array} \right)}_{M \times 1}$$

Background: Matrix Multiplication Given matrices **A** and **B**

- If **A** is $q \times r$ and **B** is $r \times s$, computing **AB** takes $O(qrs)$
- If **A** and **B** are $q \times q$, computing **AB** takes $O(q^{2.373})$
- If **A** is $q \times q$, computing A^{-1} takes $O(q^{2.373})$.

Computational Complexity of OLS:

$\mathbf{X}^T \mathbf{X}$	$O(M^2 N)$
$(\quad)^{-1}$	$O(M^{2.373})$
$\mathbf{X}^T \mathbf{Y}$	$O(MN)$
$(\quad)^{-1}(\quad)$	$O(M^2)$
total	$O(M^2 N + M^{2.373})$

Linear in # of examples, N
Polynomial in # of features, M

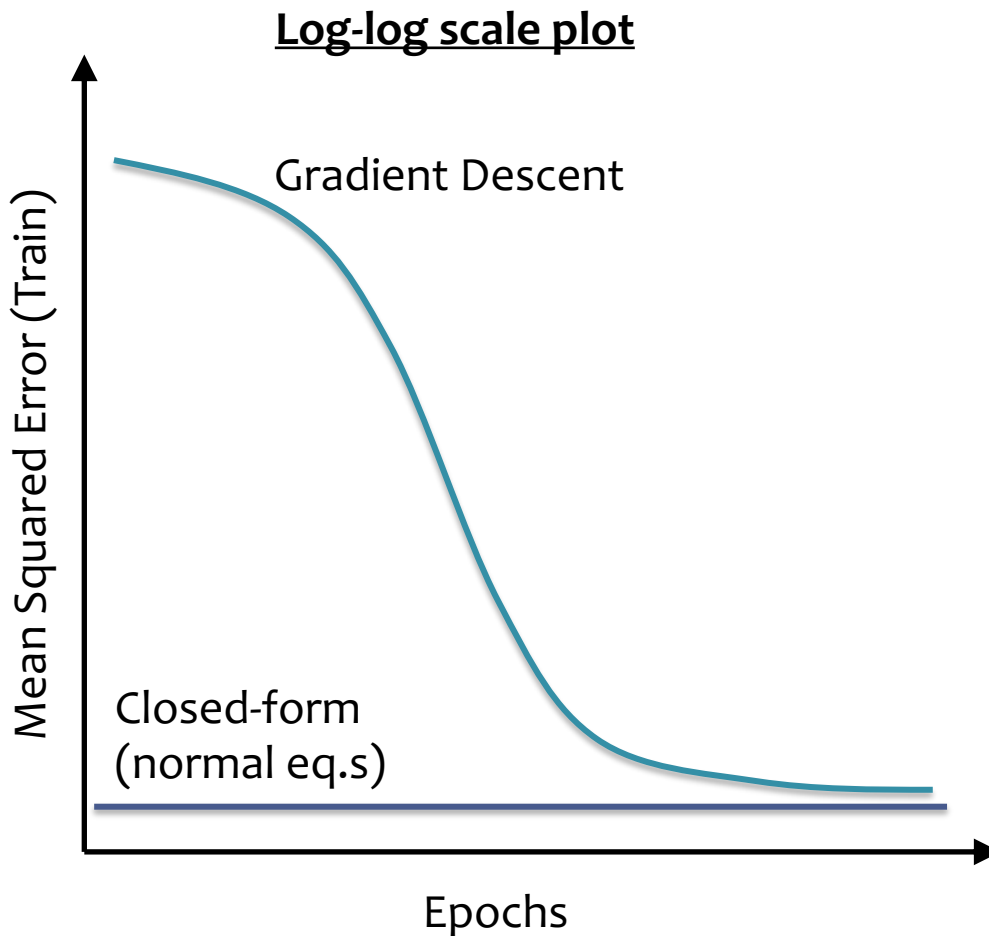


Gradient Descent

Cases to consider gradient descent:

1. What if we **can not** find a closed-form solution?
2. What if we **can**, but it's inefficient to compute?
3. What if we **can**, but it's numerically unstable to compute?

Empirical Convergence



- Def: an **epoch** is a single pass through the training data
- 1. For GD, only **one update** per epoch
- 2. For SGD, **N updates** per epoch
 $N = (\# \text{ train examples})$

- SGD reduces MSE much more rapidly than GD
- For GD / SGD, training MSE is initially large due to uninformed initialization

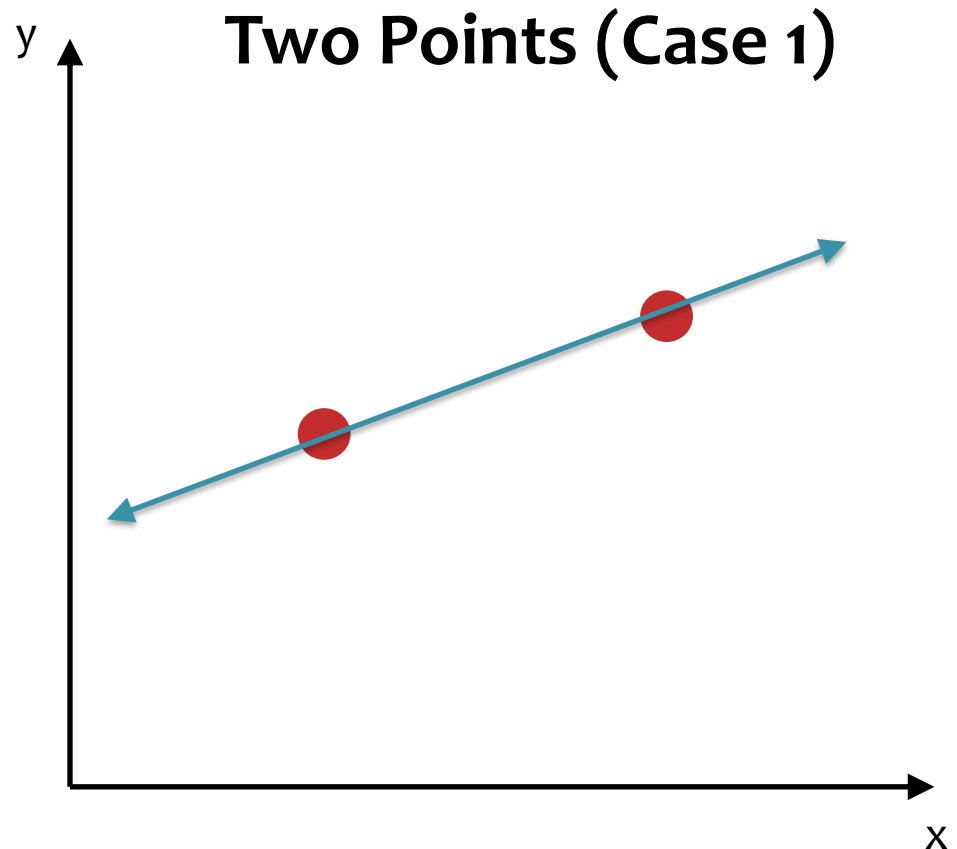
LINEAR REGRESSION: SOLUTION UNIQUENESS

Linear Regression: Uniqueness

Question:

Consider a 1D linear regression model trained to minimize MSE.

How many solutions (i.e. sets of parameters w, b) are there for the given dataset?

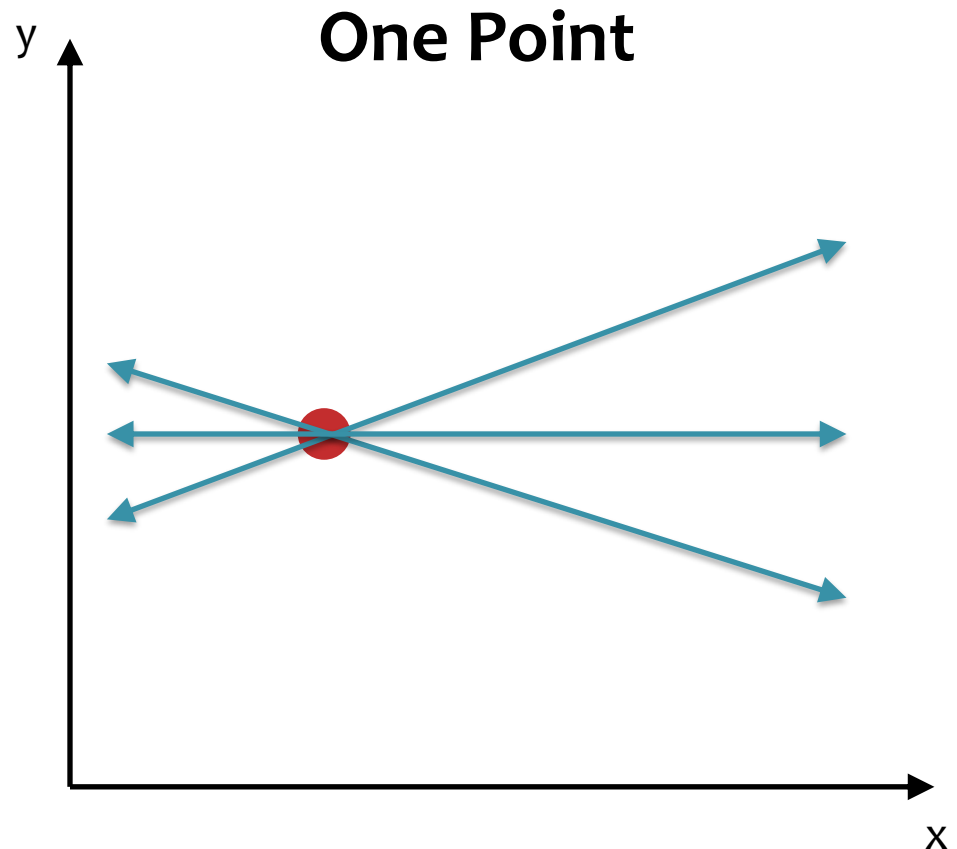


Linear Regression: Uniqueness

Question:

Consider a 1D linear regression model trained to minimize MSE.

How many solutions (i.e. sets of parameters w, b) are there for the given dataset?

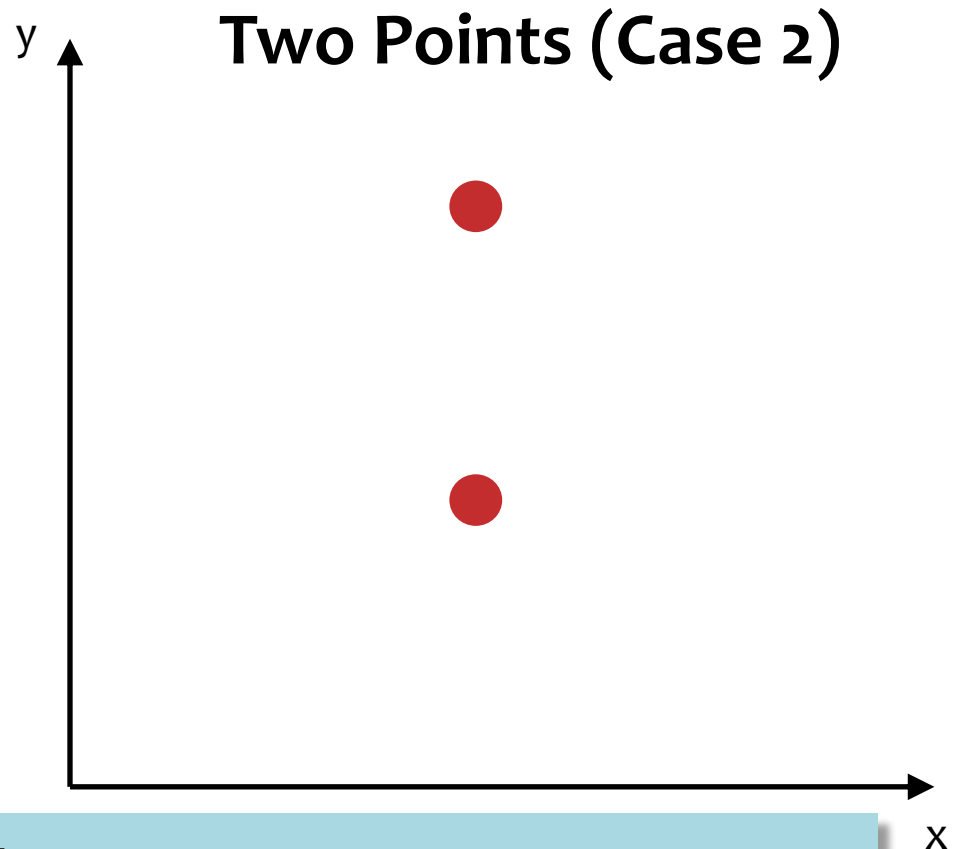


Linear Regression: Uniqueness

Question:

Consider a 1D linear regression model trained to minimize MSE.

How many solutions (i.e. sets of parameters w, b) are there for the given dataset?



Answer:

A: 0 B: 1 C: 2 D: $+\infty$

Question 2

A

B

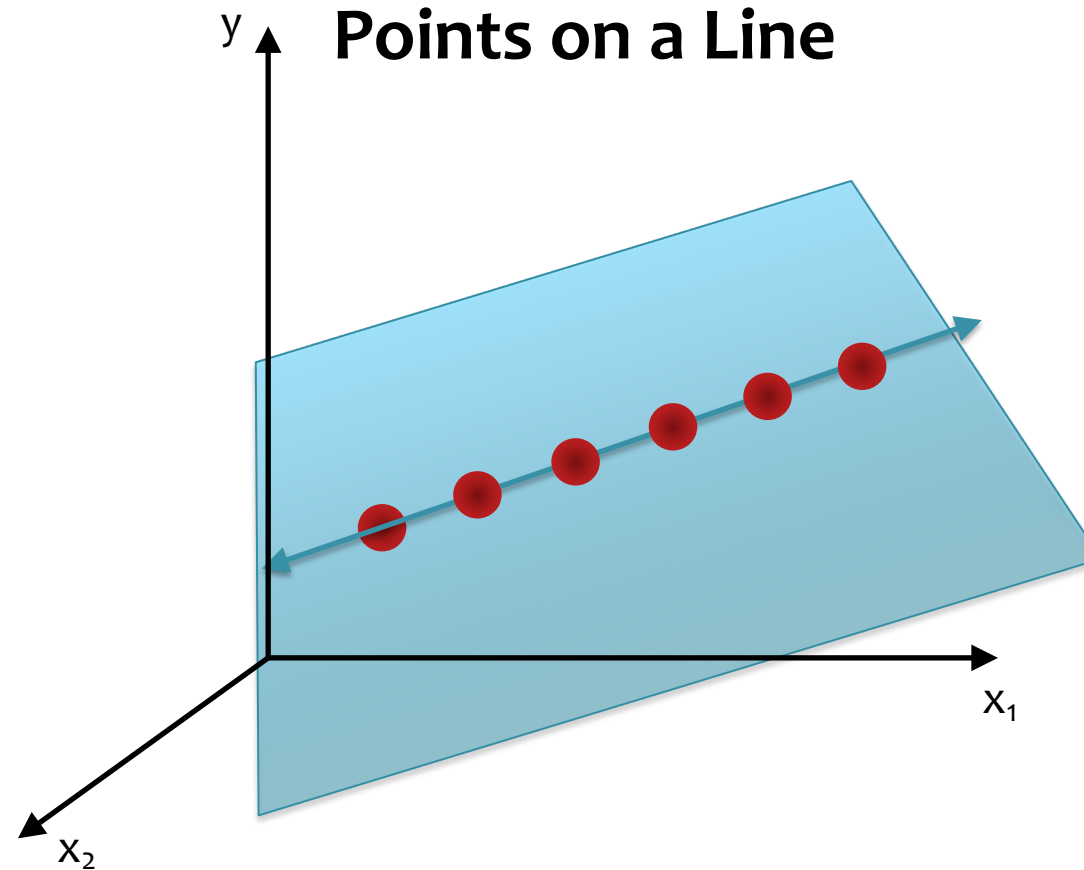
C

D

Linear Regression: Uniqueness

Question:

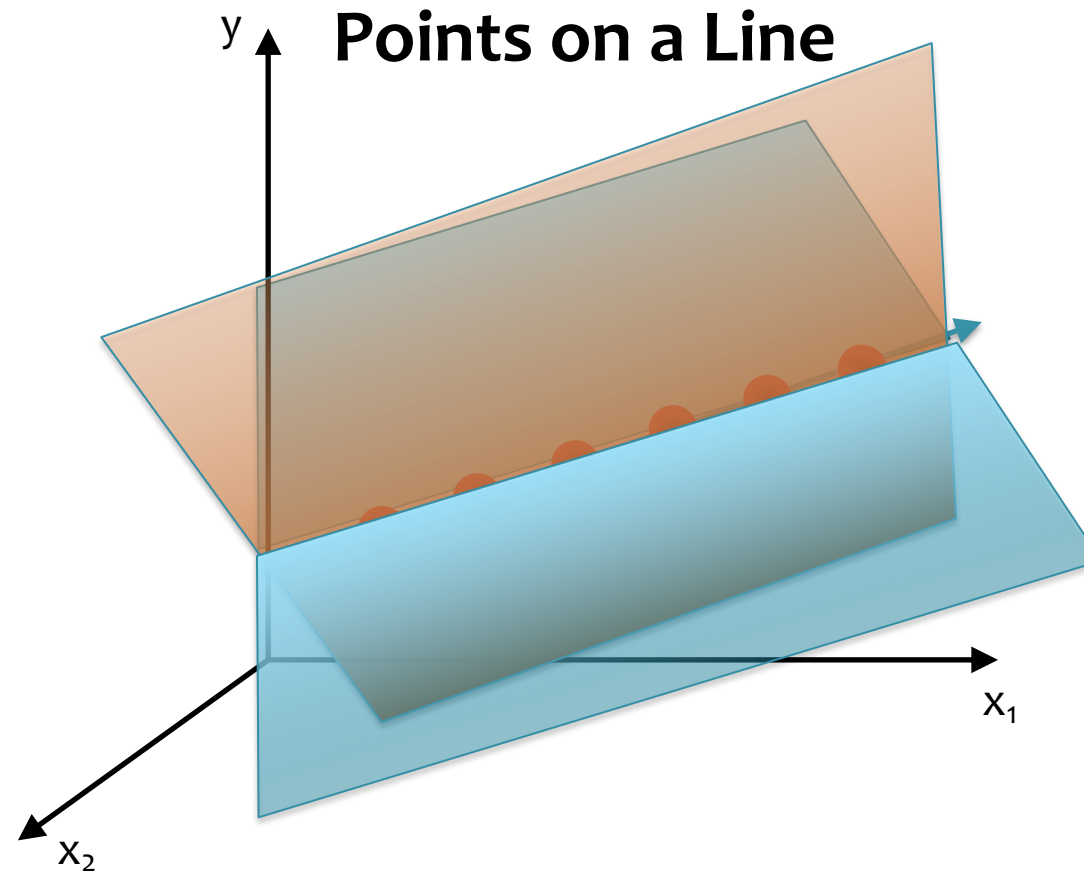
- Consider a **2D** linear regression model trained to minimize MSE
- How many solutions (i.e. sets of parameters w_1, w_2, b) are there for the given dataset?



Linear Regression: Uniqueness

Question:

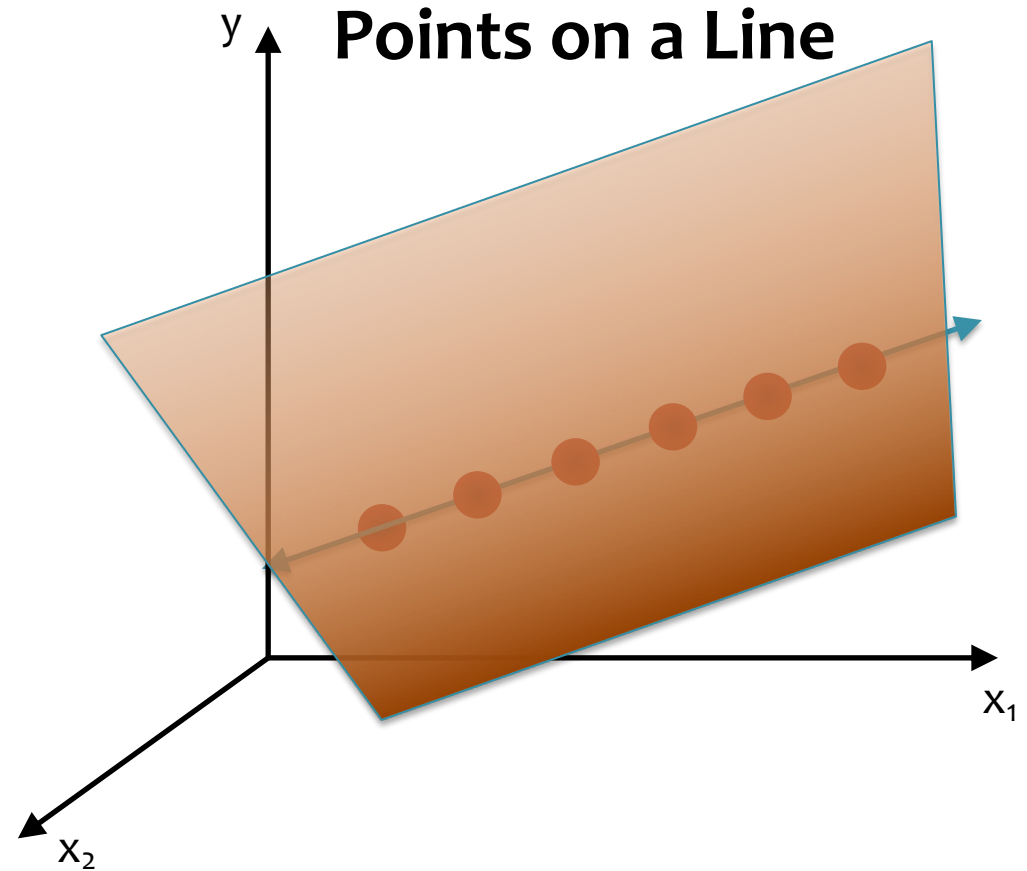
- Consider a **2D** linear regression model trained to minimize *MSE*
- How many solutions (i.e. sets of parameters w_1 , w_2 , b) are there for the given dataset?



Linear Regression: Uniqueness

Question:

- Consider a **2D** linear regression model trained to minimize MSE
- How many solutions (i.e. sets of parameters w_1, w_2, b) are there for the given dataset?



Linear Regression: Uniqueness

To solve the Ordinary Least Squares problem we compute:

$$\begin{aligned}\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} &= \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - (\boldsymbol{\theta}^T \mathbf{x}^{(i)}))^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{Y})\end{aligned}$$

These geometric intuitions align with the linear algebraic intuitions we can derive from the normal equations.

1. If $(\mathbf{X}^T \mathbf{X})$ is **invertible**, then there is exactly one solution.
2. If $(\mathbf{X}^T \mathbf{X})$ is **not invertible**, then there are either no solutions or infinitely many solutions.


Linear Regression: Uniqueness

To solve the Ordinary Least Squares problem we compute:

$$\begin{aligned}\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} &= \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - (\boldsymbol{\theta}^T \mathbf{x}^{(i)}))^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{Y})\end{aligned}$$

These geometric intuitions align with the linear algebraic intuitions we can derive from the normal equations.

1. If $(\mathbf{X}^T \mathbf{X})$ is invertible, then there is exactly one solution.
2. If $(\mathbf{X}^T \mathbf{X})$ is not invertible, there are no solutions or infinitely many solutions.



Invertability of $(\mathbf{X}^T \mathbf{X})$ is equivalent to \mathbf{X} being full rank. That is, there is no feature that is a linear combination of the other features.

Solving Linear Regression

Question:

True or False: If Mean Squared Error (i.e. $\frac{1}{N} \sum_{i=1}^N (y^{(i)} - h(\mathbf{x}^{(i)}))^2$) has a unique minimizer (i.e. argmin), then Mean Absolute Error (i.e. $\frac{1}{N} \sum_{i=1}^N |y^{(i)} - h(\mathbf{x}^{(i)})|$) must also have a unique minimizer.

Answer:

Question 3

A

B

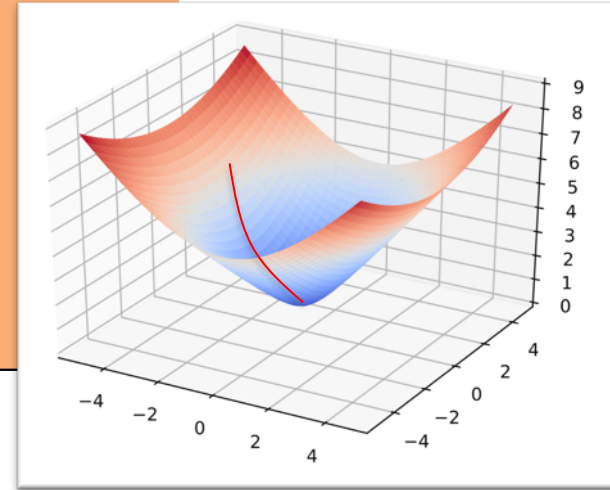
C

OPTIMIZATION METHOD #3: STOCHASTIC GRADIENT DESCENT

Gradient Descent

Algorithm 1 Gradient Descent

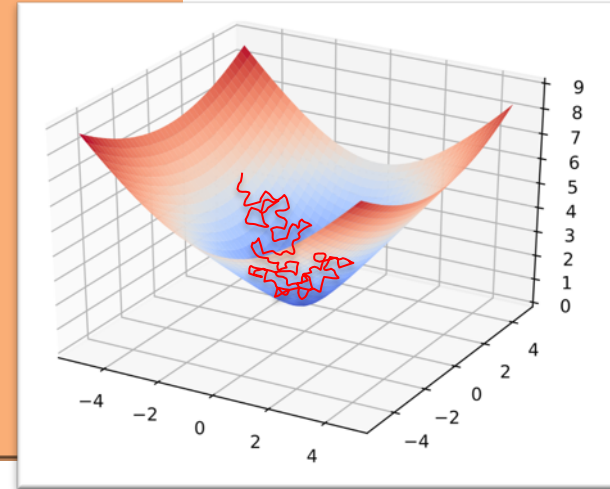
```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$ 
5:   return  $\theta$ 
```



Stochastic Gradient Descent (SGD)

Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $i \sim \text{Uniform}(\{1, 2, \dots, N\})$ 
5:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



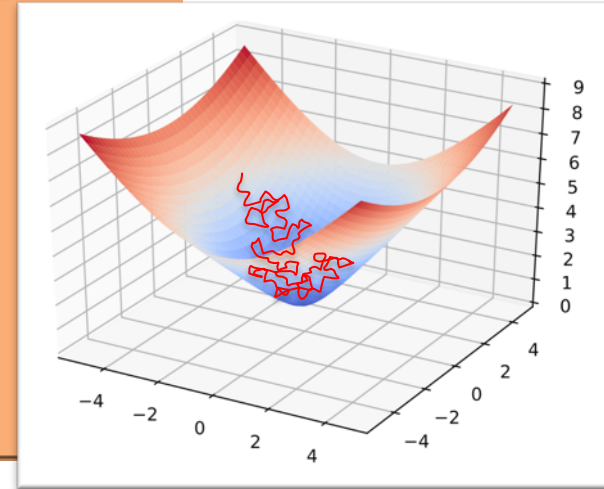
Per-example objective: $J^{(i)}(\theta)$

$$\text{Original objective: } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

Stochastic Gradient Descent (SGD)

Algorithm 2 Stochastic Gradient Descent (SGD)

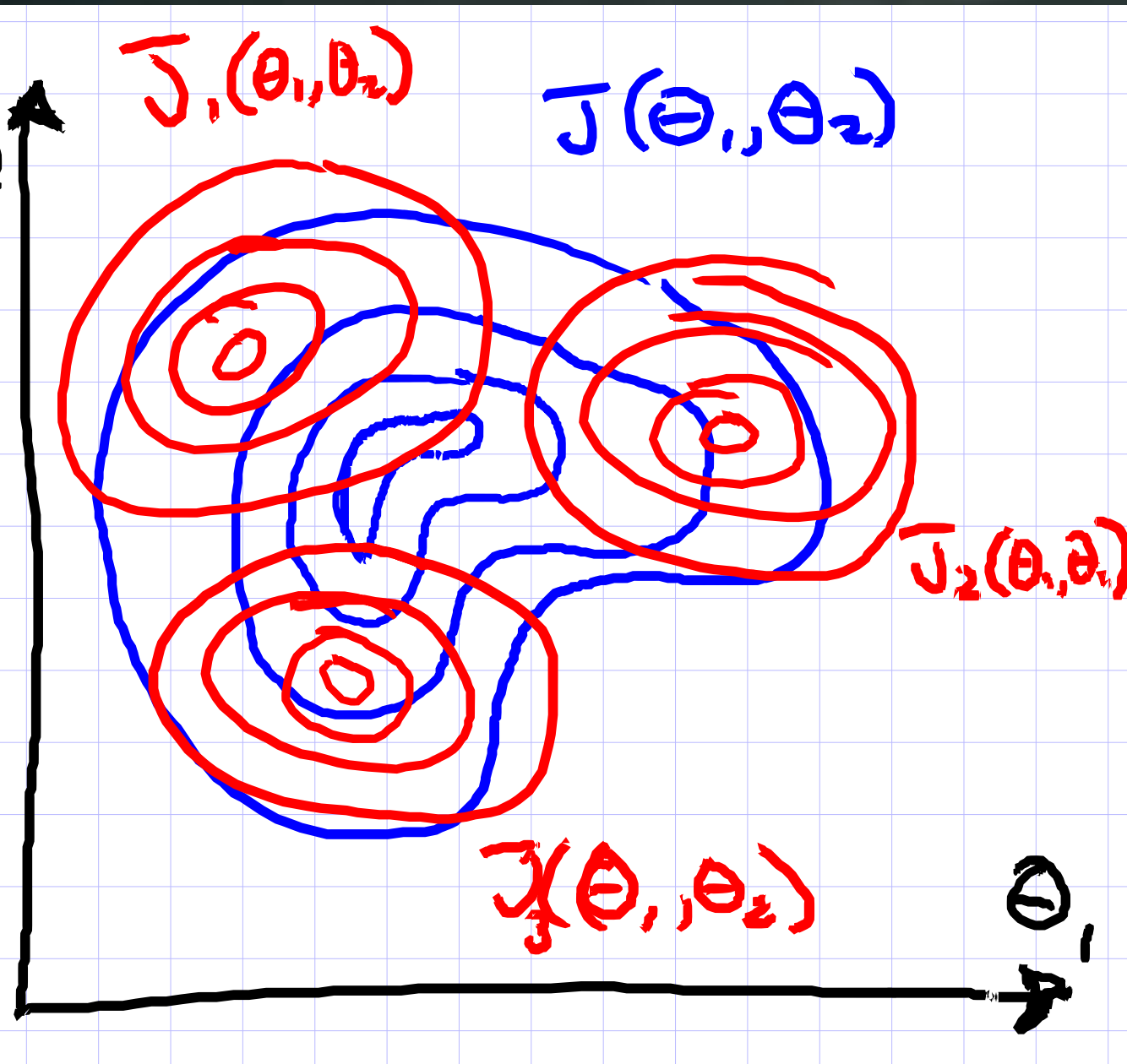
```
1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:        $\theta \leftarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



Per-example objective: $J^{(i)}(\theta)$

$$\text{Original objective: } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

In practice, it is common to implement SGD using sampling **without** replacement (i.e. $\text{shuffle}(\{1, 2, \dots, N\})$), even though most of the theory is for sampling **with** replacement (i.e. $\text{Uniform}(\{1, 2, \dots, N\})$).



Expectations of Gradients

$$\frac{dJ(\vec{\theta})}{d\theta_j} = \frac{1}{N} \sum_{i=1}^N \frac{d}{d\theta_j} (J_i(\vec{\theta}))$$
$$\nabla J(\vec{\theta}) = \begin{bmatrix} \vdots \\ \text{jth} \\ \vdots \end{bmatrix} = \frac{1}{N} \sum_{i=1}^N \nabla J_i(\vec{\theta})$$

Recall: for any discrete r.v. X

$$E_X[f(x)] \triangleq \sum_x P(X=x) f(x)$$

Q: What is the expected value of a randomly chosen $\nabla J_i(\vec{\theta})$?

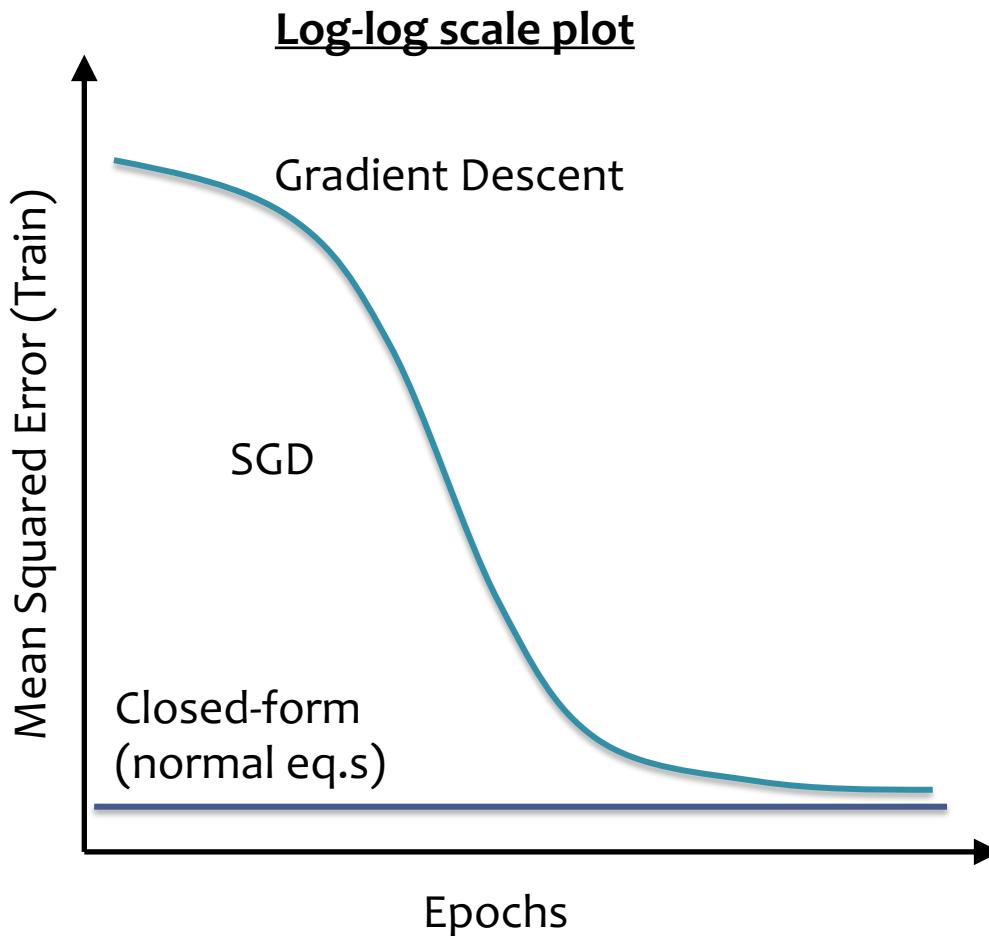
Let $I \sim \text{Uniform}(\{1, \dots, N\})$

$$\Rightarrow P(I=i) = \frac{1}{N} \text{ if } i \in \{1, \dots, N\}$$

$$E_I[\nabla J_I(\vec{\theta})] = \sum_{i=1}^N P(I=i) \nabla J_i(\vec{\theta})$$
$$= \frac{1}{N} \sum_{i=1}^N \nabla J_i(\vec{\theta})$$
$$= \nabla J(\vec{\theta})$$

LINEAR REGRESSION: PRACTICALITIES

Empirical Convergence



- Def: an **epoch** is a single pass through the training data
- 1. For GD, only **one update** per epoch
- 2. For SGD, **N updates** per epoch
 $N = (\# \text{ train examples})$

- SGD reduces MSE much more rapidly than GD
- For GD / SGD, training MSE is initially large due to uninformed initialization

Convergence of Optimizers

Convergence Analysis:

Def: Convergence is when $J(\vec{\theta}) - J(\vec{\theta}^*) < \epsilon$

↖ true unknown min

Methods	Steps to Converge	Computation per iteration
Newton's Method	$O(\ln \ln 1/\epsilon)$	$\nabla^2 J(\theta) \leftarrow O(NM^2)$
GD	$O(\ln 1/\epsilon)$	$\nabla J(\theta) \leftarrow O(NM)$
SGD	$O(1/\epsilon)$	$\nabla J_i(\theta) \leftarrow O(M)$

not converged ↗

"almost sure" convergence ↗

lots of caveats and conditions ↗

very less computation ↗

Takeaway: SGD has much slower asymptotic convergence, but is often faster in practice.

SGD FOR LINEAR REGRESSION

Linear Regression as Function Approximation

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$$

where $\mathbf{x} \in \mathbb{R}^M$ and $y \in \mathbb{R}$

1. Assume \mathcal{D} generated as:

$$\begin{aligned}\mathbf{x}^{(i)} &\sim p^*(\cdot) \\ y^{(i)} &= h^*(\mathbf{x}^{(i)})\end{aligned}$$

2. Choose hypothesis space, \mathcal{H} :
all linear functions in M -dimensional space

$$\mathcal{H} = \{h_{\boldsymbol{\theta}} : h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}, \boldsymbol{\theta} \in \mathbb{R}^M\}$$

3. Choose an objective function:
mean squared error (MSE)

$$\begin{aligned}J(\boldsymbol{\theta}) &= \frac{1}{N} \sum_{i=1}^N e_i^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})\right)^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)}\right)^2\end{aligned}$$

4. Solve the unconstrained optimization problem via favorite method:

- gradient descent
- closed form
- stochastic gradient descent
- ...

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$$

5. Test time: given a new \mathbf{x} , make prediction \hat{y}

$$\hat{y} = h_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}^T \mathbf{x}$$

Gradient Calculation for Linear Regression

Derivative of $J^{(i)}(\theta)$:

$$\begin{aligned} \frac{d}{d\theta_k} J^{(i)}(\theta) &= \frac{d}{d\theta_k} \frac{1}{2} (\theta^T \mathbf{x}^{(i)} - y^{(i)})^2 \\ &= \frac{1}{2} \frac{d}{d\theta_k} (\theta^T \mathbf{x}^{(i)} - y^{(i)})^2 \\ &= (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \frac{d}{d\theta_k} (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \\ &= (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \frac{d}{d\theta_k} \left(\sum_{j=1}^K \theta_j x_j^{(i)} - y^{(i)} \right) \\ &= (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)} \end{aligned}$$

Gradient of $J^{(i)}(\theta)$ [used by SGD]

$$\begin{aligned} \nabla_{\theta} J^{(i)}(\theta) &= \begin{bmatrix} \frac{d}{d\theta_1} J^{(i)}(\theta) \\ \frac{d}{d\theta_2} J^{(i)}(\theta) \\ \vdots \\ \frac{d}{d\theta_M} J^{(i)}(\theta) \end{bmatrix} = \begin{bmatrix} (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_1^{(i)} \\ (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_2^{(i)} \\ \vdots \\ (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_N^{(i)} \end{bmatrix} \\ &= (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \end{aligned}$$

Derivative of $J(\theta)$:

$$\begin{aligned} \frac{d}{d\theta_k} J(\theta) &= \sum_{i=1}^N \frac{d}{d\theta_k} J^{(i)}(\theta) \\ &= \sum_{i=1}^N (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)} \end{aligned}$$

Gradient of $J(\theta)$ [used by Gradient Descent]

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \begin{bmatrix} \frac{d}{d\theta_1} J(\theta) \\ \frac{d}{d\theta_2} J(\theta) \\ \vdots \\ \frac{d}{d\theta_M} J(\theta) \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_1^{(i)} \\ \sum_{i=1}^N (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_2^{(i)} \\ \vdots \\ \sum_{i=1}^N (\theta^T \mathbf{x}^{(i)} - y^{(i)}) x_N^{(i)} \end{bmatrix} \\ &= \sum_{i=1}^N (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)} \end{aligned}$$

SGD for Linear Regression

SGD applied to Linear Regression is called the “Least Mean Squares” algorithm

Algorithm 1 Least Mean Squares (LMS)

```
1: procedure LMS( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$  ▷ Initialize parameters
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:        $\mathbf{g} \leftarrow (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$  ▷ Compute gradient
6:        $\theta \leftarrow \theta - \gamma \mathbf{g}$  ▷ Update parameters
7:   return  $\theta$ 
```

Optimization Objectives

You should be able to...

- Apply gradient descent to optimize a function
- Apply stochastic gradient descent (SGD) to optimize a function
- Apply knowledge of zero derivatives to identify a closed-form solution (if one exists) to an optimization problem
- Distinguish between convex, concave, and nonconvex functions
- Obtain the gradient (and Hessian) of a (twice) differentiable function

Linear Regression Objectives

You should be able to...

- Design k-NN Regression and Decision Tree Regression
- Implement learning for Linear Regression using three optimization techniques: (1) closed form, (2) gradient descent, (3) stochastic gradient descent
- Choose a Linear Regression optimization technique that is appropriate for a particular dataset by analyzing the tradeoff of computational complexity vs. convergence speed
- Distinguish the three sources of error identified by the bias-variance decomposition: bias, variance, and irreducible error.