# 10-301/601: Introduction to Machine Learning Lecture 8 – Optimization for Machine Learning

Henry Chai & Matt Gormley

9/25/23

## Exam 1 Logistics

- Exam 1 on 2/19 (next Monday!) from 7 PM – 9 PM

- Location & Seats: You all will be split across multiple (large) rooms.

  - Everyone will have an assigned seat

  - Please watch Piazza carefully for more details

  - If you have exam accommodations through ODR, they will be proctoring your exam on our behalf; **you are responsible for submitting the exam proctoring request through your student portal.**

# Exam 1 Logistics

- Format of questions:
  - Multiple choice
  - True / False (with justification)
  - Derivations
  - Short answers
  - Drawing & Interpreting figures
  - Implementing algorithms on paper
- No electronic devices (you won't need them!)
- You are allowed to bring one letter-size sheet of notes; you can put *whatever* you want on *both sides*

# Exam 1 Topics

- Covered material: Lectures 1 – 7
  - Foundations
    - Probability, Linear Algebra, Geometry, Calculus
    - Optimization
  - Important Concepts
    - Overfitting
    - Model selection / Hyperparameter optimization
  - Decision Trees
  - $k$-NN
  - Perceptron
  - Regression
    - Decision Tree and $k$-NN Regression
    - Linear Regression

# Exam 1 Preparation

- Review the exam practice problems (released 2/12 on the course website, under Coursework)

- Attend the dedicated exam 1 review OH (in lieu of recitation on 2/16)

- Review HWs 1 - 3

- Consider whether you have achieved the "learning objectives" for each lecture / section

- Write your one-page cheat sheet (back and front)

# Exam 1 Tips

- Solve the easy problems first

- If a problem seems extremely complicated, you might be missing something

- If you make an assumption, write it down

- Don't leave any answer blank
  - If you look at a question and don't know the answer:
    - just start trying things
    - consider multiple approaches
    - imagine arguing for some answer and see if you like it

# 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:

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

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

$$\vec{\theta} = (\theta_1, \cdots, \theta_M)$$

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

$$\hookrightarrow \mathbb{R}^M$$

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

$$\text{MSE} = J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} e_i^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left( y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right)^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left( y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)} \right)^2$$

4. Solve the unconstrained optimization problem via favorite method:

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

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

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

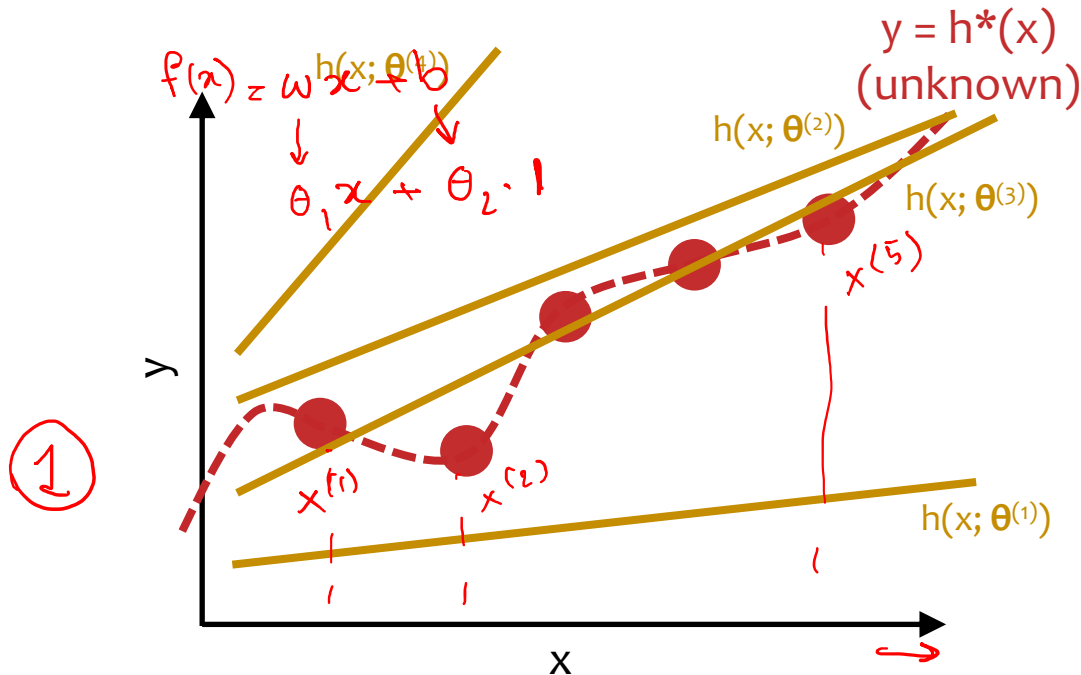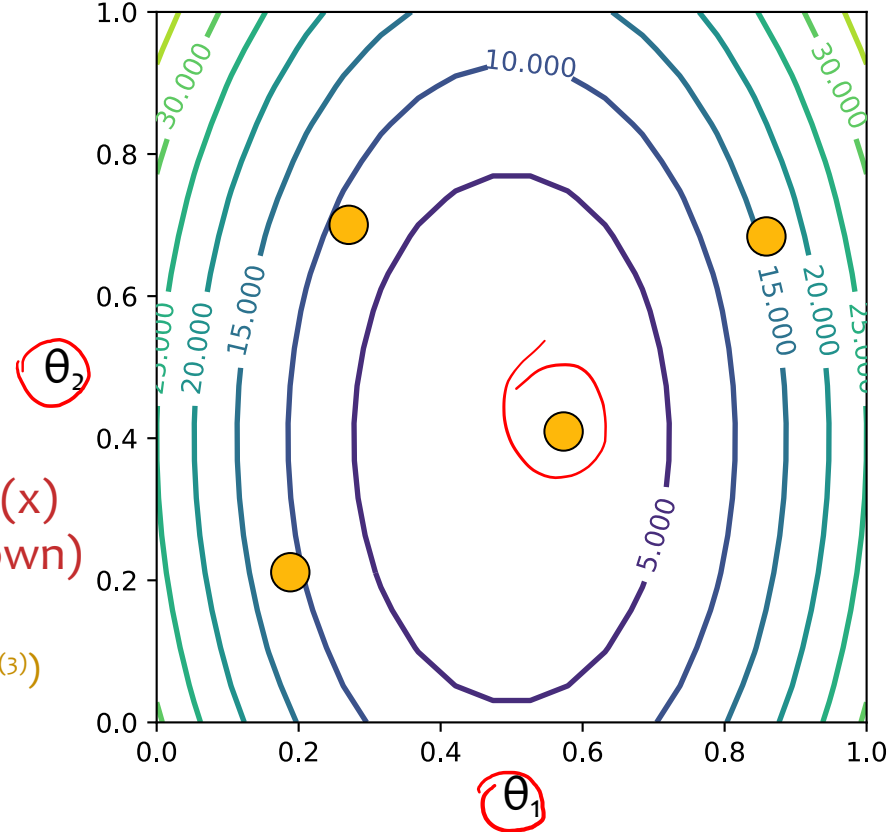$$\hat{y} = h_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}^T \mathbf{x}$$

# Linear Regression by Rand. Guessing

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

**Optimization Method #0:
Random Guessing**

1. Pick a random **θ**

2. Evaluate J(**θ**)

3. Repeat steps 1 and 2 many times

4. Return **θ** that gives smallest J(**θ**)

y = h*(x)
(unknown)

h(x; θ⁽¹⁾)

h(x; θ⁽²⁾)

h(x; θ⁽³⁾)

h(x; θ⁽¹⁾)

$f(x) = w x + b$

$\theta_1 x + \theta_2 \cdot 1$

x⁽⁵⁾

x⁽¹⁾   x⁽²⁾

①

y

x

θ₂

θ₁

| t | θ₁ | θ₂ | J(θ₁, θ₂) |
|---|-----|-----|-----------|
| 1 | 0.2 | 0.2 | 10.4 |
| 2 | 0.3 | 0.7 | 7.2 |
| 3 | 0.6 | 0.4 | 1.0 |
| 4 | 0.9 | 0.7 | 16.2 |

8

Gradients

$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2)$

# Gradients



$\min J(\theta)$
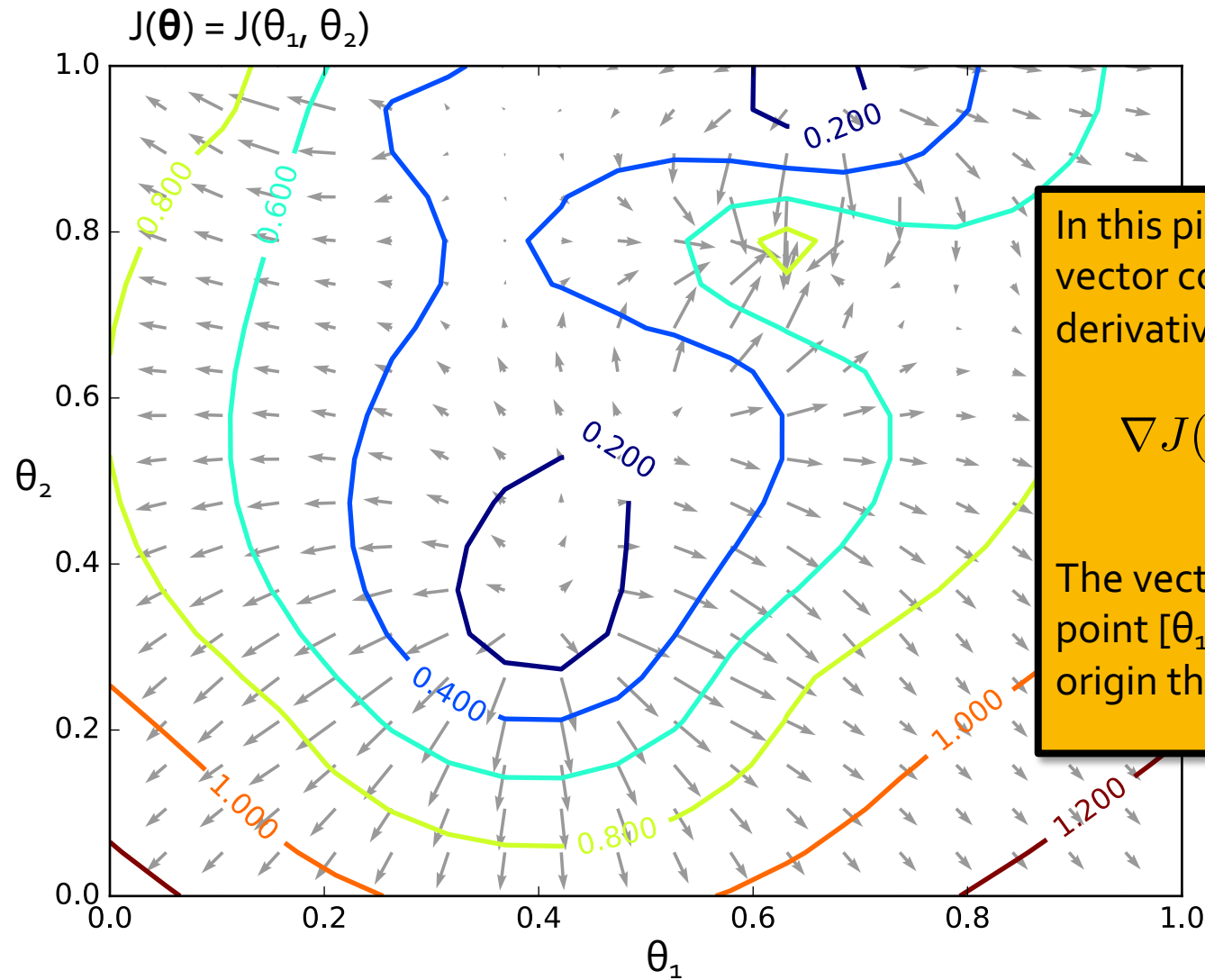
$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2)$

$\vec{\theta}'$

$\nabla J(\vec{\theta}') = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_m} \end{bmatrix}$

$\theta_2$

$\theta_1$

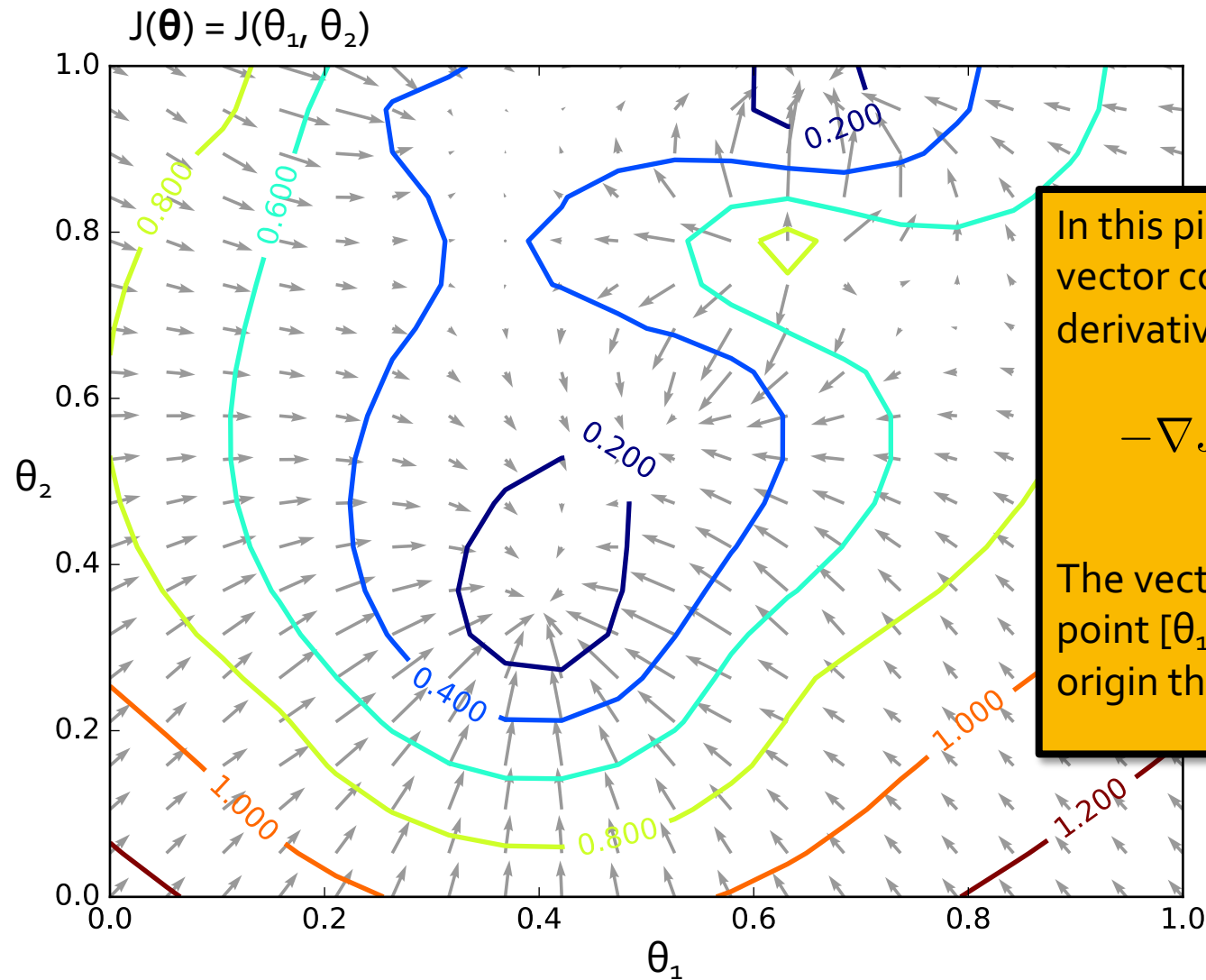These are the **gradients** that
Gradient **Ascent** would follow.

# Gradients
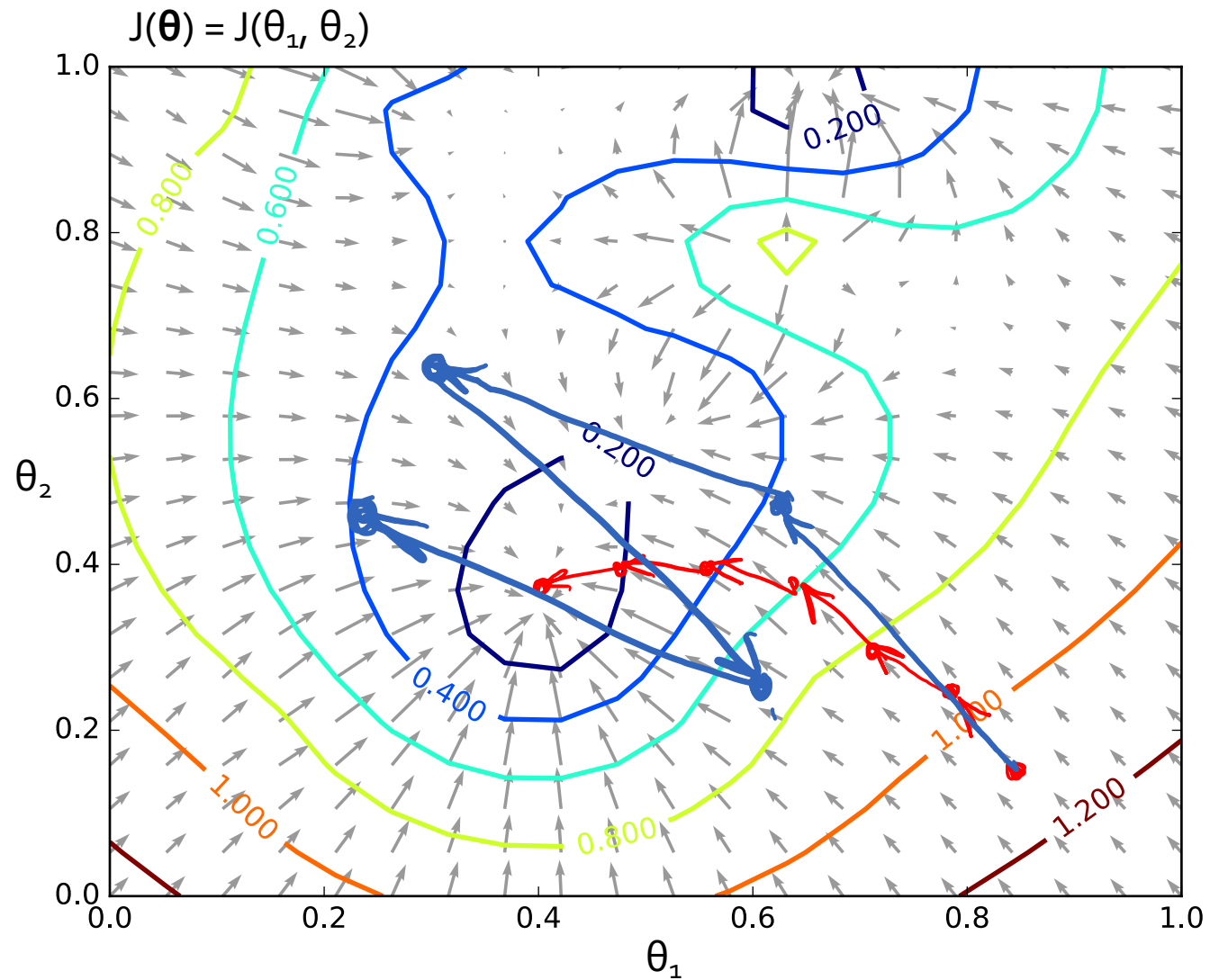
$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2)$



In this picture, each arrow is a 2D vector consisting of two partial derivatives.

$$\nabla J(\theta_1, \theta_2) = \begin{bmatrix} \dfrac{\partial J}{\partial \theta_1} \\[2ex] \dfrac{\partial J}{\partial \theta_2} \end{bmatrix}$$

The vector is evaluated at the point $[\theta_1, \theta_2]^\mathsf{T}$ and plotted with its origin there as well.

These are the **gradients** that
Gradient **Ascent** would follow.

# (Negative) Gradients

$J(\mathbf{\theta}) = J(\theta_1, \theta_2)$



In this picture, each arrow is a 2D vector consisting of two partial derivatives.

$$-\nabla J(\theta_1, \theta_2) = \begin{bmatrix} -\frac{\partial J}{\partial \theta_1} \\ -\frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

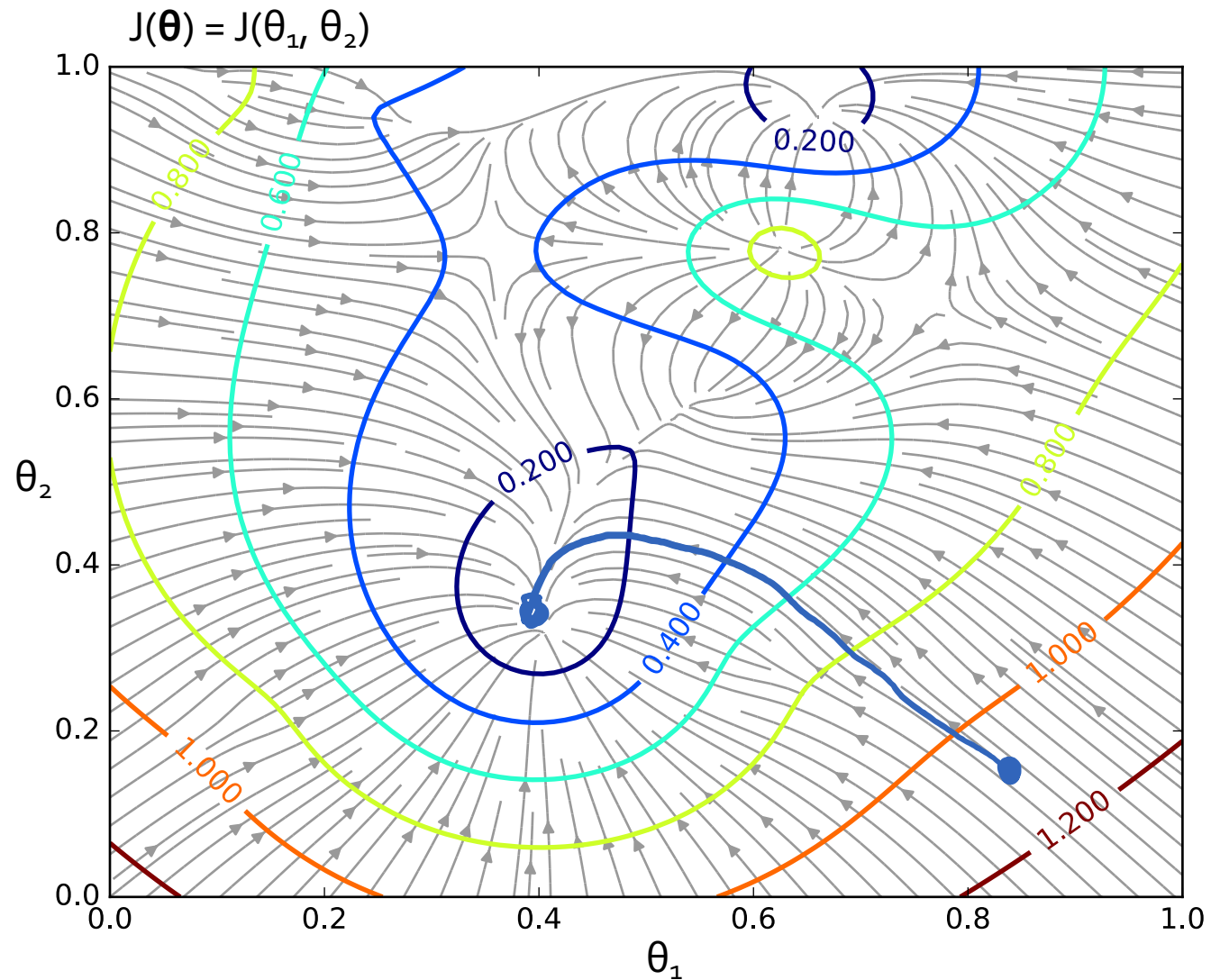The vector is evaluated at the point $[\theta_1, \theta_2]^T$ and plotted with its origin there as well.

These are the **negative** gradients that Gradient **Descent** would follow.

# (Negative) Gradients



$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2)$

These are the **negative** gradients that
Gradient **Descent** would follow.

# (Negative) Gradient *Pa*



$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2)$

Shown are the **paths** that Gradient Descent
would follow if it were making **infinitesimally
small steps**.

- Gradient descent for linear regression repeatedly takes steps opposite the gradient of the objective function

## Recall: Gradient Descent for Linear Regression

| **Algorithm 1** GD for Linear Regression |
|---|

1: **procedure** GDLR($\mathcal{D}, \boldsymbol{\theta}^{(0)}$)

2:     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$      ▷ Initialize parameters

3:     **while** not converged **do**

4:     $\nabla J(\theta)$   $\mathbf{g} \leftarrow \sum_{i=1}^{N}(\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})\mathbf{x}^{(i)}$    ▷ Compute gradient

5:     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \mathbf{g}$    ▷ Update parameters

6:     **return** $\boldsymbol{\theta}$

step size
learning rate

$$\nabla J = \vec{0}$$

$$||\nabla J|| = 0$$

# Gradient Calculation for Linear Regression

$i \rightarrow$ traing instance $i \in 1, \dots, N$

$k \rightarrow$ feature $k \in 1, \dots, M$

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

④

$$\frac{d}{d\theta_k} J^{(i)}(\boldsymbol{\theta}) = \frac{d}{d\theta_k} \frac{1}{2} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2} \frac{d}{d\theta_k} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2$$

$$= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \frac{d}{d\theta_k} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})$$

$$= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \frac{d}{d\theta_k} \left( \sum_{j=1}^{K} \theta_j x_j^{(i)} - y^{(i)} \right)$$

$$= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)}$$

① $\text{MSE} = J(\vec{\theta}) = \frac{1}{2N} \sum_{i=1}^{N} \left( y^{(i)} - \vec{\theta} \cdot \vec{x}^{(i)} \right)^2$

$i = 1, \dots N : J(\vec{\theta}) = \frac{1}{2N} \sum_{i=1}^{N} J^{(i)}(\theta) \quad J^{(i)}(\theta)$

**Derivative of $J(\boldsymbol{\theta})$:**

③

$$\frac{d}{d\theta_k} J(\boldsymbol{\theta}) = \sum_{i=1}^{N} \frac{d}{d\theta_k} J^{(i)}(\boldsymbol{\theta})$$

$$= \sum_{i=1}^{N} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)}$$

② **Gradient of $J(\boldsymbol{\theta})$**          [used by Gradient Descent]

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

$k = 1, \dots, M$

$$= \sum_{i=1}^{N} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$$
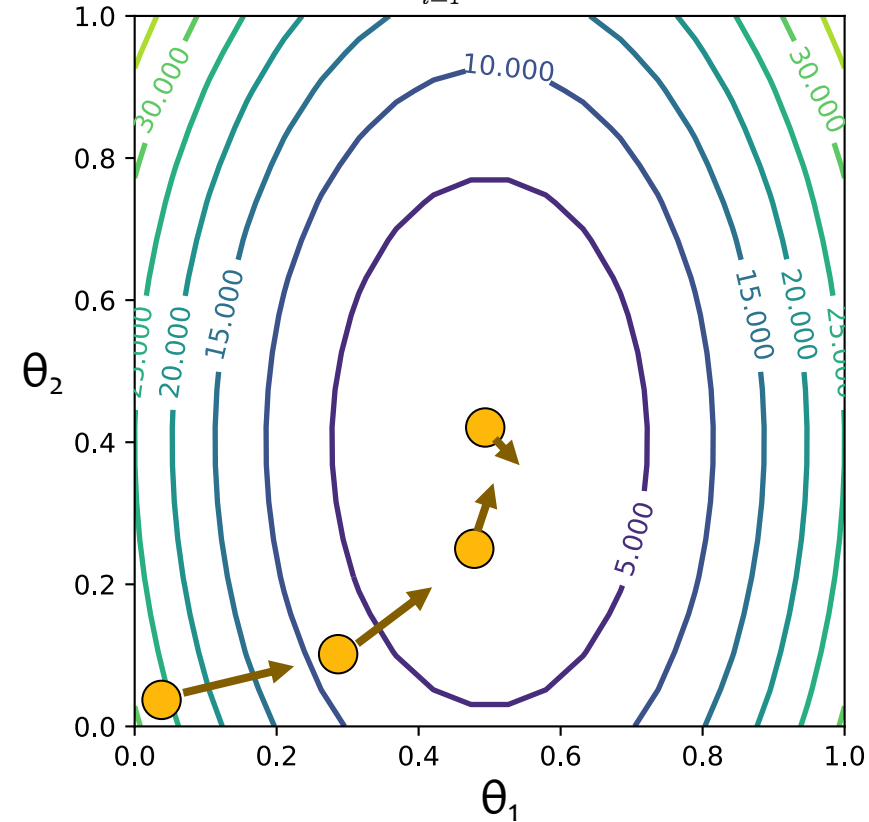
# Linear Regression by Gradient Desc.

**Optimization Method #1: Gradient Descent**

1. Pick a random **θ**

2. Repeat:
   a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
   b. Step opposite gradient

3. Return **θ** that gives smallest $J(\boldsymbol{\theta})$

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



| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|------|------|------|
| 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.
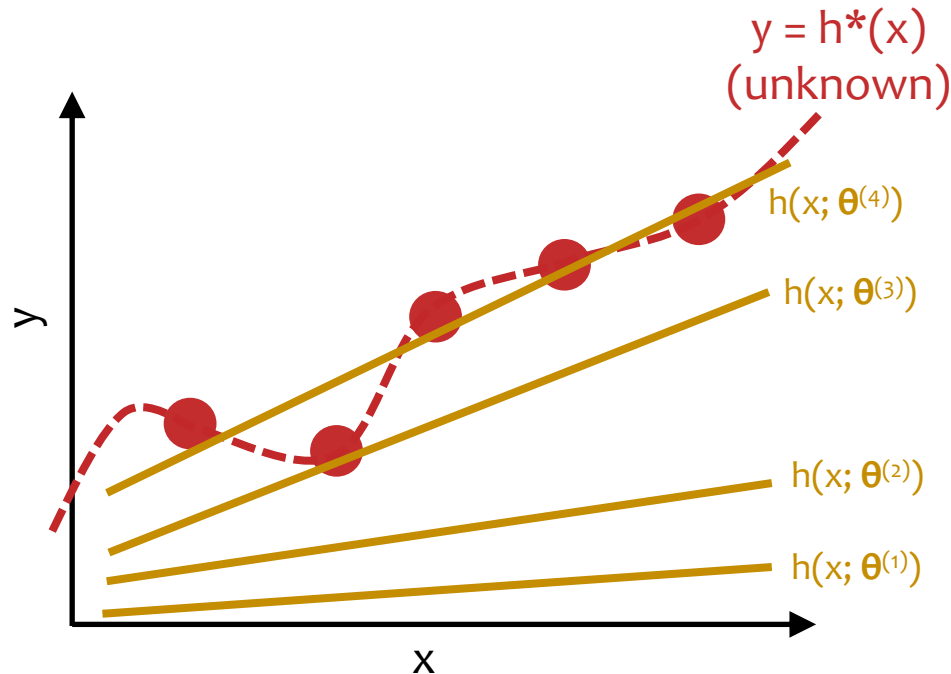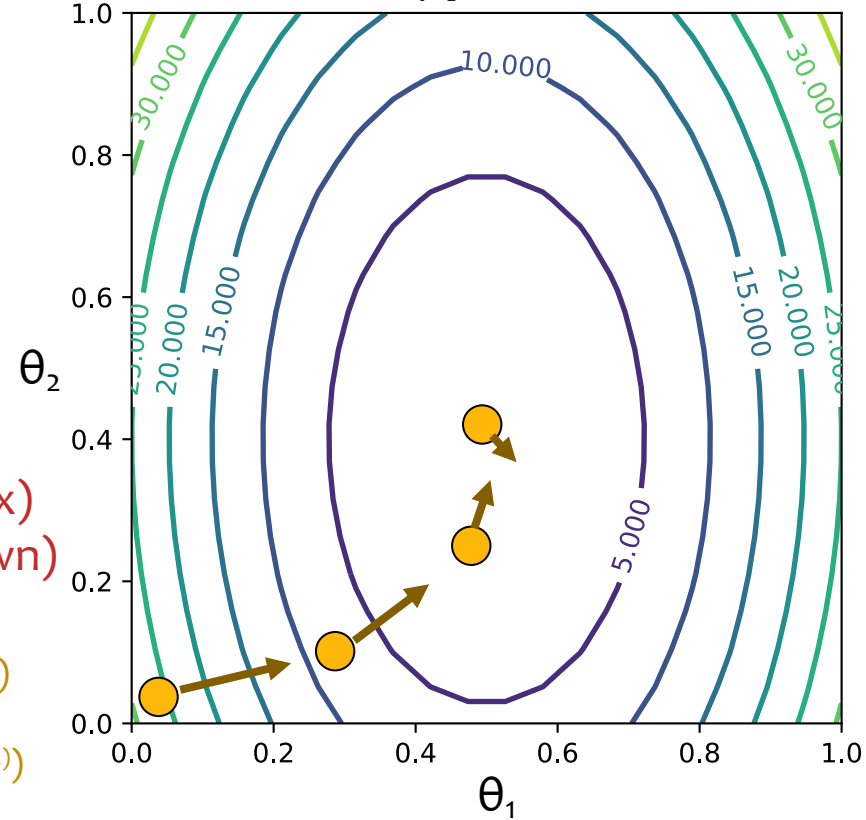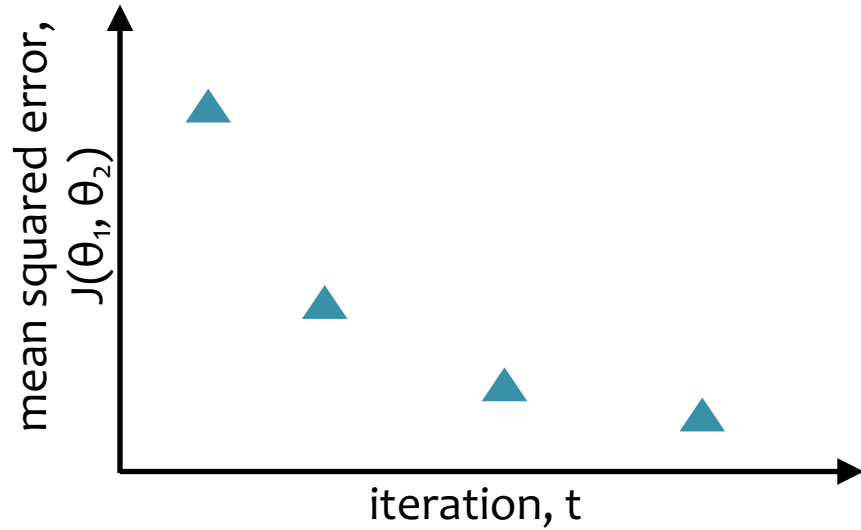
Optimization Method #1:
**Gradient Descent**

1. Pick a random $\theta$

2. Repeat:
   a. Evaluate gradient $\nabla J(\theta)$
   b. Step opposite gradient

3. Return $\theta$ that gives smallest $J(\theta)$



y = h*(x)
(unknown)

$h(x; \theta^{(4)})$

$h(x; \theta^{(3)})$

$h(x; \theta^{(2)})$

$h(x; \theta^{(1)})$

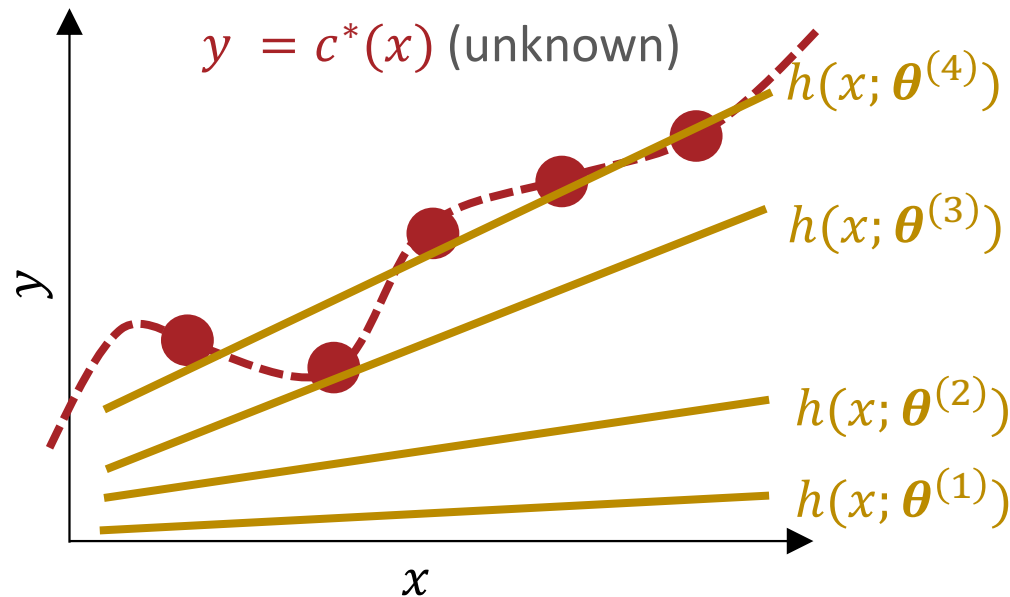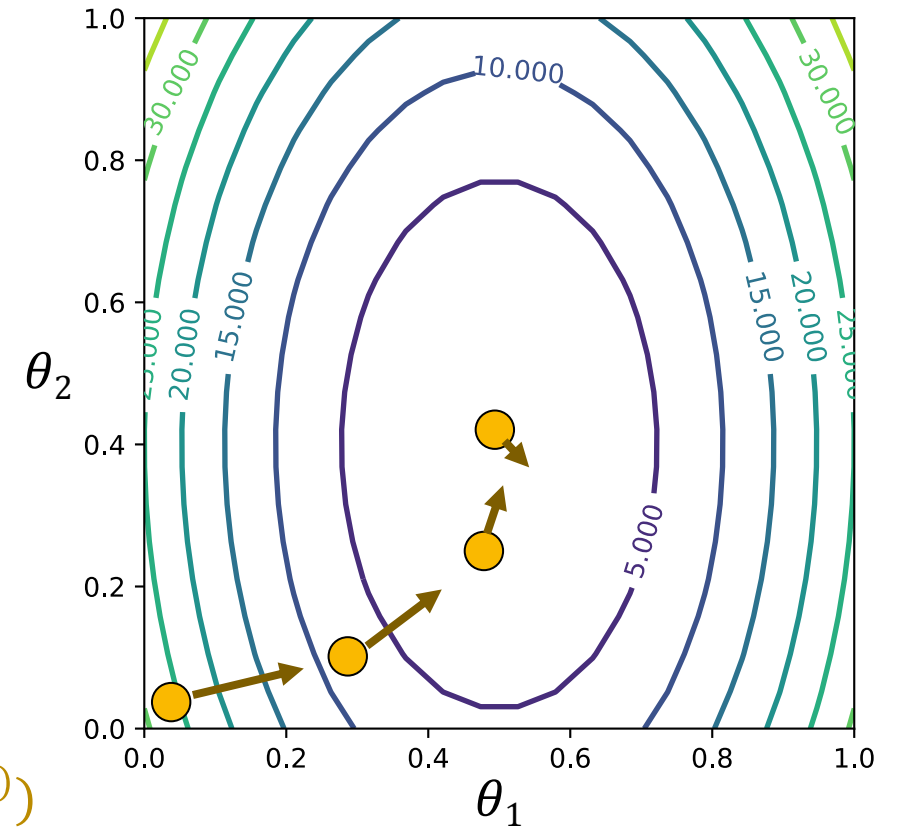| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|---|---|---|
| 1 | 0.01 | 0.02 | 25.2 |
| 2 | 0.30 | 0.12 | 8.7 |
| 3 | 0.51 | 0.30 | 1.5 |
| 4 | 0.59 | 0.43 | 0.2 |

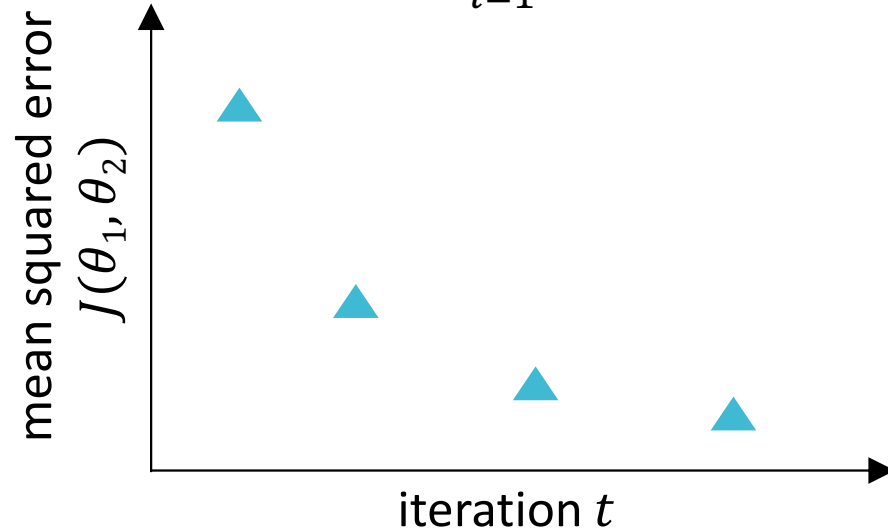# Linear Regression by Gradient Desc.

**Optimization Method #1: Gradient Descent**

1. Pick a random **θ**

2. Repeat:
   a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
   b. Step opposite gradient

3. Return **θ** that gives smallest $J(\boldsymbol{\theta})$

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



y = h*(x) (unknown)

h(x; **θ**⁽⁴⁾)

h(x; **θ**⁽³⁾)

h(x; **θ**⁽²⁾)

h(x; **θ**⁽¹⁾)
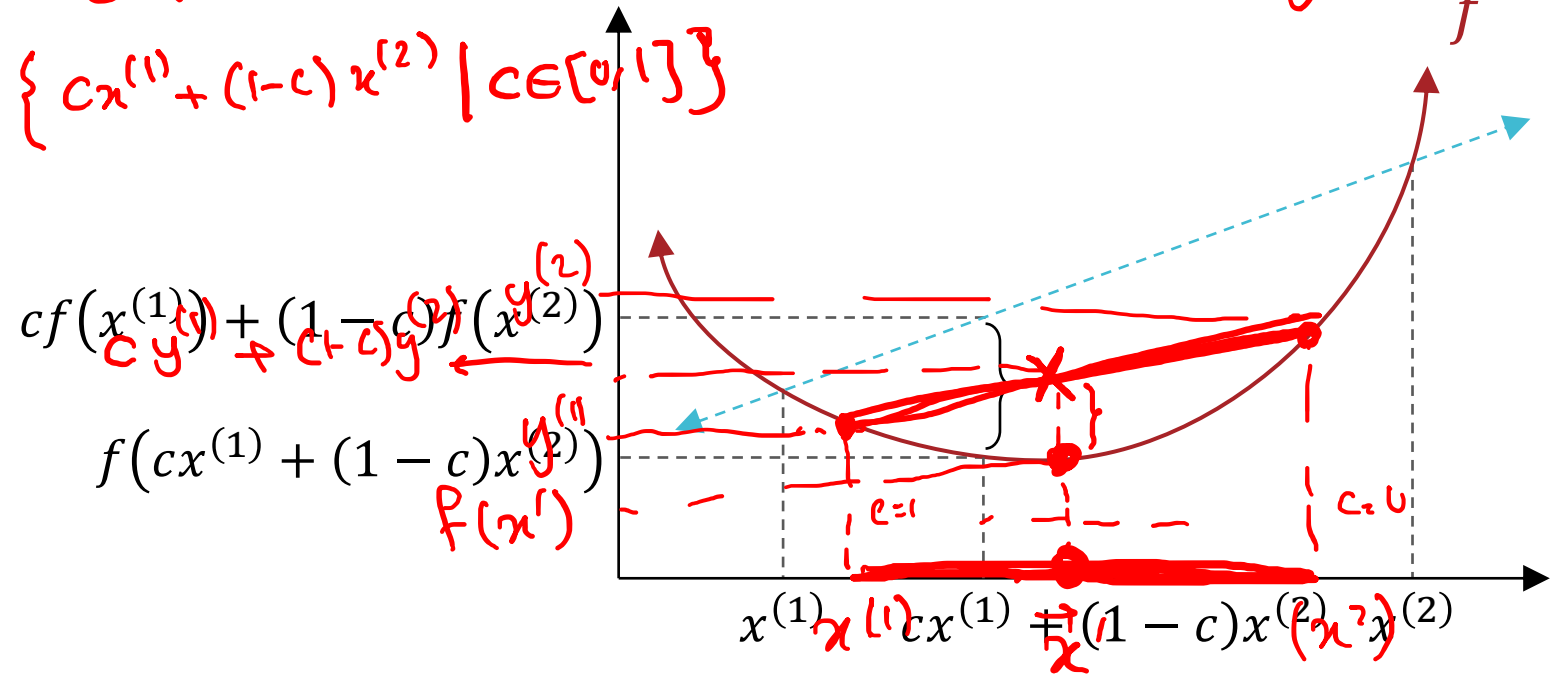
| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|------|------|------|
| 1 | 0.01 | 0.02 | 25.2 |
| 2 | 0.30 | 0.12 | 8.7 |
| 3 | 0.51 | 0.30 | 1.5 |
| 4 | 0.59 | 0.43 | 0.2 |

# Linear Regression by Gradient Desc.



| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|------|------|---------|
| 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.

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



| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|---|---|---|
| 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 |

# Why Gradient Descent for Linear Regression?

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



$y = c^*(x)$ (unknown)

$h(x; \boldsymbol{\theta}^{(4)})$

$h(x; \boldsymbol{\theta}^{(3)})$

$h(x; \boldsymbol{\theta}^{(2)})$

$h(x; \boldsymbol{\theta}^{(1)})$

| $t$ | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|---|---|---|
| 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 |

# Convexity

- A function $f: \mathbb{R}^D \to \mathbb{R}$ is          convex if
  $\forall \, \boldsymbol{x}^{(1)} \in \mathbb{R}^D, \boldsymbol{x}^{(2)} \in \mathbb{R}^D$ and $0 \leq c \leq 1$

$$* \quad f\left(c\boldsymbol{x}^{(1)} + (1-c)\boldsymbol{x}^{(2)}\right) \leq cf\left(\boldsymbol{x}^{(1)}\right) + (1-c)f\left(\boldsymbol{x}^{(2)}\right)$$

Convex combination of $x^{(1)}, x^{(2)}$     $y^{(1)}$     $y^{(2)}$

$$\sim \left\{ cx^{(1)} + (1-c)x^{(2)} \mid c \in [0,1] \right\}$$

$f$

$cf\left(x^{(1)}\right) + (1-c)f\left(x^{(2)}\right)$     $c\,y^{(2)}$
$c\,y + (1-c)\,y$

$f\left(cx^{(1)} + (1-c)x^{(2)}\right)$     $y^{(1)}$
$f(x')$     $c=1$     $c=0$

$x^{(1)}$     $x^{(1)}$     $cx^{(1)} + (1-c)x^{(2)}$     $x^{(2)}$     $x^{(2)}$
$x'$

# Convexity

- A function $f: \mathbb{R}^D \to \mathbb{R}$ is convex if
  $\forall\, \boldsymbol{x}^{(1)} \in \mathbb{R}^D, \boldsymbol{x}^{(2)} \in \mathbb{R}^D$ and $0 \leq c \leq 1$

$$f\left(c\boldsymbol{x}^{(1)} + (1-c)\boldsymbol{x}^{(2)}\right) \leq cf\left(\boldsymbol{x}^{(1)}\right) + (1-c)f\left(\boldsymbol{x}^{(2)}\right)$$

# Convexity

- A function $f: \mathbb{R}^D \to \mathbb{R}$ is *strictly* convex if

  $\forall\, \boldsymbol{x}^{(1)} \in \mathbb{R}^D, \boldsymbol{x}^{(2)} \in \mathbb{R}^D$ and $0 < c < 1$

  $$f\left(c\boldsymbol{x}^{(1)} + (1-c)\boldsymbol{x}^{(2)}\right) < cf\left(\boldsymbol{x}^{(1)}\right) + (1-c)f\left(\boldsymbol{x}^{(2)}\right)$$

# Convexity

Convex functions

Non-convex functions

# Convexity

Given a function $f: \mathbb{R}^D \to \mathbb{R}$

- $\boldsymbol{x}^*$ is a global minimum iff
$f(\boldsymbol{x}^*) \leq f(\boldsymbol{x}) \; \forall \; \boldsymbol{x} \in \mathbb{R}^D$

- $\boldsymbol{x}^*$ is a local minimum iff
$\exists \; \epsilon$ s.t. $f(\boldsymbol{x}^*) \leq f(\boldsymbol{x}) \; \forall$
$\boldsymbol{x}$ s.t. $\|\boldsymbol{x} - \boldsymbol{x}^*\|_2 < \epsilon$

# Convexity



Convex functions:

Each local minimum is a global minimum!



Non-convex functions:

A local minimum may or may not be a global minimum…

# Convexity

Strictly convex functions:

There exists a unique global minimum!

Non-convex functions:
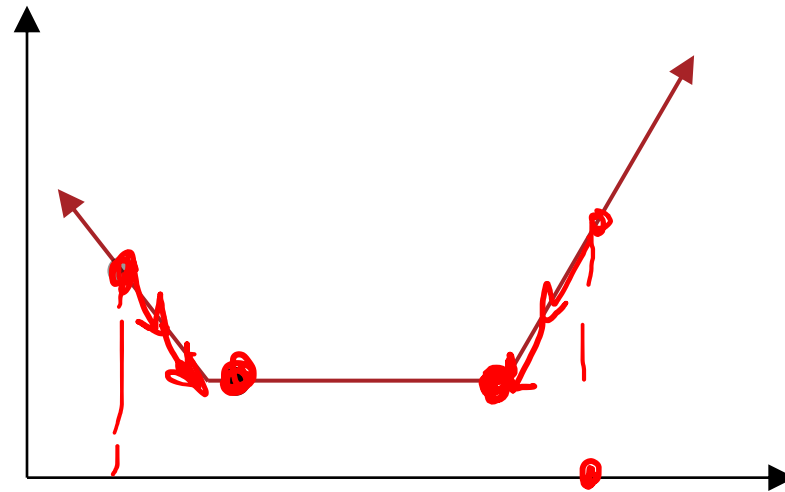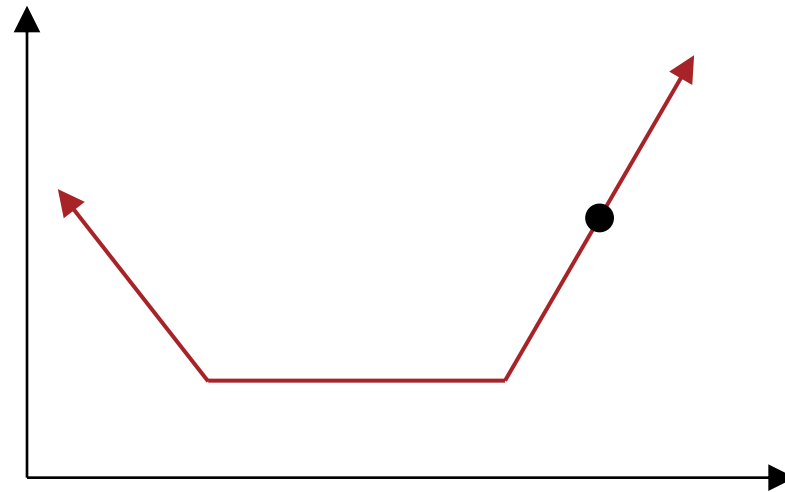
A local minimum may or may not be a global minimum…
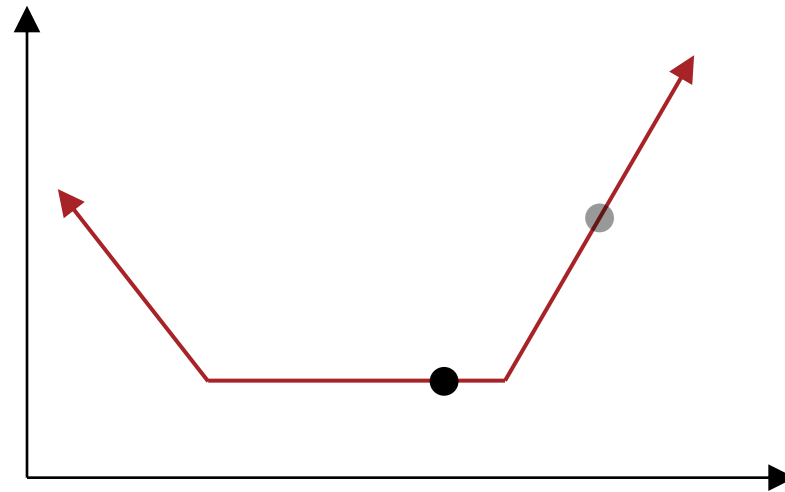
# Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
  - Works great if the objective function is convex!
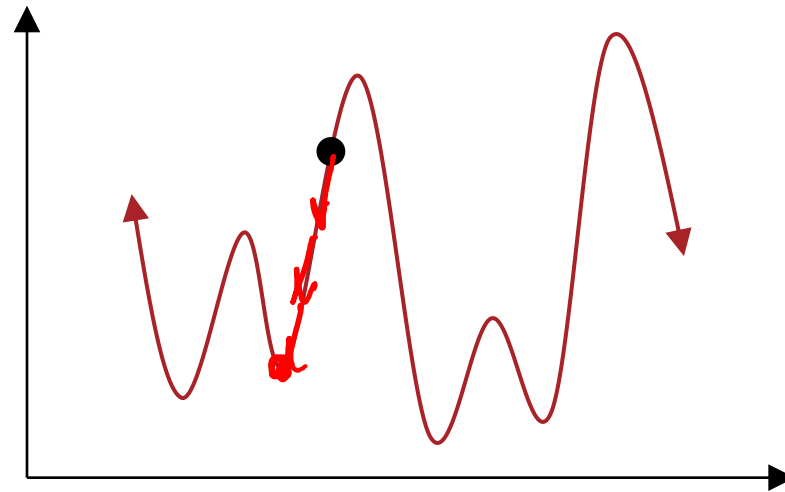
# Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
  - Works great if the objective function is convex!

# Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
  - Works great if the objective function is convex!
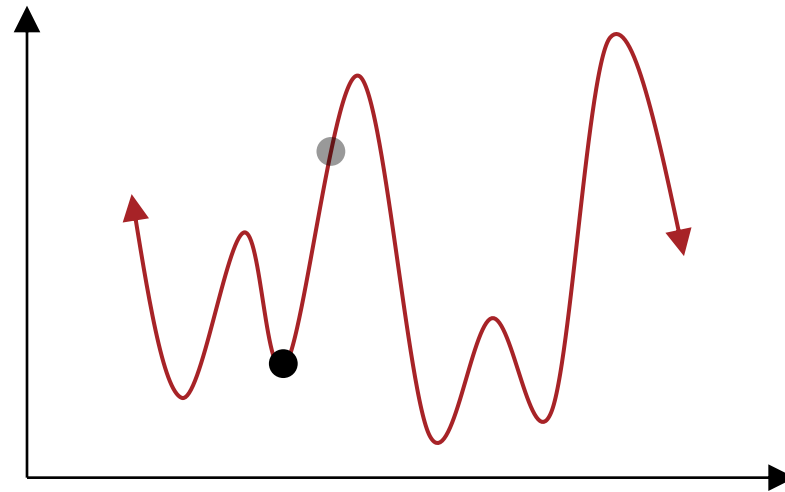
# Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
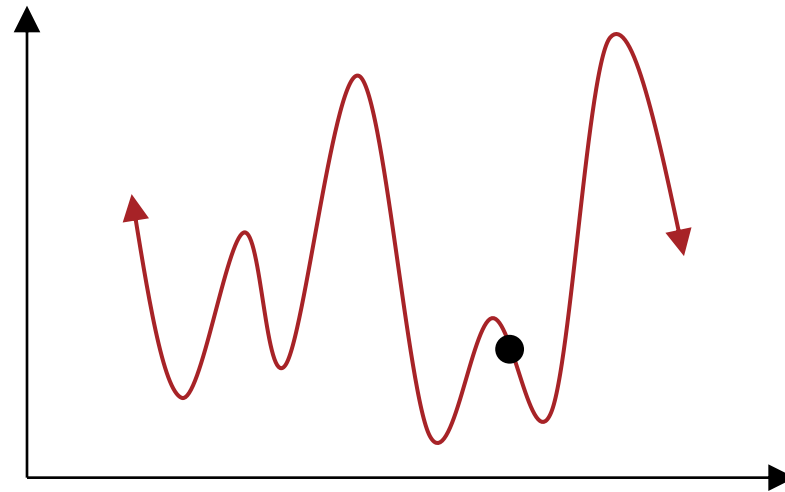  - Works great if the objective function is convex!

# Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
  - Not ideal if the objective function is non-convex...
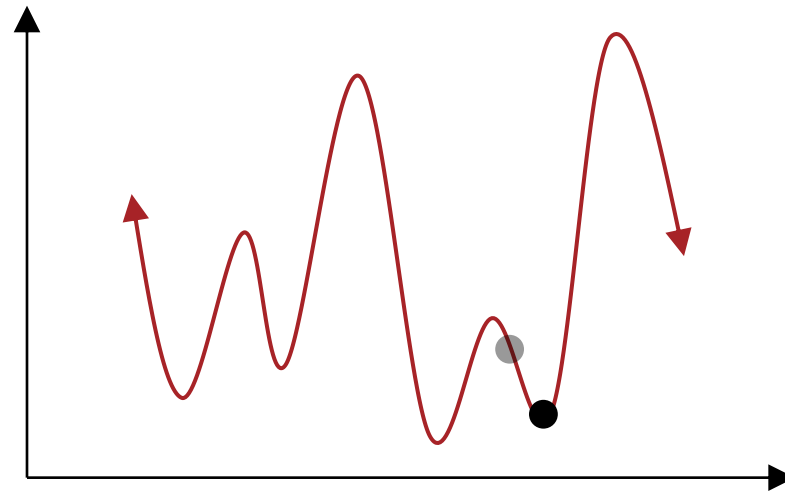
# Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
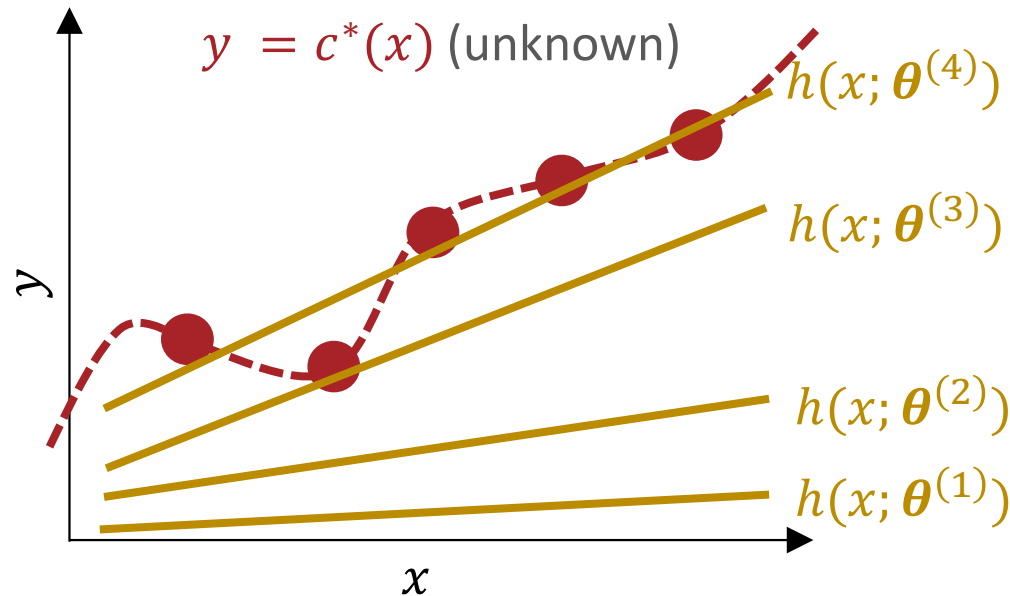  - Not ideal if the objective function is non-convex…

# Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
  - Not ideal if the objective function is non-convex…
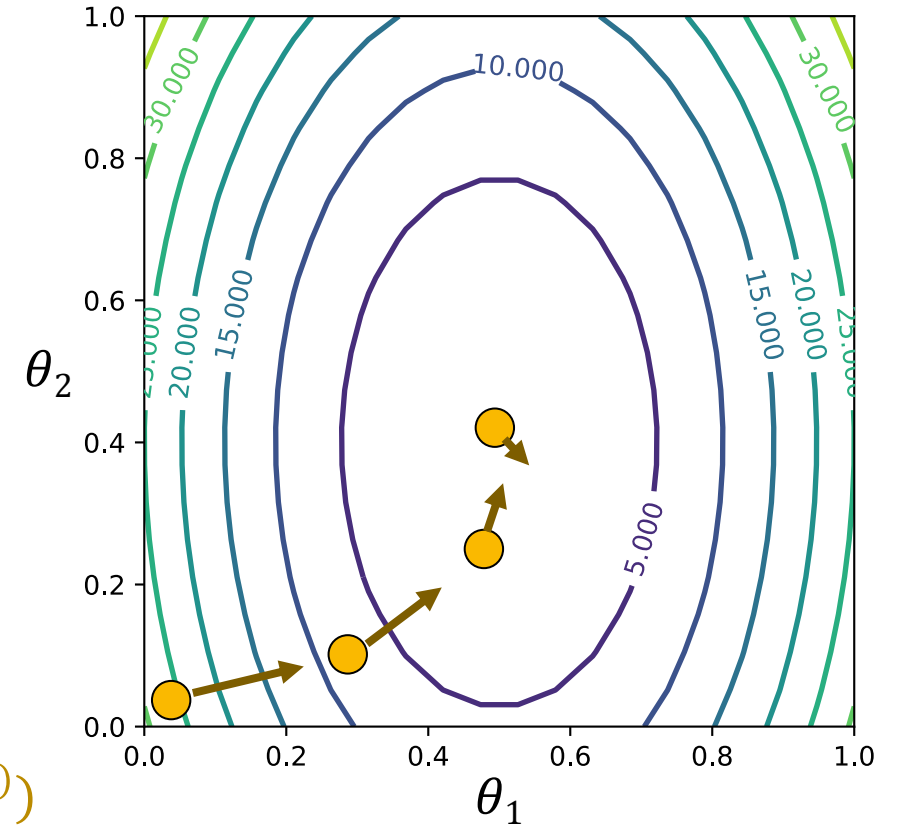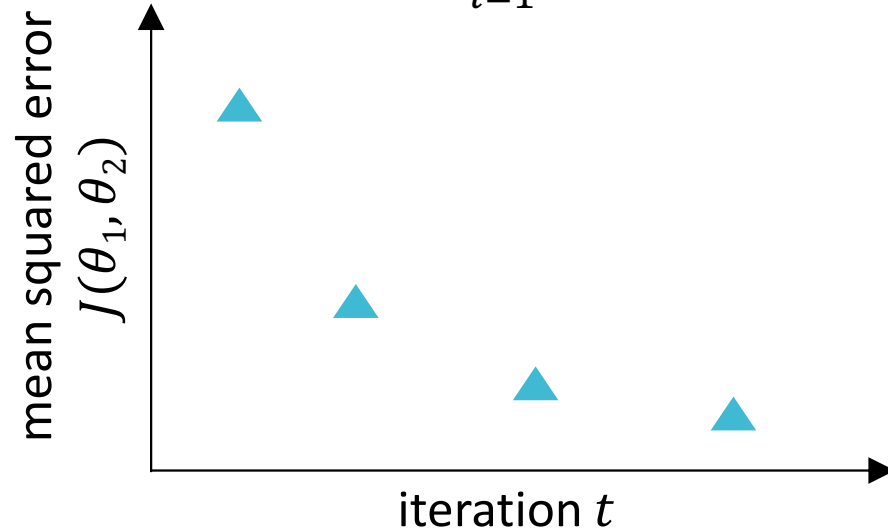
# Gradient Descent & Convexity

- Gradient descent is a local optimization algorithm – it will converge to a local minimum (if it converges)
  - Not ideal if the objective function is non-convex…
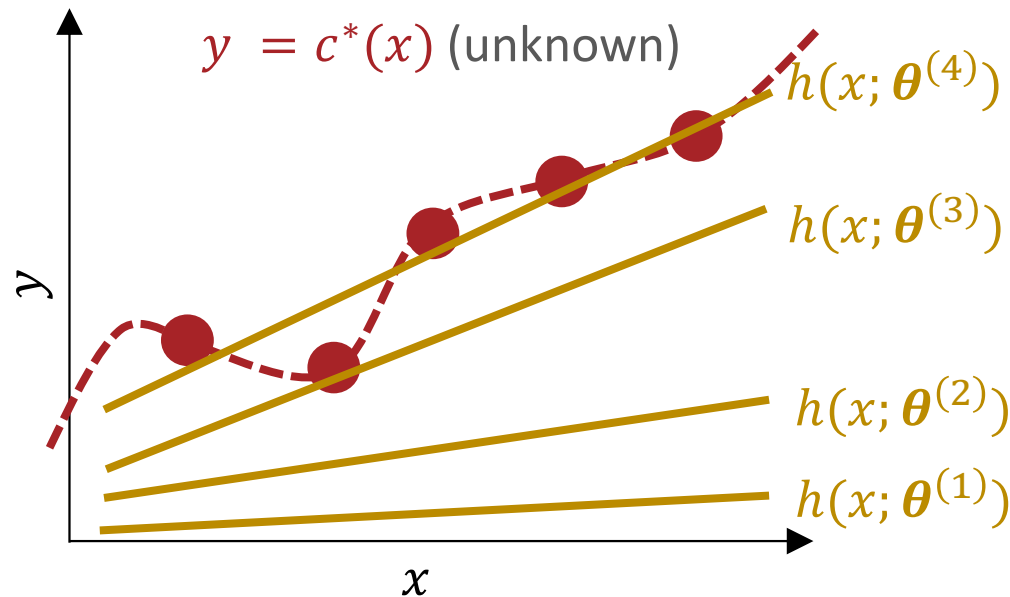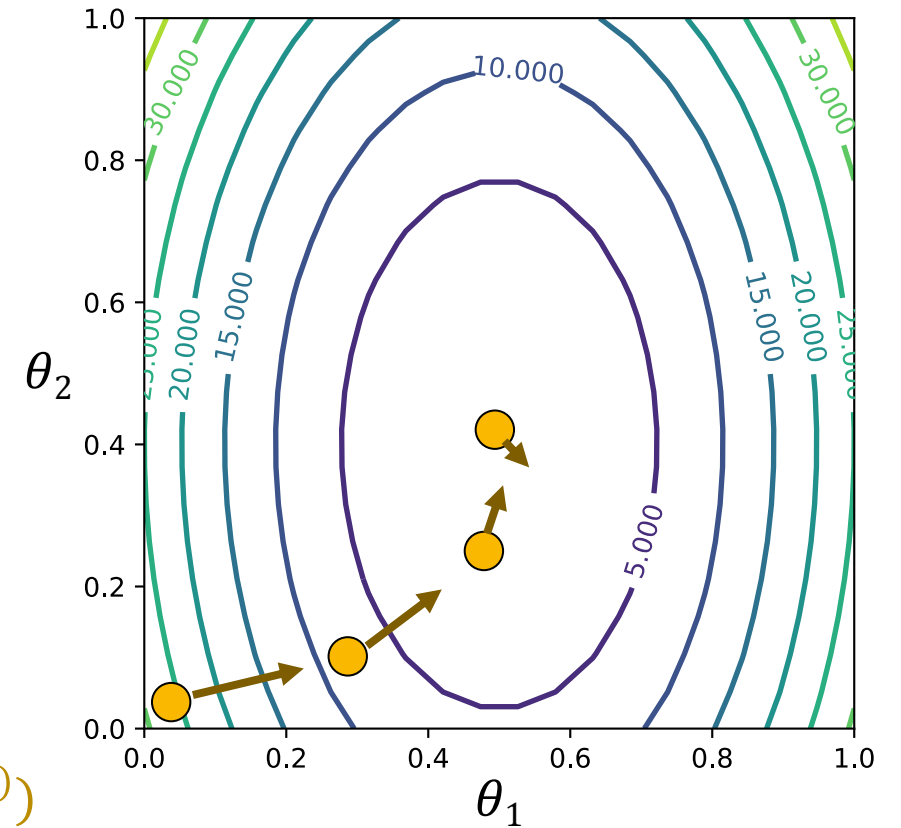
# Why Gradient Descent for Linear Regression?

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



| $t$ | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|-----|-----|-----|-----|
| 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 |

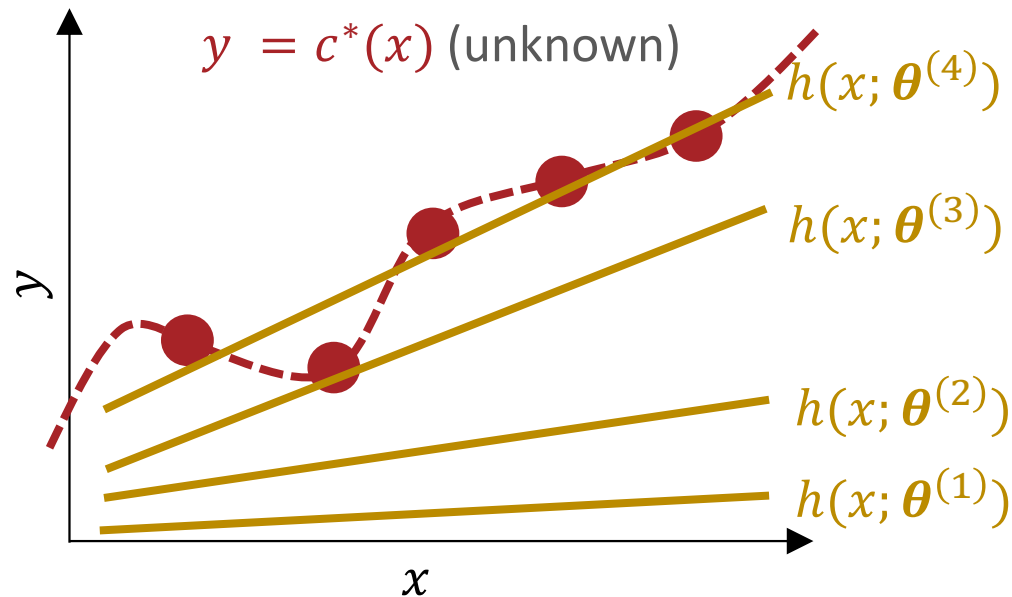The mean squared error is convex (but not always strictly convex)

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



mean squared error $J(\theta_1, \theta_2)$

iteration $t$



$\theta_2$

$\theta_1$



$y = c^*(x)$ (unknown)

$h(x; \boldsymbol{\theta}^{(4)})$

$h(x; \boldsymbol{\theta}^{(3)})$

$h(x; \boldsymbol{\theta}^{(2)})$

$h(x; \boldsymbol{\theta}^{(1)})$

$y$

$x$

| $t$ | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|-----|-----|-----|-----|
| 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 |

Okay, fine but couldn't we do something simpler?

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

mean squared error $J(\theta_1, \theta_2)$

iteration $t$

$y = c^*(x)$ (unknown)

$h(x; \boldsymbol{\theta}^{(4)})$

$h(x; \boldsymbol{\theta}^{(3)})$

$h(x; \boldsymbol{\theta}^{(2)})$

$h(x; \boldsymbol{\theta}^{(1)})$

$y$

$x$

$\theta_2$

$\theta_1$

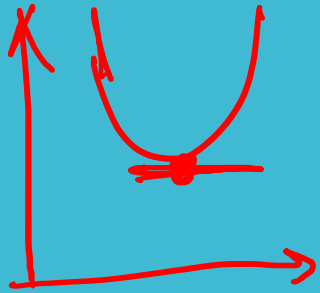| $t$ | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|---|---|---|
| 1 | 0.01 | 0.02 | 25.2 |
| 2 | 0.30 | 0.12 | 8.7 |
| 3 | 0.51 | 0.30 | 1.5 |
| 4 | 0.59 | 0.43 | 0.2 |

# Closed Form Optimization

- Idea: find the *critical points* of the objective function, specifically the ones where $\nabla J(\theta) = \mathbf{0}$ (the vector of all zeros), and check if any of them are local minima

- Notation: given training data $\mathcal{D} = \left\{ \left( \boldsymbol{x}^{(n)}, y^{(n)} \right) \right\}_{n=1}^{N}$

- $X = \begin{bmatrix} 1 & \boldsymbol{x}^{(1)^T} \\ 1 & \boldsymbol{x}^{(2)^T} \\ \vdots & \vdots \\ 1 & \boldsymbol{x}^{(N)^T} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_D^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_D^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \cdots & x_D^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times}$ 

is the *design matrix*

- $\boldsymbol{y} = \left[ y^{(1)}, \ldots, y^{(N)} \right]^T \in \mathbb{R}^N$ is the *target vector*

## Minimizing the Mean Squared Error

$$\text{MSE}: J(\boldsymbol{\theta}) = \frac{1}{N}\sum_{i=1}^{N}\frac{1}{2}\left(y^{(i)} - \boldsymbol{\theta}^T \boldsymbol{x}^{(i)}\right)^2 = \frac{1}{2N}\sum_{i=1}^{N}\left(y^{(i)} - \boldsymbol{\theta}^T x^{(i)}\right)^2$$

$$= \frac{1}{2N}\left(\vec{X}\vec{\theta} - \vec{y}\right)^T\left(\vec{X}\vec{\theta} - \vec{y}\right)$$

$$\nabla_{\theta} J(\theta) = \frac{1}{2N}\left(2\vec{X}^T\vec{X}\vec{\theta} - 2\vec{X}^T\vec{y}\right) = \vec{0}$$

$$\Leftrightarrow \vec{X}^T\vec{X}\vec{\theta} - \vec{X}^T\vec{y} = 0$$

Hessian of $J$ must be PSD.

$$\Leftrightarrow \vec{X}^T\vec{X}\vec{\theta} = \vec{X}^T\vec{y}$$

$$(X^T X)^{-1} X^T X \vec{\theta} = \boxed{(X^T X)^{-1}\vec{X}^T\vec{y}}$$

# Closed Form Optimization

$$\widehat{\boldsymbol{\theta}} = (X^T X)^{-1} X^T \boldsymbol{y}$$

$y = c^*(x)$ (unknown)

$h(x; \widehat{\boldsymbol{\theta}})$



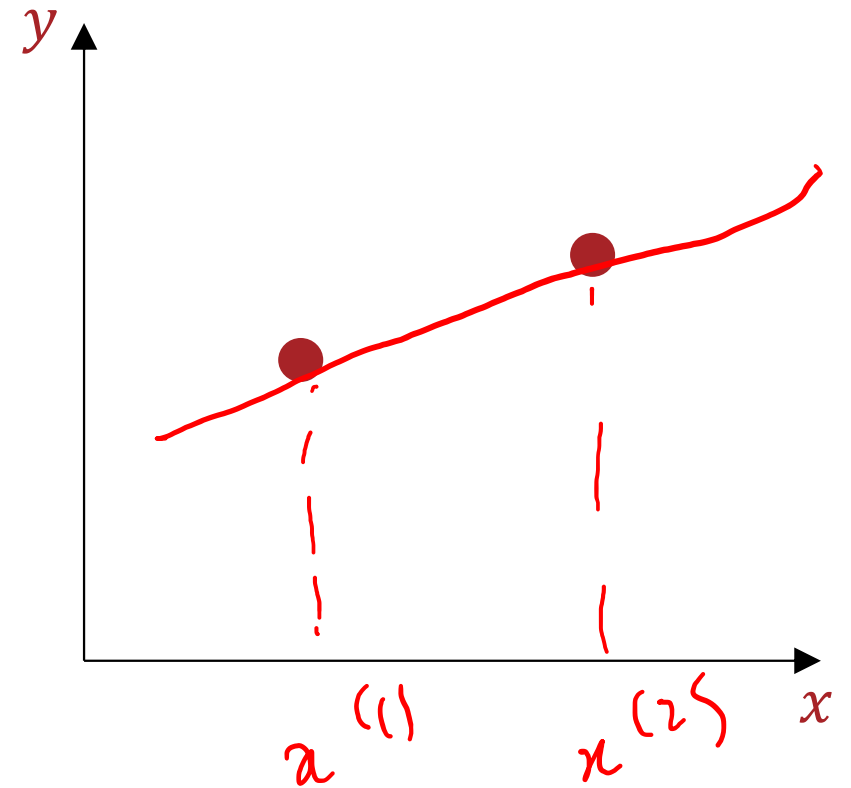| $t$ | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|-----|------------|------------|-------------------------|
| 1   | 0.59       | 0.43       | 0.2                     |

$$\widehat{\boldsymbol{\theta}} = (X^T X)^{-1} X^T \boldsymbol{y}$$

1. Is $X^T X$ invertible?

## Closed Form Solution

2. If so, how computationally expensive is inverting $X^T X$?

# Linear Regression: Uniqueness

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters $\theta$) are there for the given dataset?
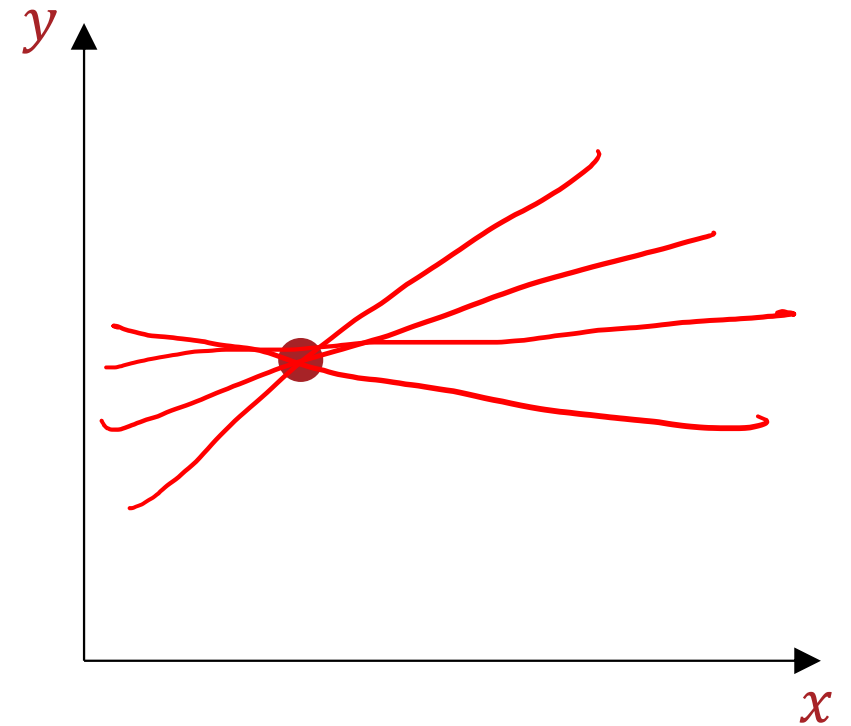
# Linear Regression: Uniqueness

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters $\theta$) are there for the given dataset?
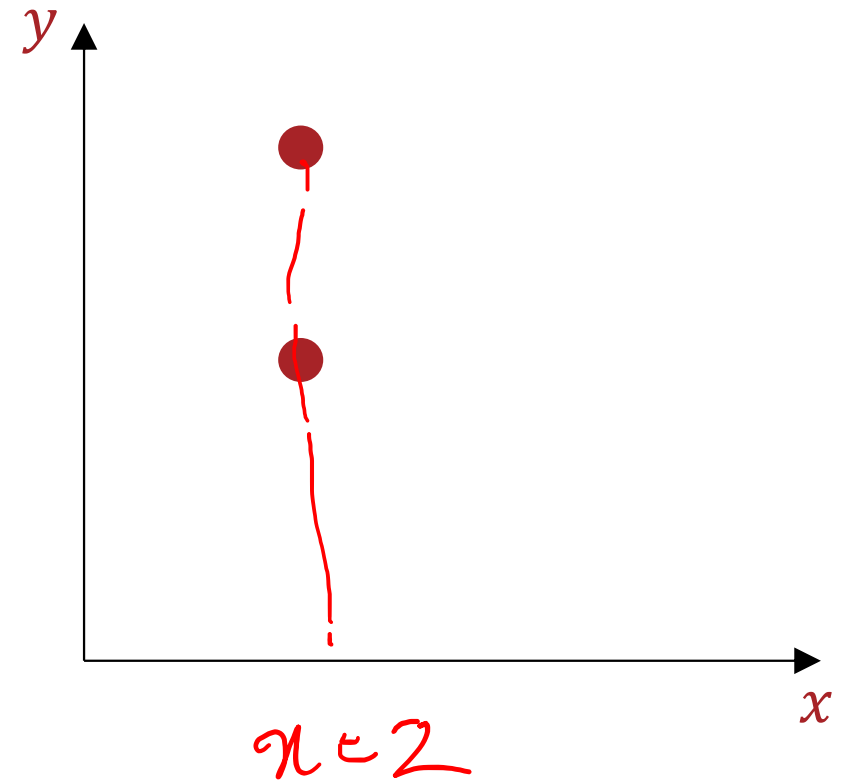
# Linear Regression: Uniqueness

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters $\theta$) are there for the given dataset?

Poll Question 3

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters $\theta$) are there for the given dataset?

A. -1 **(TOXIC)**     B. 0     C. 1     D. 2     E. ∞

# Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters $\theta$) are there for the given dataset?
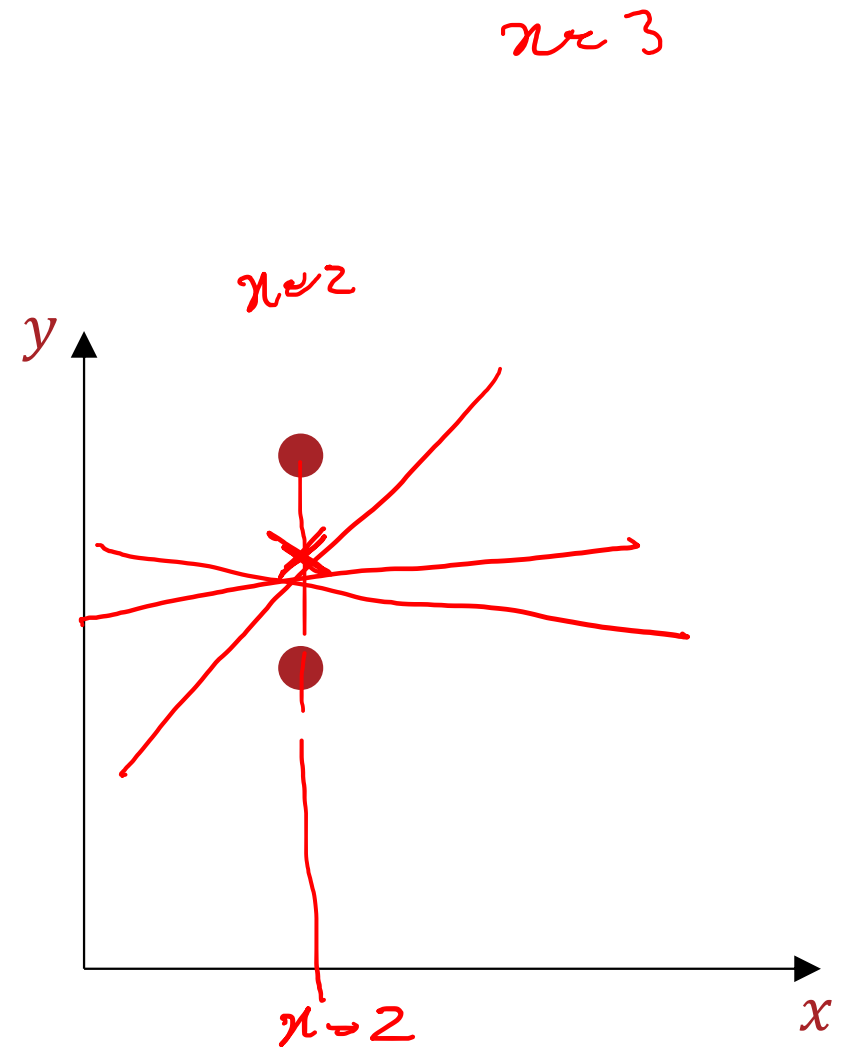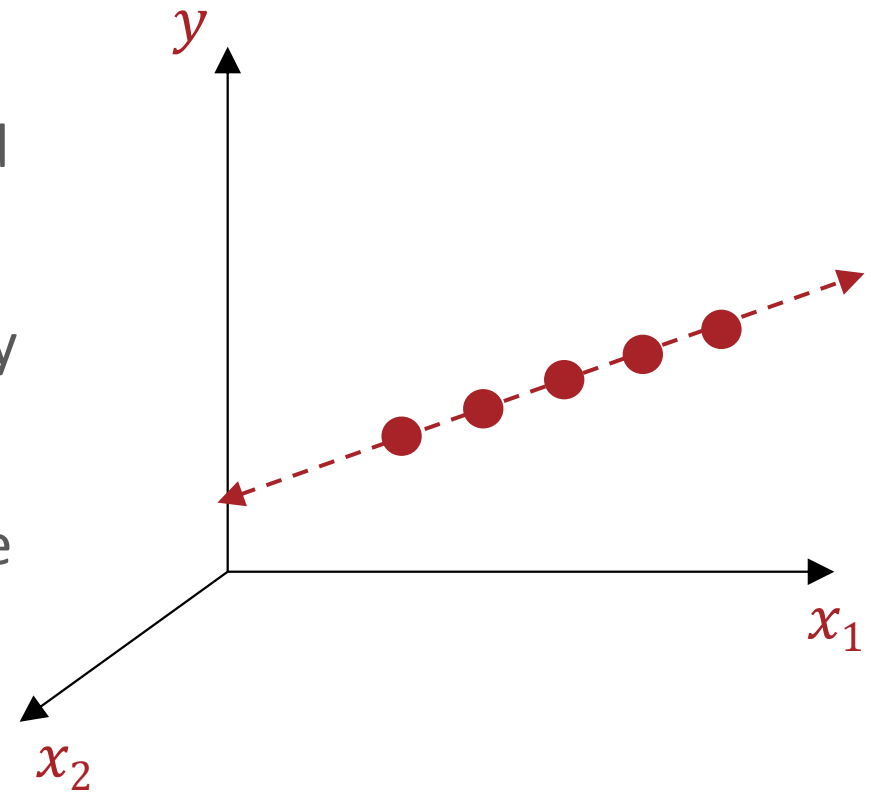
# Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters $\theta$) are there for the given dataset?
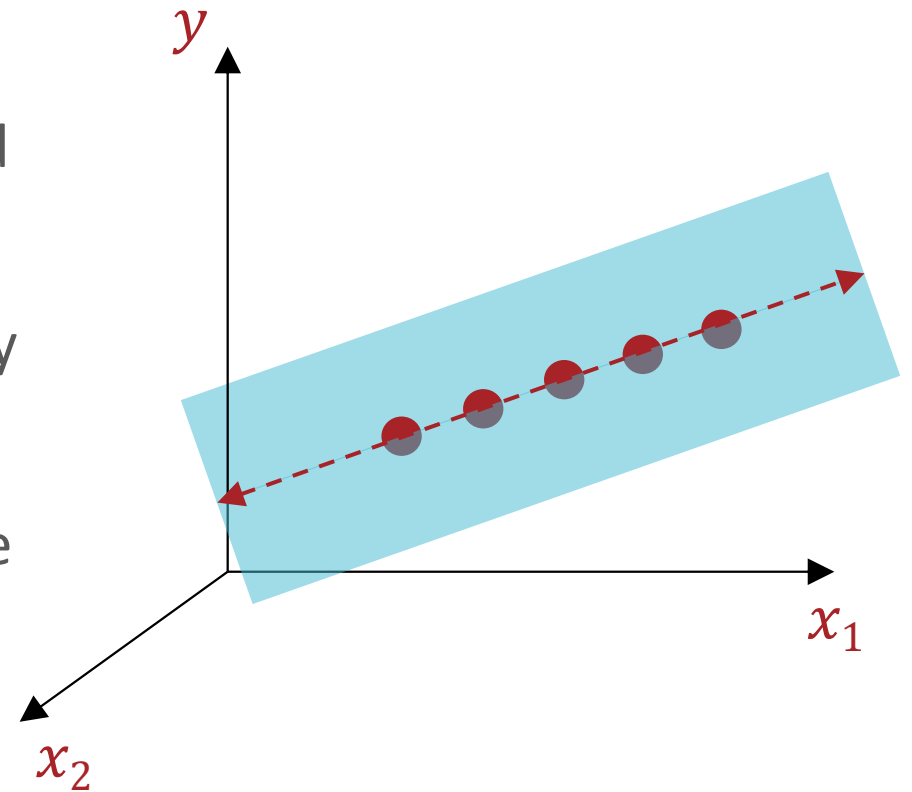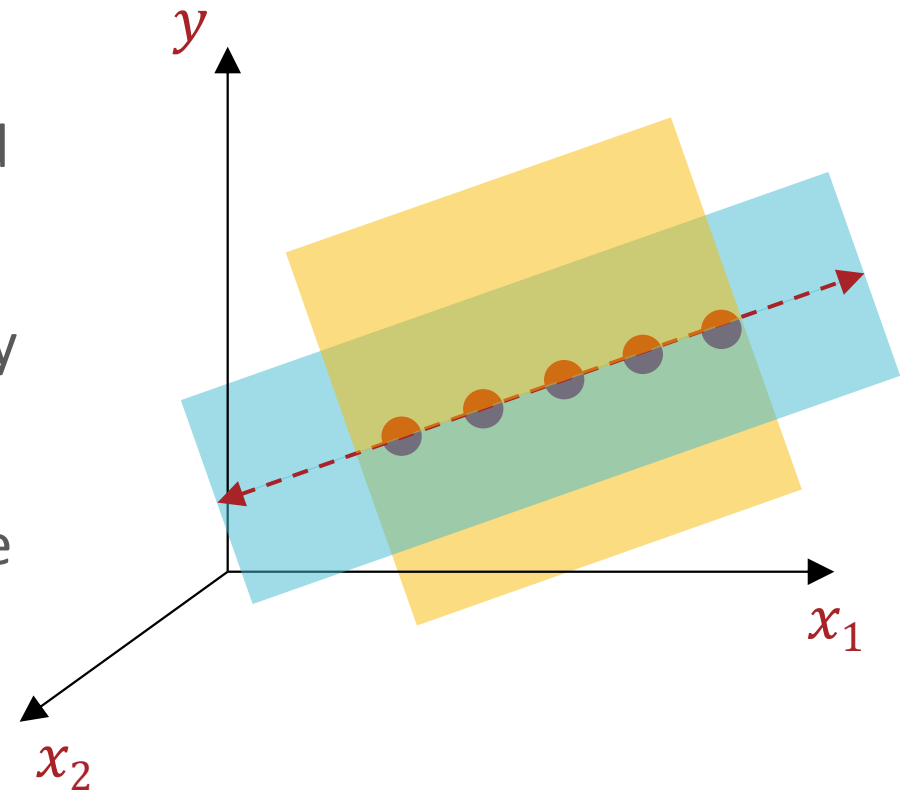
# Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters $\theta$) are there for the given dataset?

$$\widehat{\boldsymbol{\theta}} = (X^T X)^{-1} X^T \boldsymbol{y}$$

# Closed Form Solution

1. Is $X^T X$ invertible?

2. If so, how computationally expensive is inverting $X^T X$?

## Closed Form Solution

$$\widehat{\boldsymbol{\theta}} = (X^T X)^{-1} X^T \boldsymbol{y}$$

1. Is $X^T X$ invertible?
   - When $N \gg D + 1$, $X^T X$ is (almost always) full rank and therefore, invertible!
   - If $X^T X$ is not invertible (occurs when one of the features is a linear combination of the others), then there are infinitely many solutions
2. If so, how computationally expensive is inverting $X^T X$?
   - $X^T X \in \mathbb{R}^{D+1 \times D+1}$ so inverting $X^T X$ takes $O(D^3)$ time…
     - Computing $X^T X$ takes $O(ND^2)$ time
   - Can use gradient descent to (potentially) speed things up when $N$ and $D$ are large!

# Linear Regression Learning Objectives

You should be able to…

- Design k-NN Regression and Decision Tree Regression
- Implement learning for Linear Regression using gradient descent or closed form optimization
- Choose a Linear Regression optimization technique that is appropriate for a particular dataset by analyzing the tradeoff of computational complexity vs. convergence speed
- Identify situations where least squares regression has exactly one solution or infinitely many solutions