



Announcements

Assignments

- HW3
 - Mon, 9/28, 11:59 pm

Midterm 1

- Mon, 10/5
- See Piazza for details
- Fill out swap-section / conflict form by Friday

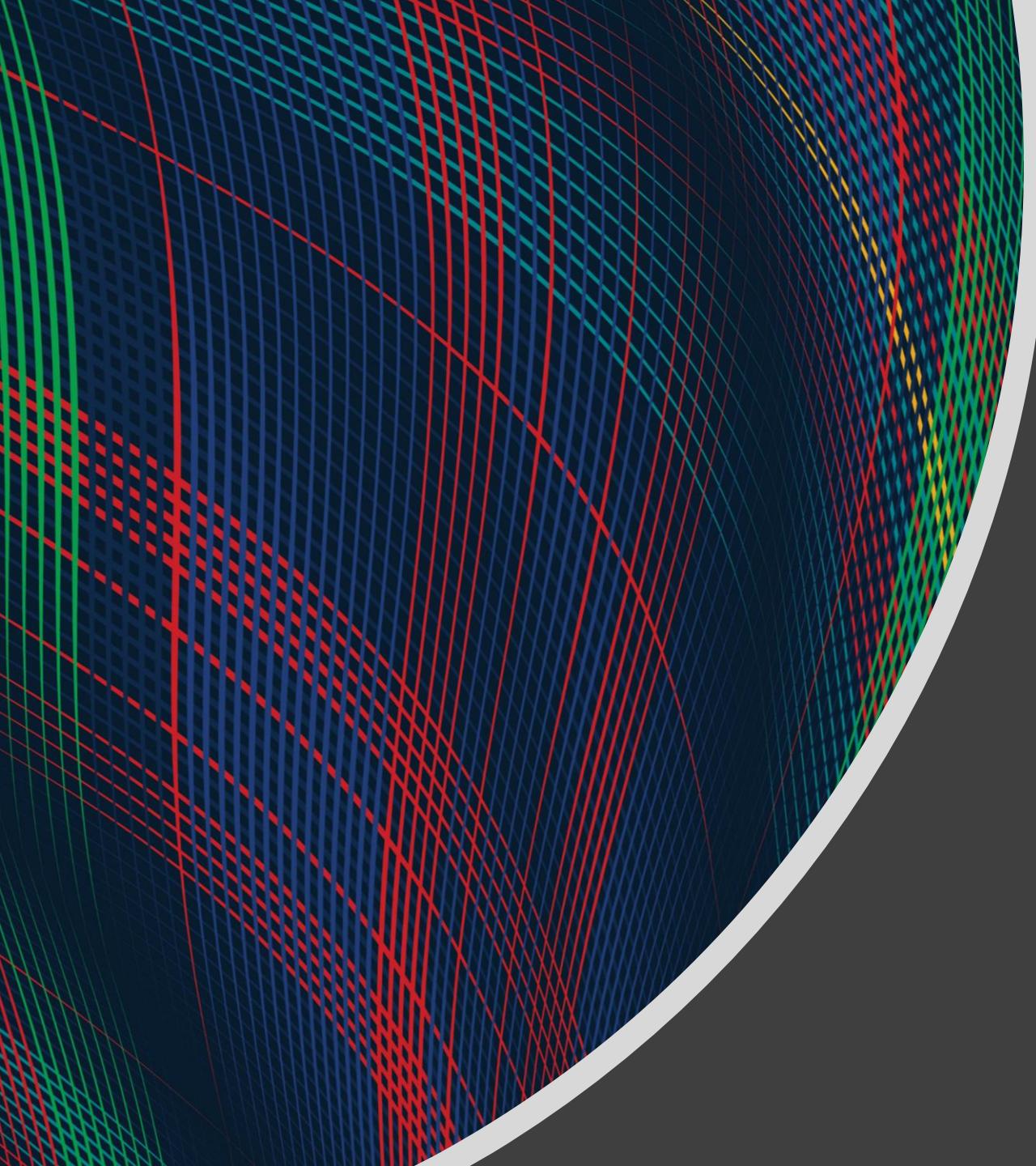
Plan

Last time

- Regression
- Linear regression
- Optimization for linear regression

Today

- Optimization for linear regression
 - Linear and convex function
 - (Batch) Gradient descent
 - Closed-form solution
 - Stochastic gradient descent 



Introduction to Machine Learning

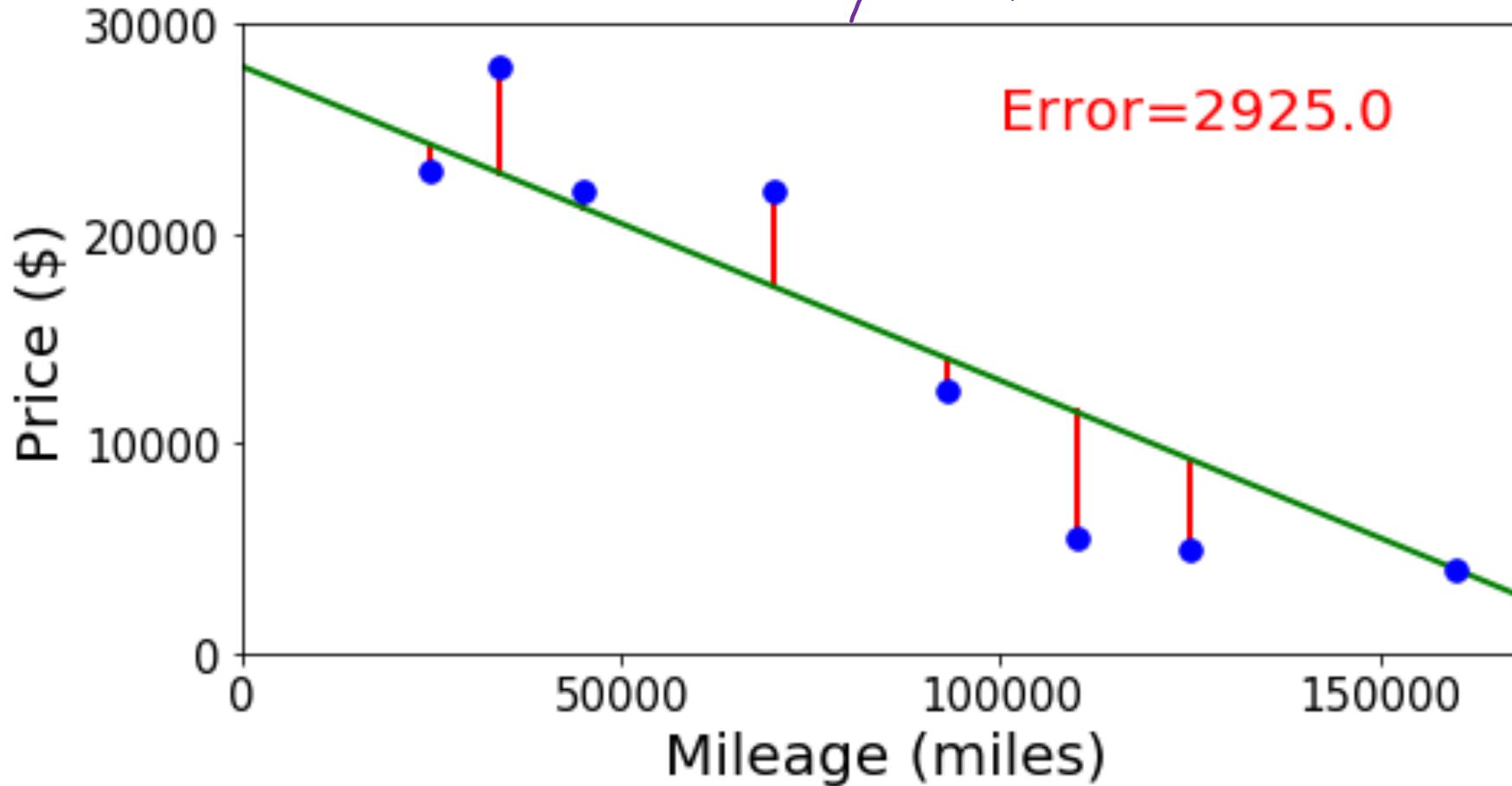
Linear Regression and Optimization

Instructor: Pat Virtue

Linear Regression

Selling my car

$$y = mx + b$$
$$y = w x + b$$
$$y = \omega_1 x + \omega_0$$
$$y = \theta_1 x + \theta_0$$



Linear in Higher Dimensions

$$1\text{-D} \quad y = w_0 x + b$$
$$2\text{-D} \quad y = w_0 x_1 + w_1 x_2 + b$$

What are these linear shapes called for 1-D, 2-D, 3-D, M-D input?

$$x \in \mathbb{R}$$

$$x \in \mathbb{R}^2$$

$$x \in \mathbb{R}^3$$

$$x \in \mathbb{R}^M$$

$$\rightarrow y = w^T x + b$$

line

plane

hyperplane

hyperplane

$$w^T x + b = 0$$

point

line

plane

hyperplane

$$w^T x + b \geq 0$$

halfline

halfplane

halfspace

halfspace

Linear Function

Linear function

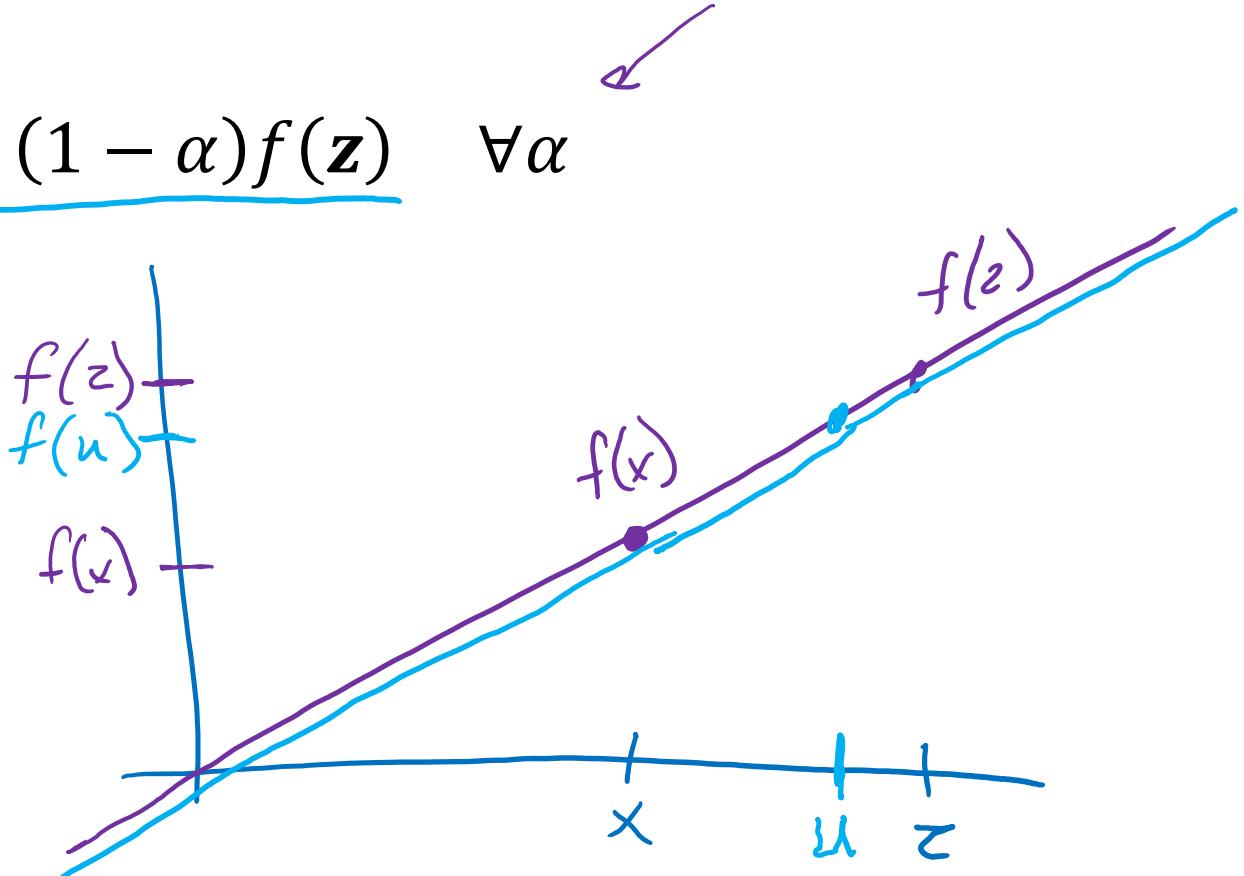
If $f(x)$ is linear, then:

- $f(x + z) = f(x) + f(z)$

- $f(\underline{\alpha}x) = \alpha f(x) \quad \forall \alpha$

→ ■ $f(\alpha x + (1 - \alpha)z) = \alpha f(x) + (1 - \alpha)f(z)$ $\forall \alpha$

$\alpha = 0.25$



Piazza Poll 1

Based on the following definition of a linear, is the equation for a line,
 $y = wx + b$, linear? Example: $y = 3x + 5$

$$y = \underbrace{wx}_\text{linear} + b$$

affine

$$y = \underbrace{\vec{w}^T \vec{x}}_\text{linear} + b$$

affine

$f(x)$ is linear if and only if:

- $f(x+z) = f(x) + f(z)$ and
- $f(\alpha x) = \alpha f(x) \quad \forall \alpha \leftarrow$

$$\alpha = 7$$

$$\begin{aligned}f(7x) &= 3(7x) + 5 \\&= 21x + 5 \\&\neq 7(3x + 5)\end{aligned}$$

Yes

$$62 \rightarrow 32\%$$

No

$$32 \rightarrow 67\%$$



Linear Regression

Linear algebra formulation

$$y = \vec{w}^T \vec{x}_{\text{orig}} + b$$

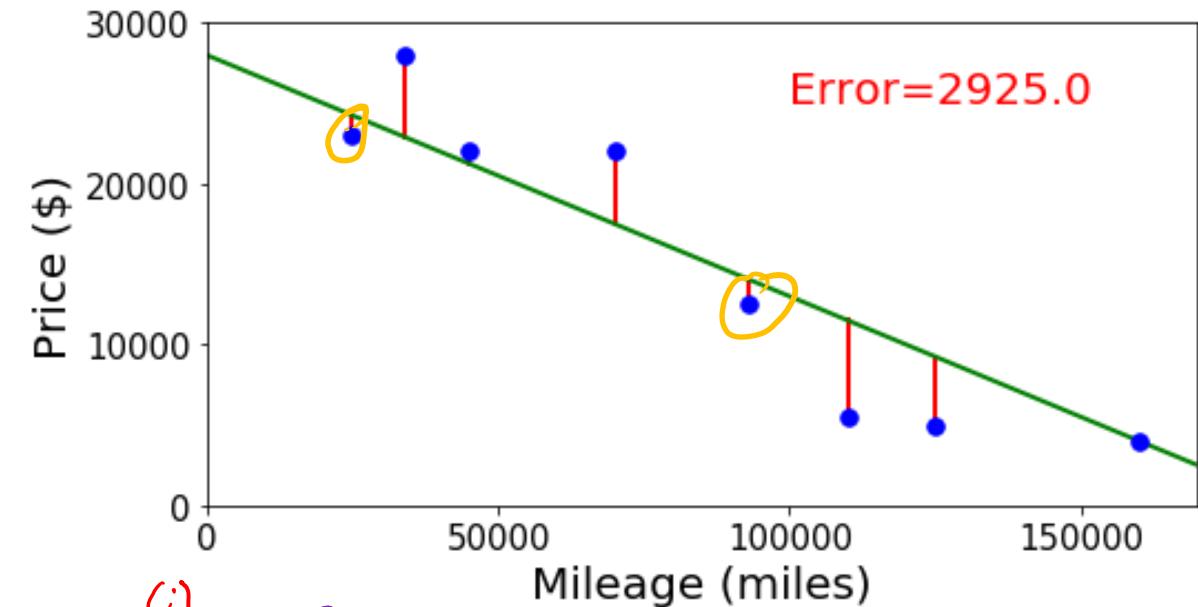
$$y^{(i)} = \vec{w}^T \vec{x}_{\text{orig}}^{(i)} + b$$

$$\underline{y^{(i)} = \vec{\theta}^T \vec{x} = \vec{x}^T \vec{\theta}}$$

$$\rightarrow \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} = \begin{bmatrix} -x^{(1)T} \\ \vdots \\ -x^{(N)T} \end{bmatrix} \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_N \end{bmatrix}$$

$$\rightarrow \vec{y} = \vec{X} \vec{\theta}$$

Design Matrix \vec{X}



$$\vec{x}_{\text{orig}}^{(i)} \in \mathbb{R}^2$$

$$\vec{x}_{\text{orig}}^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \end{bmatrix}$$

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix}$$

$$\vec{x}^{(i)} \in \mathbb{R}^3$$

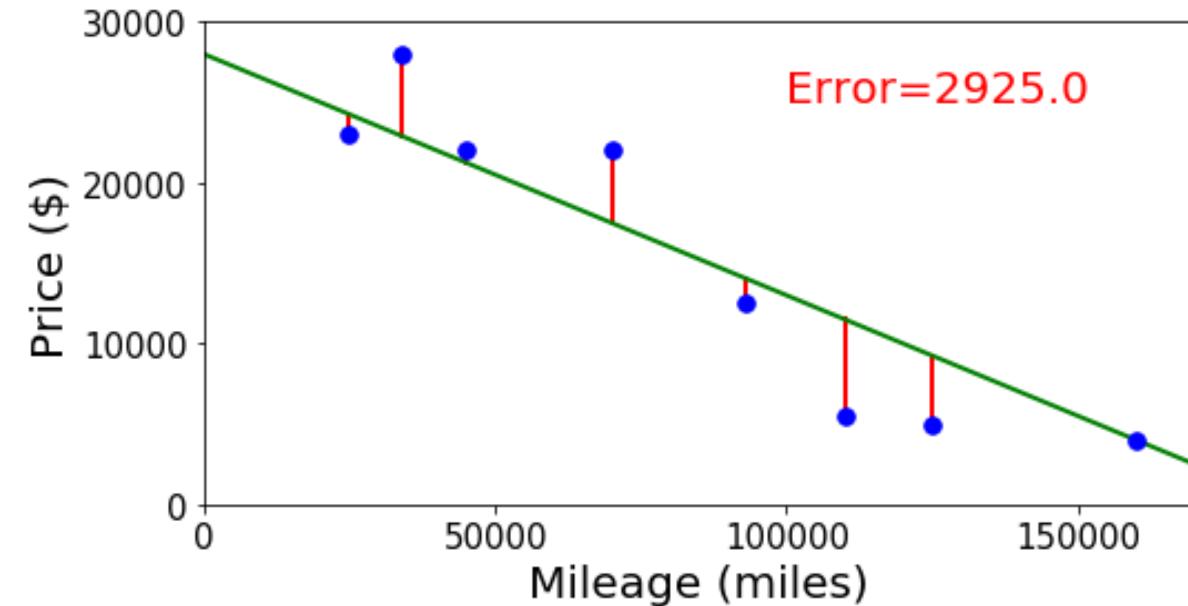
$$\vec{x}^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ x_2^{(i)} \end{bmatrix}$$

$$\vec{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$

Linear Regression

Error and objectives

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$
$$\hat{y}^{(i)} = w x^{(i)} + b$$



$$J(w_1, w_2, b) = \frac{1}{n} \sum (y^{(i)} - \hat{y}^{(i)})^2$$
$$\hat{y}^{(i)} = w_1 x_1^{(i)} + w_2 x_2^{(i)} + b$$

$$J(w_1, \dots, w_M, b) =$$
$$\hat{y}^{(i)} = \sum_{j=1}^M w_j x_j^{(i)} + b$$

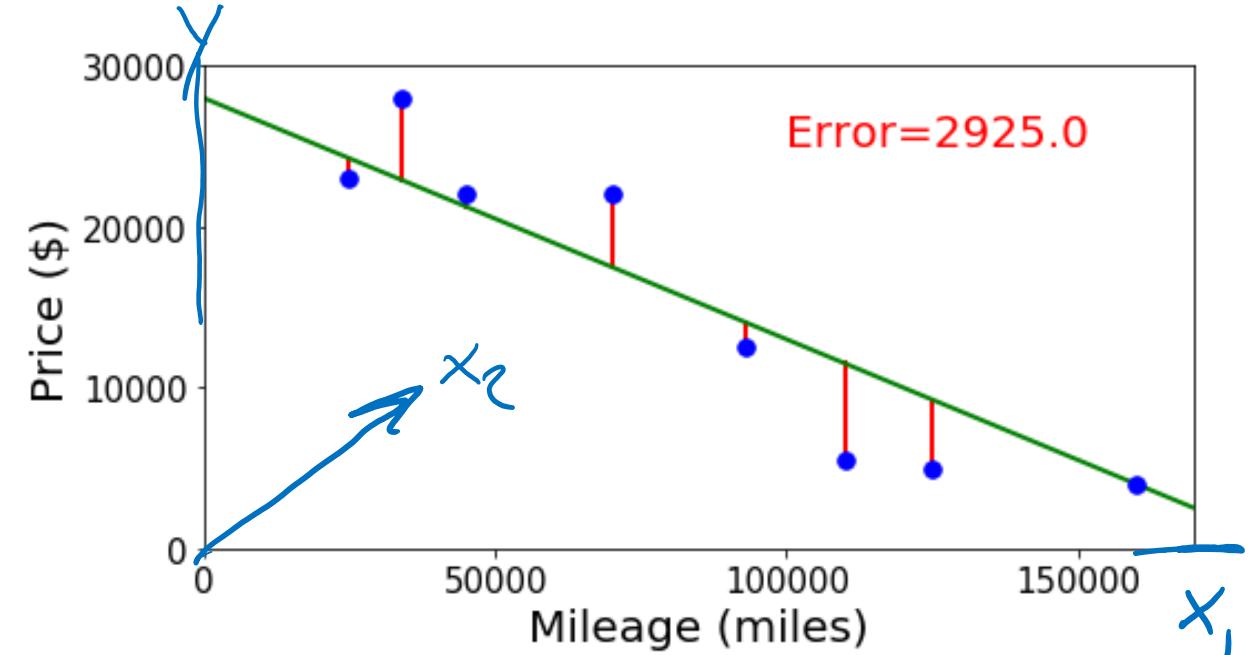
Linear Regression

Linear algebra formulation

$$\vec{\theta} = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} \quad \vec{x}^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ x_2^{(i)} \end{bmatrix}$$

$$\vec{X} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ & \vdots & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} \end{bmatrix}$$

$$J(\vec{\theta}) = \frac{1}{N} \left\| \vec{y} - \vec{y} \right\|_2^2 = \frac{1}{N} \left\| \vec{y} - \vec{X} \vec{\theta} \right\|_2^2$$



$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

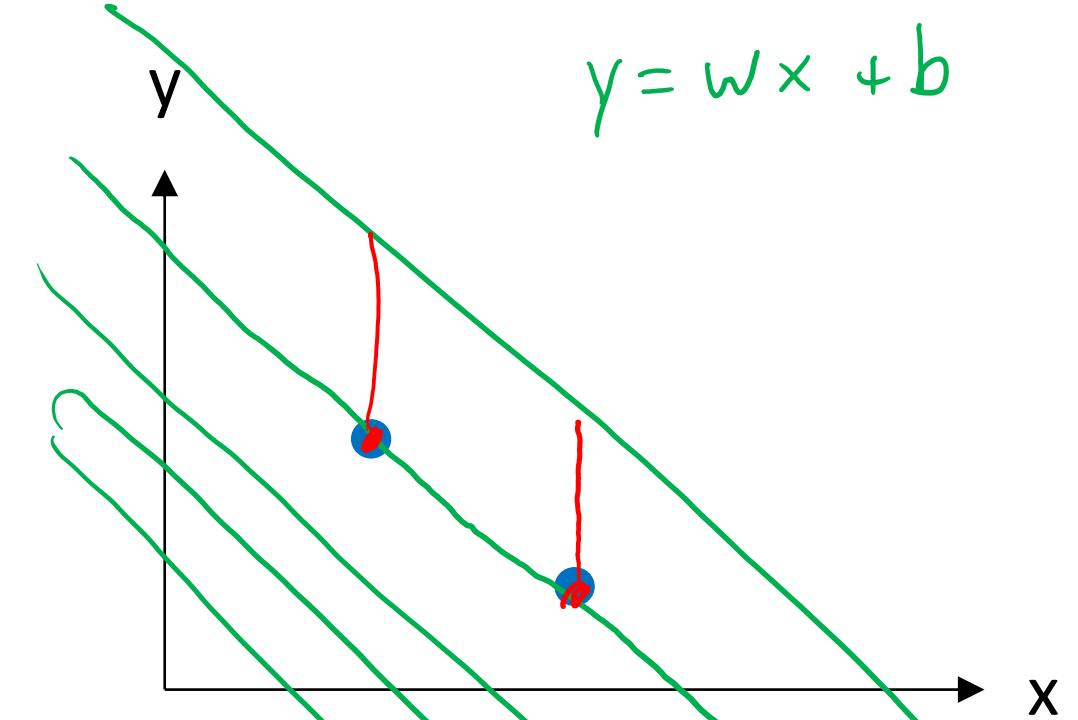
ℓ_2 -norm squared

$$\left\| \vec{z} \right\|_2^2 = \sum_{i=1}^N (z_i)^2$$

Previous Piazza Poll

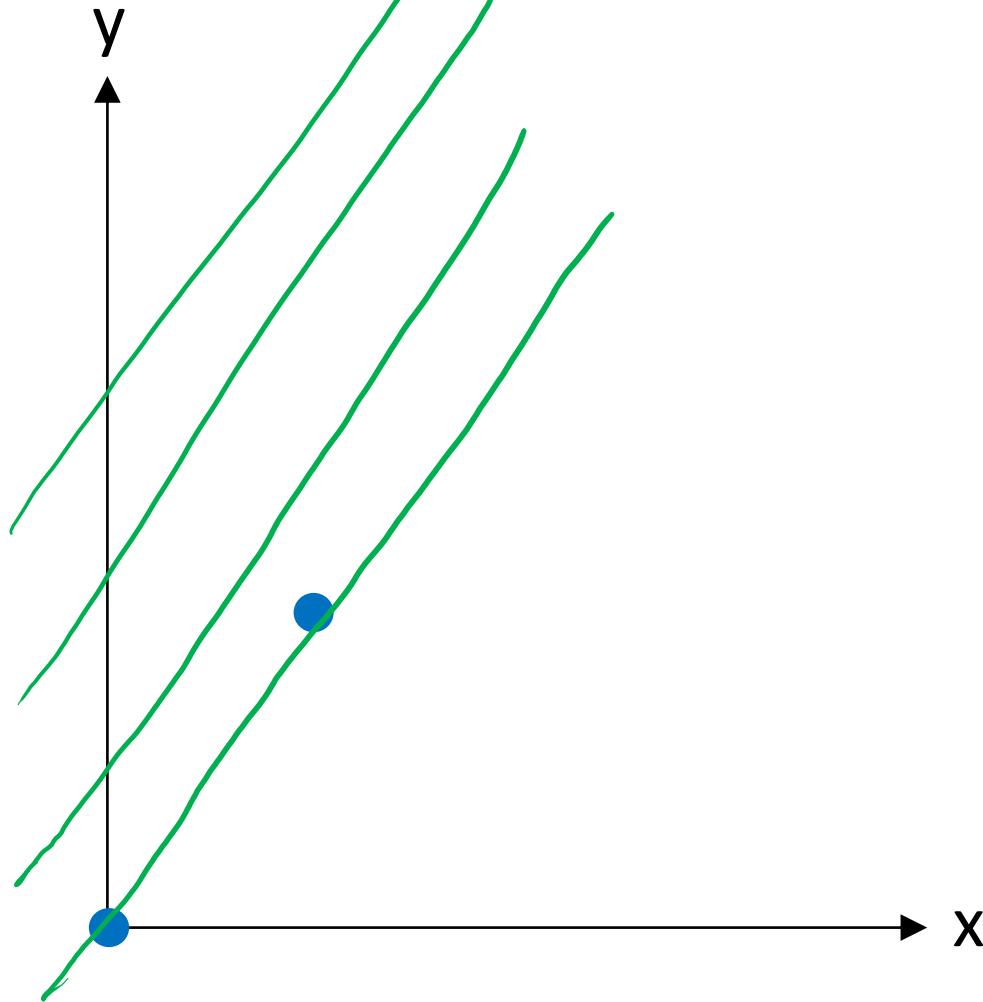
For fixed data and fixed slope, w , what shape do we get by plotting MSE objective vs intercept, b ?

- A. Line
- B. Plane
- C. Half-plane
- D. Convex Parabola (U-shape)
- E. Concave parabola (up-side-down U)
- F. None of the above

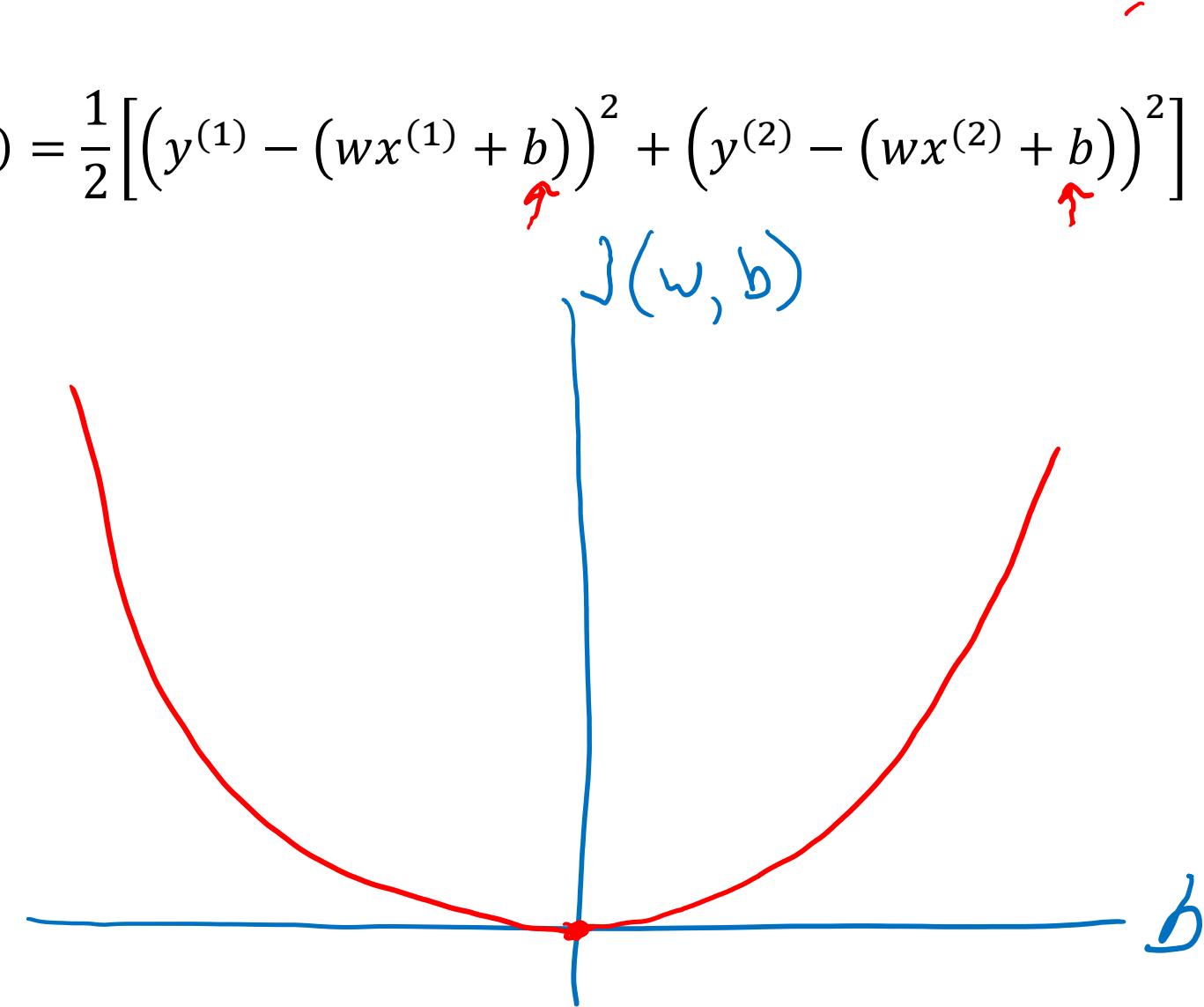


Linear Regression

Optimizing the objective



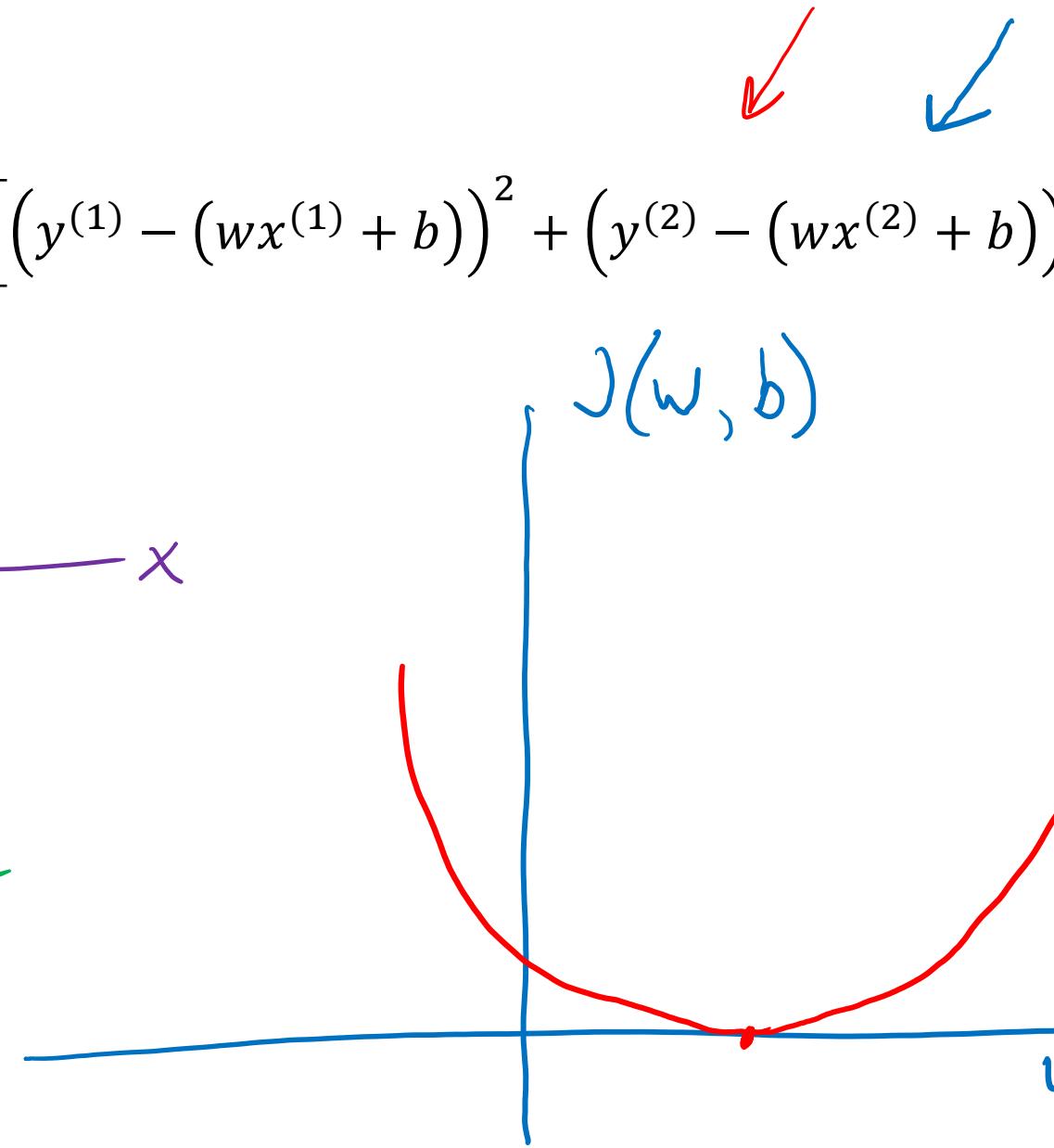
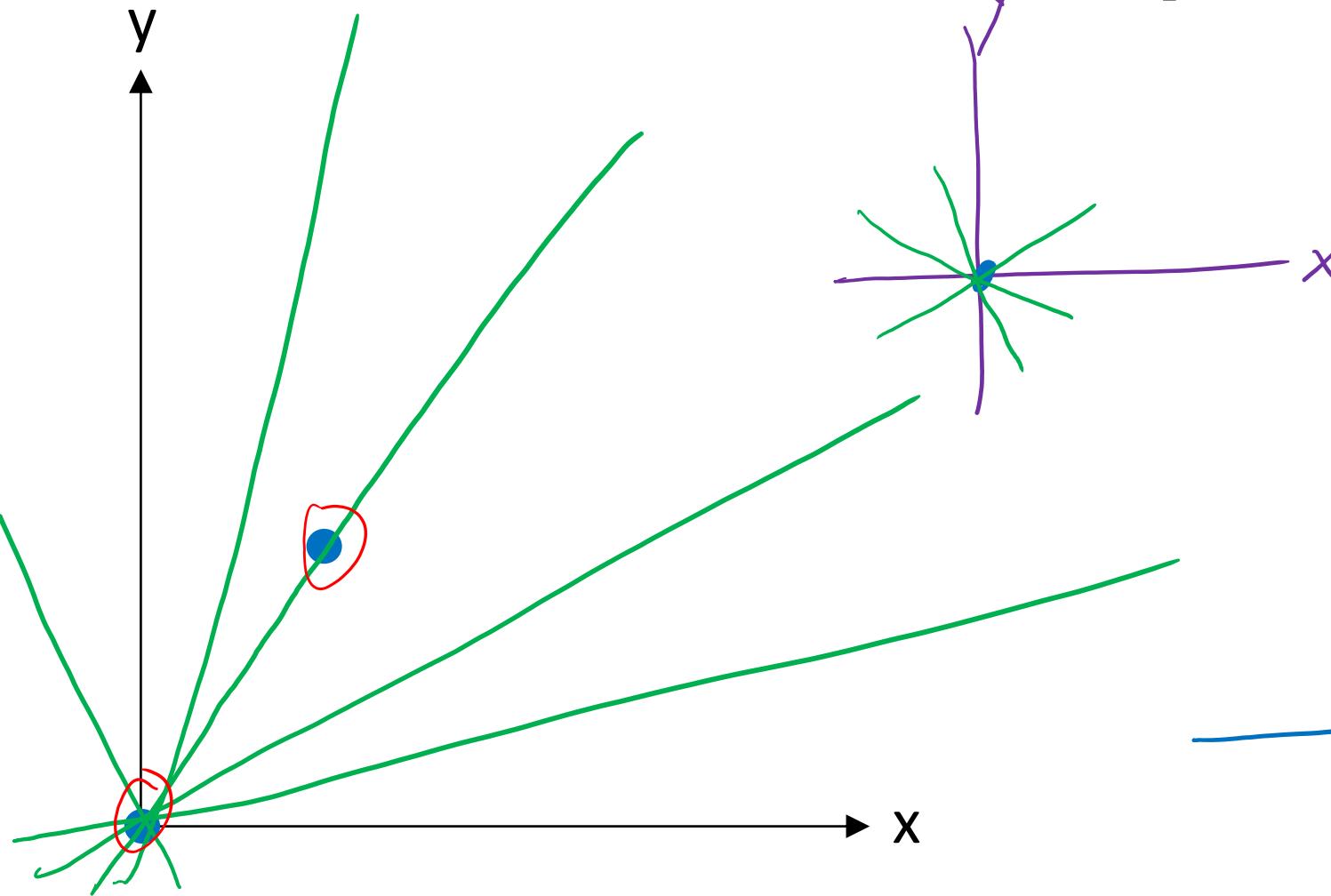
$$J(w, b) = \frac{1}{2} \left[(y^{(1)} - (wx^{(1)} + b))^2 + (y^{(2)} - (wx^{(2)} + b))^2 \right]$$



Linear Regression

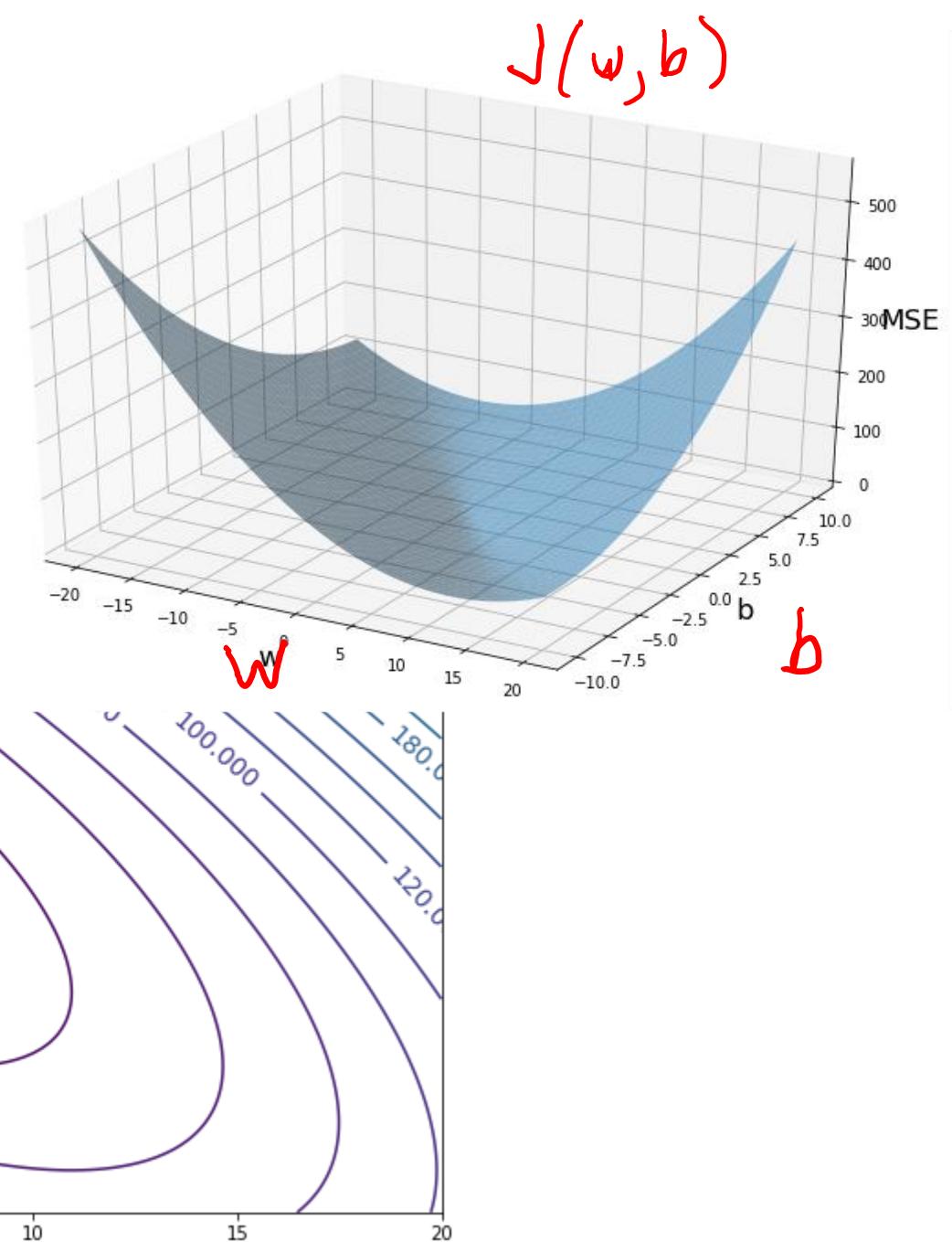
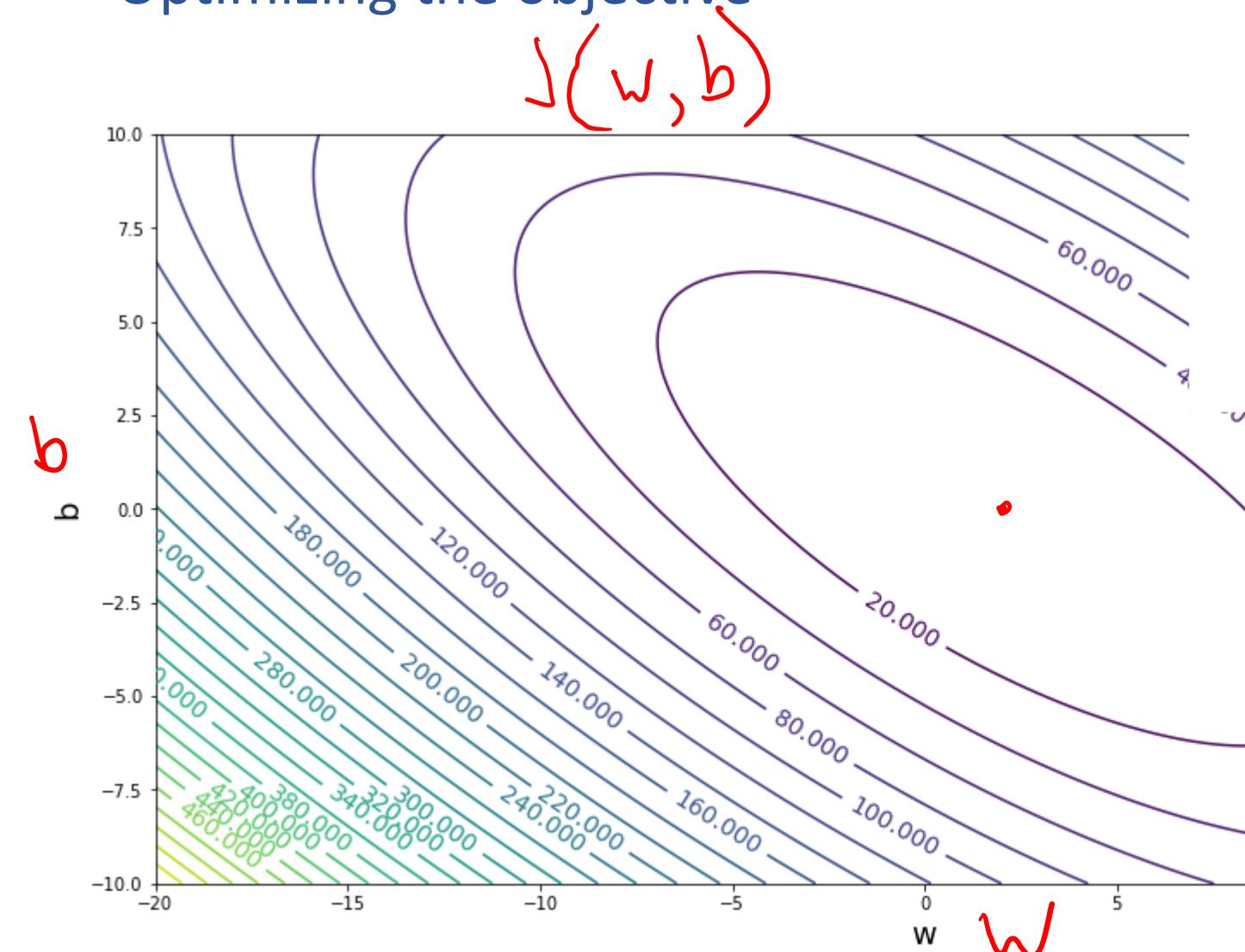
Optimizing the objective

$$J(w, b) = \frac{1}{2} \left[(y^{(1)} - (wx^{(1)} + b))^2 + (y^{(2)} - (wx^{(2)} + b))^2 \right]$$



Linear Regression

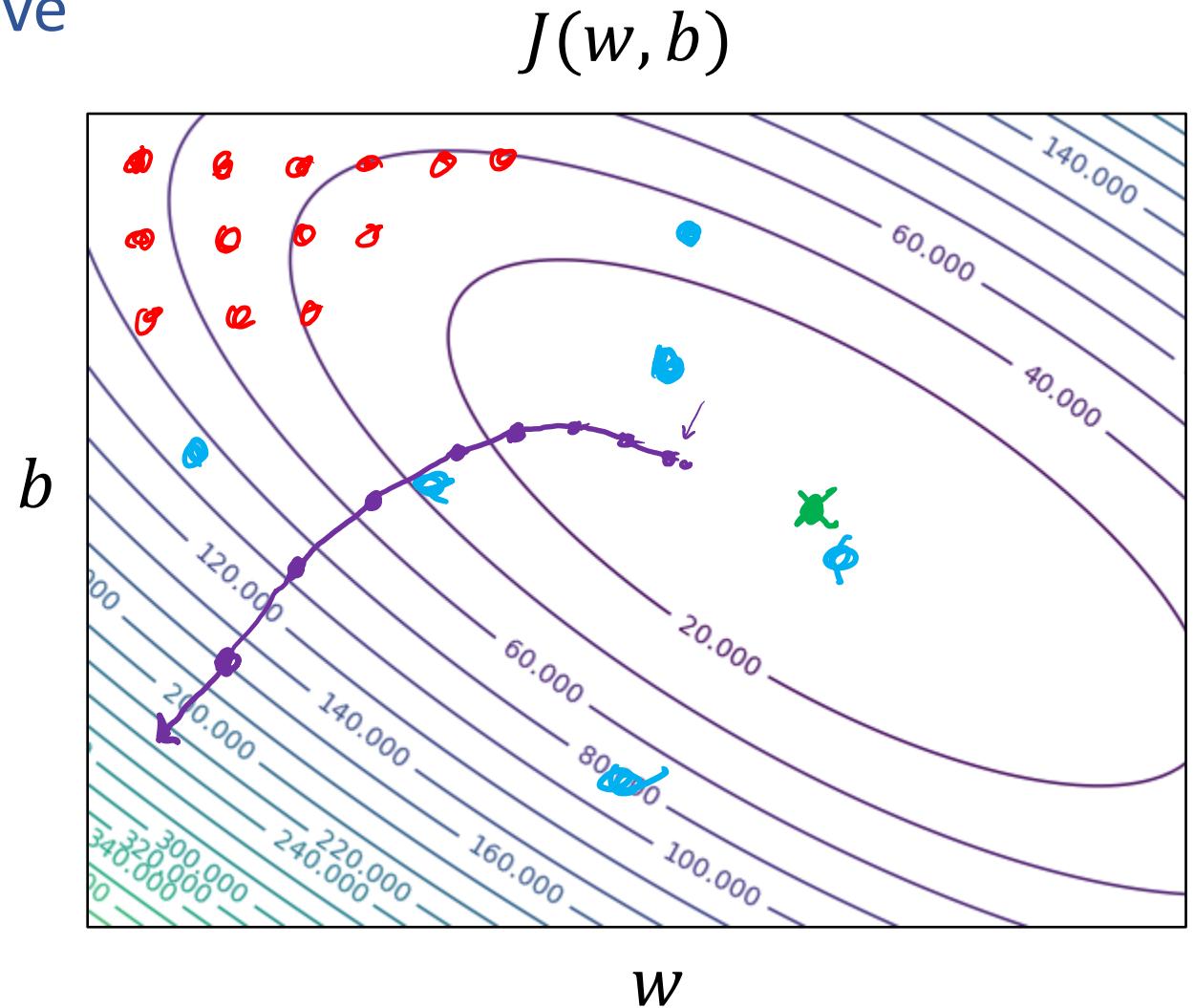
Optimizing the objective



Linear Regression

Methods for optimizing the objective

- Grid search
- Random search
- Closed-form solution
- (Batch) Gradient descent
- Stochastic gradient descent



Optimization

Linear function

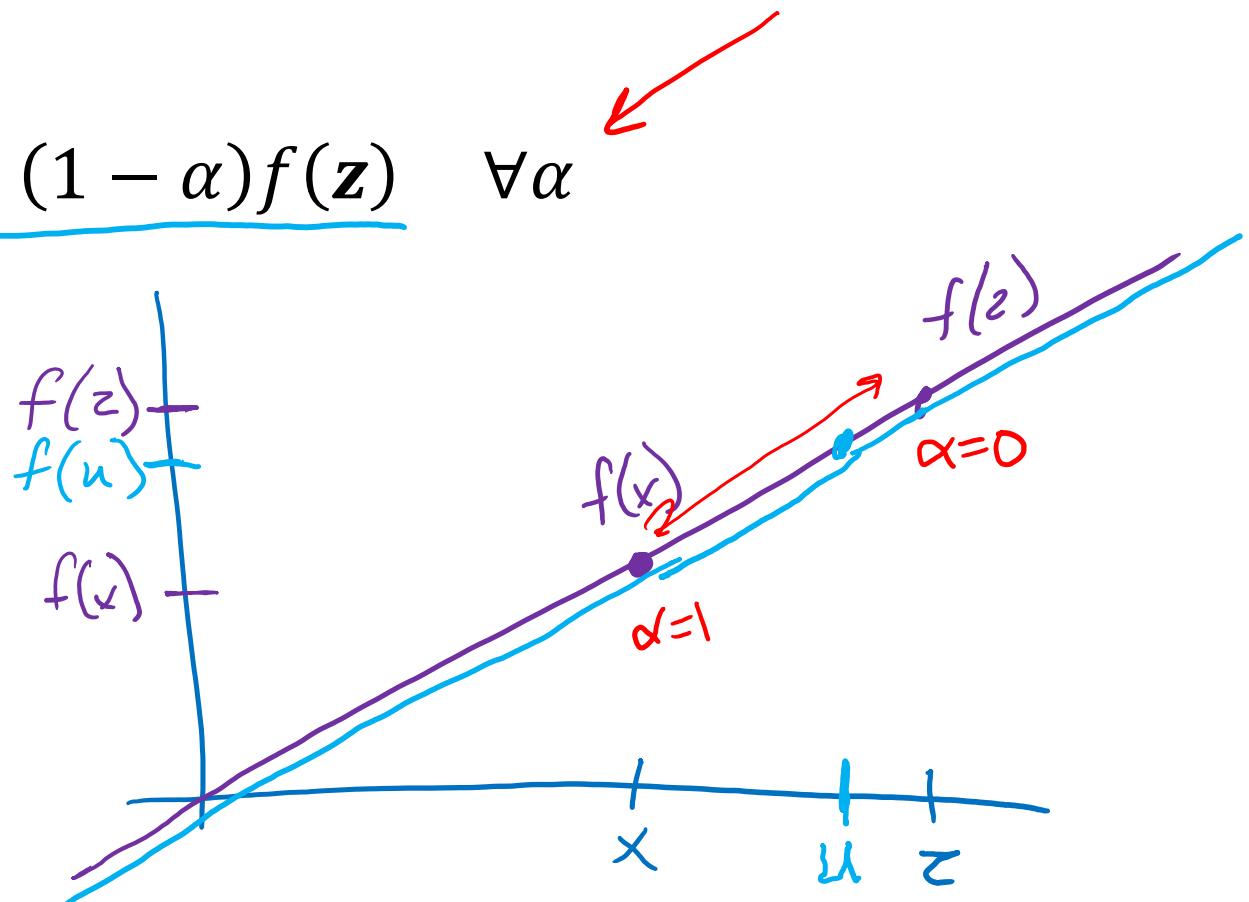
If $f(x)$ is linear, then:

- $f(x + z) = f(x) + f(z)$

- $f(\underline{\alpha}x) = \alpha f(x) \quad \forall \alpha$

→ ■ $f(\alpha x + (1 - \alpha)z) = \alpha f(x) + (1 - \alpha)f(z)$ $\forall \alpha$

$\alpha = 0.25$



Optimization

Convex function

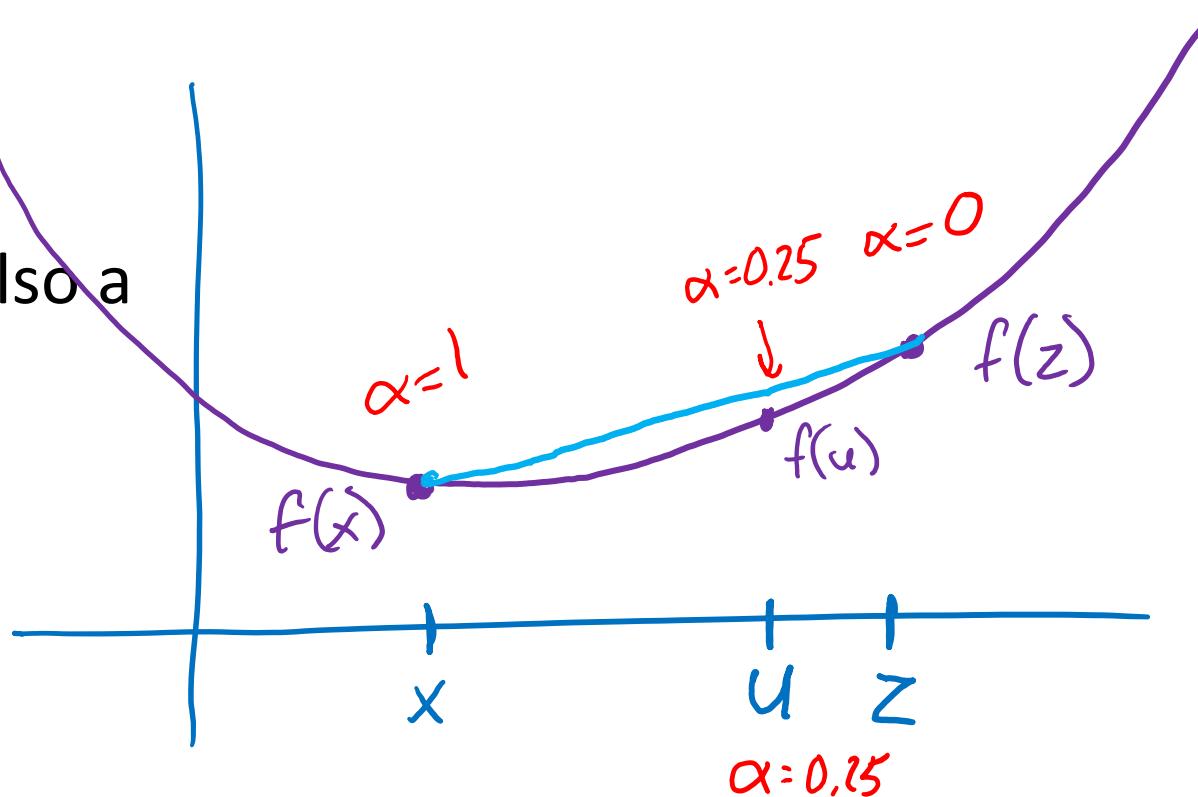
If $f(x)$ is convex, then:

- $f(\alpha x + (1 - \alpha)z) \leq \alpha f(x) + (1 - \alpha)f(z)$ $\forall 0 \leq \alpha \leq 1$

Convex optimization

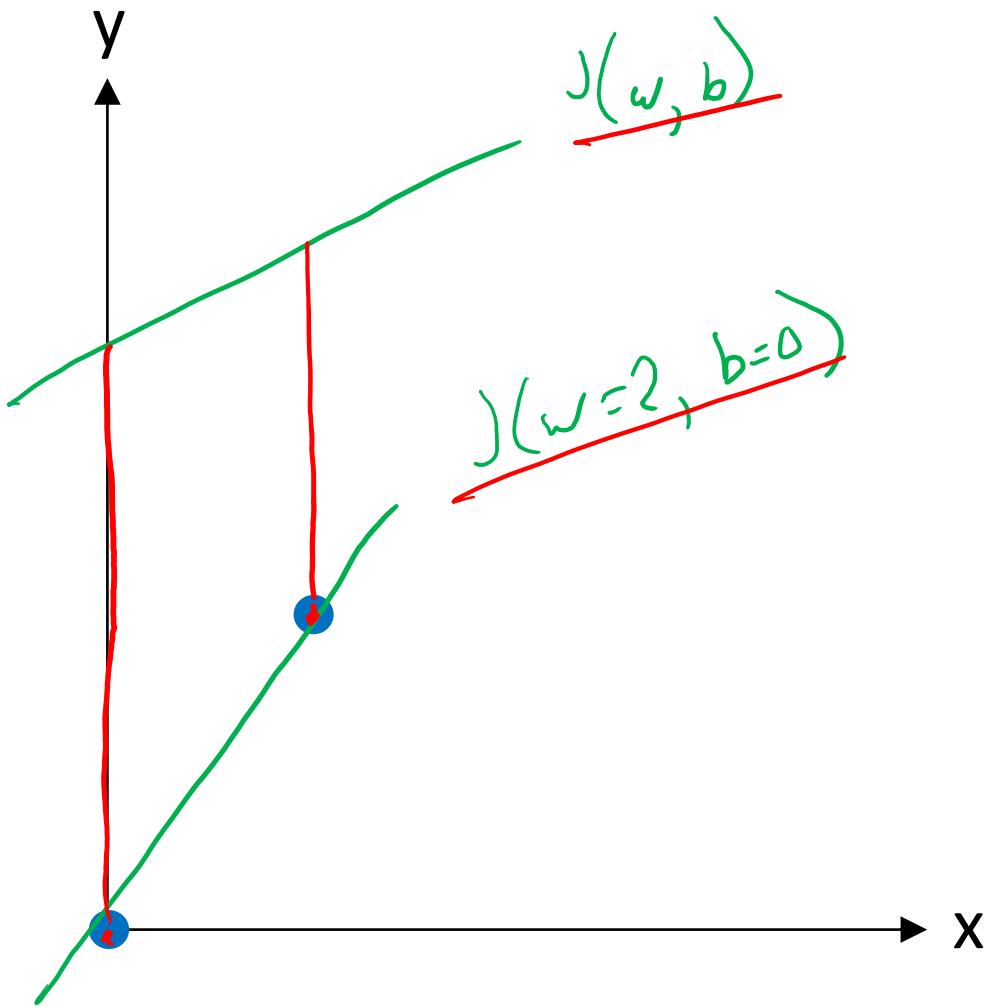
If $f(x)$ is convex, then:

- Every local minimum is also a global minimum ☺



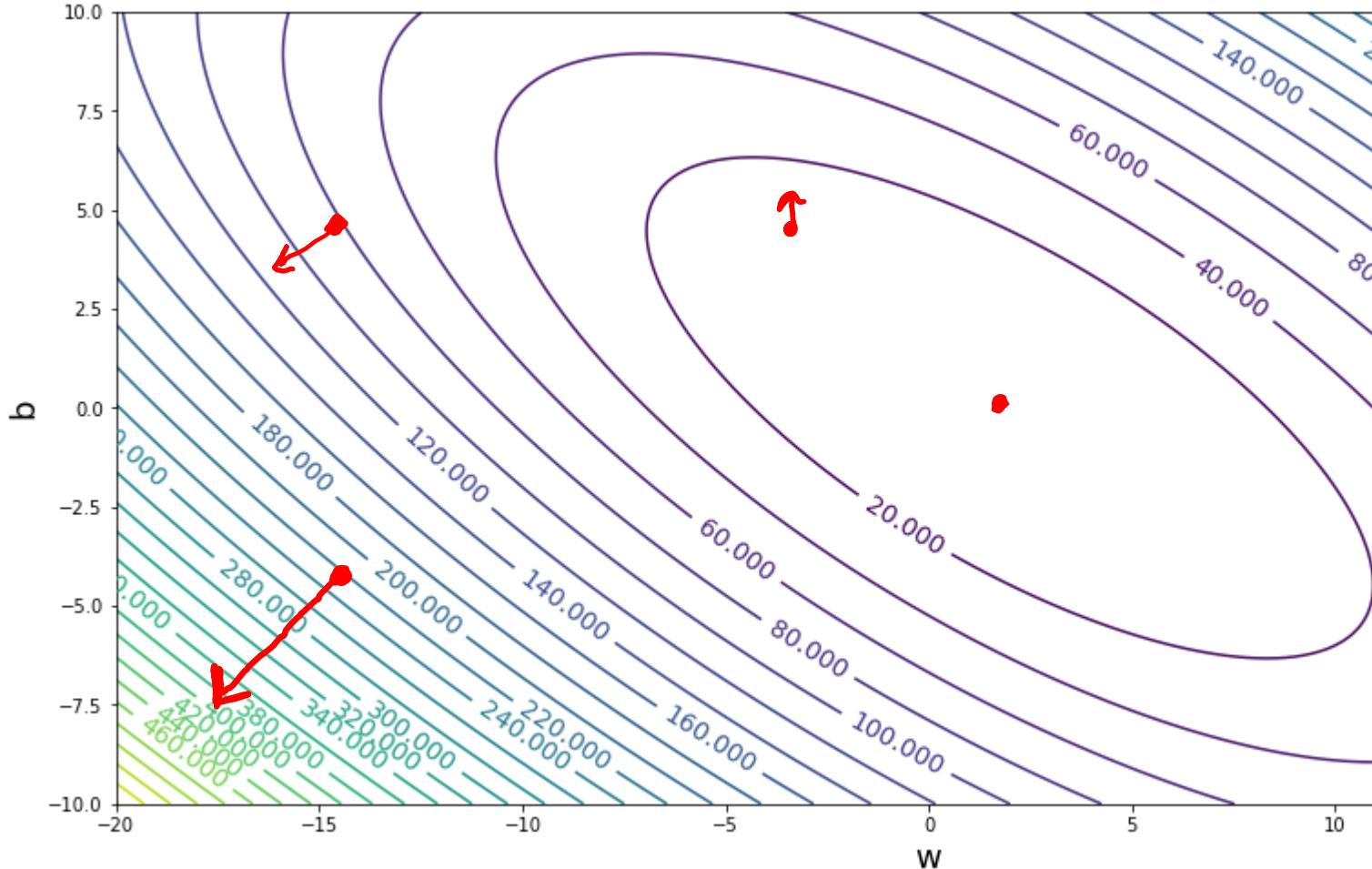
Linear Regression

Optimizing the objective

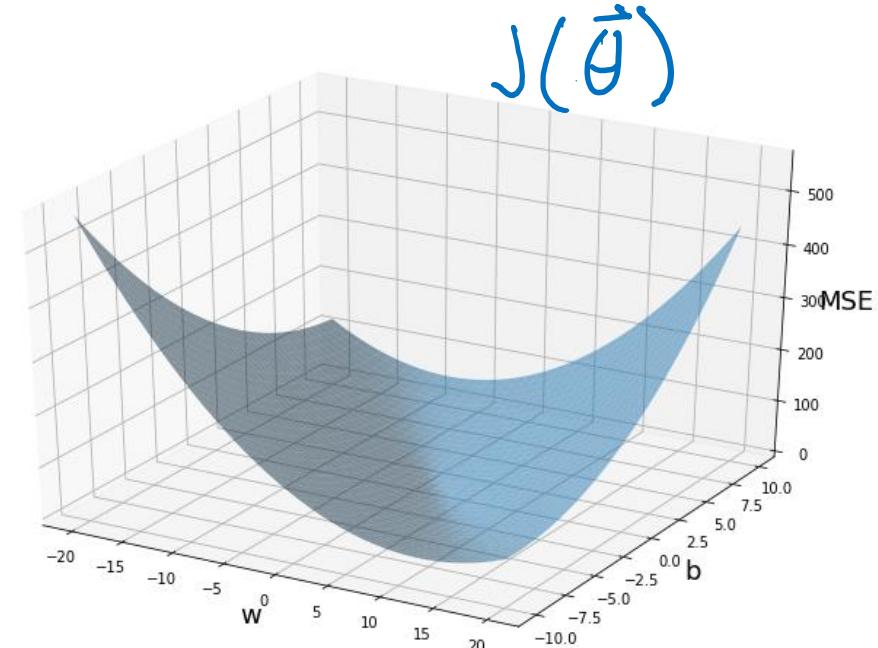


Optimization

Gradients



$$\vec{\theta} = \begin{bmatrix} b \\ w \end{bmatrix}$$



$$\nabla_{\vec{\theta}} J(\vec{\theta}) = \begin{bmatrix} \frac{\partial J}{\partial b} \\ \frac{\partial J}{\partial w} \end{bmatrix}$$

Optimization

Gradients

function $f: \mathbb{R}^M \rightarrow \mathbb{R}$

gradient $\nabla f: \mathbb{R}^M \rightarrow \mathbb{R}^M$

Input $\vec{z} \in \mathbb{R}^M$

$f(\vec{z})$

$\nabla f(\vec{z})$

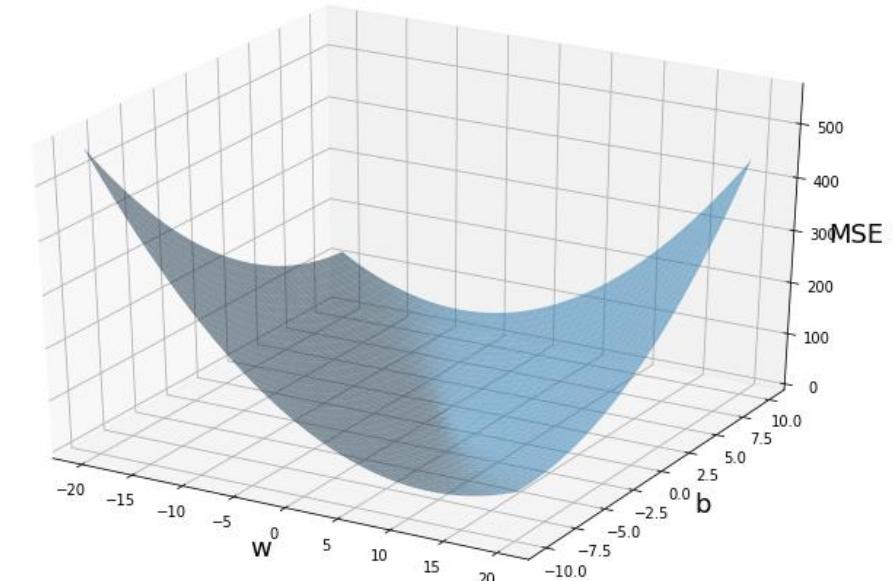
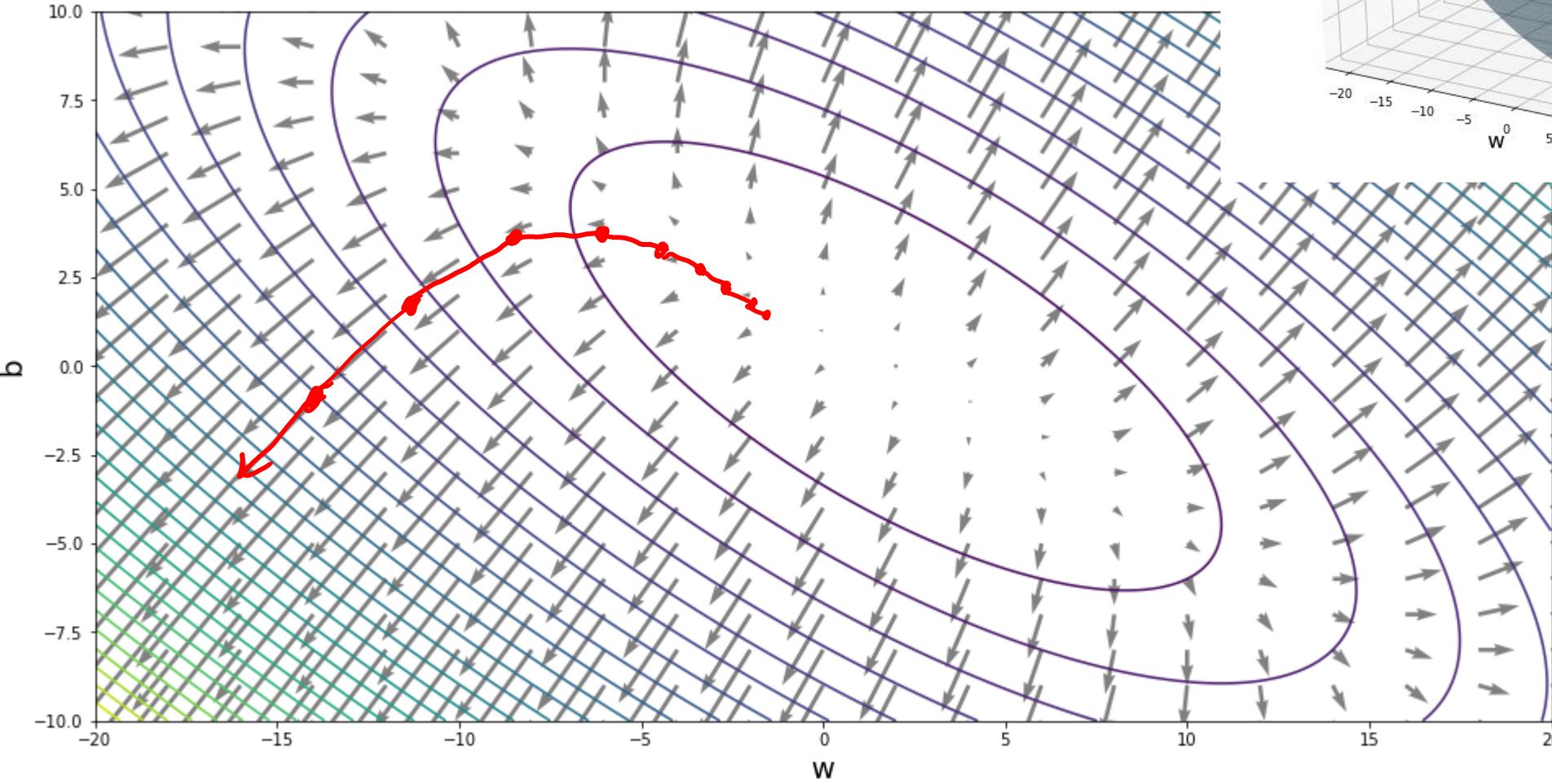
$$z = \begin{bmatrix} z_1 \\ \vdots \\ z_M \end{bmatrix}$$

$$\nabla_z f(z) = \begin{bmatrix} \frac{\partial f}{\partial z_1} \\ \frac{\partial f}{\partial z_2} \\ \vdots \\ \frac{\partial f}{\partial z_M} \end{bmatrix}$$

$$\nabla_{\vec{z}} g(\vec{z}, \vec{u})$$

Optimization

Gradients



Optimization

Gradient descent

Choose learning rate

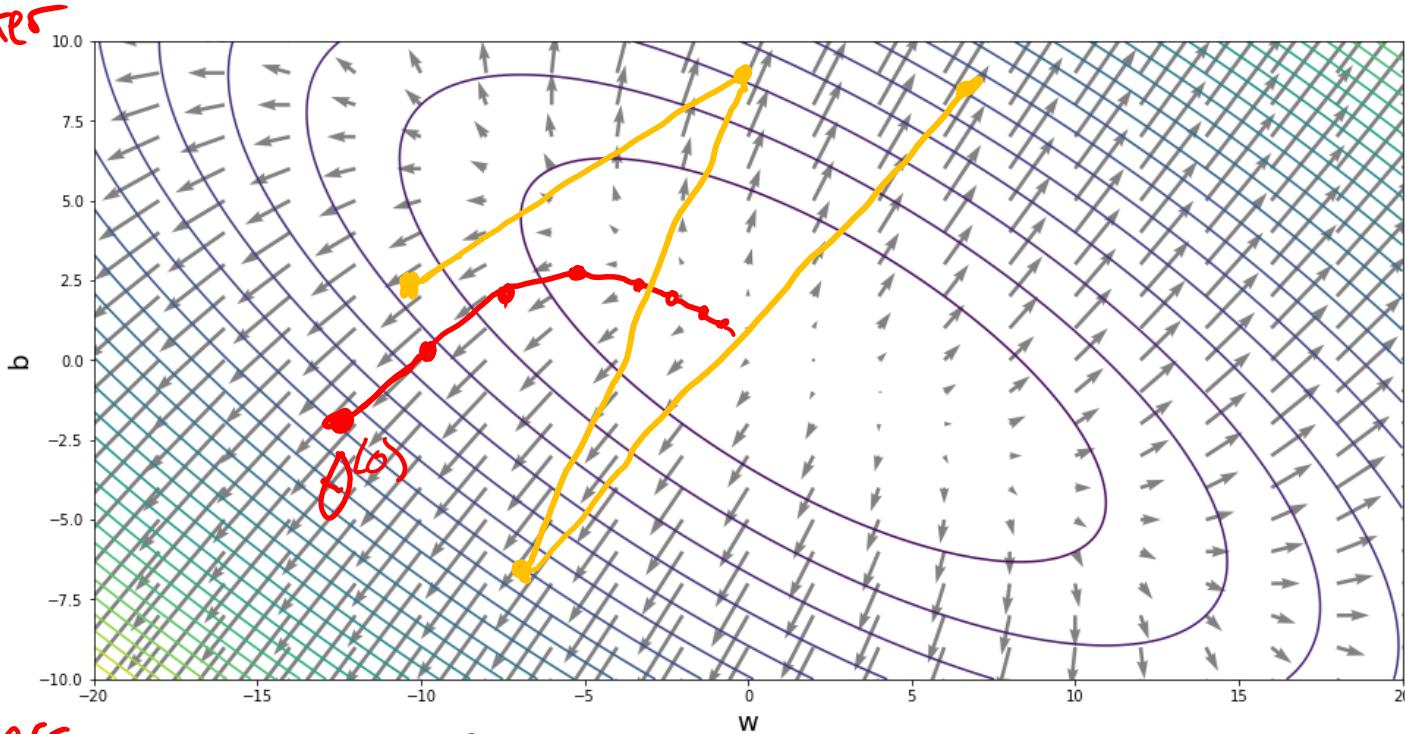
$$\alpha > 0$$

Initial $\vec{\theta}^{(0)}$ $b^{(0)}$ $w^{(0)}$

parameters

Loop

$$\vec{\theta}^{(t+1)} \leftarrow \vec{\theta}^{(t)} - \alpha \nabla J(\theta^{(t)})$$



$$\left\{ \begin{array}{l} b^{(t+1)} \leftarrow b^{(t)} - \alpha \frac{\partial J}{\partial b}(b^{(t)}, w^{(t)}) \\ w^{(t+1)} \leftarrow w^{(t)} - \alpha \frac{\partial J}{\partial w}(b^{(t)}, w^{(t)}) \end{array} \right.$$

Linear Regression

$$\|\vec{z}\|_2^2 = \sum_{i=1}^N z_i^2 = \sum_{i=1}^N z_i z_i = \vec{z}^T \vec{z}$$

Expanding objective before computing gradient

$$\begin{aligned} J(\theta) &= \frac{1}{N} \|y - X\theta\|_2^2 \\ &= \frac{1}{N} (y - X\theta)^T (y - X\theta) \\ &= \frac{1}{N} (y^T - \theta^T X^T)(y - X\theta) \\ &= \frac{1}{N} (y^T y - \underline{\theta^T X^T y} - \underline{y^T X\theta} + \theta^T X^T X\theta) \\ &= \frac{1}{N} (y^T y - 2\theta^T X^T y + \theta^T X^T X\theta) \end{aligned}$$

these two are the same

Linear Regression

Gradient of objective with respect to parameters

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{N} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \\ &= \frac{1}{N} (\mathbf{y}^T \mathbf{y} - 2\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta}) \end{aligned}$$

$$\begin{aligned} \nabla J(\boldsymbol{\theta}) &= \frac{1}{N} (0 - \cancel{2\mathbf{X}^T \mathbf{y}} + \cancel{2\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X}}) \leftarrow \\ &= \frac{1}{N} (0 - 2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \boldsymbol{\theta}) \\ &= \frac{2}{N} (-\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \boldsymbol{\theta}) \end{aligned}$$

Dimension mismatch

$$\frac{\partial \mathbf{z}^T \mathbf{u}}{\partial \mathbf{z}} = \mathbf{u}$$

-- or --

$$\frac{\partial \mathbf{z}^T \mathbf{u}}{\partial \mathbf{z}} = \mathbf{u}^T$$

$$\frac{\partial \mathbf{z}^T \mathbf{A} \mathbf{z}}{\partial \mathbf{z}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{z}$$

-- or --

$$\frac{\partial \mathbf{z}^T \mathbf{A} \mathbf{z}}{\partial \mathbf{z}} = \mathbf{z}^T (\mathbf{A} + \mathbf{A}^T)$$

Linear Regression

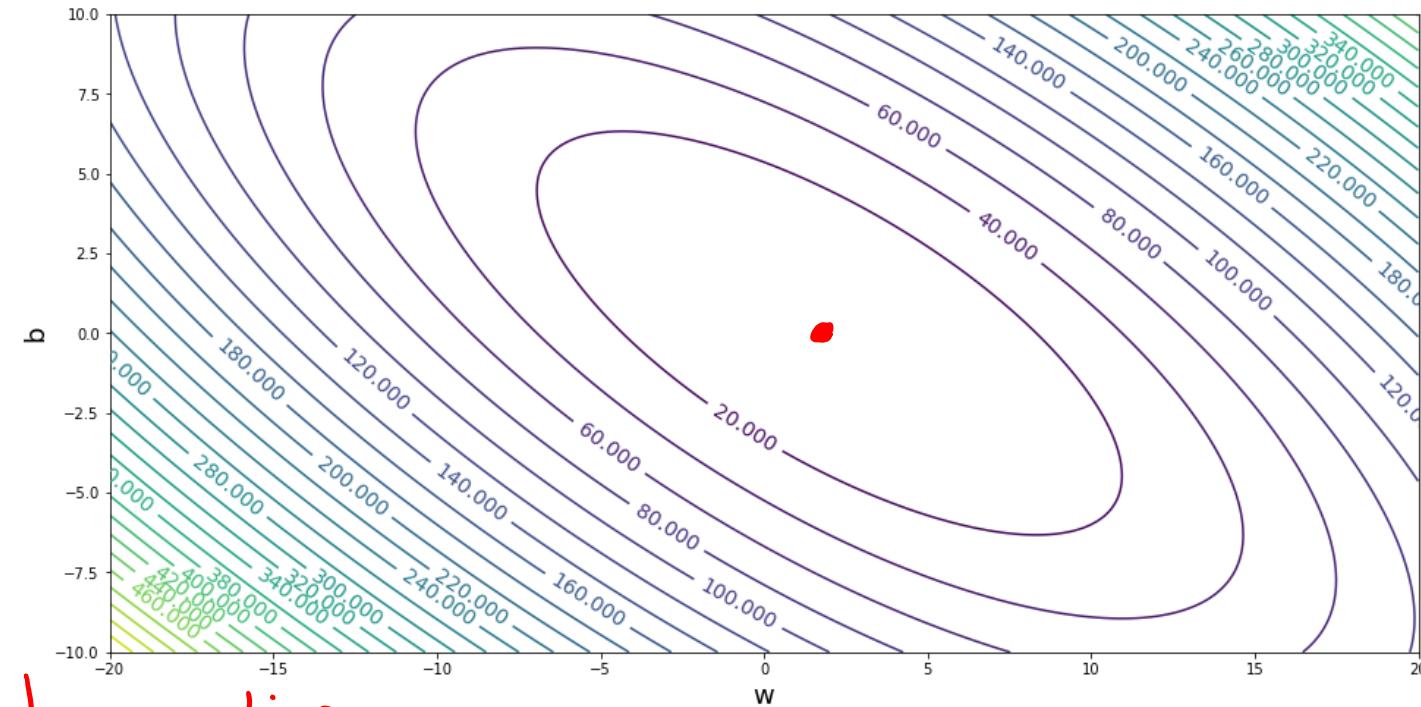
Closed-form solution

$$\nabla J(\theta) = \frac{2}{N} (-X^T y + X^T X \theta)$$

$$\nabla J(\theta) = 0$$

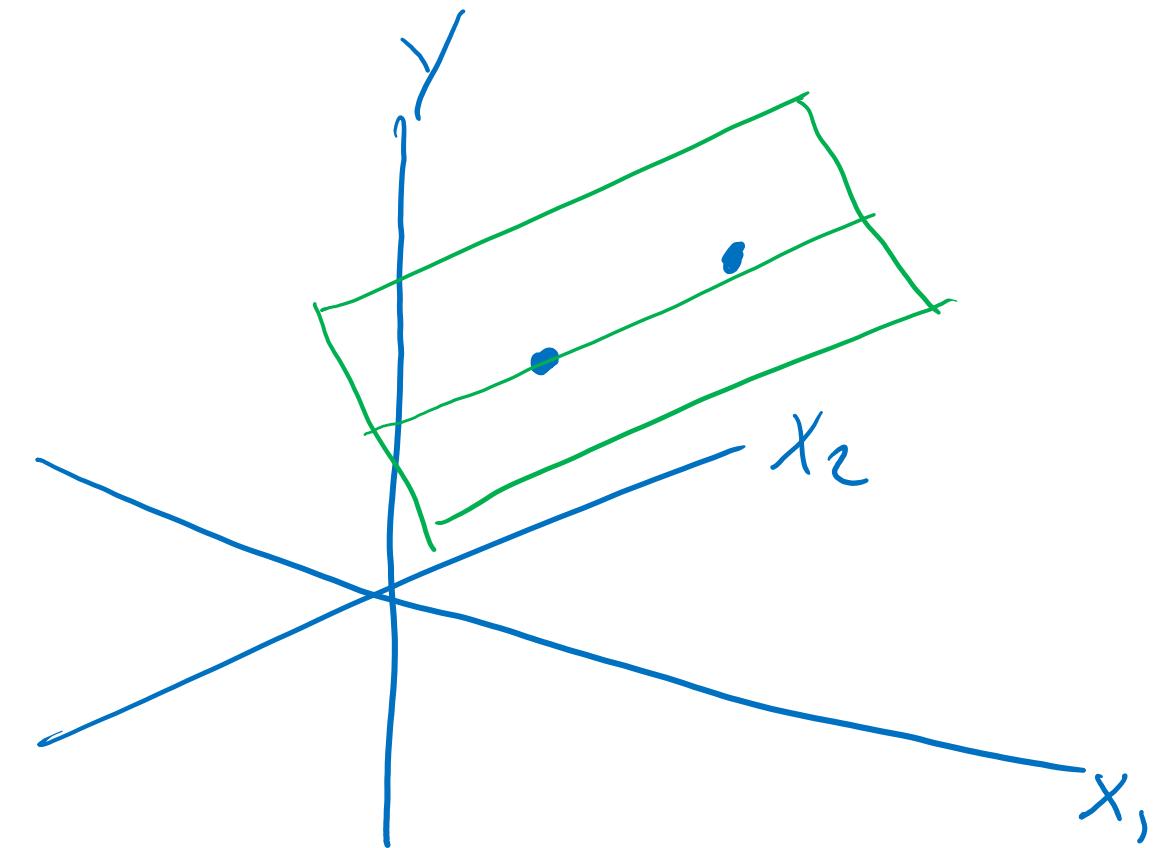
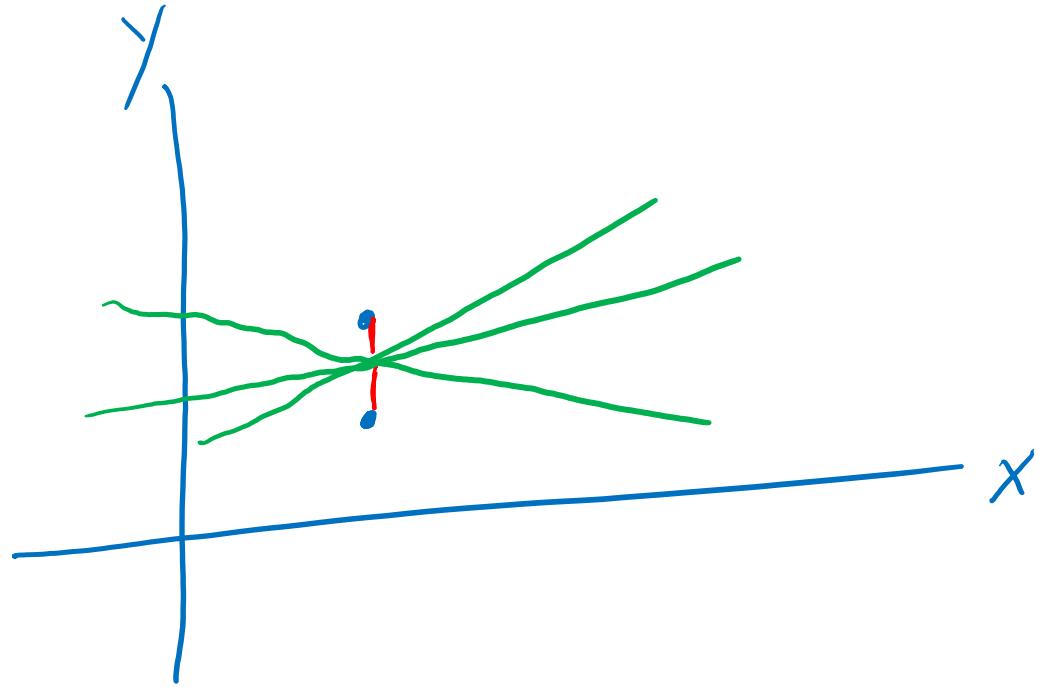
$$X^T X \theta = X^T y \leftarrow \text{Normal equation}$$

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



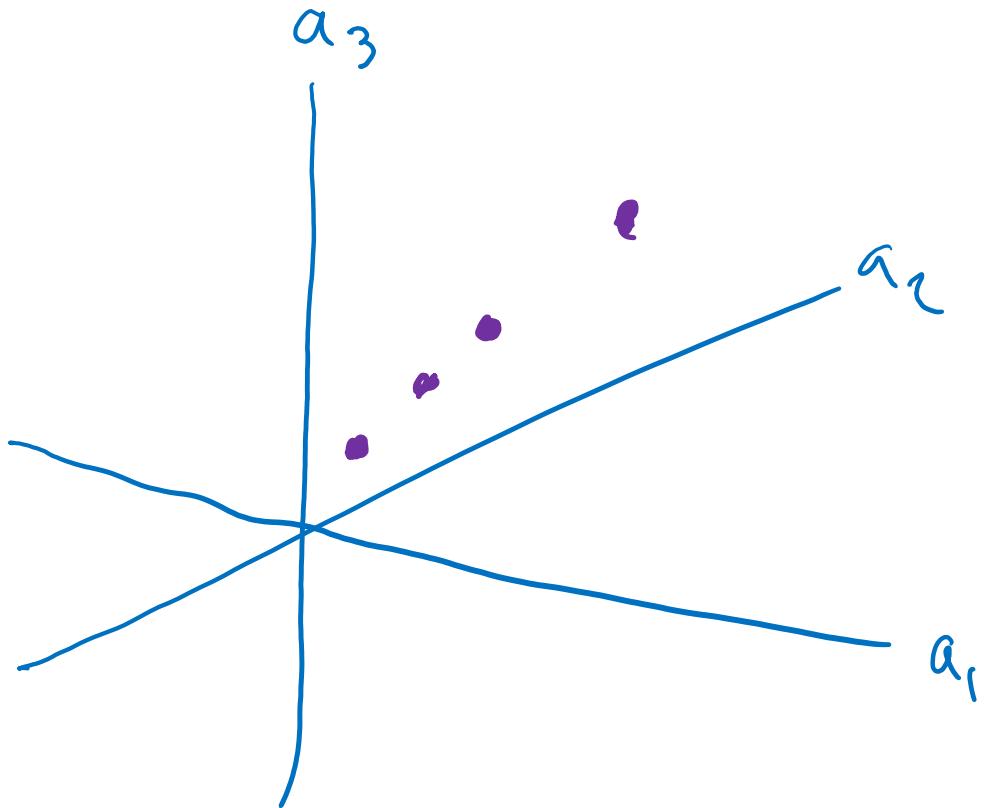
Linear Regression

Number of solutions



A Note on Matrix Rank

Underlying dimensionality of the data



$$A = \begin{bmatrix} a_1 & a_2 & a_3 \\ 1 & 1 & 1 \\ 2 & 2 & 2 \\ 5 & 5 & 5 \\ 3 & 3 & 3 \end{bmatrix}$$

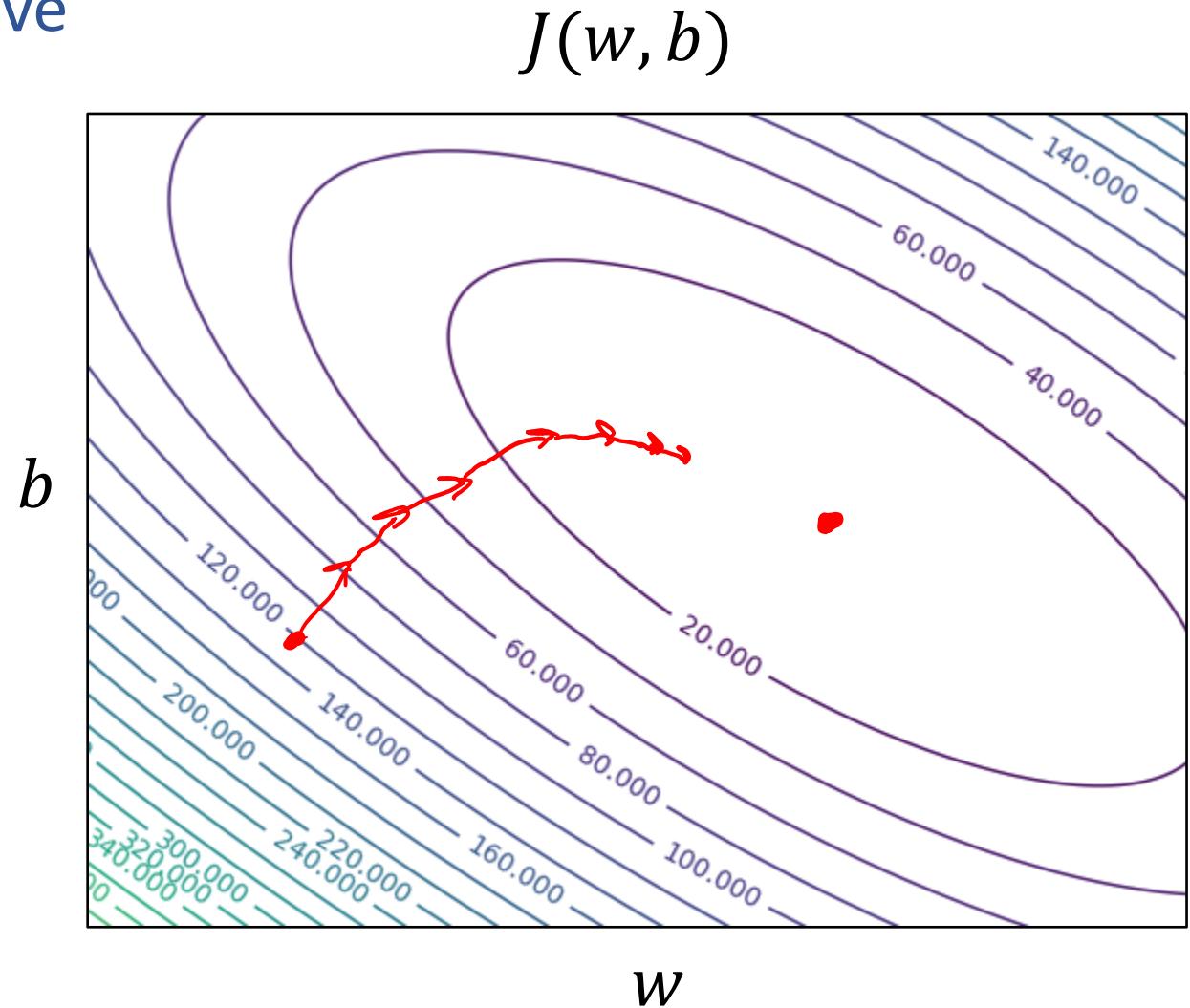
$$\text{Rank}(A) = 1$$

Linear Regression



Methods for optimizing the objective

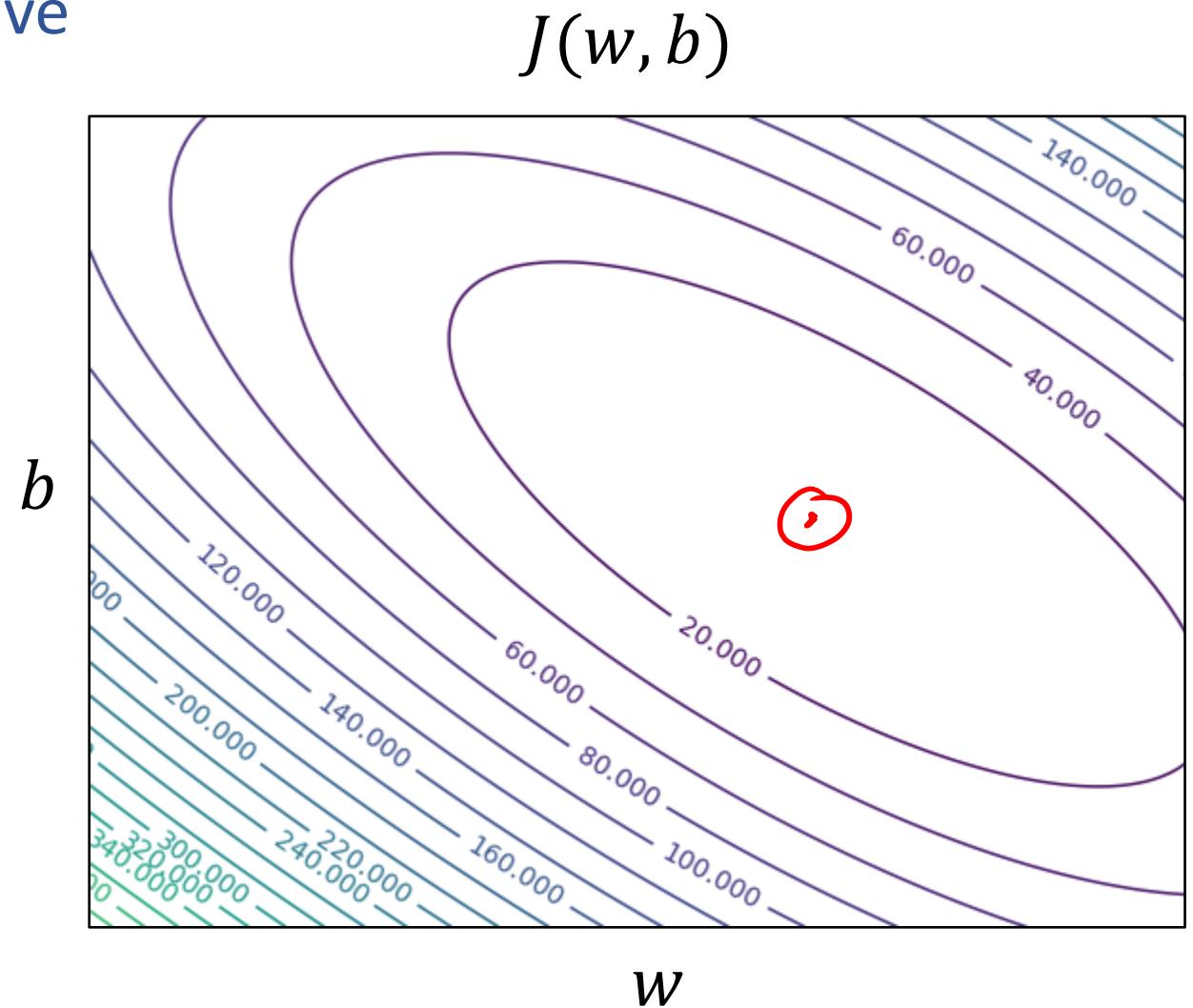
- Grid search
- Random search
- Closed-form solution
- (Batch) Gradient descent
- Stochastic gradient descent



Linear Regression

Methods for optimizing the objective

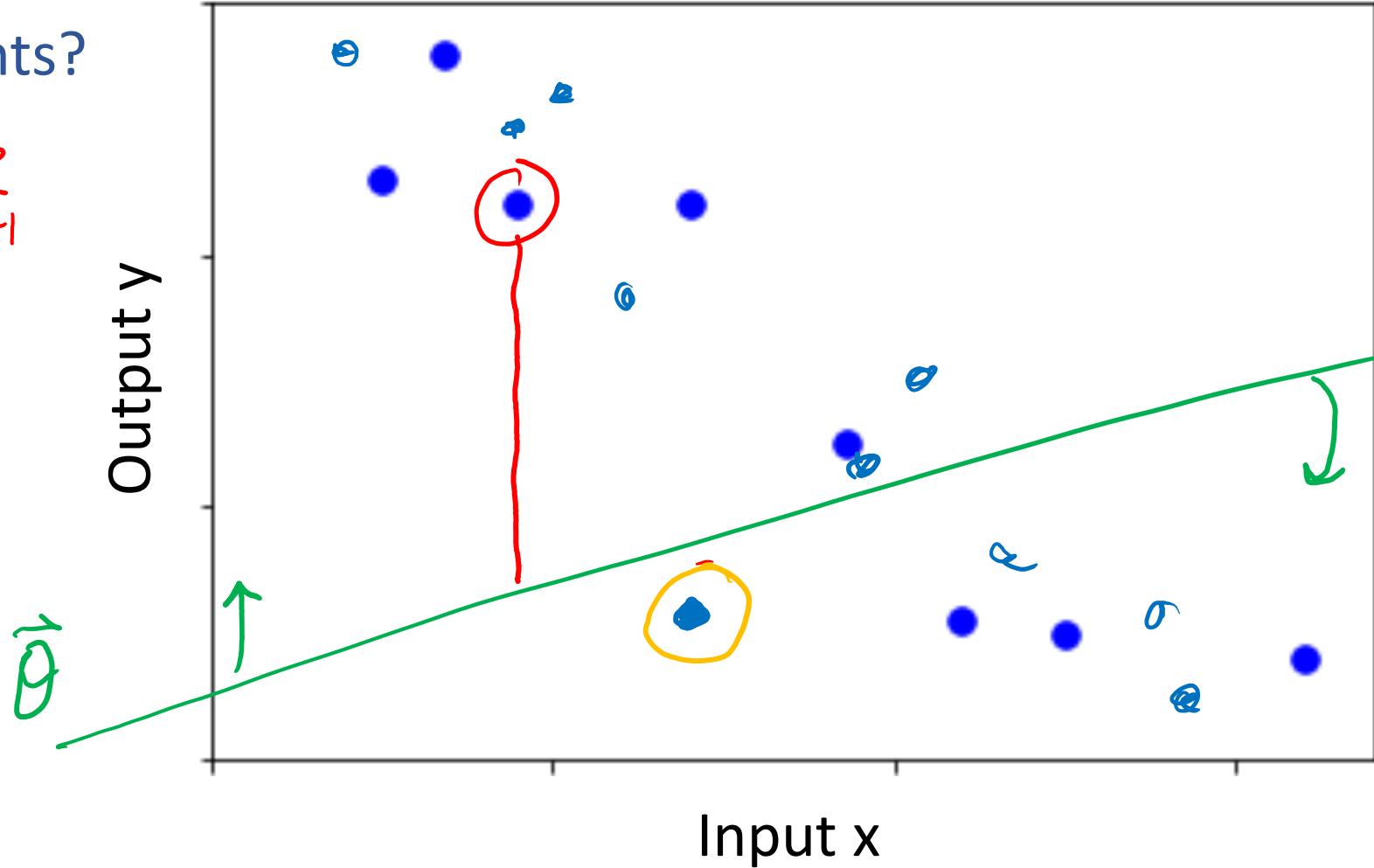
- Grid search
- Random search
- Closed-form solution
- (Batch) Gradient descent
- Stochastic gradient descent



Linear Regression Gradient Descent

What happens in gradient descent when we have $N=1,000,000$ training points?

$$\nabla_{\theta} J = \left[\begin{array}{c} \frac{\partial J}{\partial b} \\ \vdots \\ \frac{\partial J}{\partial w} \end{array} \right] \quad \sum_{j=1}^N$$

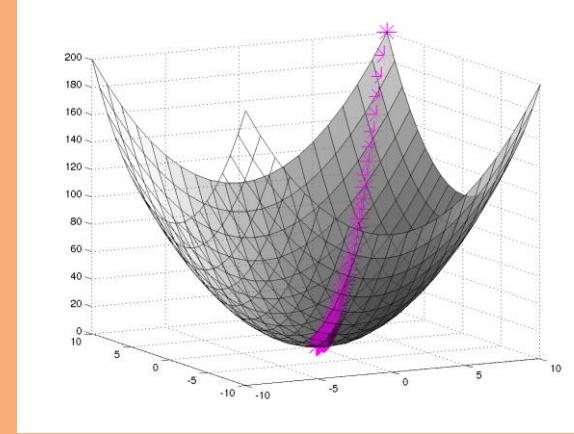


$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (\gamma^{(i)} - \theta^T x^{(i)})^2$$

(Batch) Gradient Descent

Algorithm 1 Gradient Descent

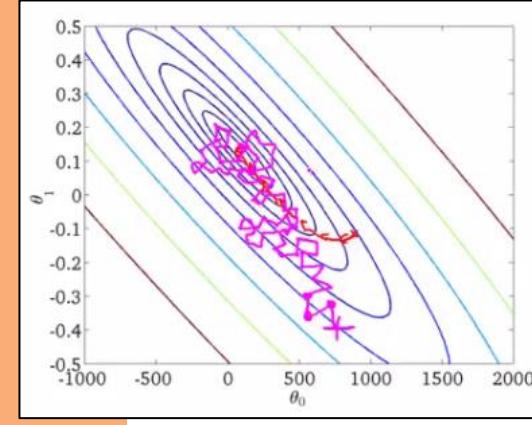
```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )  
2:    $\theta \leftarrow \theta^{(0)}$   
3:   while not converged do  
4:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$   
5:   return  $\theta$ 
```



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:      $\theta \leftarrow \theta - \lambda \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```

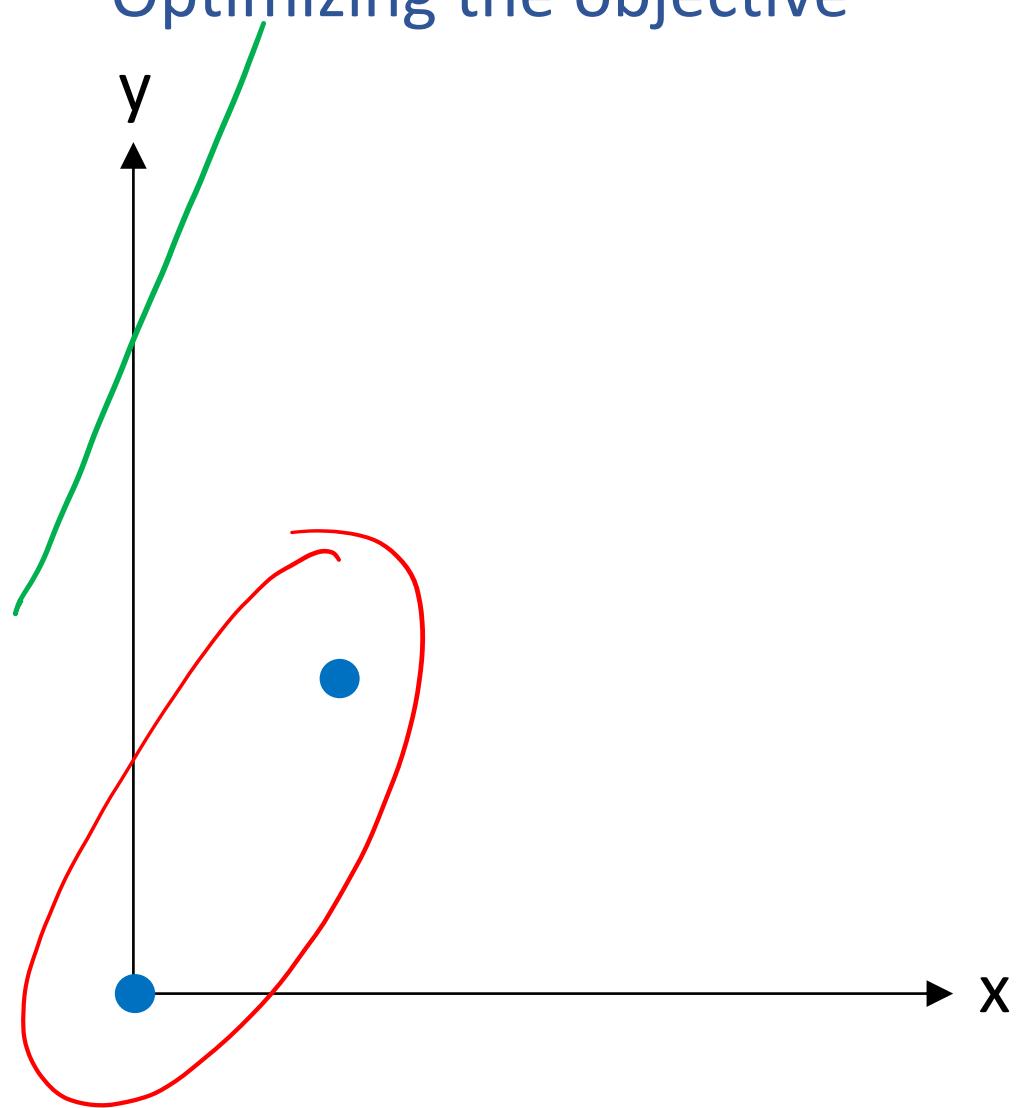


We need a per-example objective:

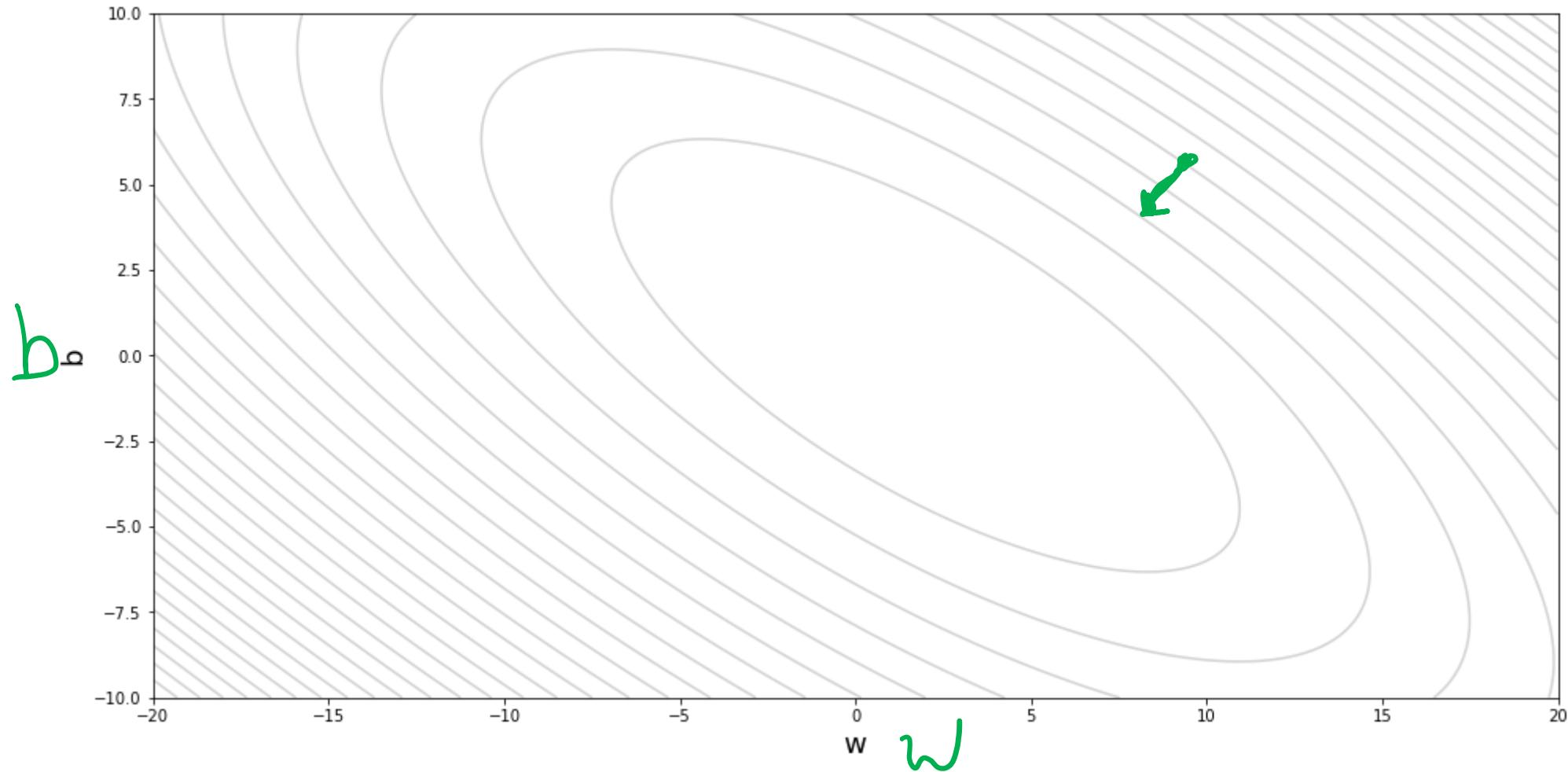
$$\text{Let } \underline{J(\theta)} = \sum_{i=1}^N J^{(i)}(\theta)$$

Linear Regression

Optimizing the objective

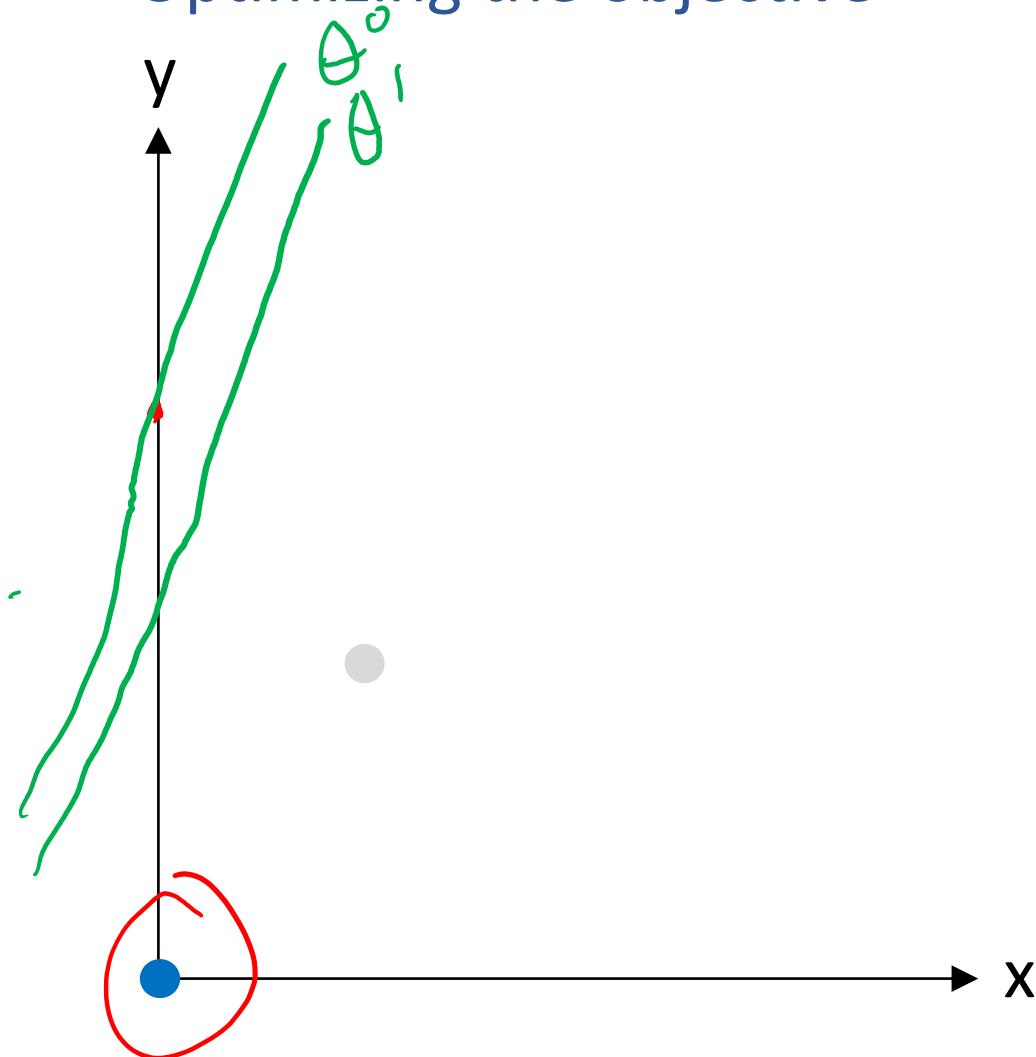


Stochastic Gradient Descent

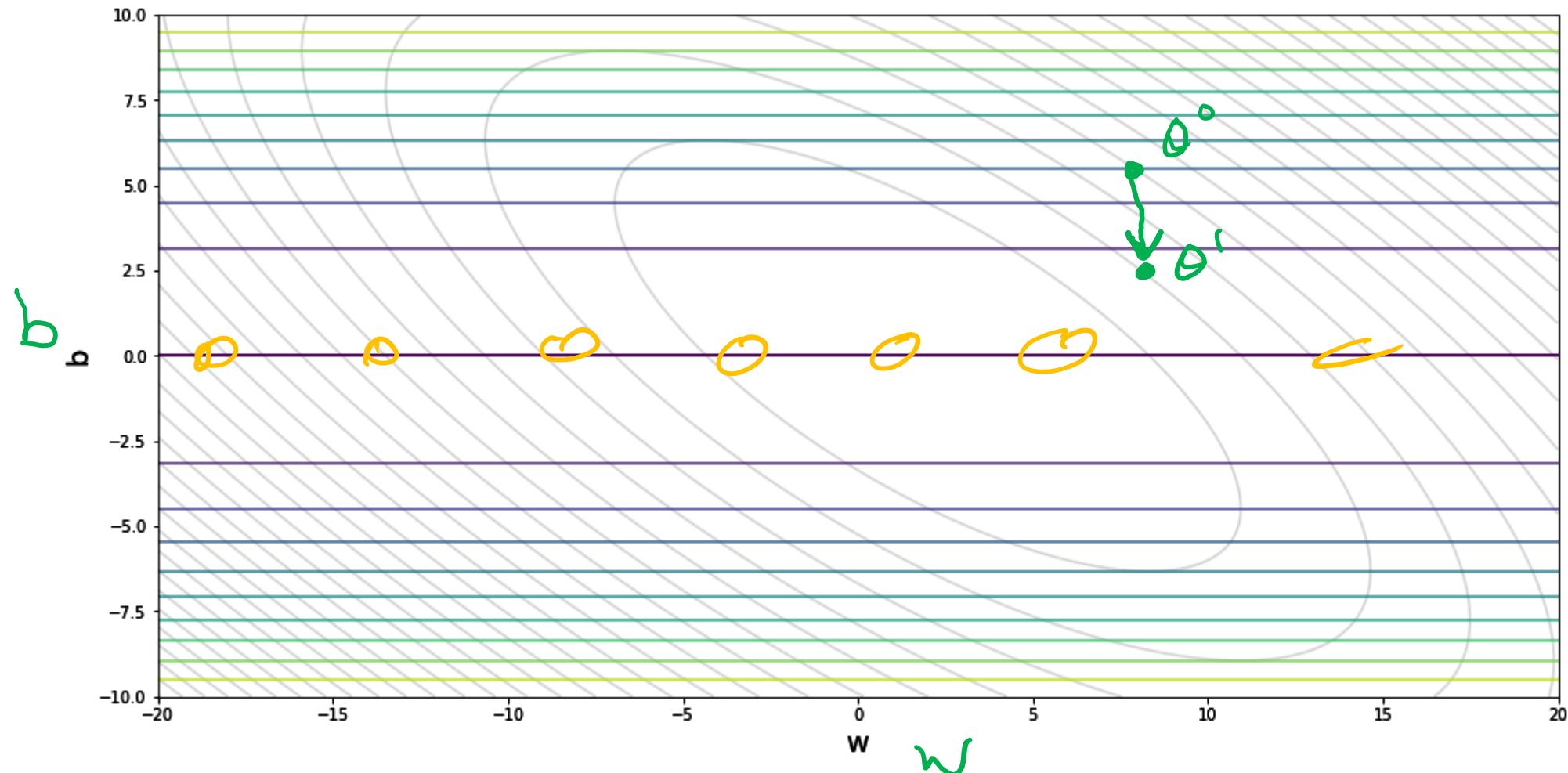


Linear Regression

Optimizing the objective

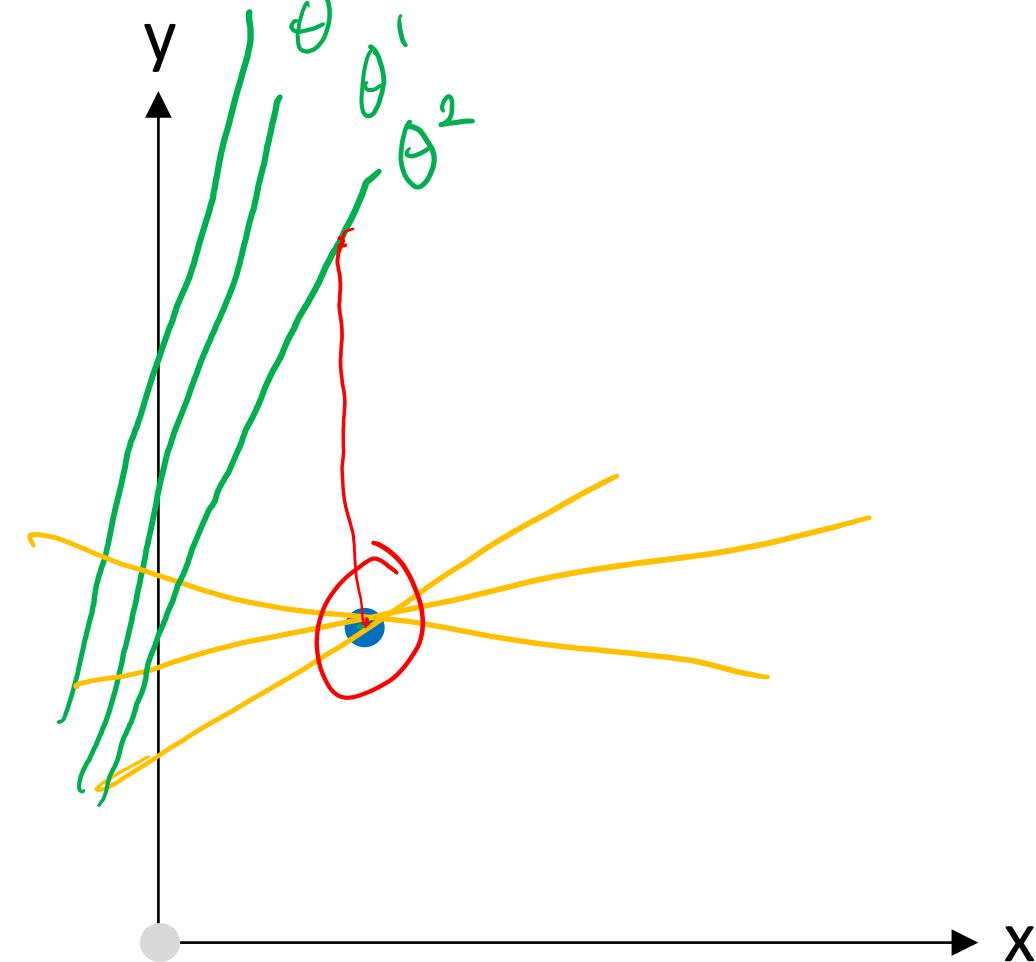


Stochastic Gradient Descent

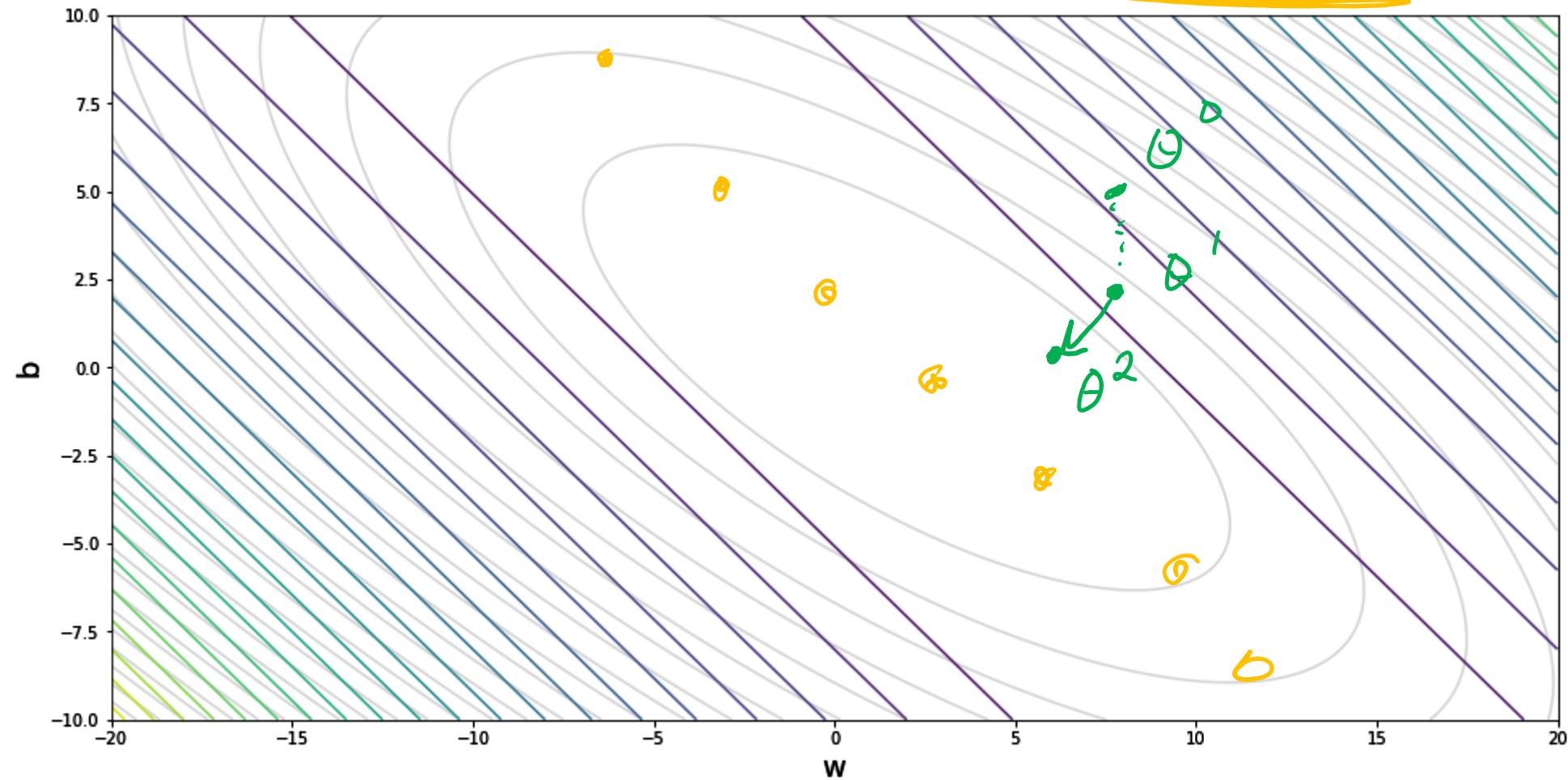


Linear Regression

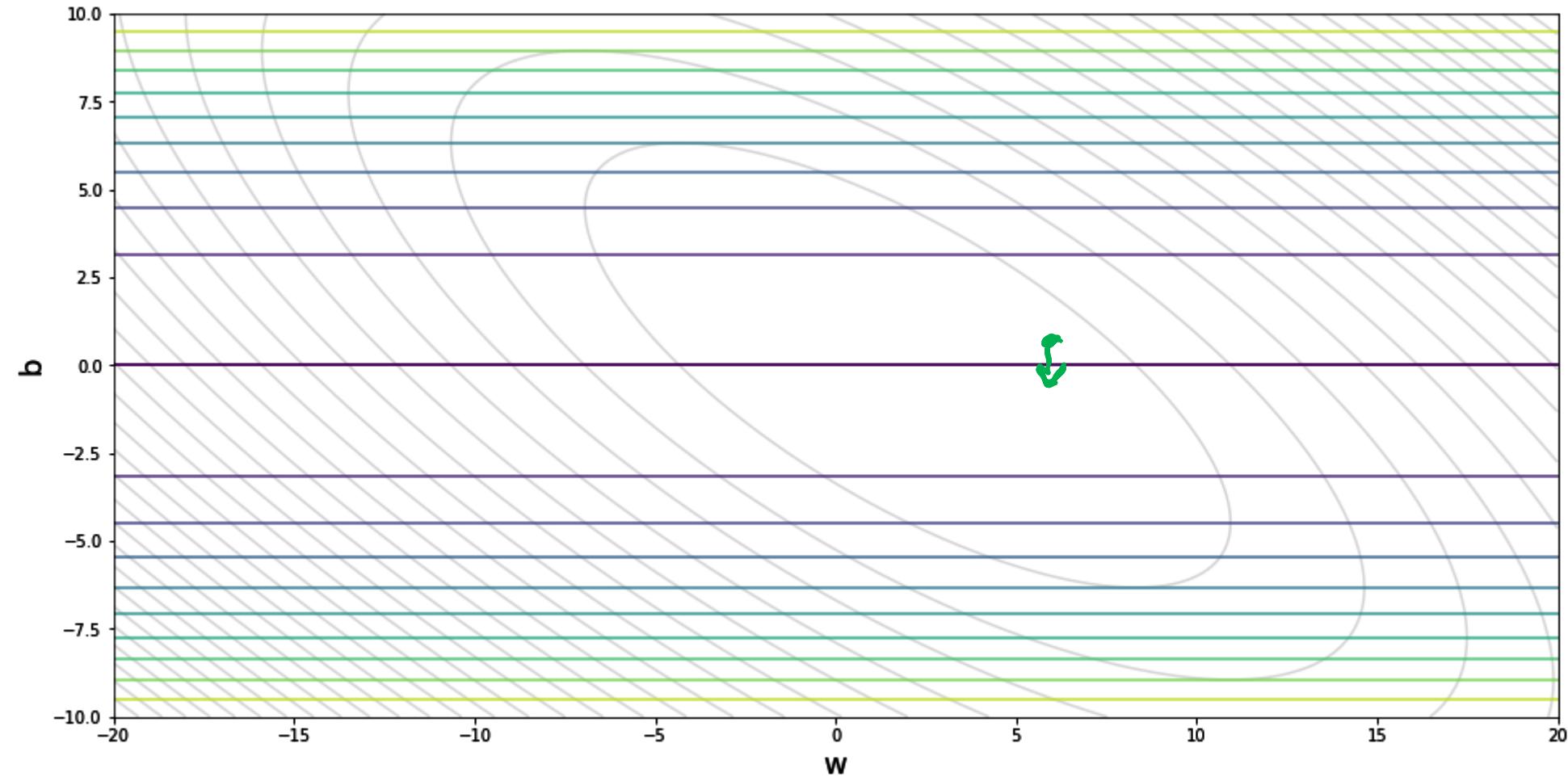
Optimizing the objective



Stochastic Gradient Descent



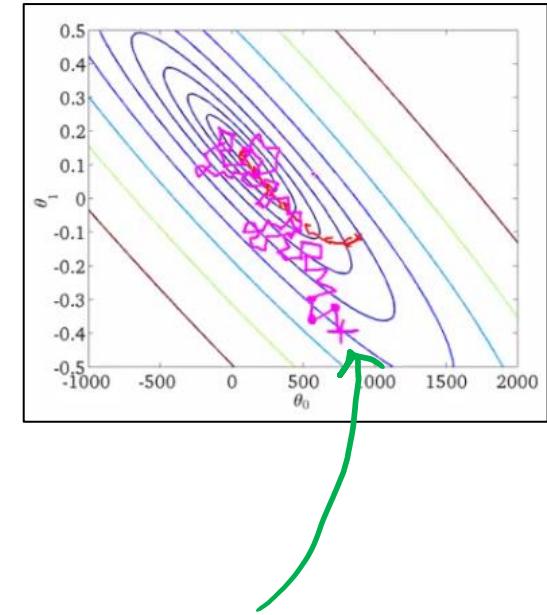
Stochastic Gradient Descent



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:      $\theta \leftarrow \theta - \lambda \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



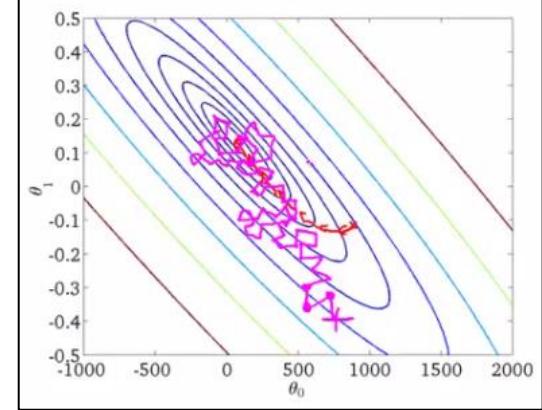
We need a per-example objective:

$$\text{Let } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

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  $\theta$ 
```

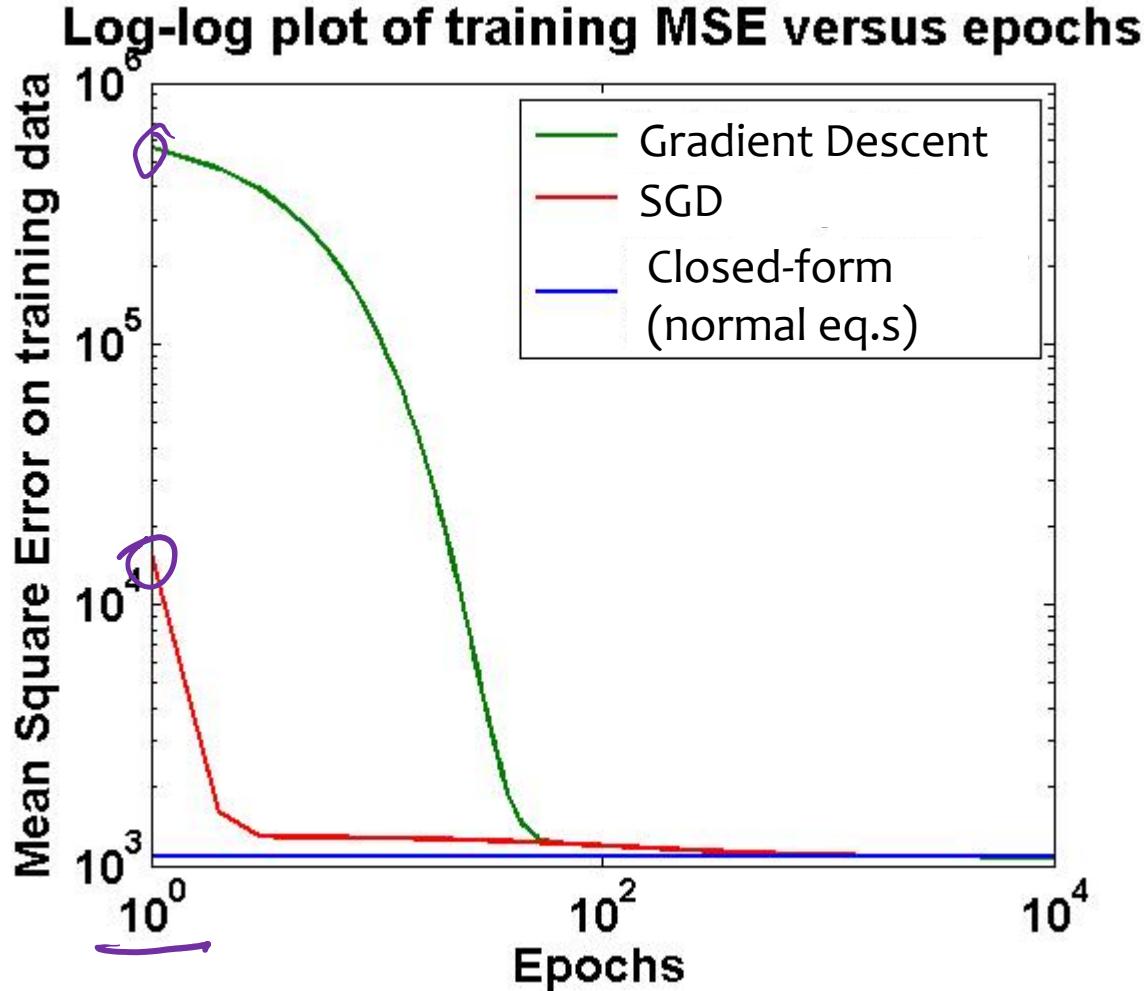


We need a per-example objective:

$$\text{Let } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

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. $\text{Uniform}(\{1,2,\dots,N\})$).

Convergence Curves



- 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