

# 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

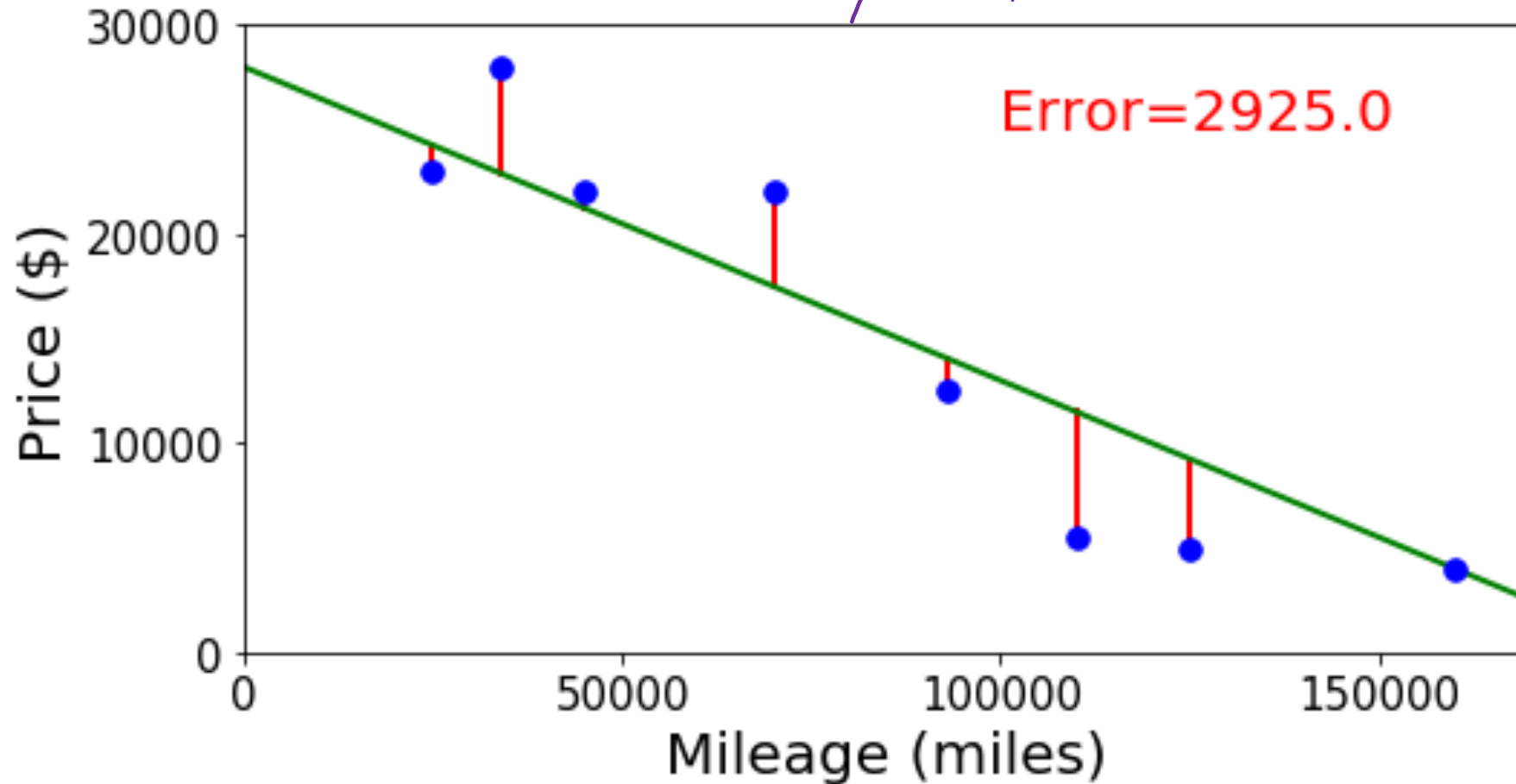
## Linear Regression and Optimization

Instructor: Pat Virtue

# Linear Regression

Selling my car

$$y = mx + b$$
$$y = wx + b$$
$$y = w_1x + w_0$$
$$y = \theta_1x + \theta_0$$



# Linear in Higher Dimensions

1-D  $y = w x + b$   
2-D  $y = w_1 x_1 + w_2 x_2 + b$

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

	$x \in \mathbb{R}$	$x \in \mathbb{R}^2$	$x \in \mathbb{R}^3$	$x \in \mathbb{R}^M$
$\rightarrow y = \mathbf{w}^T \mathbf{x} + b$	line	plane ↓	hyperplane	hyperplane
$\mathbf{w}^T \mathbf{x} + b = 0$	point	line	plane	hyperplane
$\mathbf{w}^T \mathbf{x} + b \geq 0$	halfline	halfplane	halfspace	halfspace

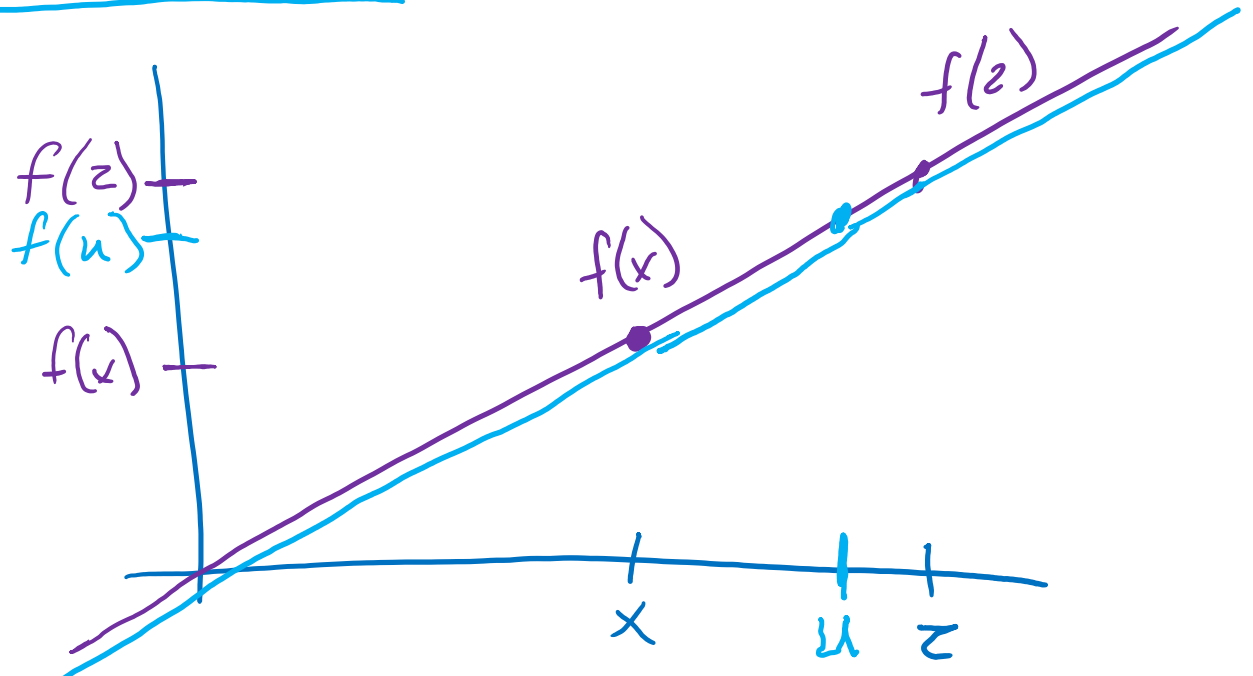
# Linear Function

## Linear function

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

- $f(x + z) = f(x) + f(z)$
- $f(\alpha x) = \alpha f(x) \quad \forall \alpha$
- ▪  $f(\alpha x + (1 - \alpha)z) = \alpha f(x) + (1 - \alpha)f(z) \quad \forall \alpha$

$$\alpha = 0.25$$



# Piazza Poll 1

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

$f(\mathbf{x})$  is linear if and only if:

- $f(\mathbf{x} + \mathbf{z}) = f(\mathbf{x}) + f(\mathbf{z})$  and
- $f(\alpha\mathbf{x}) = \alpha f(\mathbf{x}) \quad \forall \alpha$

# Linear Regression

Linear algebra formulation

$$y = \vec{w}^T x_{orig} + b$$

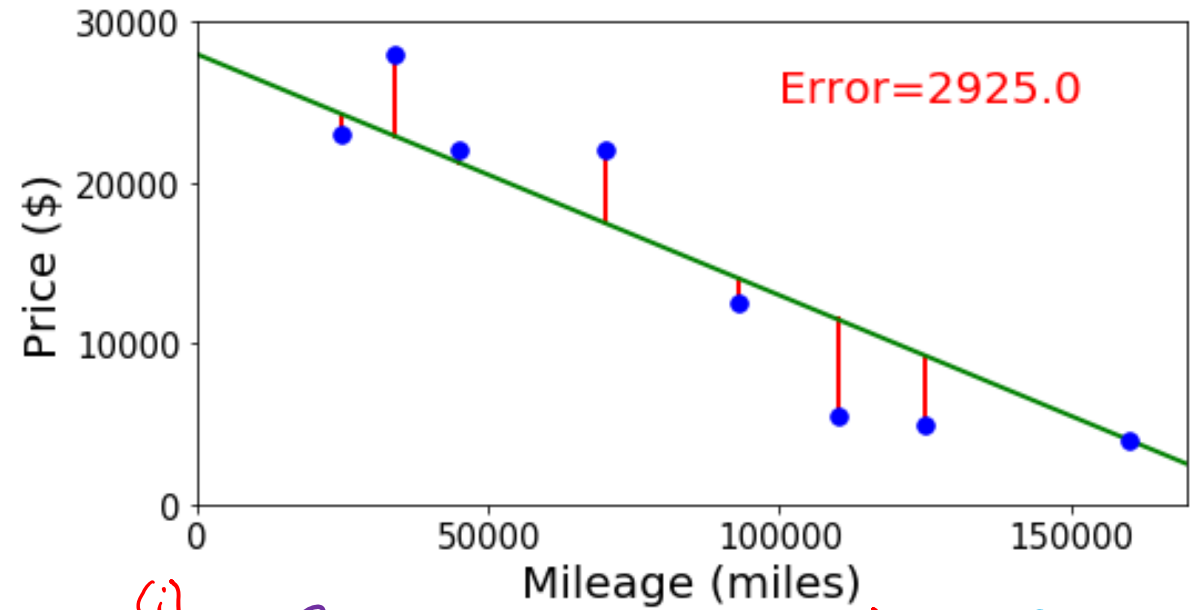
$$y^{(i)} = \vec{w}^T x_{orig}^{(i)} + b$$

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

$$\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} = \begin{bmatrix} - & x^{(1)T} & - \\ & \vdots & \\ - & x^{(n)T} & - \end{bmatrix} \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

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

Design Matrix  
X



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

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

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

b

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

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

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



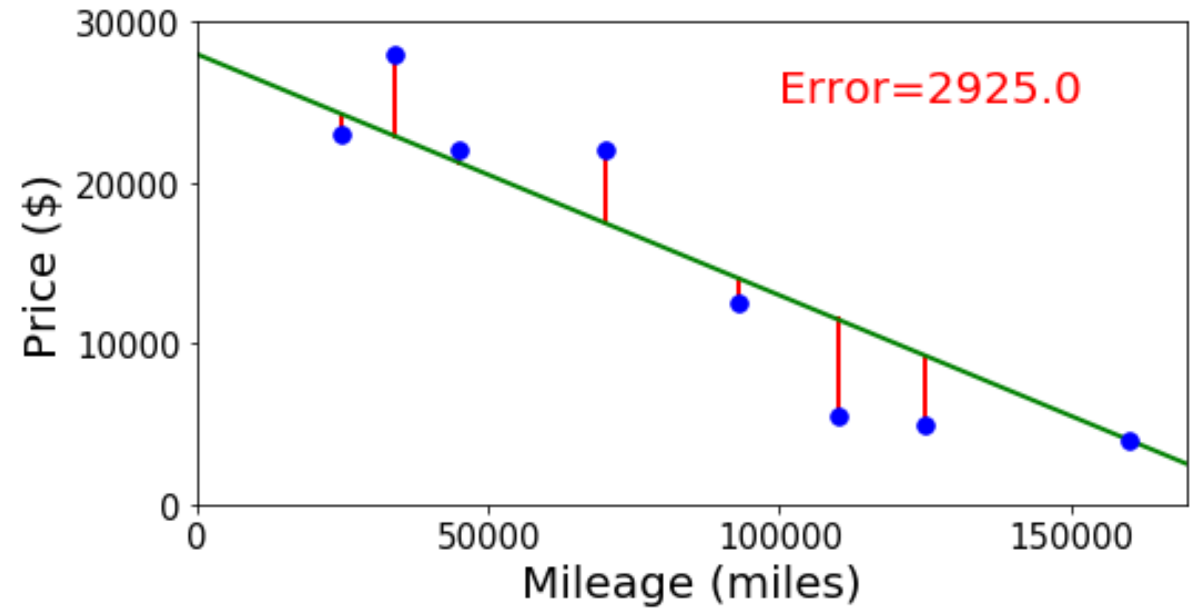
# Linear Regression

## Error and objectives

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

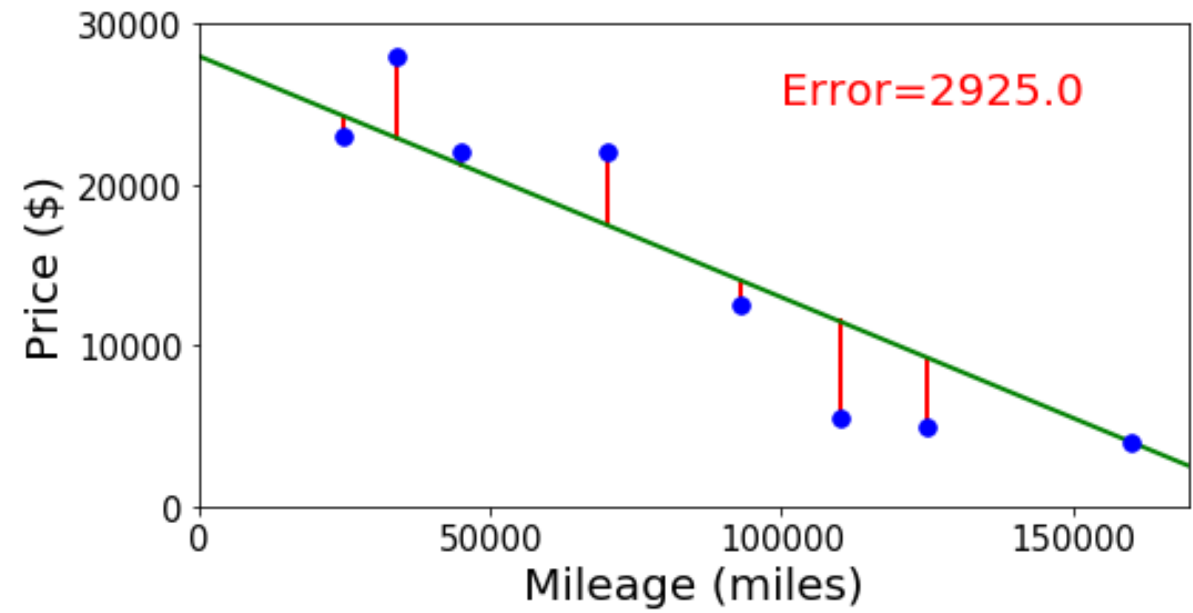
$$J(w_1, w_2, b) = \frac{1}{N} \sum (y^{(i)} - \hat{y}^{(i)})^2$$
$$\hat{y}^{(i)} = w_1 x_1^{(i)} + w_2 x_2^{(i)} + b$$

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



# Linear Regression

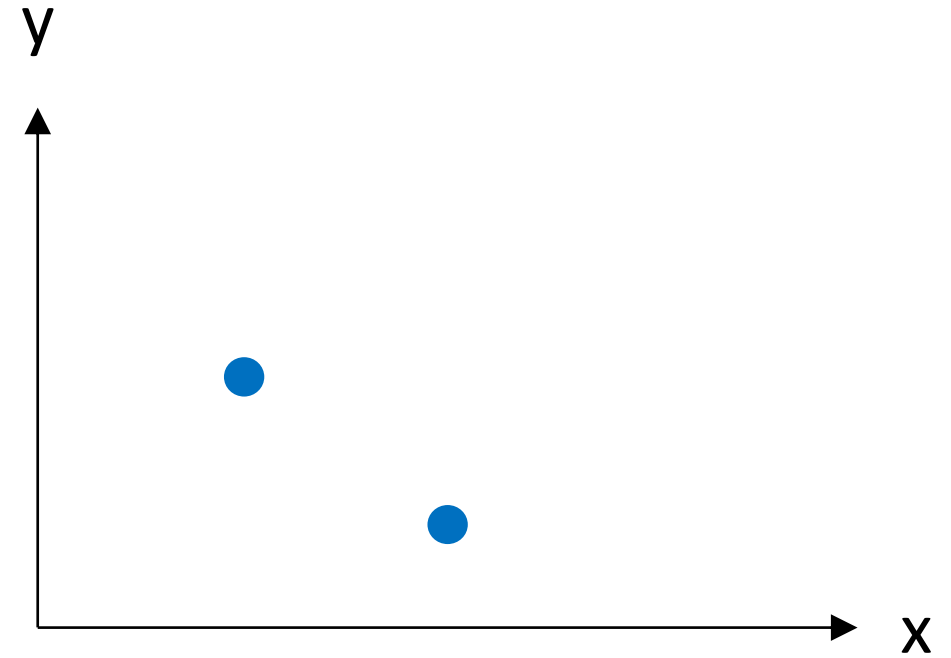
Linear algebra formulation



# Previous Piazza Poll

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

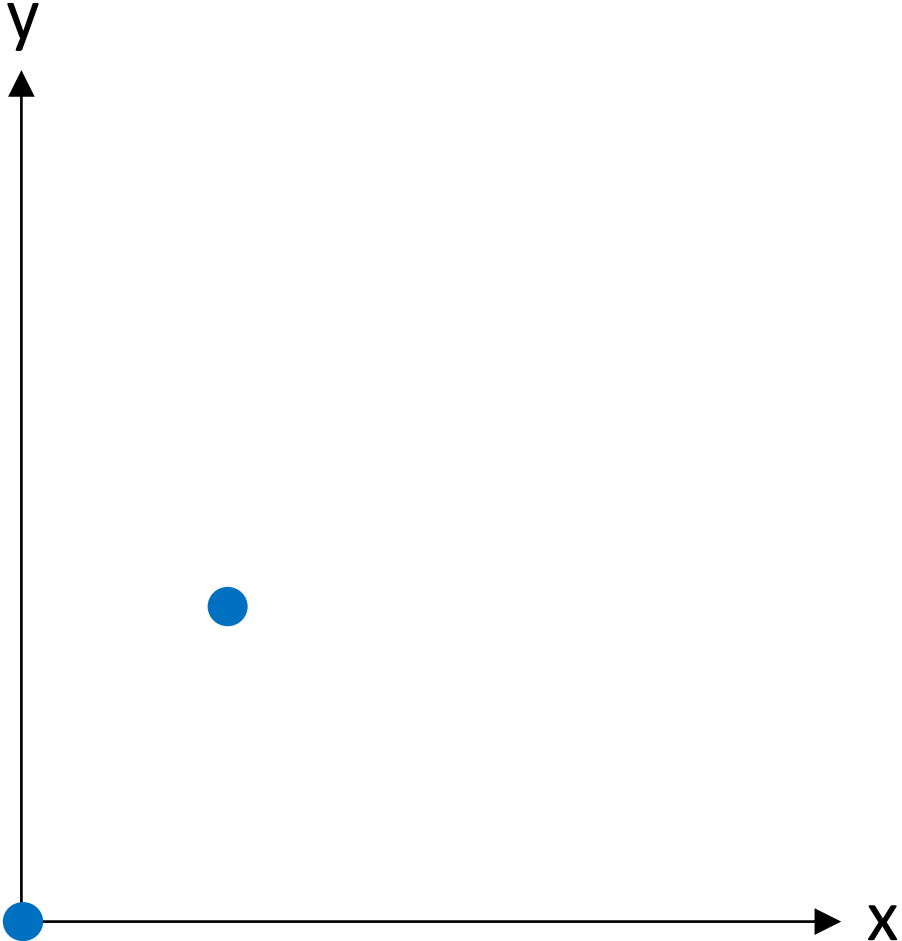
- A. Line
- B. Plane
- C. Half-plane
- D. Convex Parabola (U-shape)
- E. Concave parabola (up-side-down U)
- F. None of the above



# Linear Regression

Optimizing the objective

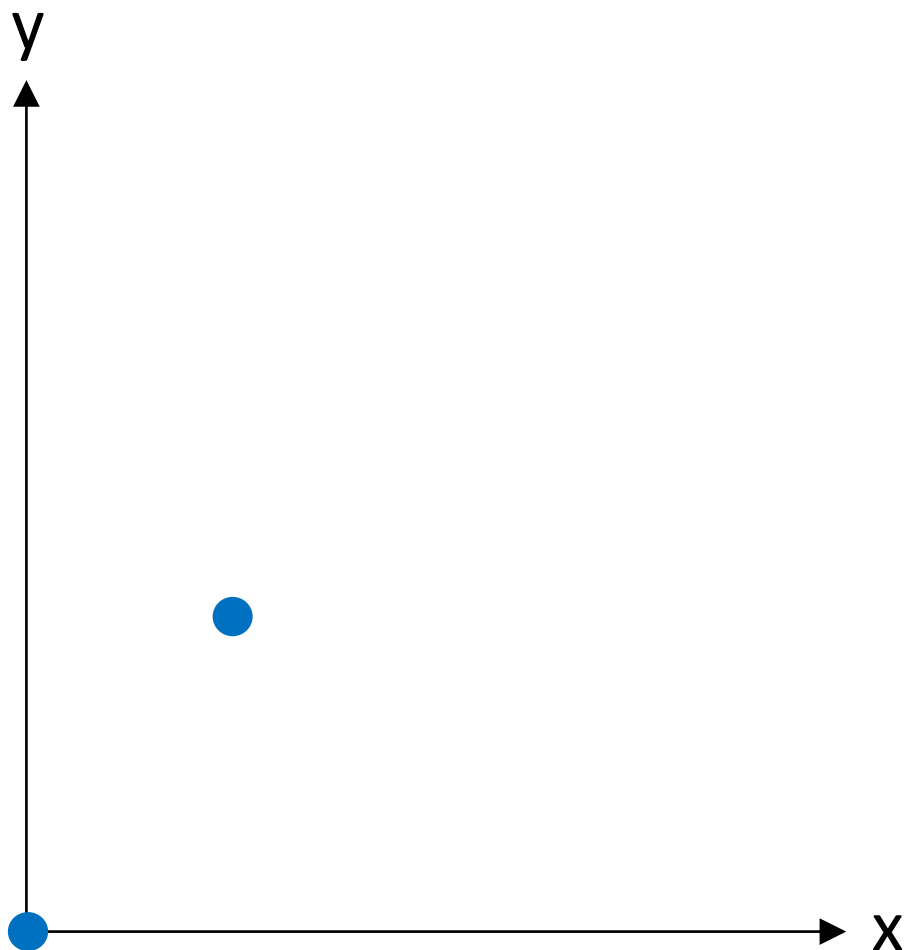
$$J(w, b) = \frac{1}{2} \left[ \left( y^{(1)} - (wx^{(1)} + b) \right)^2 + \left( y^{(2)} - (wx^{(2)} + b) \right)^2 \right]$$



# Linear Regression

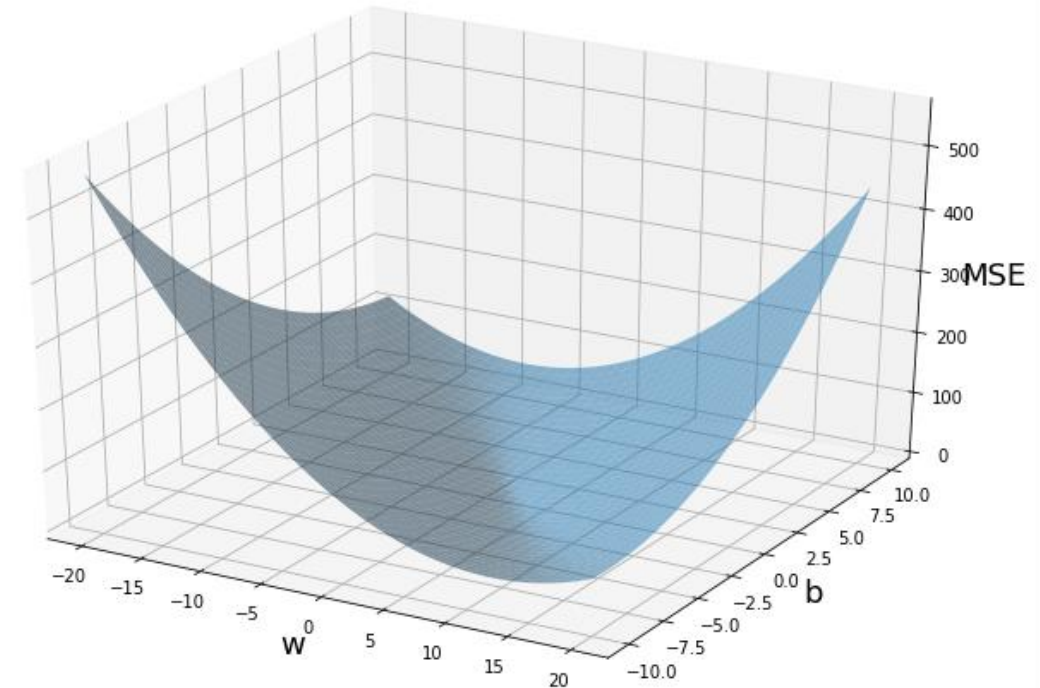
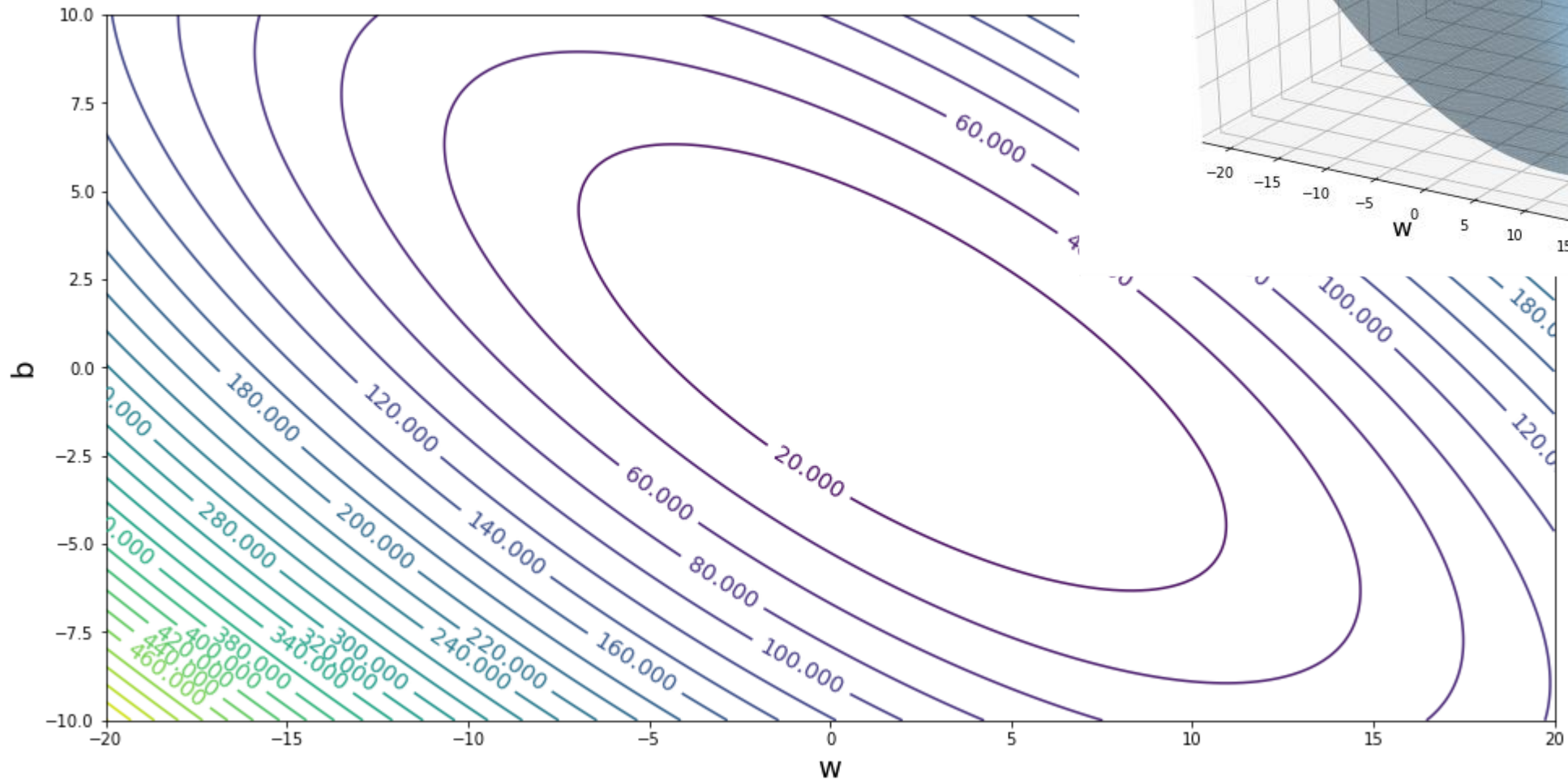
Optimizing the objective

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



# Linear Regression

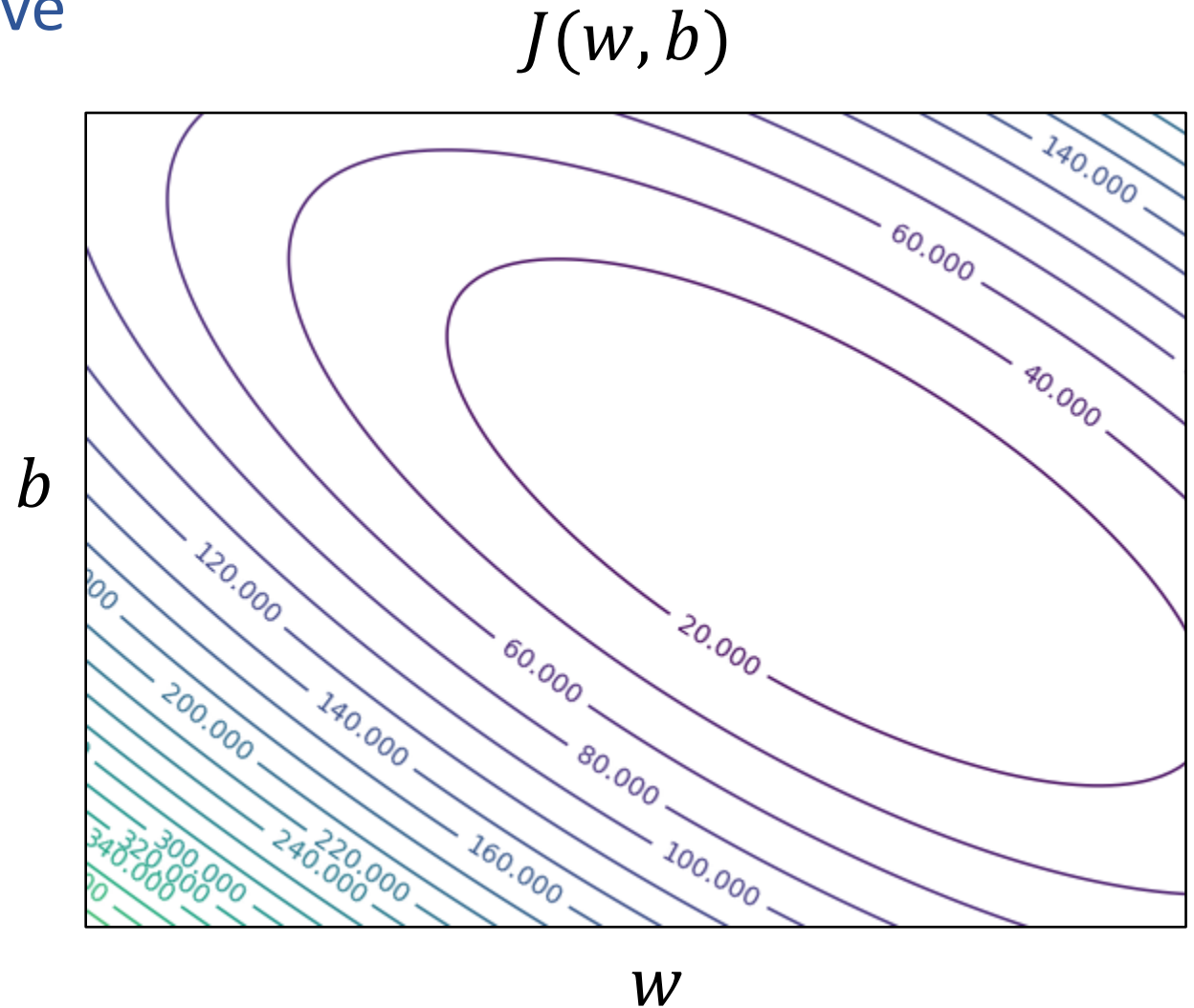
## Optimizing the objective



# Linear Regression

## Methods for optimizing the objective

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



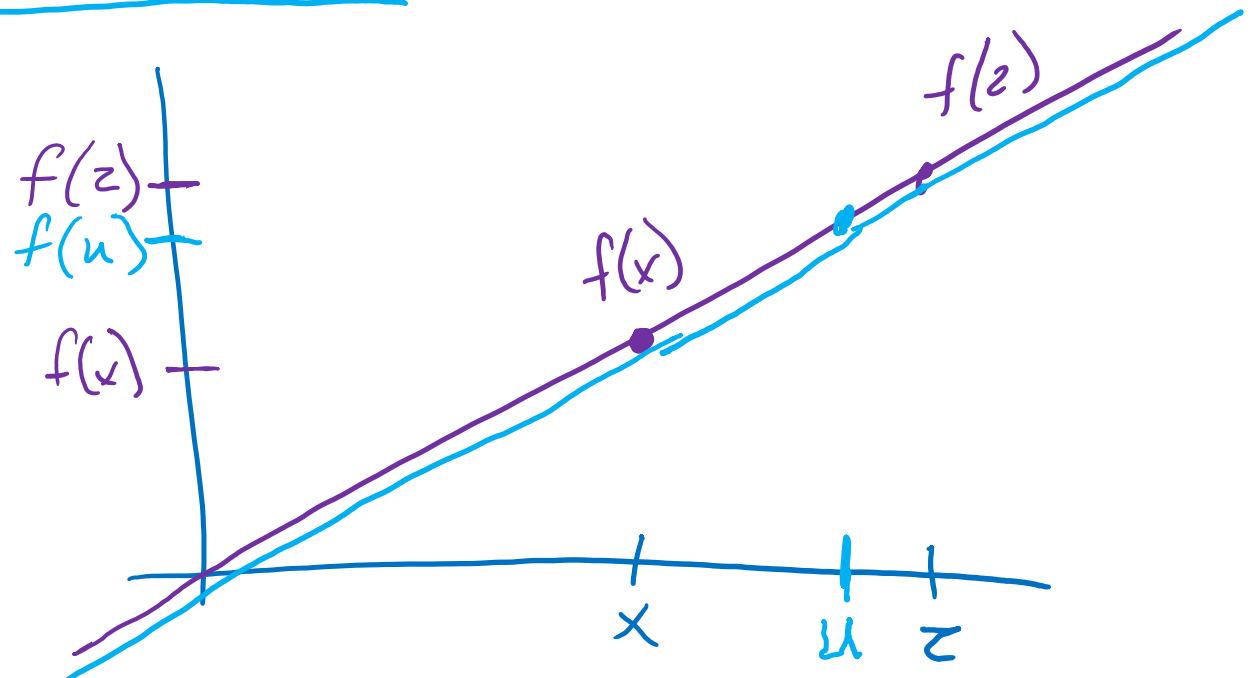
# Optimization

## Linear function

If  $f(\mathbf{x})$  is linear, then:

- $f(\mathbf{x} + \mathbf{z}) = f(\mathbf{x}) + f(\mathbf{z})$
- $f(\underline{\alpha}\mathbf{x}) = \alpha f(\mathbf{x}) \quad \forall \alpha$
- ▪  $\underline{f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{z}) = \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{z})} \quad \forall \alpha$

$$\alpha = 0.25$$





# Optimization

## Convex function

If  $f(\mathbf{x})$  is convex, then:

- $$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{z}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{z}) \quad \forall 0 \leq \alpha \leq 1$$

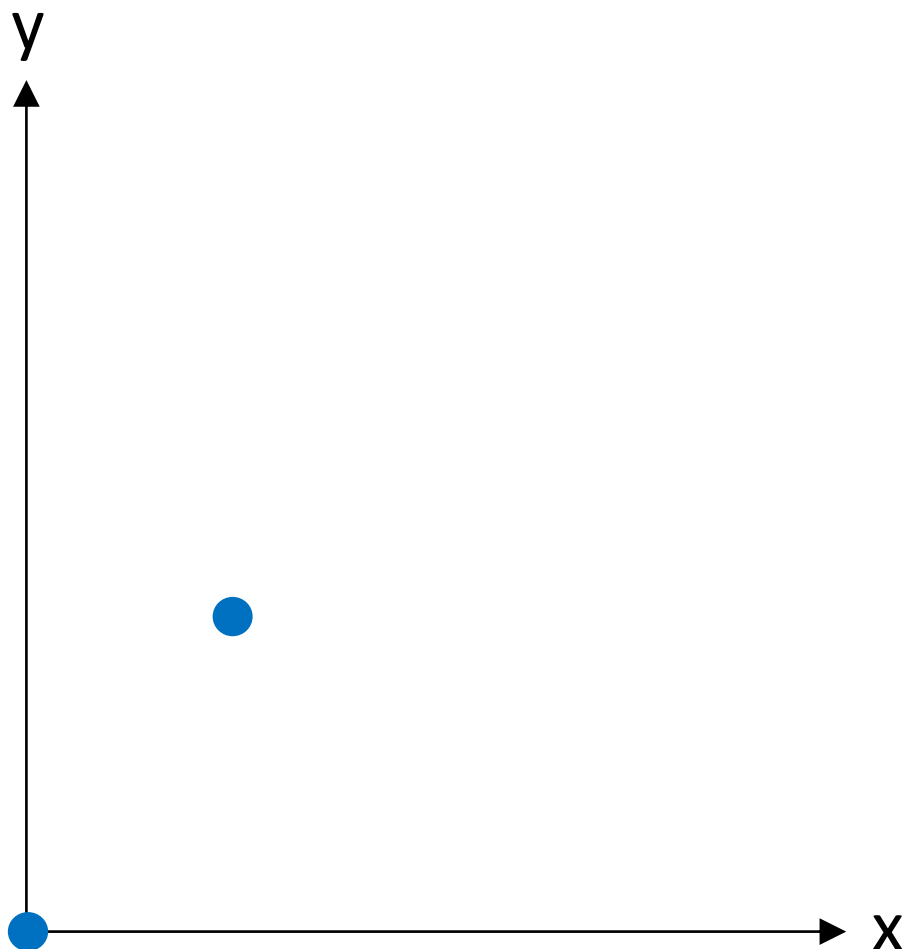
## Convex optimization

If  $f(\mathbf{x})$  is convex, then:

- Every local minimum is also a global minimum 😊

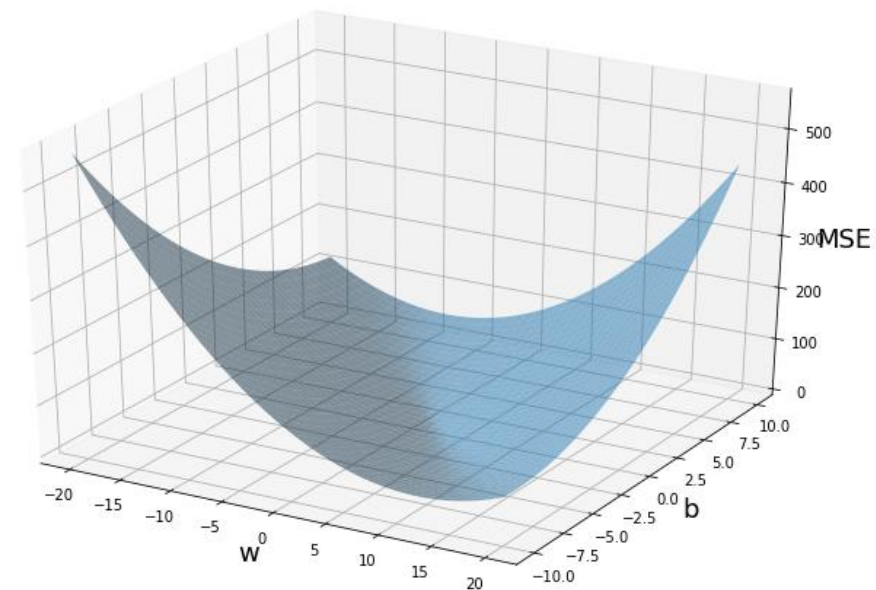
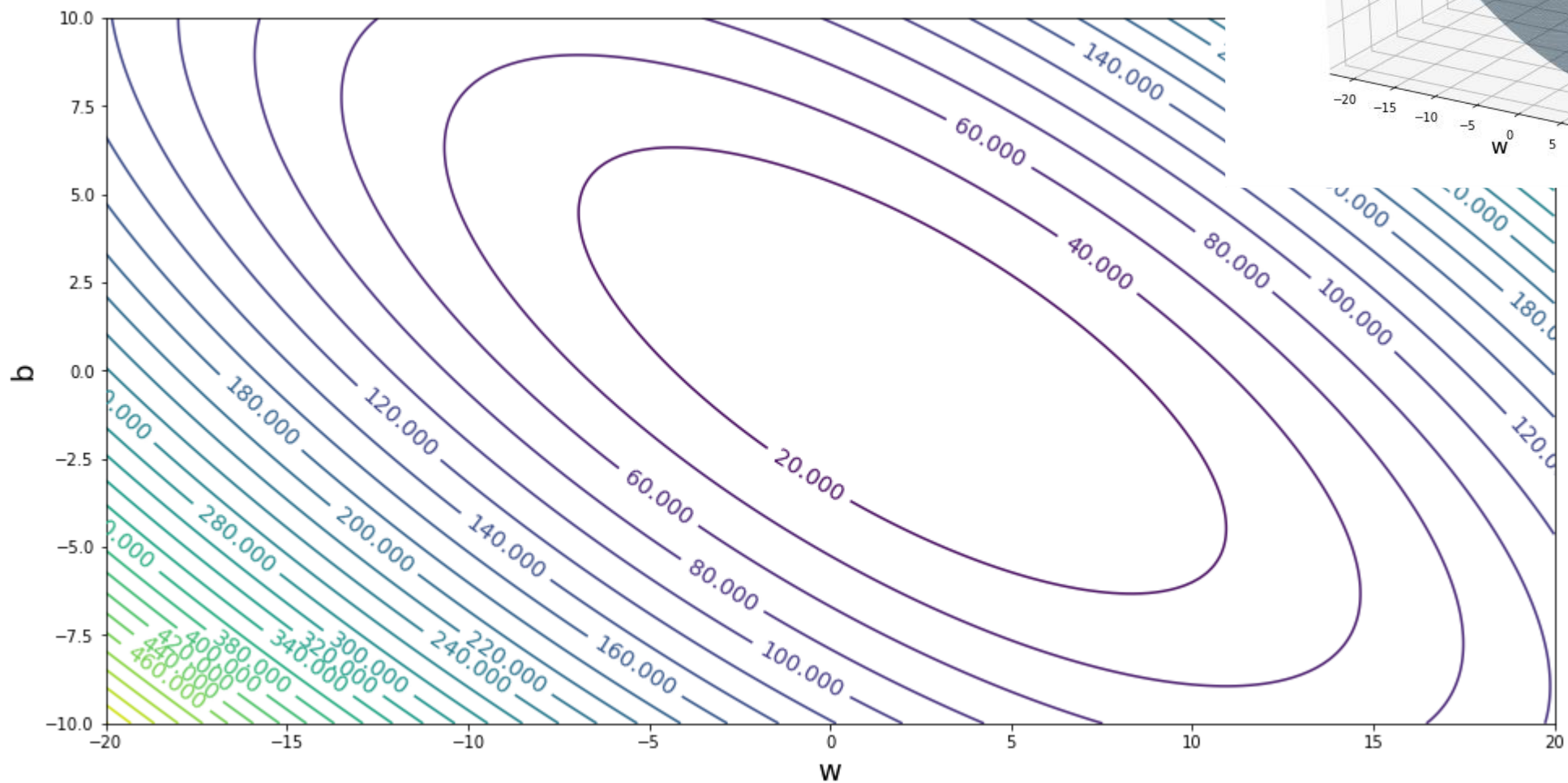
# Linear Regression

Optimizing the objective



# Optimization

## Gradients

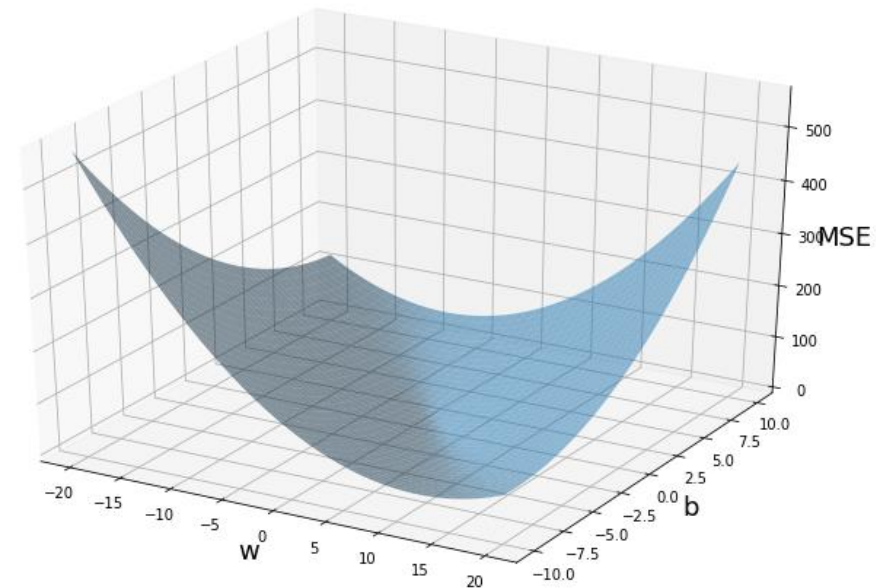
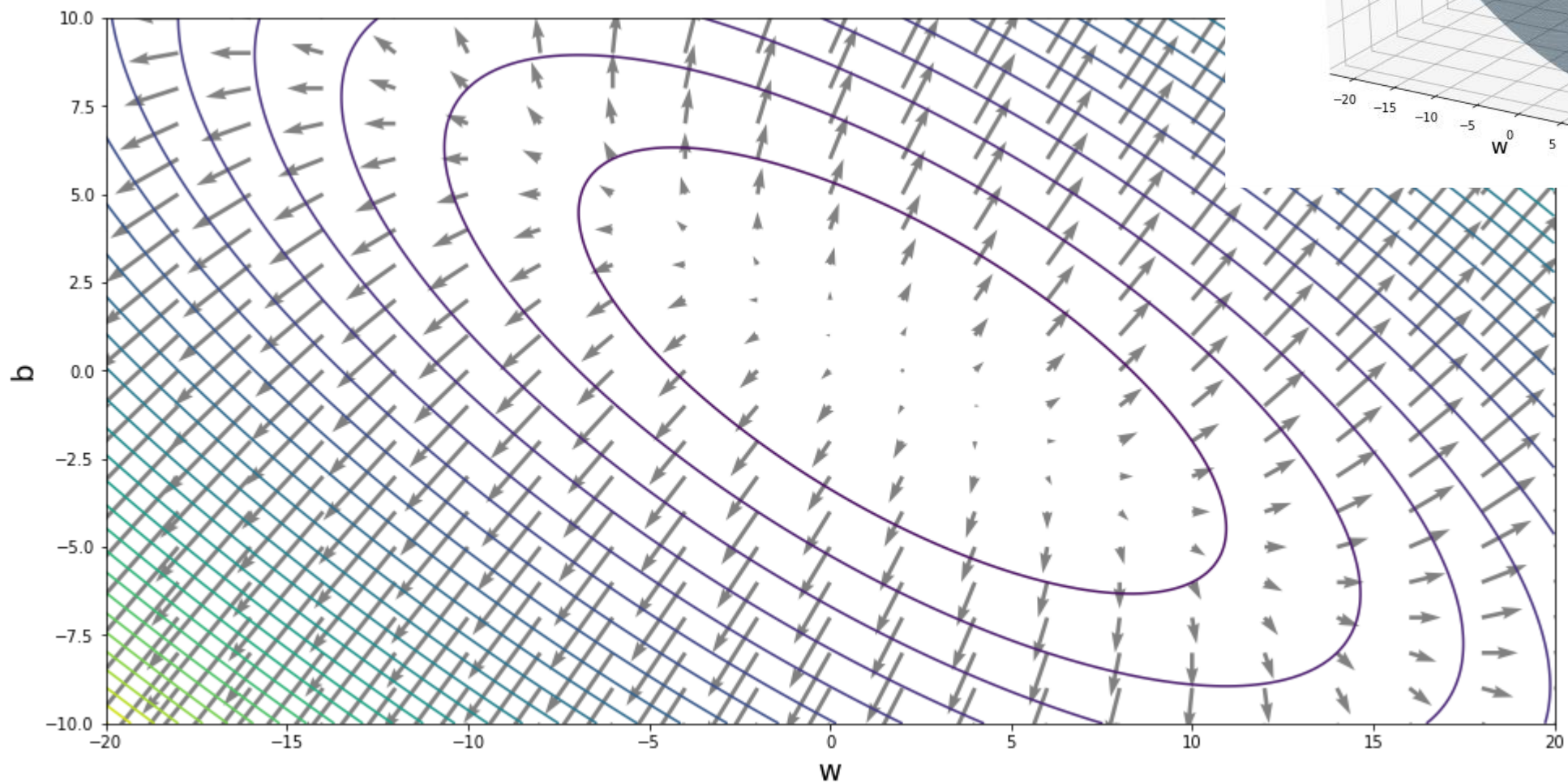


# Optimization

## Gradients

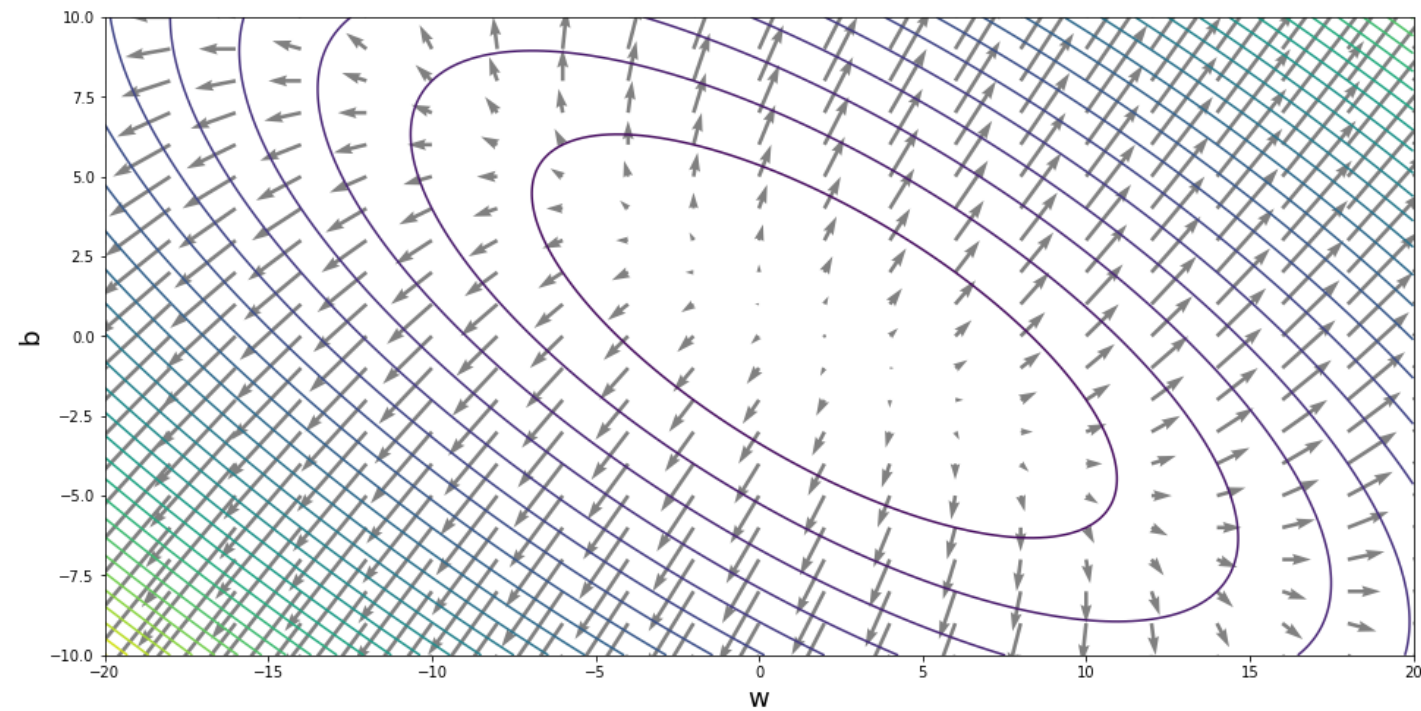
# Optimization

## Gradients



# Optimization

## Gradient descent



# Linear Regression

Expanding objective before computing gradient

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

# Linear Regression

Gradient of objective with respect to parameters

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

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

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

-- 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 --

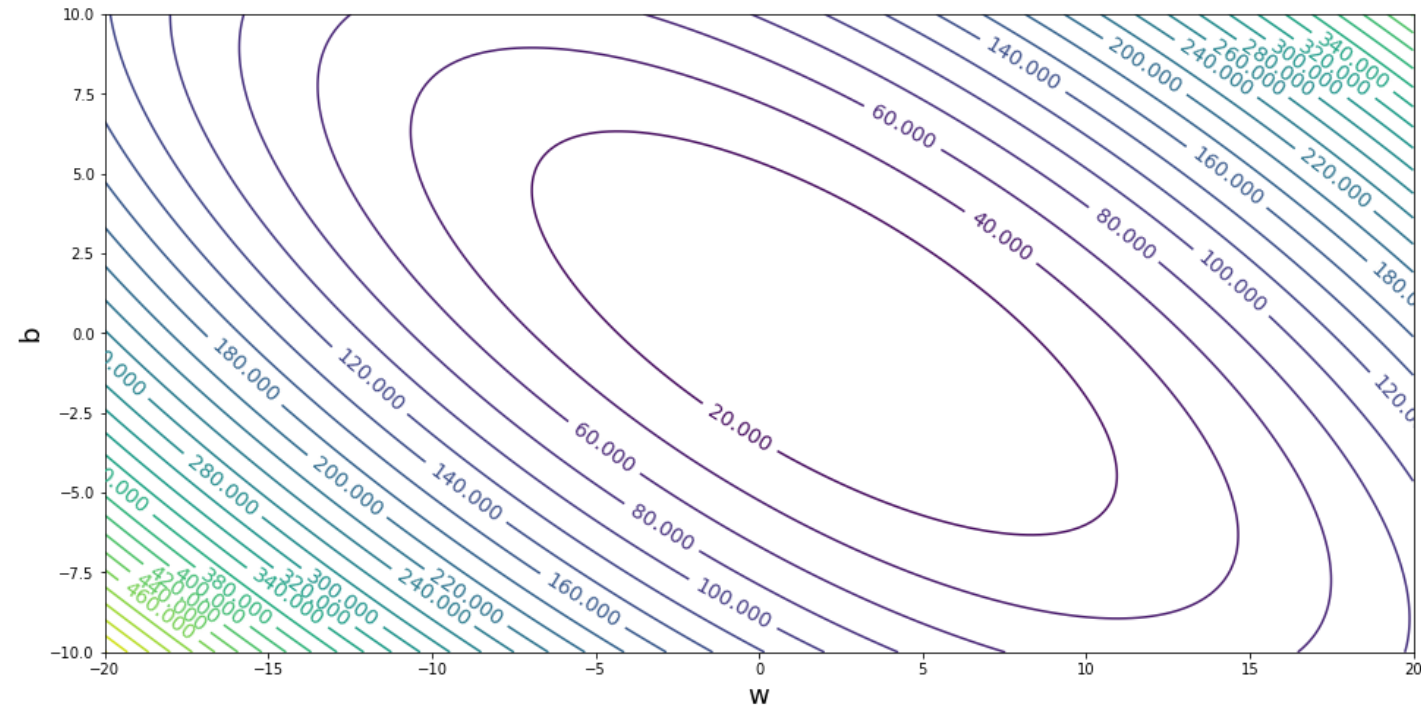
$$\frac{\partial \mathbf{z}^T \mathbf{A} \mathbf{z}}{\partial \mathbf{z}} = \mathbf{z}^T (\mathbf{A} + \mathbf{A}^T)$$



# Linear Regression

Closed-form solution

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



# Linear Regression

Number of solutions

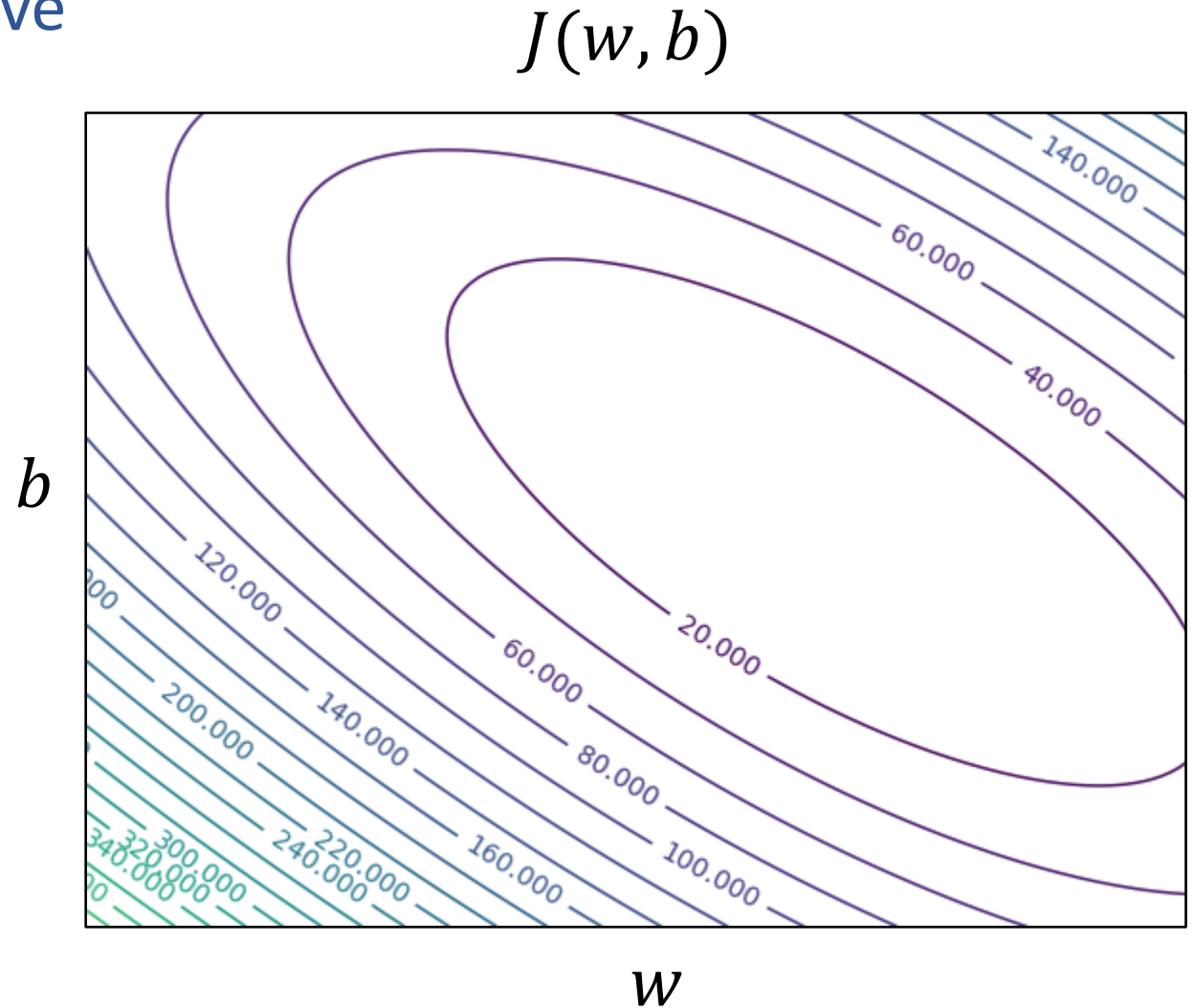
# A Note on Matrix Rank

Underlying dimensionality of the data

# Linear Regression

## Methods for optimizing the objective

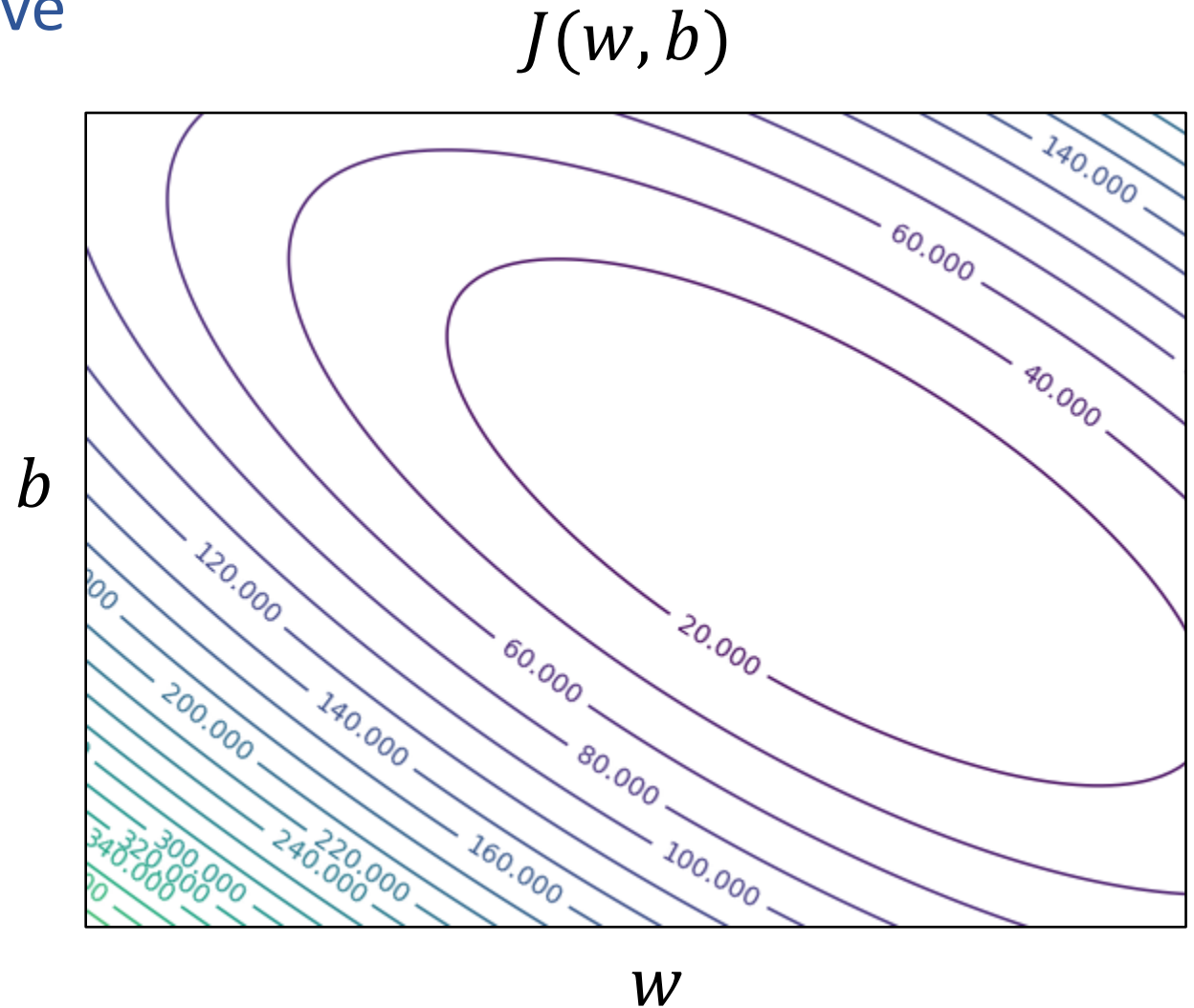
- Grid search
- Random search
- Closed-form solution
- (Batch) Gradient descent
- Stochastic gradient descent



# Linear Regression

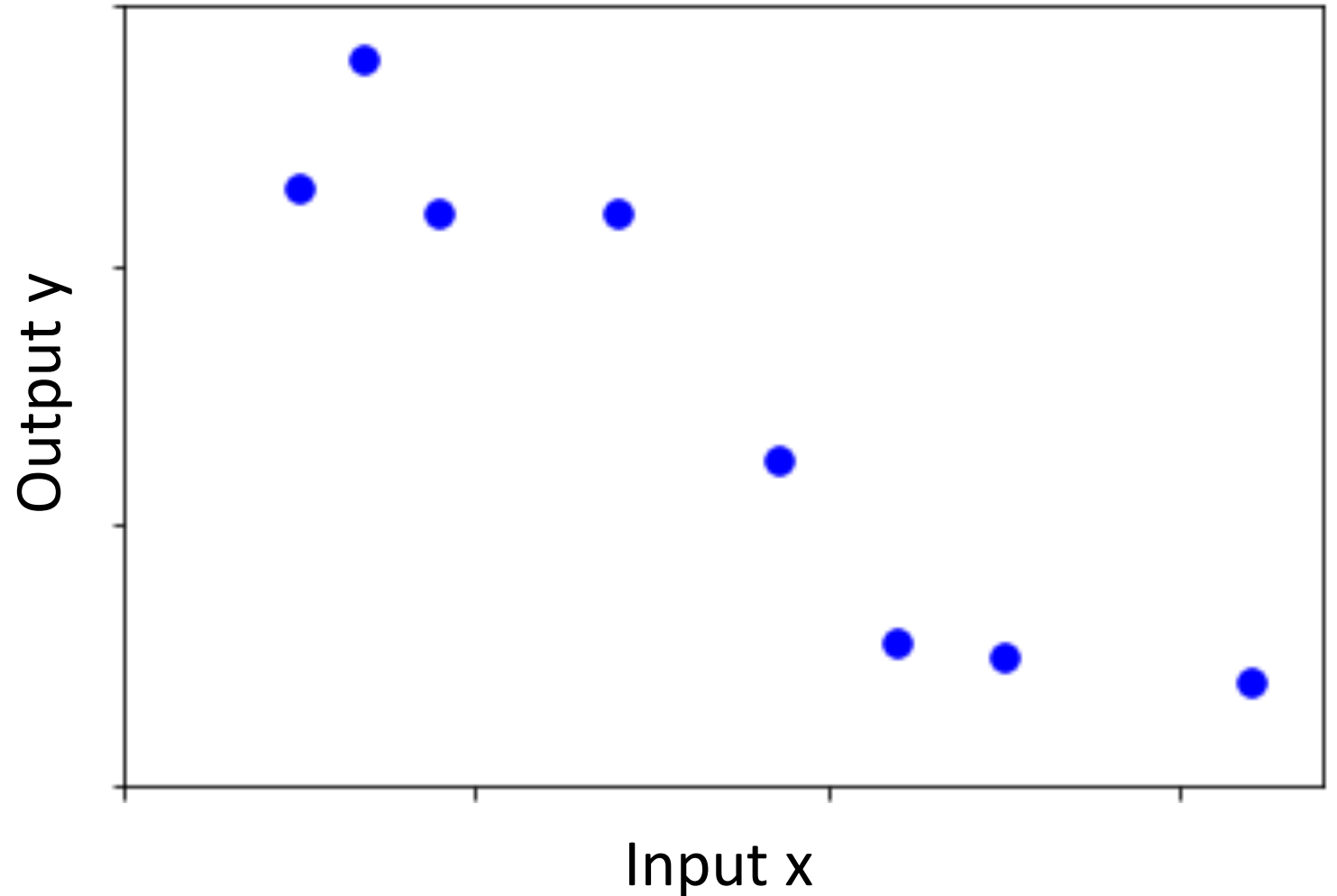
## Methods for optimizing the objective

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



# Linear Regression Gradient Descent

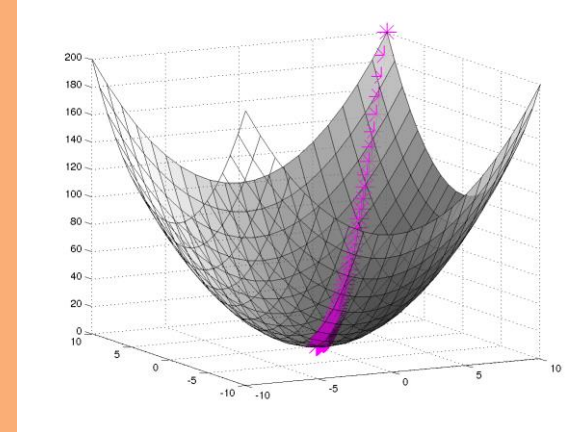
What happens in gradient descent when we have  $N=1,000,000$  training points?



# (Batch) Gradient Descent

## Algorithm 1 Gradient Descent

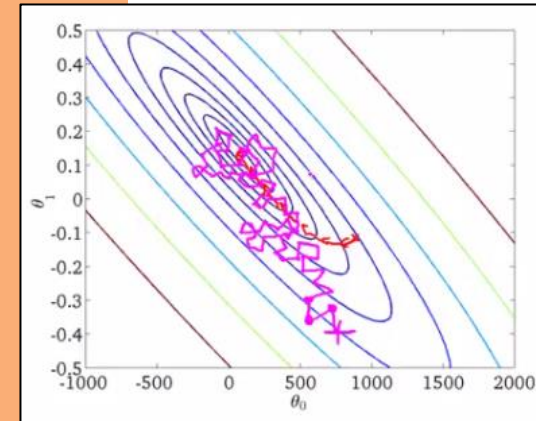
```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )  
2:    $\theta \leftarrow \theta^{(0)}$   
3:   while not converged do  
4:      $\theta \leftarrow \theta - \gamma \nabla_{\theta} J(\theta)$   
5:   return  $\theta$ 
```



# 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  $\theta$ 
```



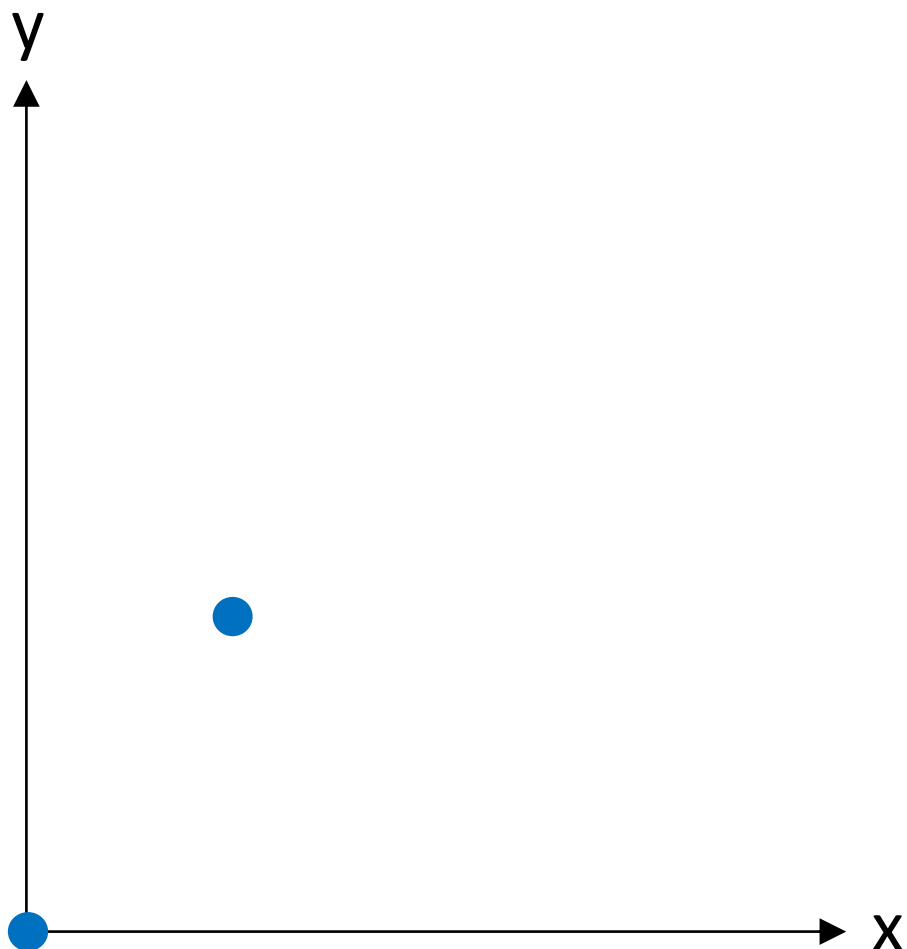
We need a per-example objective:

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

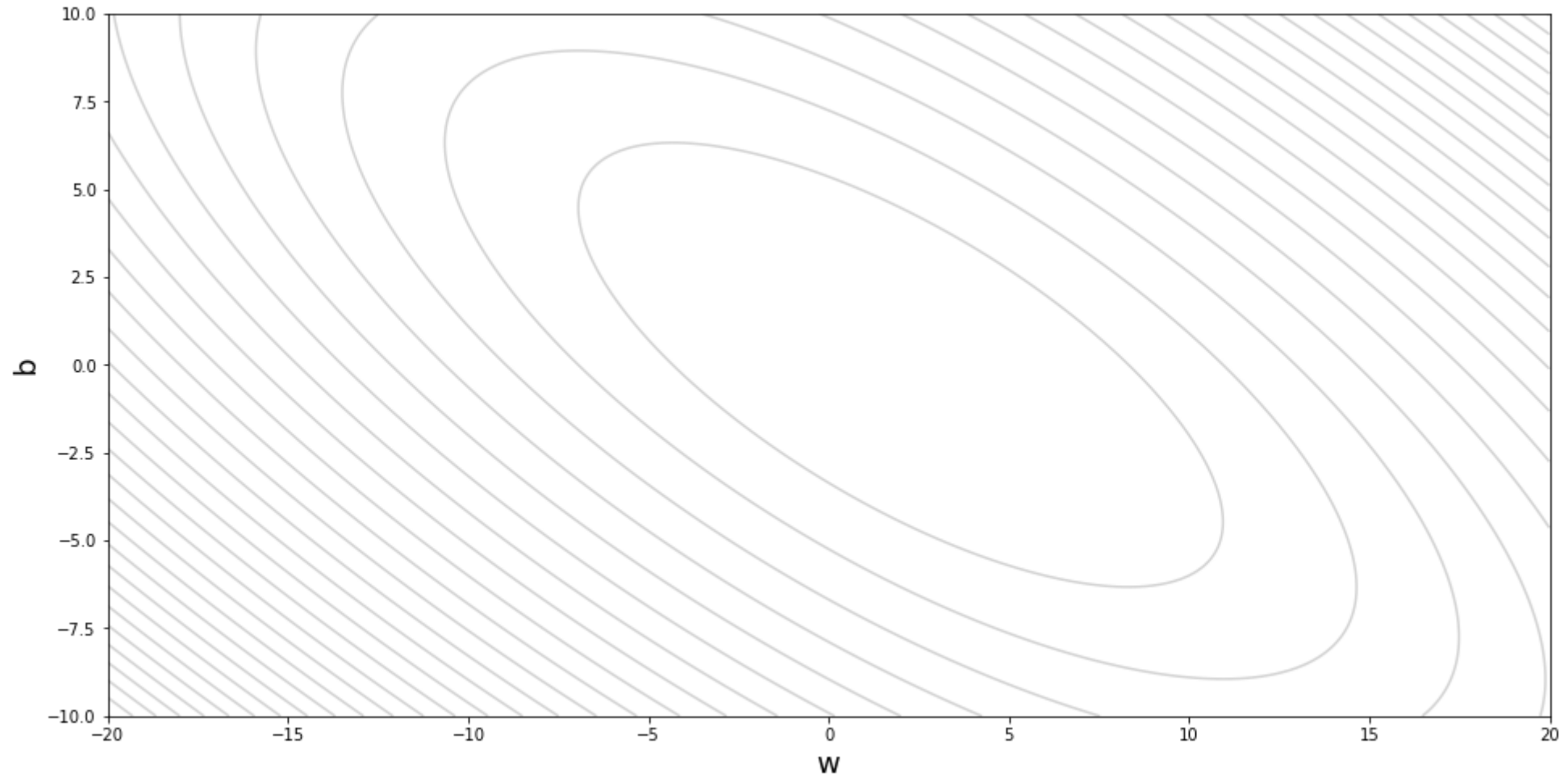


# Linear Regression

Optimizing the objective

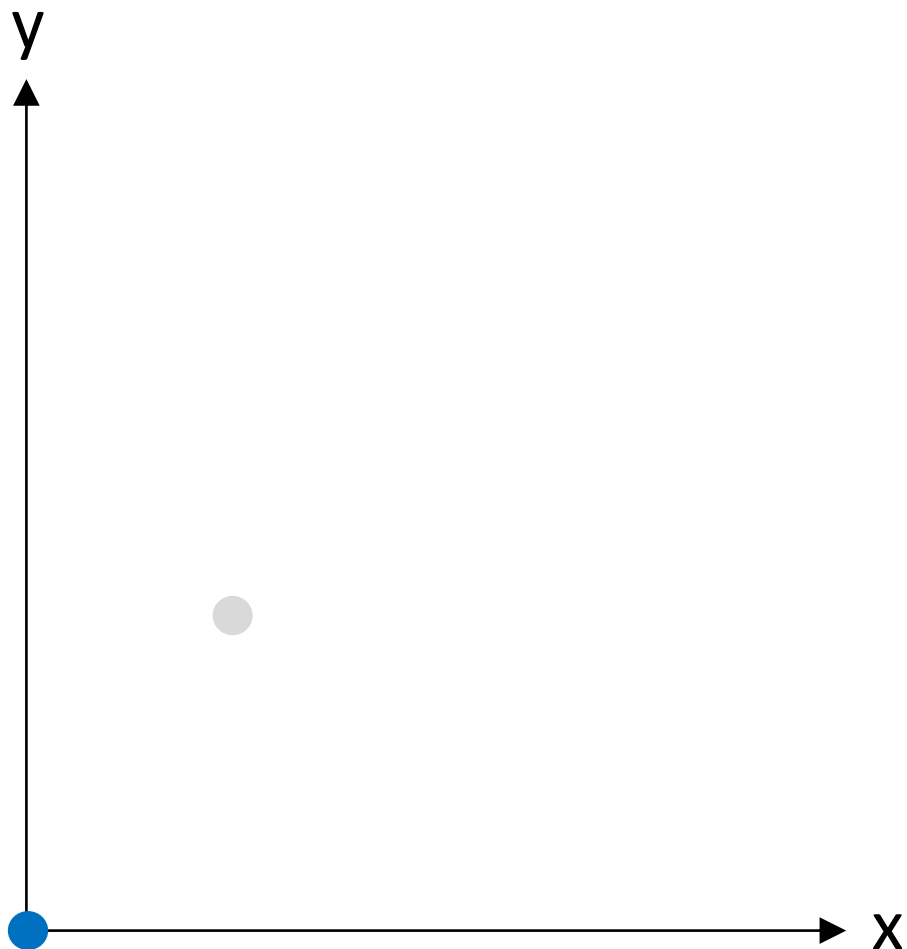


# Stochastic Gradient Descent

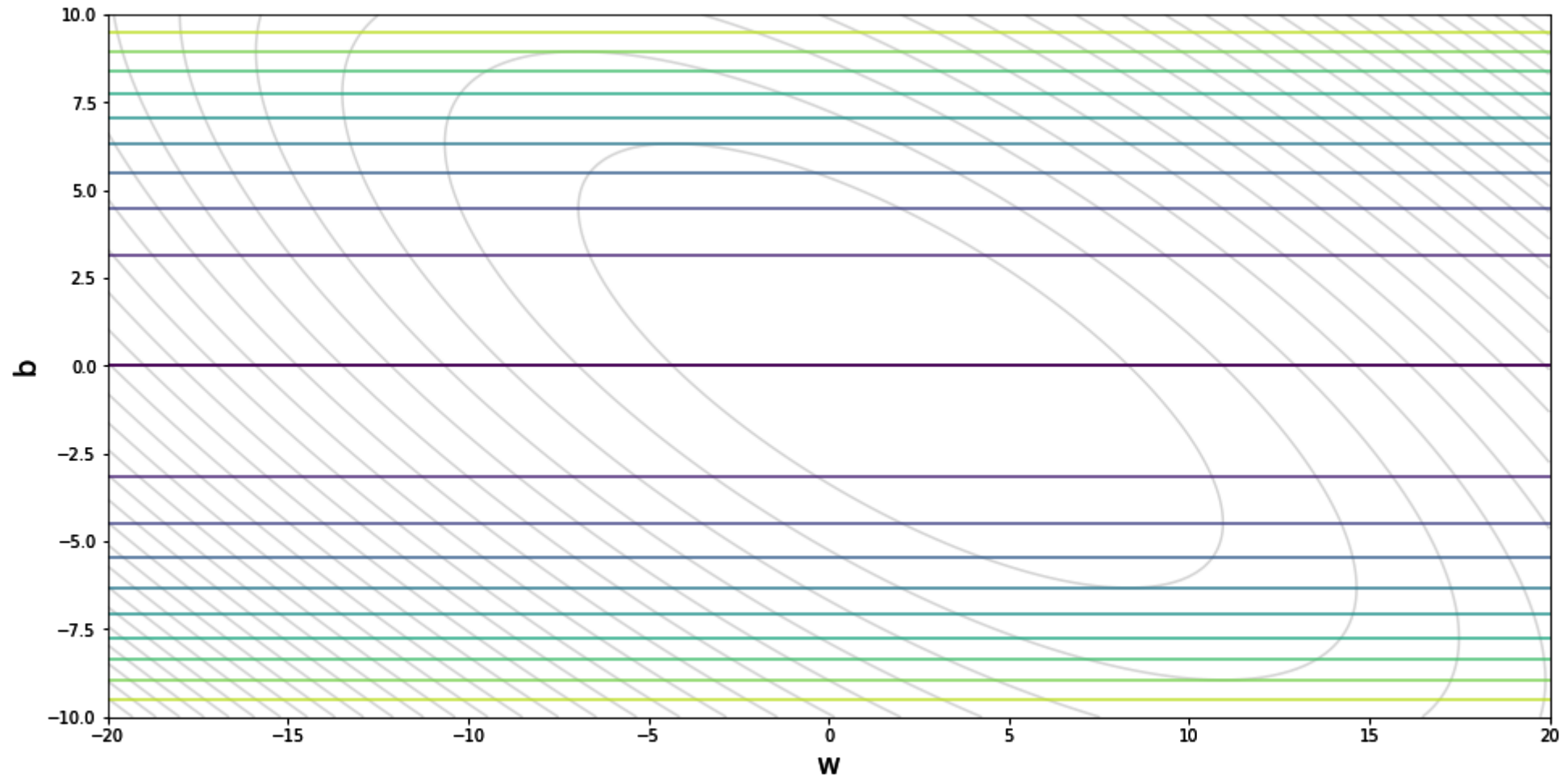


# Linear Regression

Optimizing the objective

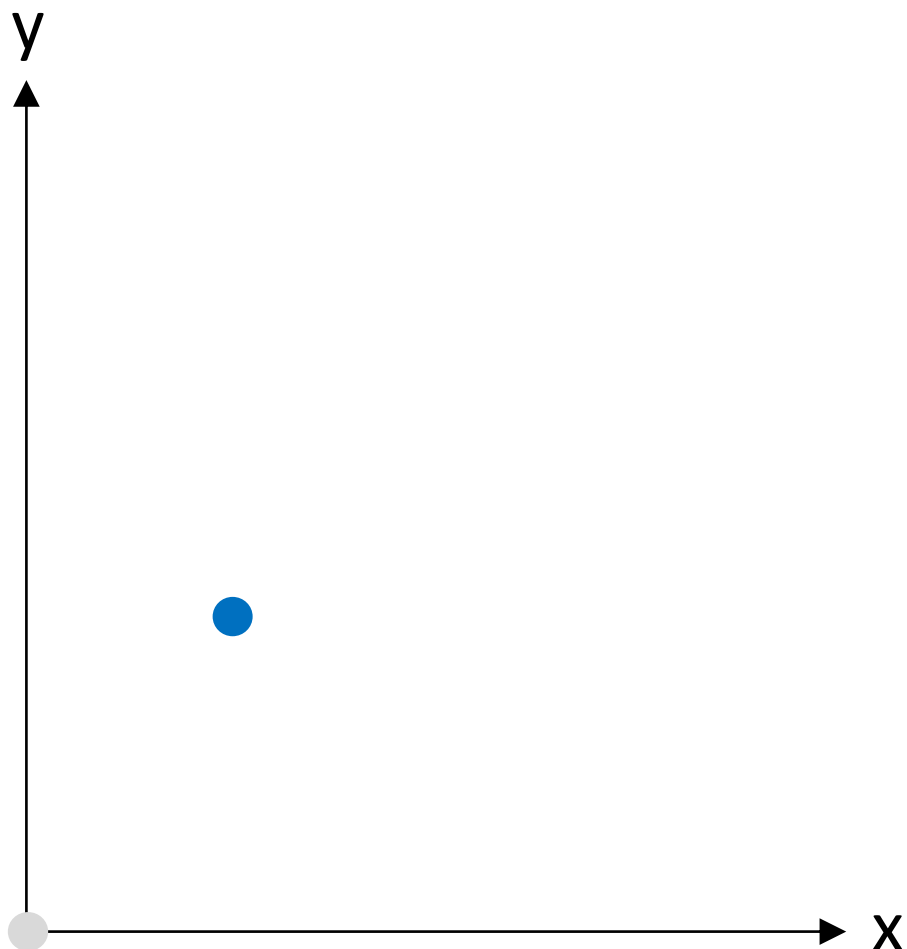


# Stochastic Gradient Descent

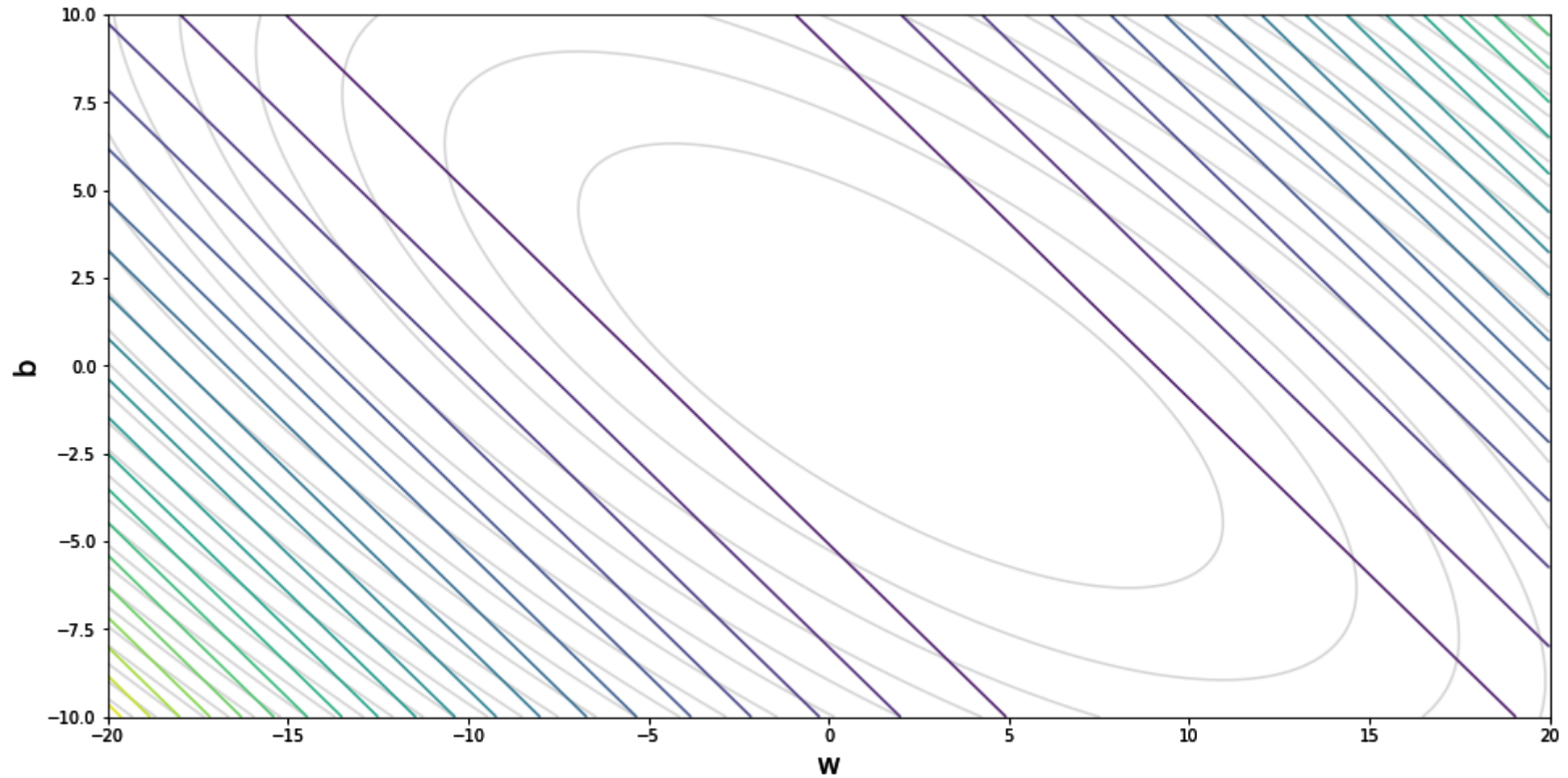


# Linear Regression

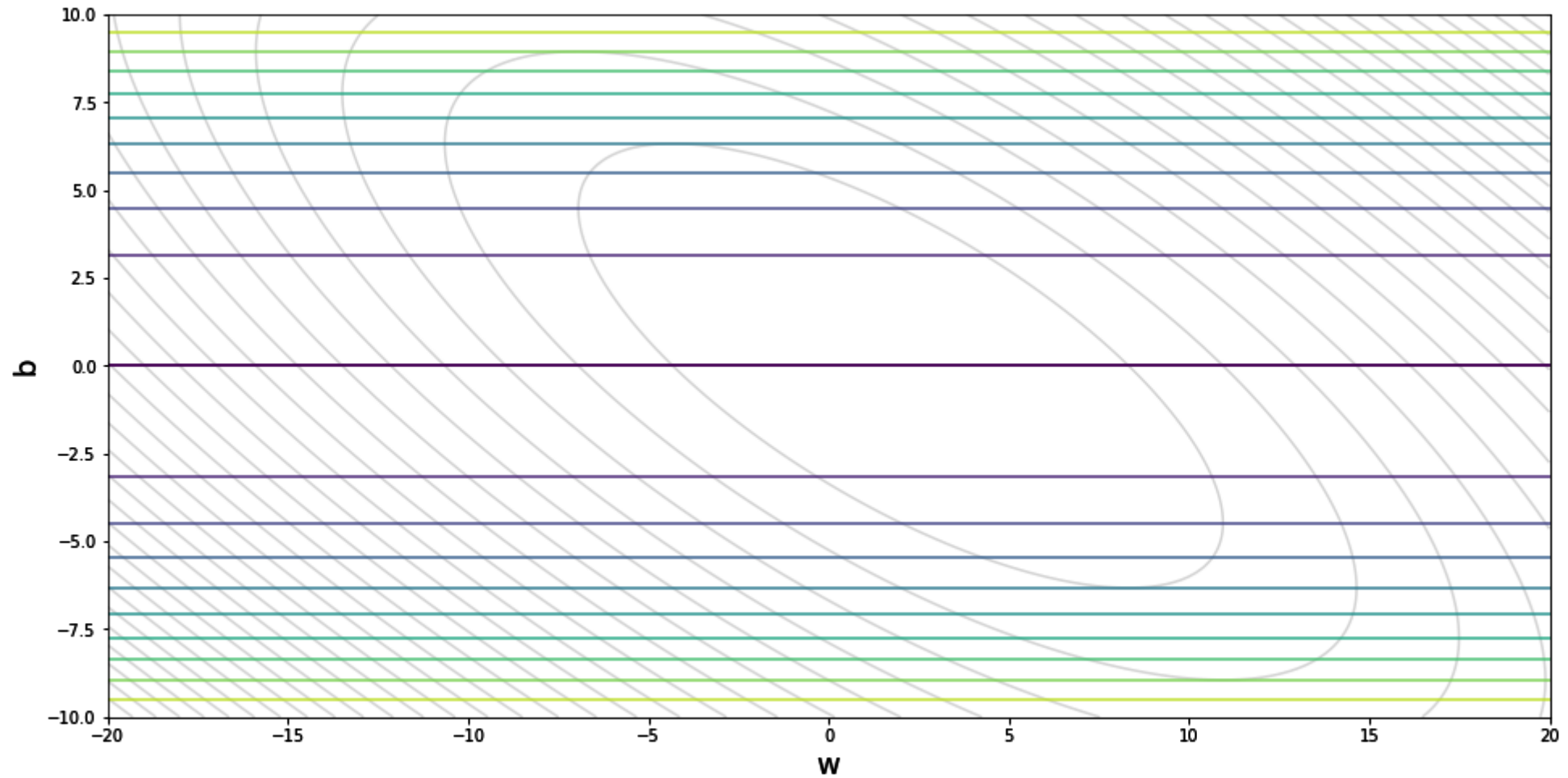
Optimizing the objective



# Stochastic Gradient Descent



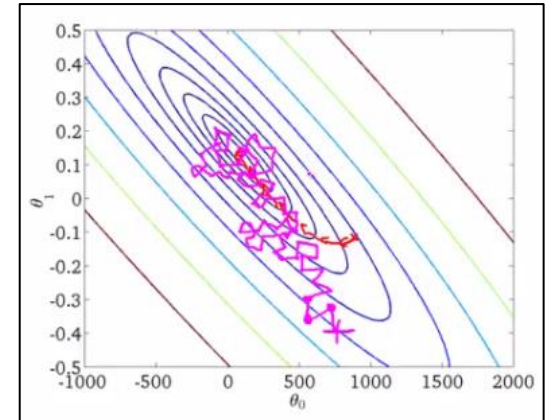
# Stochastic Gradient Descent



# 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:

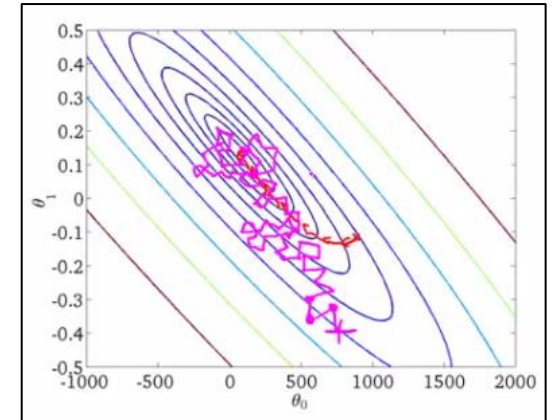
$$\text{Let } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$



# 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  $\theta$ 
```

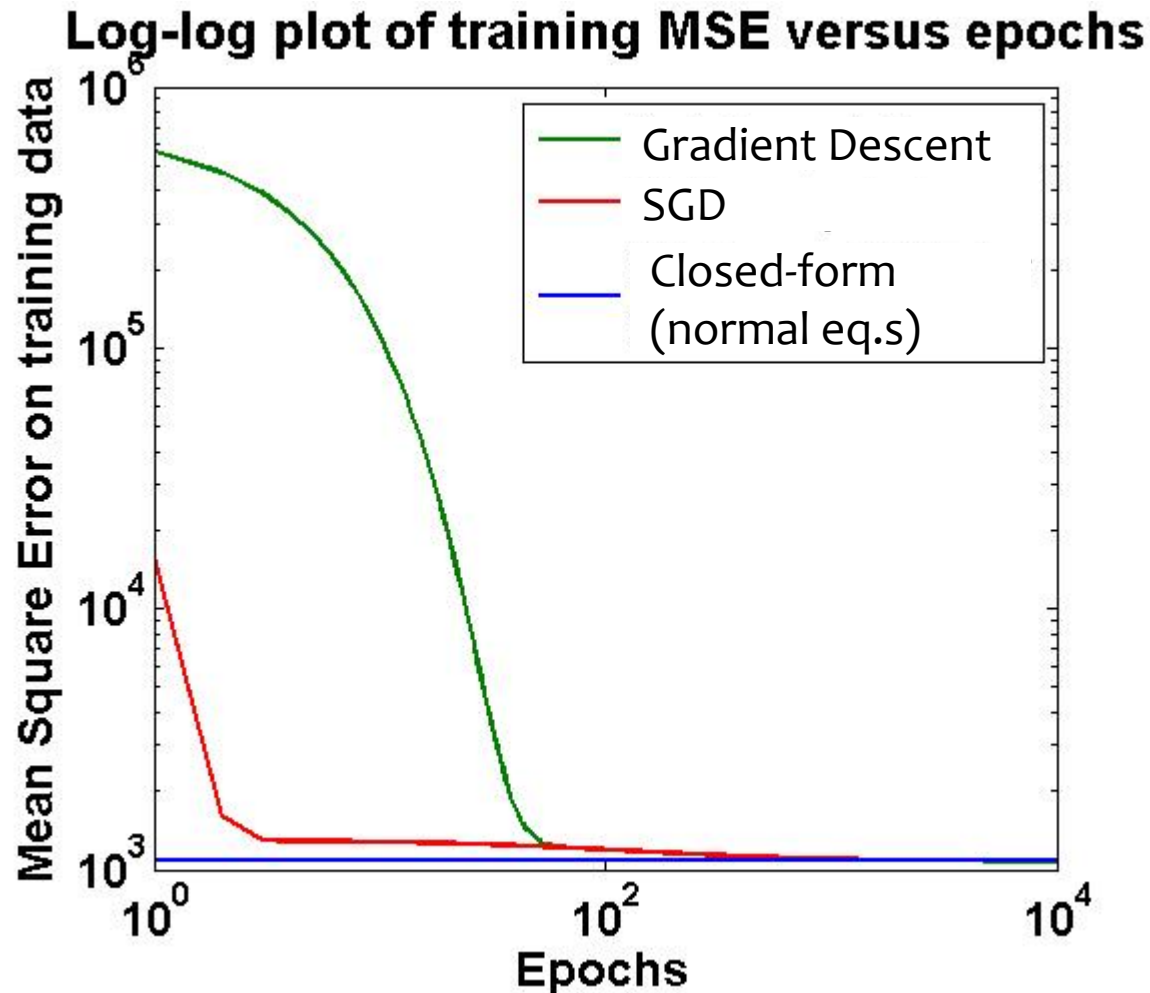


We need a per-example objective:

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

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