

Linear combination of features

$$\hat{y} = \underset{\substack{\downarrow \\ \text{bias term}}}{\Theta_0} + \Theta_1 x_1 + \dots + \Theta_d x_d = \underbrace{X \Theta}_{\substack{\downarrow \\ n \times (d+1)}} \rightsquigarrow (d+1 \times 1)$$

Training data

Now, we need define a goal or objective function: Minimize the difference between  $y_{\text{actual}}$  &  $\hat{y}_{\text{predicted}}$

$$L(\Theta) = E(\Theta) = \left( \frac{1}{N} \sum_{i=1}^N \right) (y_a - \hat{y}_p)^2 = E[(y_a - \hat{y}_p)^2]$$

$$y = X\Theta \Rightarrow \underset{(d+1 \times 1)}{\Theta} = \underbrace{(X^{-1})}_{\text{}} y \quad \underbrace{(X^{-1})}_{\text{}} = \underbrace{(X^T X)^{-1}}_{\text{}} X^T$$

$$\underset{(d+1 \times 1)}{\Theta} = (X^T X)^{-1} X^T y$$

$$\Theta^{t+1} \leftarrow \Theta^t - \alpha \frac{\partial L(\Theta)}{\partial \Theta} \Rightarrow \text{GD} : \text{In each iteration, we need to go over All data points}$$

$$\Rightarrow \text{SGD} : \text{" " " " " " " " One "}$$

$$\Rightarrow \text{BGD} : \text{" " " " " " " " a subset or batch "}$$

$$X = \begin{bmatrix} \text{sqf} \end{bmatrix} \quad Y = \begin{bmatrix} \text{Rent price} \end{bmatrix} \quad Z = \begin{bmatrix} \text{sqf} & \text{sqf}^2 & \text{sqf}^3 & \dots & \text{sqf}^d \end{bmatrix} \quad \text{overfitting}$$

$$E(\Theta) = L(\Theta) = E[(y_a - \hat{y}_p)^2] = \underbrace{\text{bias}^2 + \text{Variance}}$$

bias-variance trade off

$$X = \begin{bmatrix} h & w \end{bmatrix} \quad Z = \begin{bmatrix} h & w & h^2 & w^2 & hw & h^2w & hw^2 & \dots \end{bmatrix}$$

# Regularized Linear Regression

Mahdi Roozbahani  
Georgia Tech

# EVERY GROUP PROJECT




DOES 99%  
OF THE WORK

HAS NO IDEA  
WHAT'S GOING  
ON THE  
WHOLE TIME

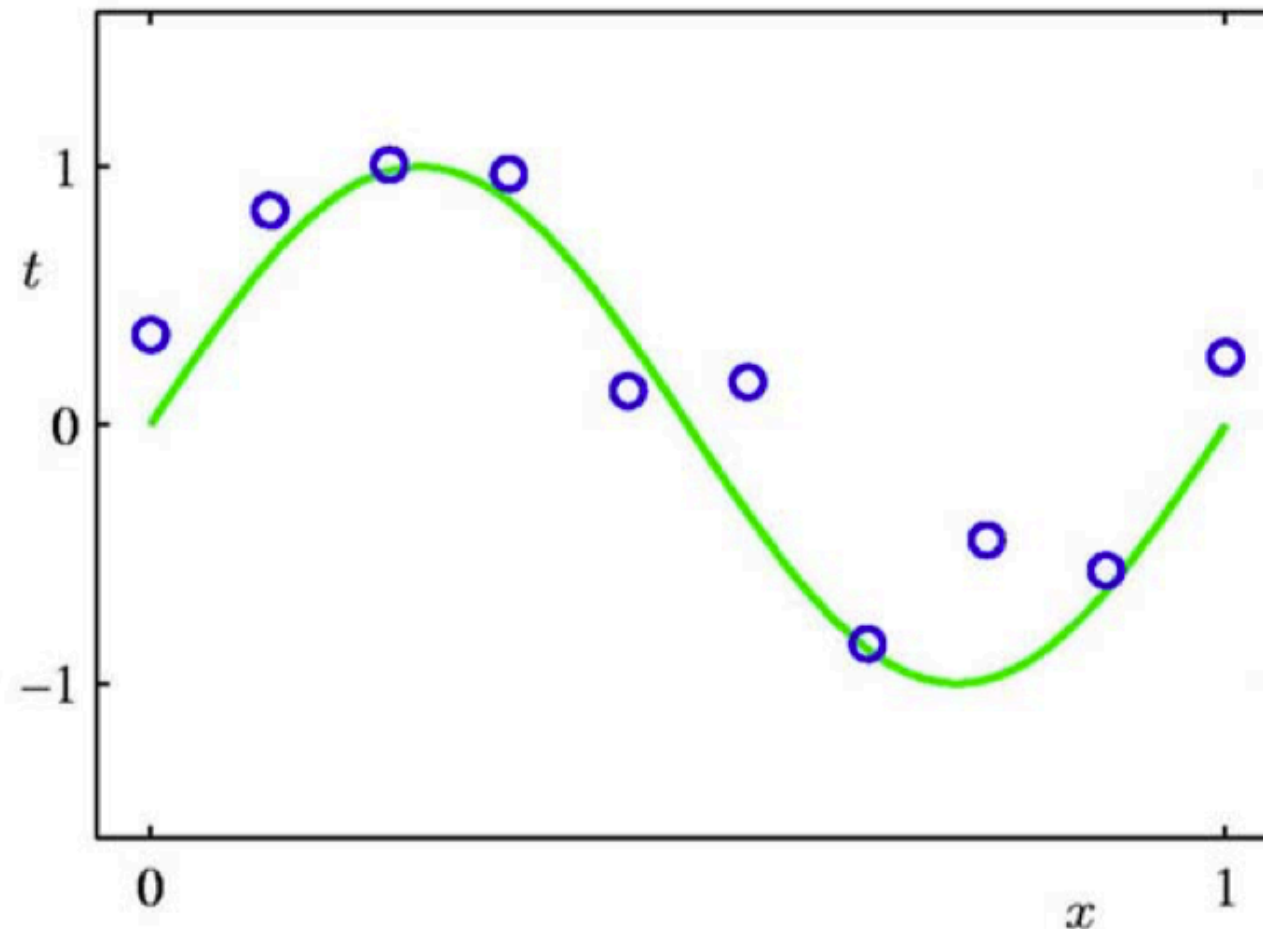
SAYS HE'S  
GOING TO  
HELP  
BUT HE'S  
NOT

DISAPPEAR  
AT THE VERY  
BEGINNING AND  
DOESN'T SHOW  
UP AGAIN TIL  
THE VERY END

# Outline

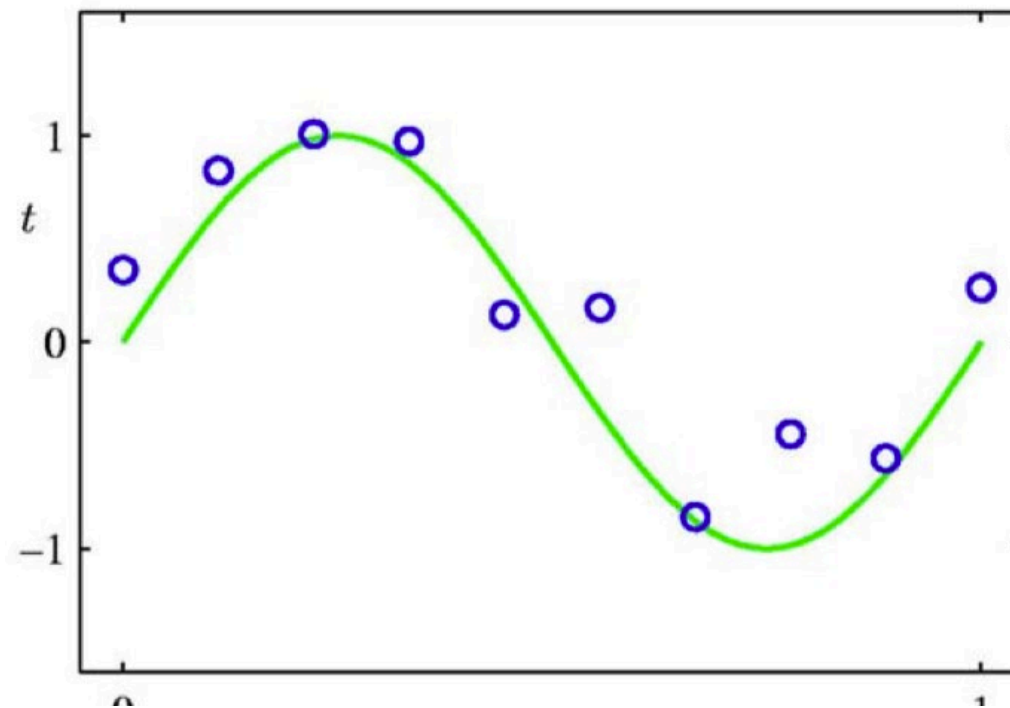
- Overfitting and regularized learning 
- Ridge regression
- Lasso regression
- Determining regularization strength

# Regression: Recap



- Suppose we are given a training set of  $N$  observations  $(x_1, \dots, x_N)$  and  $(y_1, \dots, y_N)$
- Regression problem is to estimate  $y(x)$  from this data

# Regression: Recap



- Want to fit a polynomial regression model

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_d x^d + \epsilon$$

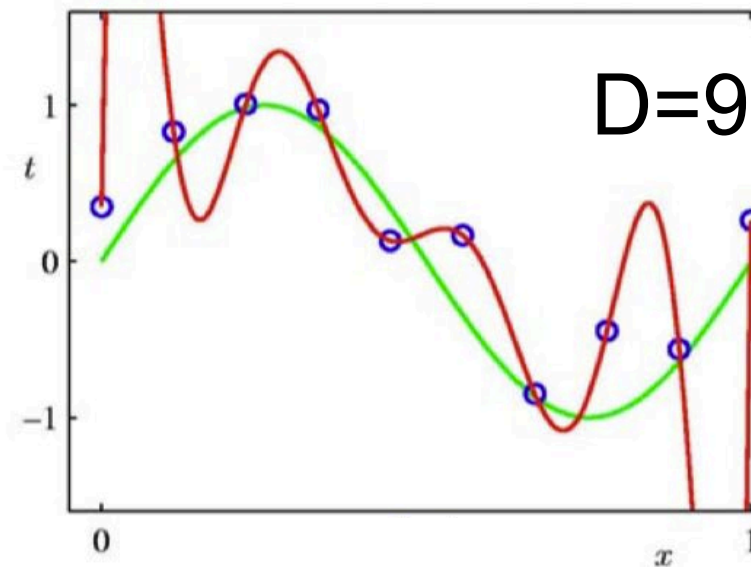
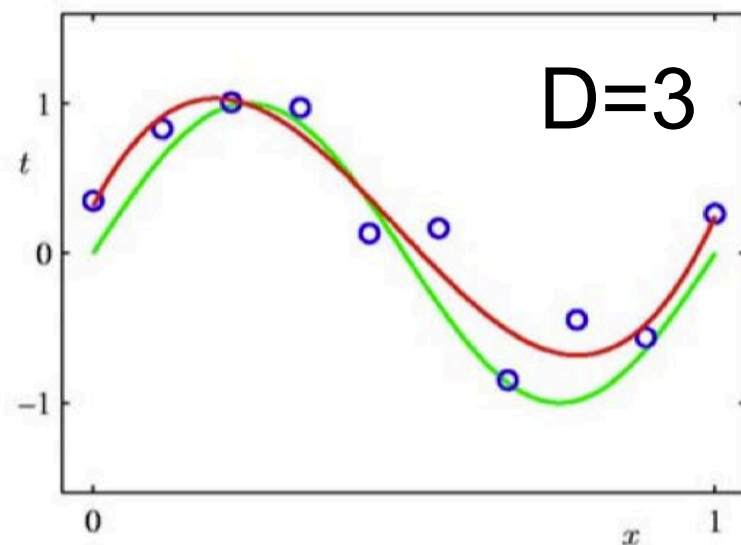
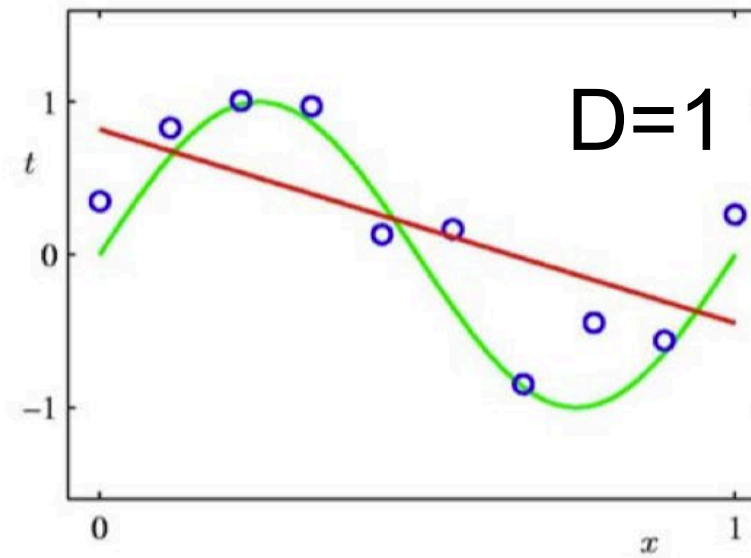
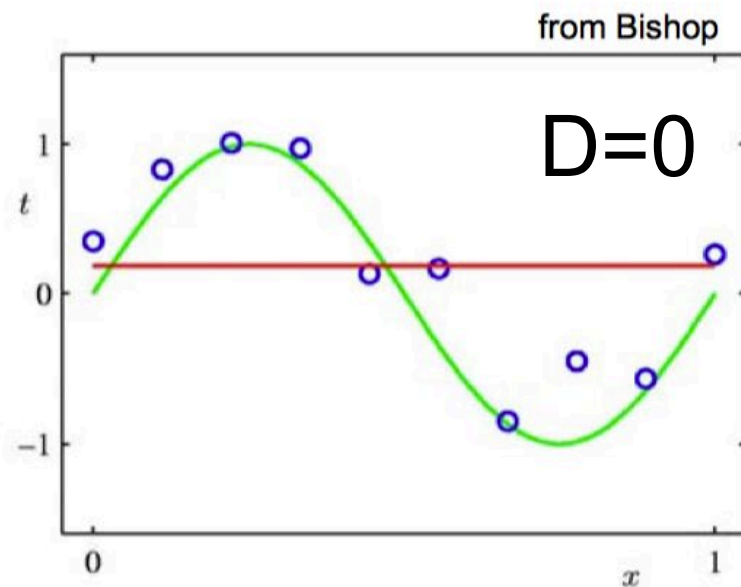
*Handwritten annotations in purple:*  
Above the equation, the terms are grouped as  $\theta_0 + \theta_1 z_1 + \theta_2 z_2 + \dots + \theta_d z_d$ .  
In the equation,  $\theta_2$  and  $x^2$  are circled, and  $\theta_d$  and  $x^d$  are also circled.

- $z = \{1, x, x^2, \dots, x^d\} \in R^d$  and  $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_d)^T$

$$y = z\theta$$

*Handwritten annotation in purple:* The equation  $y = z\theta$  is enclosed in a purple box.

# Which One is Better?

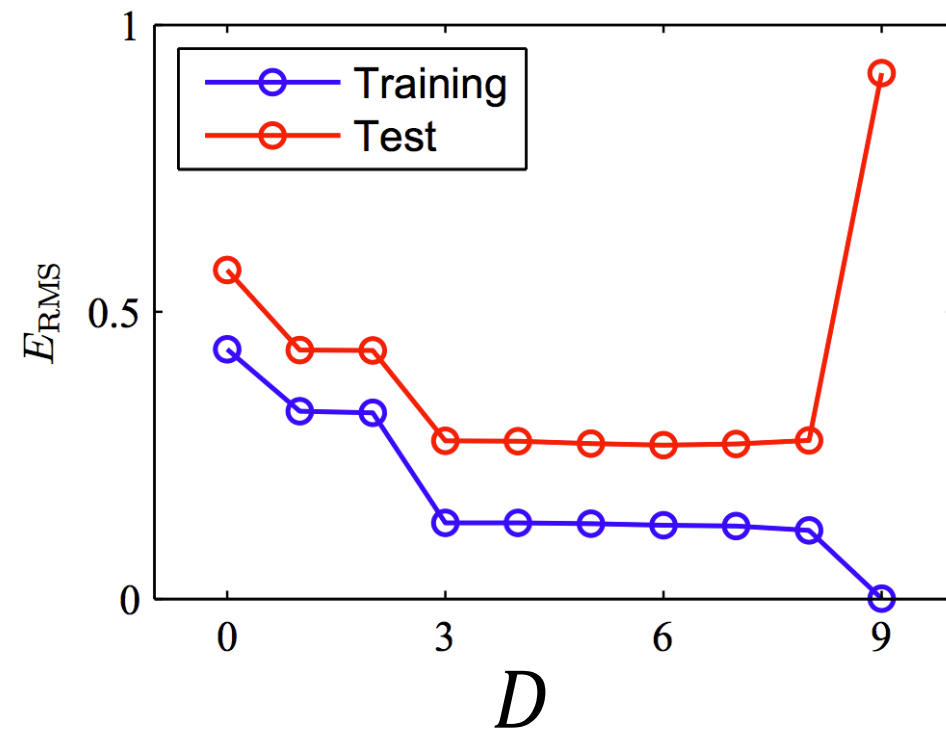
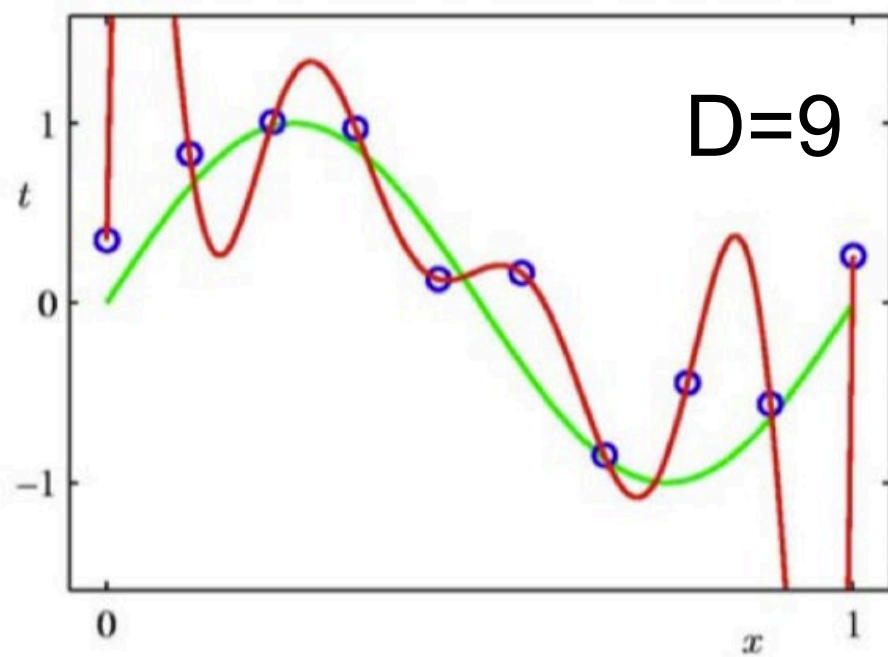


- Can we increase the maximal polynomial degree to very large, such that the curve passes through all training points?

No, this can lead to **overfitting**!

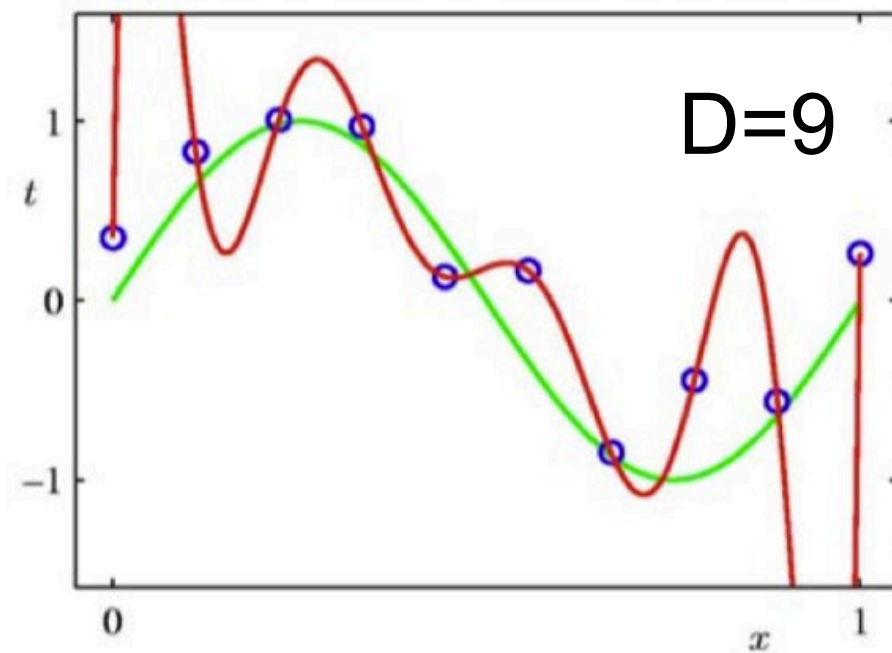


# The Overfitting Problem



- The training error is very low, but the error on test set is large.
- The model captures not only patterns but also noisy nuisances in the training data.

# The Overfitting Problem



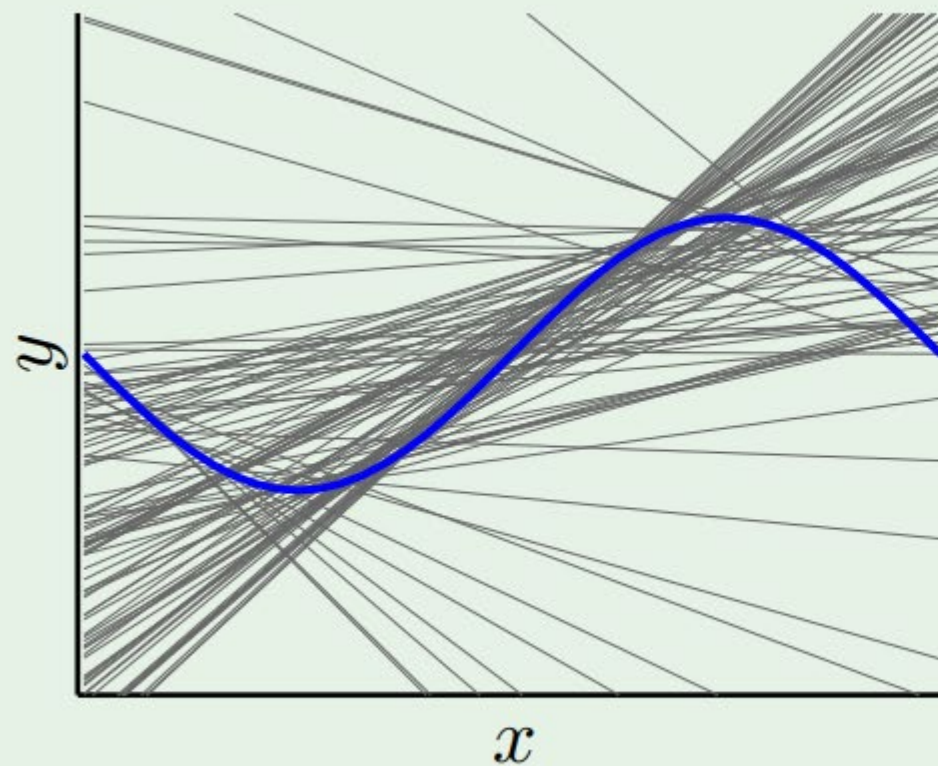
- In regression, overfitting is often associated with large Weights  
(severe oscillation)
- How can we address overfitting?

$\Theta$

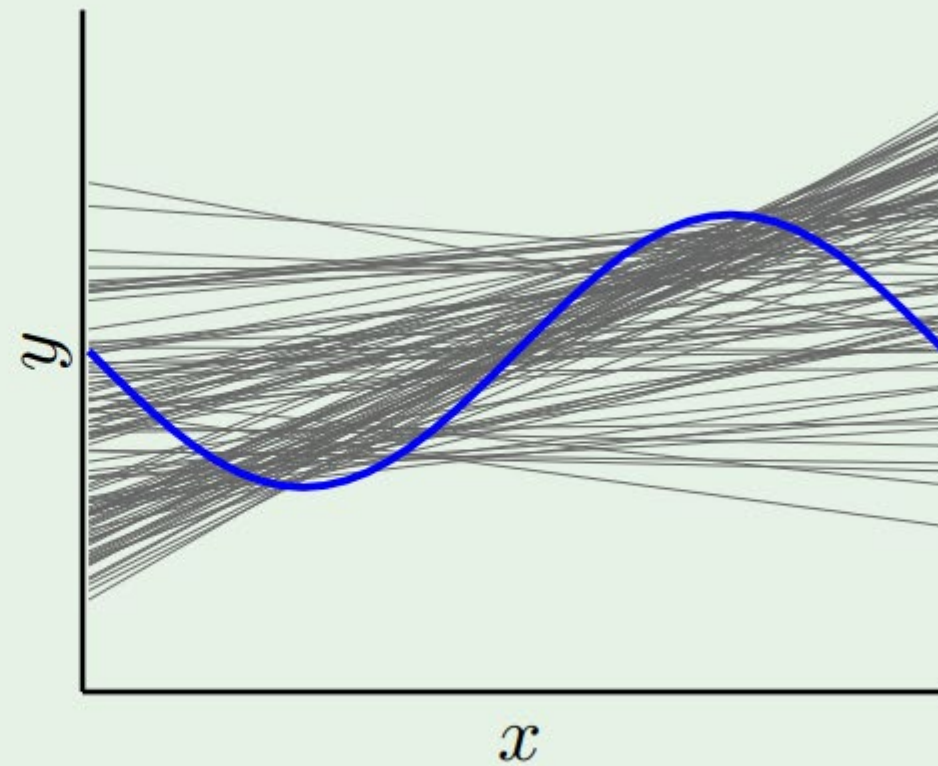
Parameters

# Regularization

(smart way to cure overfitting disease )



without regularization



with regularization

Put a brake on fitting

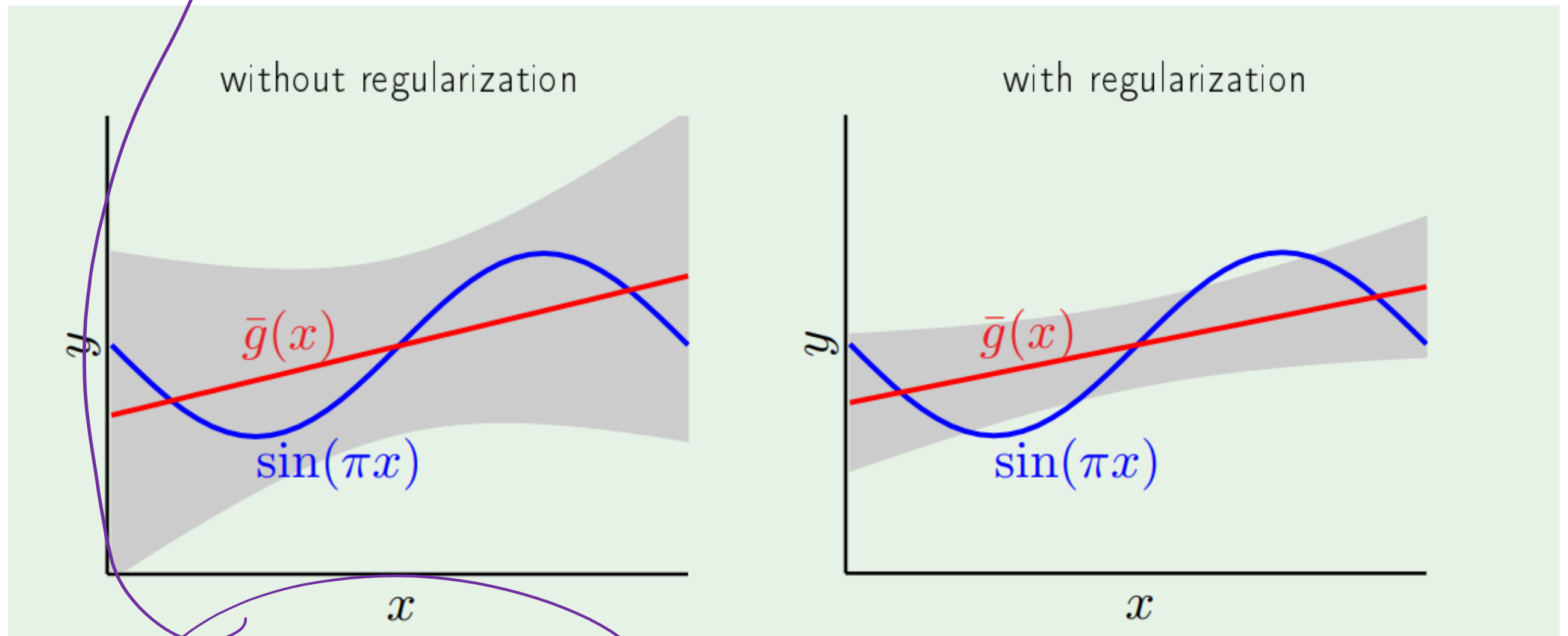


Fit a linear line on sinusoidal with just two data points

# Who is the winner?

$$E(\Theta) = L(\Theta) = 0.21^2 + 1.69$$

$\bar{g}(x)$ : average over all lines



bias=0.21; var=1.69

bias=0.23; var=0.33

# Polynomial Model

Want to fit a polynomial regression model

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \cdots + \theta_d x^d + \epsilon$$

Let's rewrite it as:

$$y = \theta_0 + \theta_1 z_1 + \theta_2 z_2 + \cdots + \theta_d z_d + \epsilon = \mathbf{z}\boldsymbol{\theta}$$

# Regularizing is just constraining the weights ( $\theta$ )

For example: let's do a **hard** constraining

$$y = \theta_0 + \theta_1 z_1 + \theta_2 z_2 + \cdots + \theta_d z_d$$

subject to

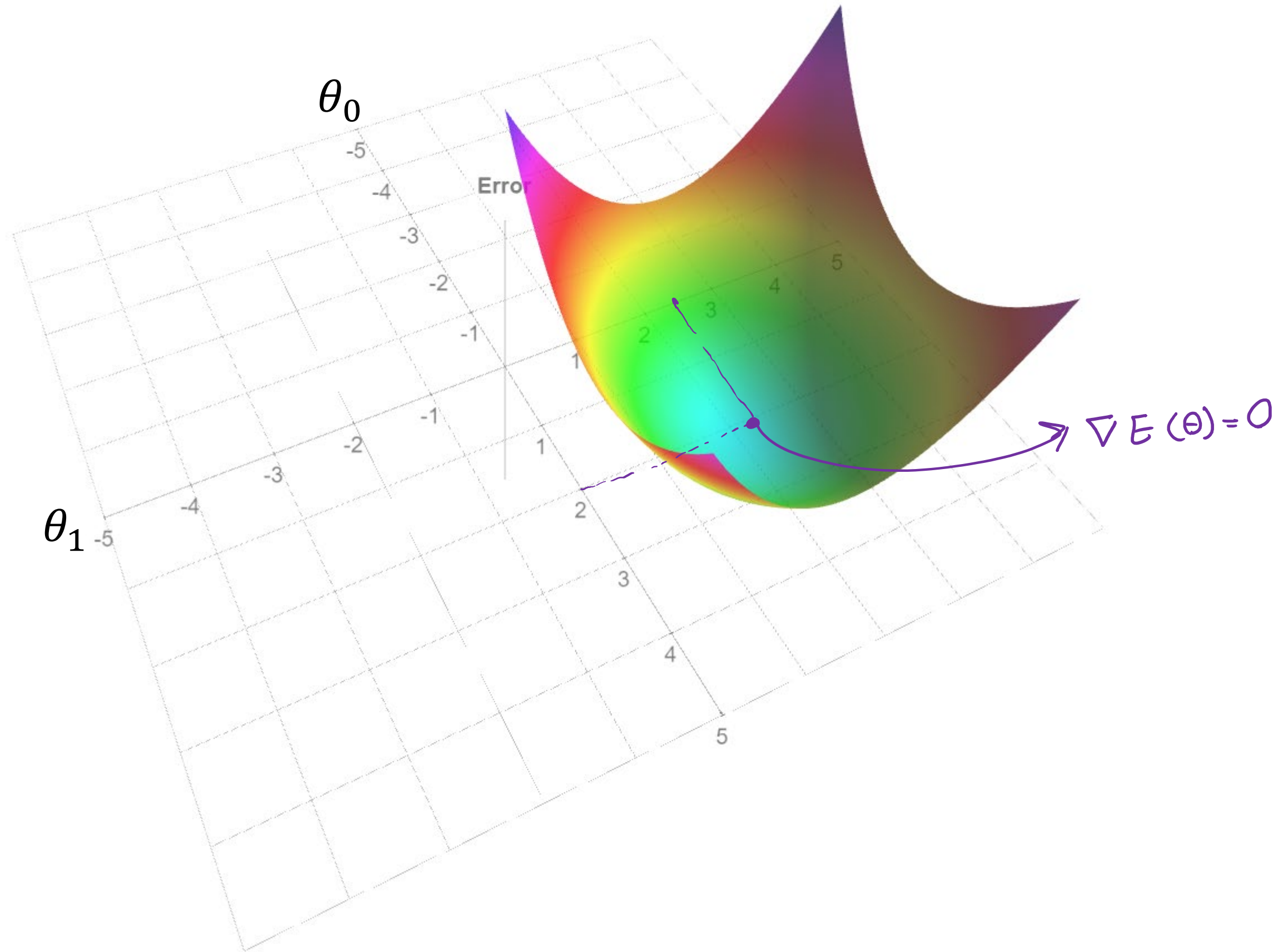
$$\theta_d = 0 \text{ for } d > 2$$



$$y = \theta_0 + \theta_1 z_1 + \theta_2 z_2 + 0 + \cdots + 0$$

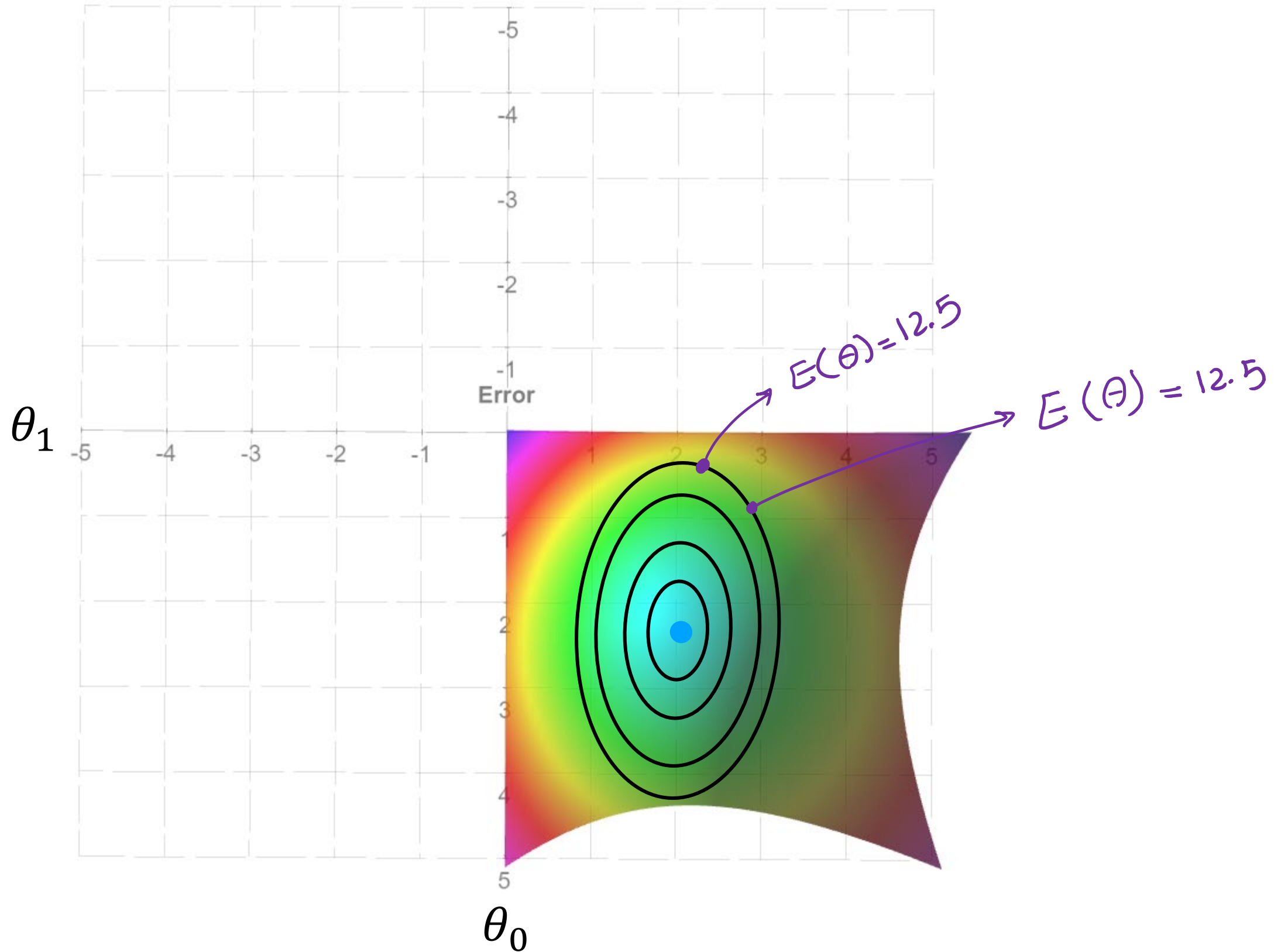
$$E(\theta) = \frac{1}{N} \sum_{i=1}^n (y^i - z_i \theta)^2$$

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



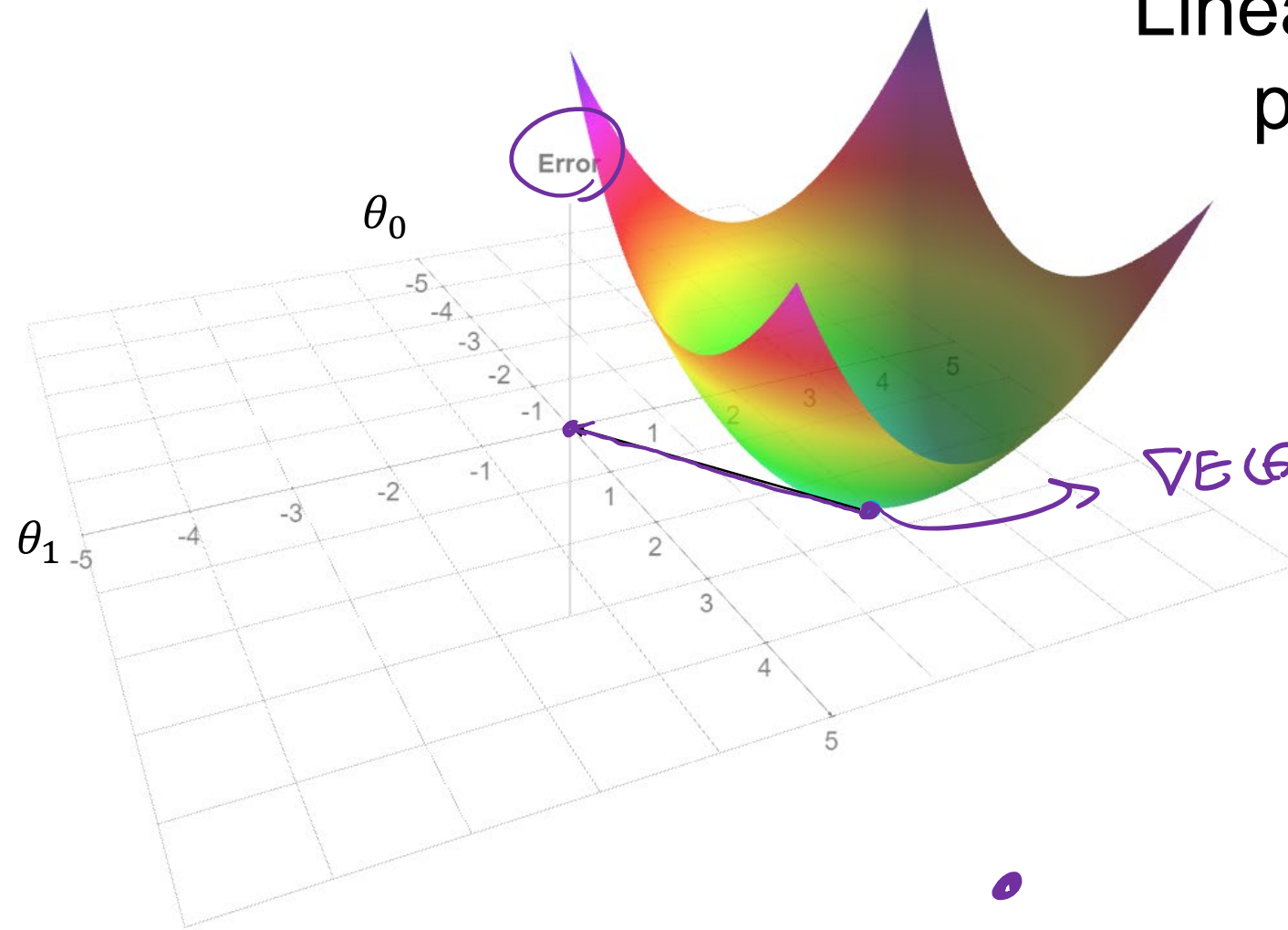


# Project the same graph on x-y using contour plot



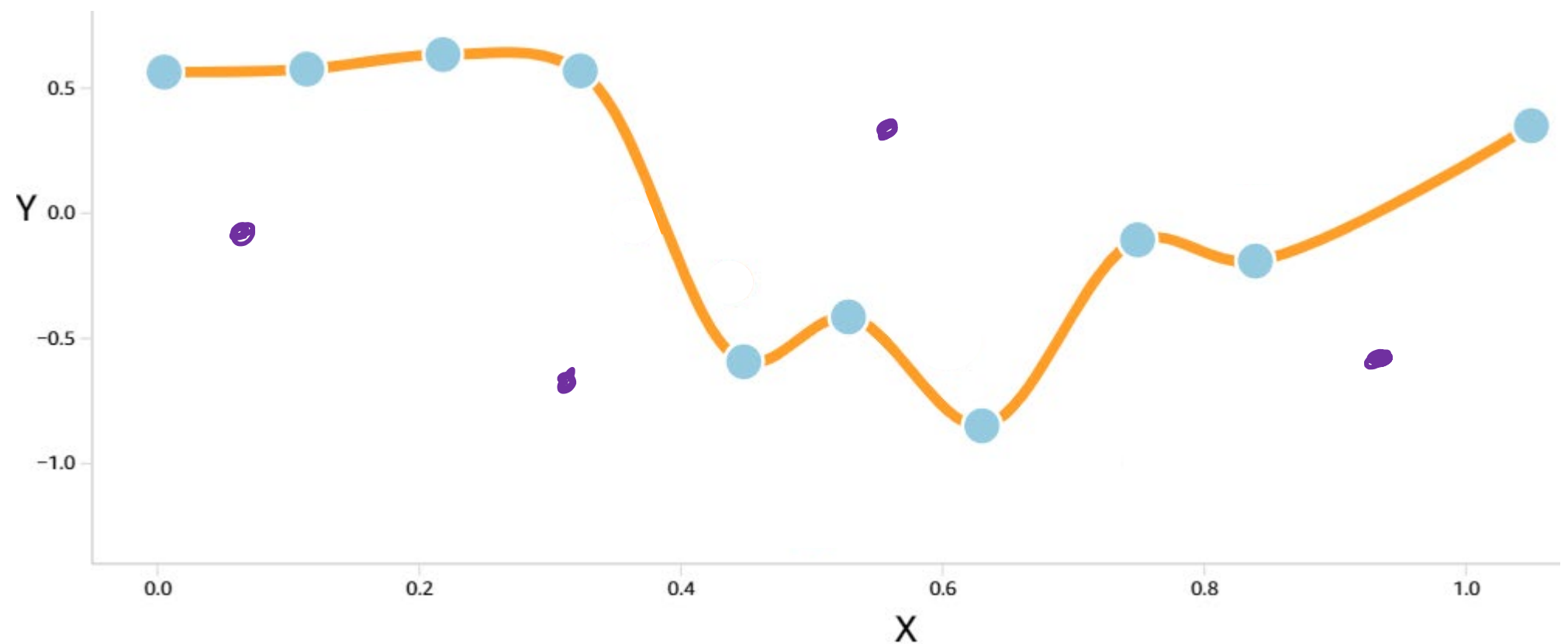


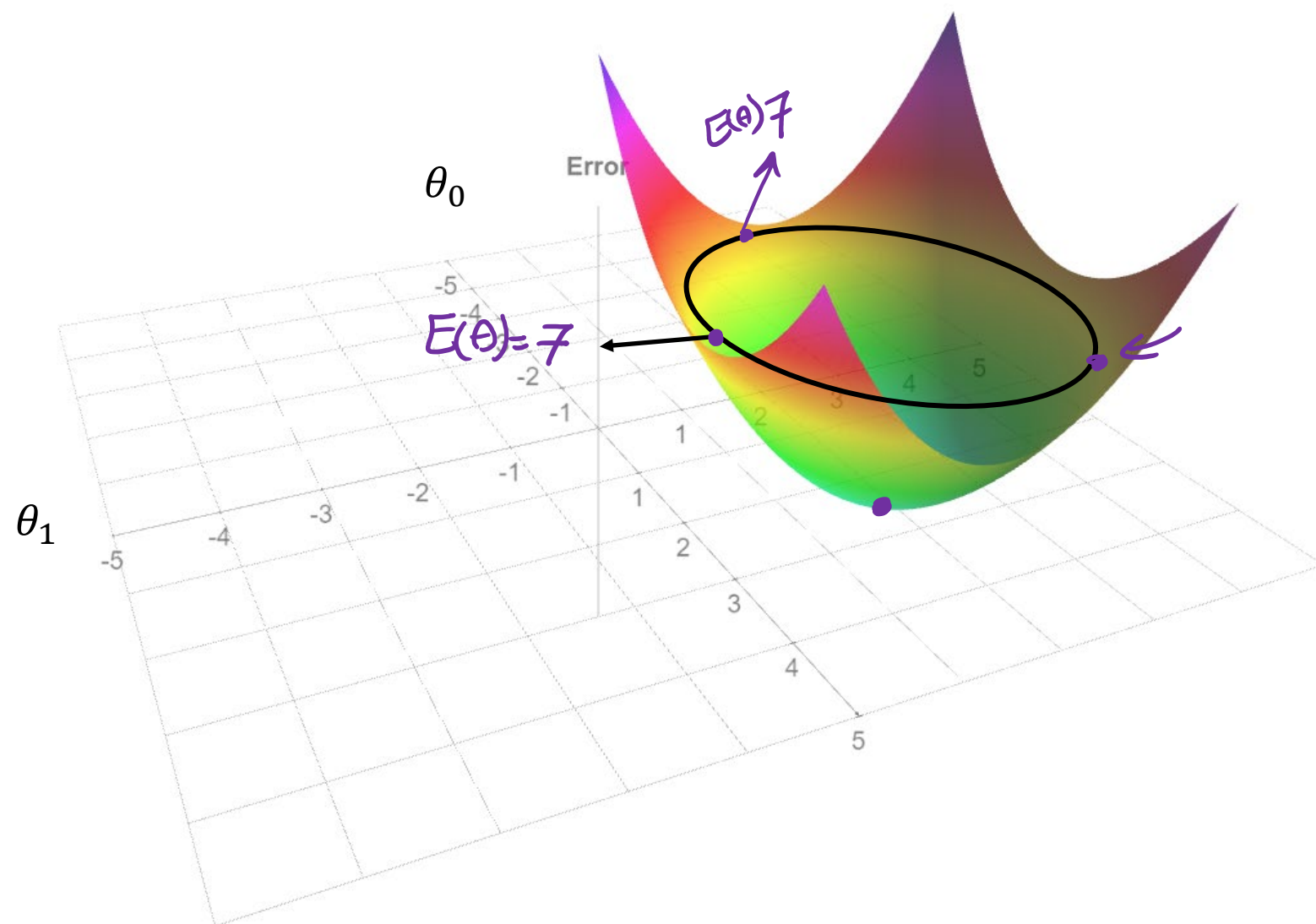
# Linear regression with a very high polynomial degree solution



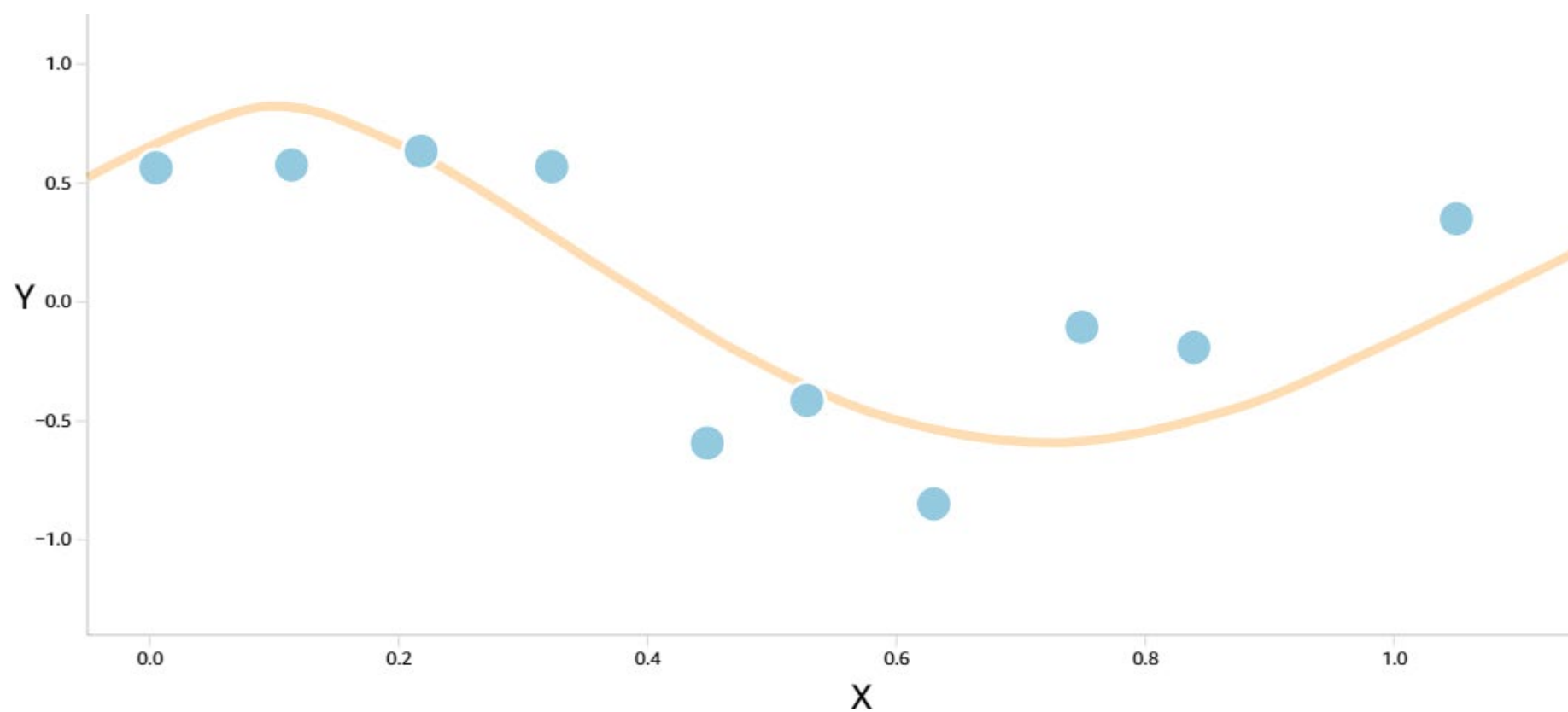
$$E(\theta) = 0$$

$$E(\theta) = \text{bias}^2 + \text{Variance}$$





$$E(\theta) = \text{bias}^2 + \text{variance}$$



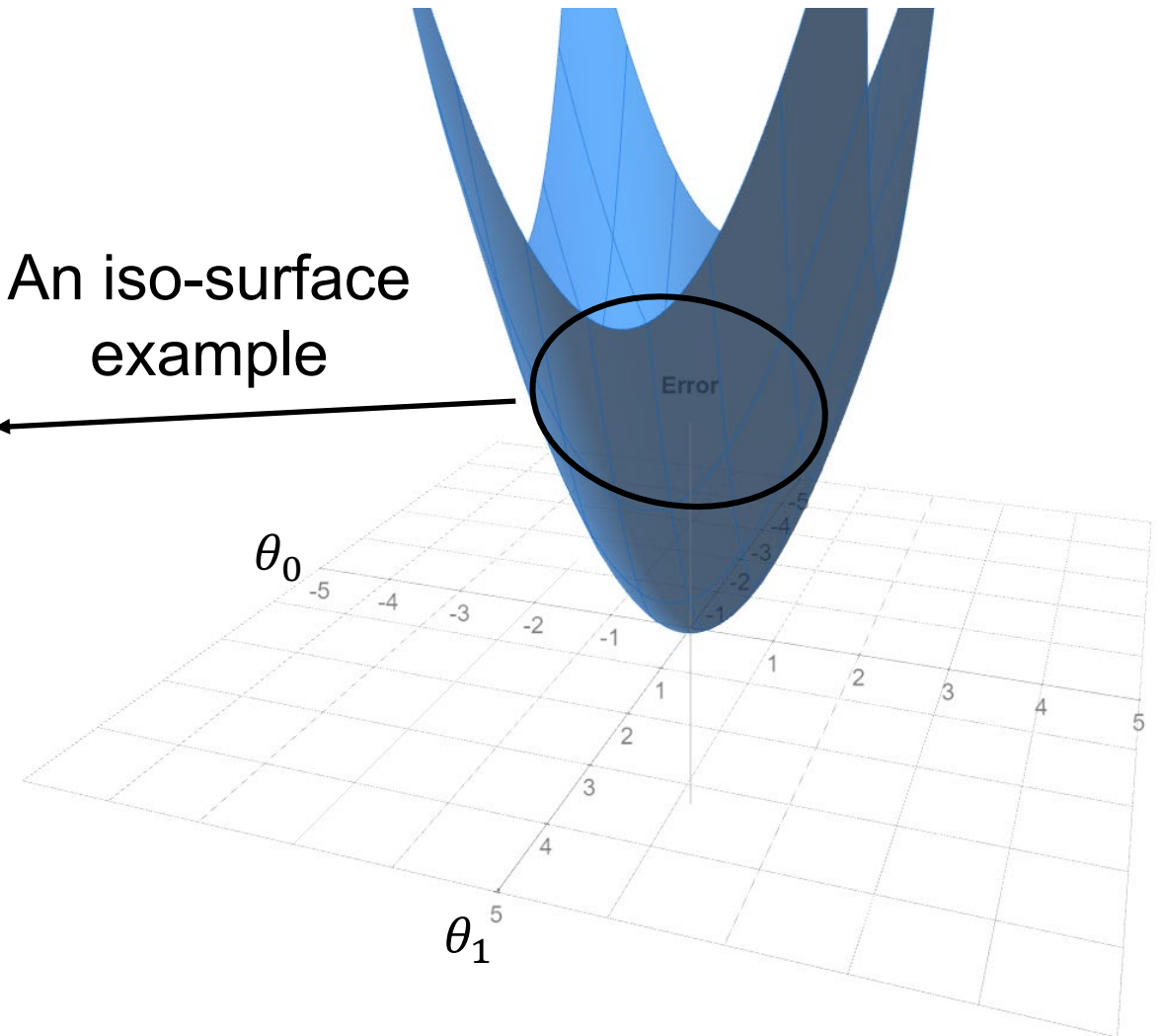
How can we get an optimal solution with a positive error for a model that overfits?

We need to introduce a constraint

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \quad \theta^T \theta = \theta_0^2 + \theta_1^2$$

$$g(\theta) = \theta_0^2 + \theta_1^2 = \theta^T \theta = C$$

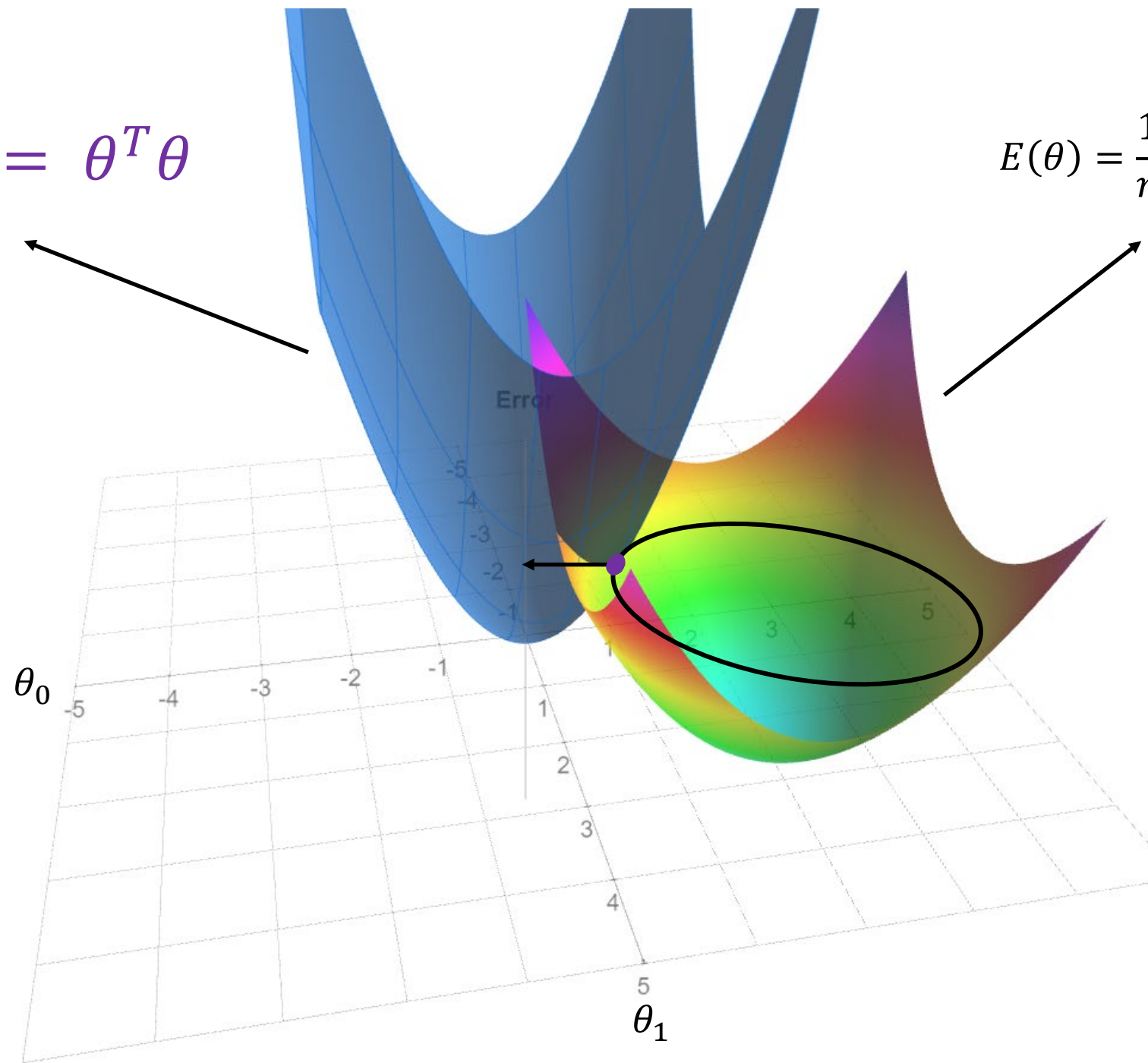
An iso-surface  
example



Error function together with a  
new introduced constraint

$$g(\theta) = \theta_0^2 + \theta_1^2 = \theta^T \theta$$

$$E(\theta) = \frac{1}{n} \sum_{i=1}^n (y^i - z_i \theta)^2$$



Let's define the Lagrange function

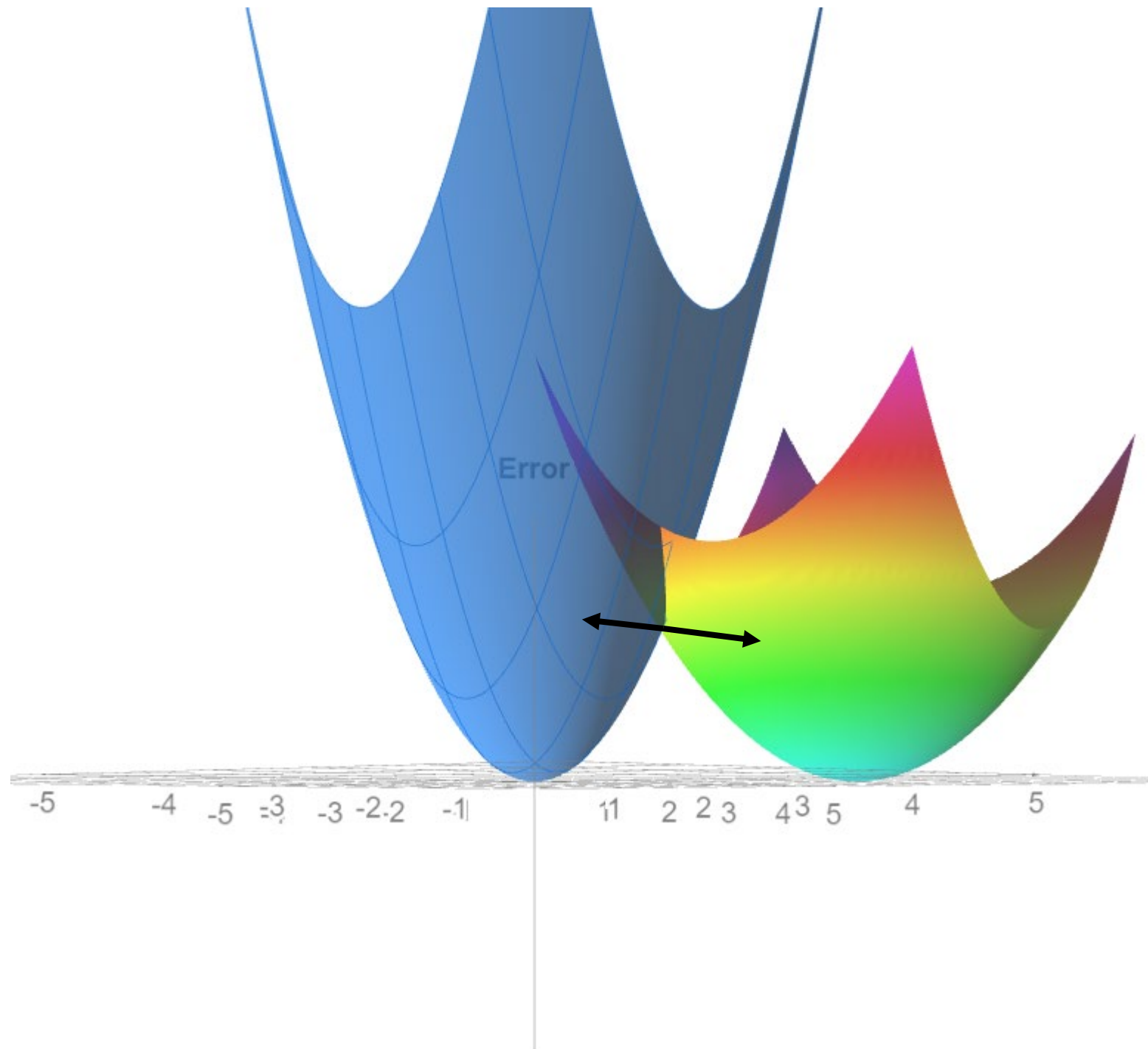
$$L(\theta, \lambda) = E(\theta) + \lambda g(\theta)$$

$$L(\theta, \lambda) = E(\theta) + \lambda \theta^T \theta$$

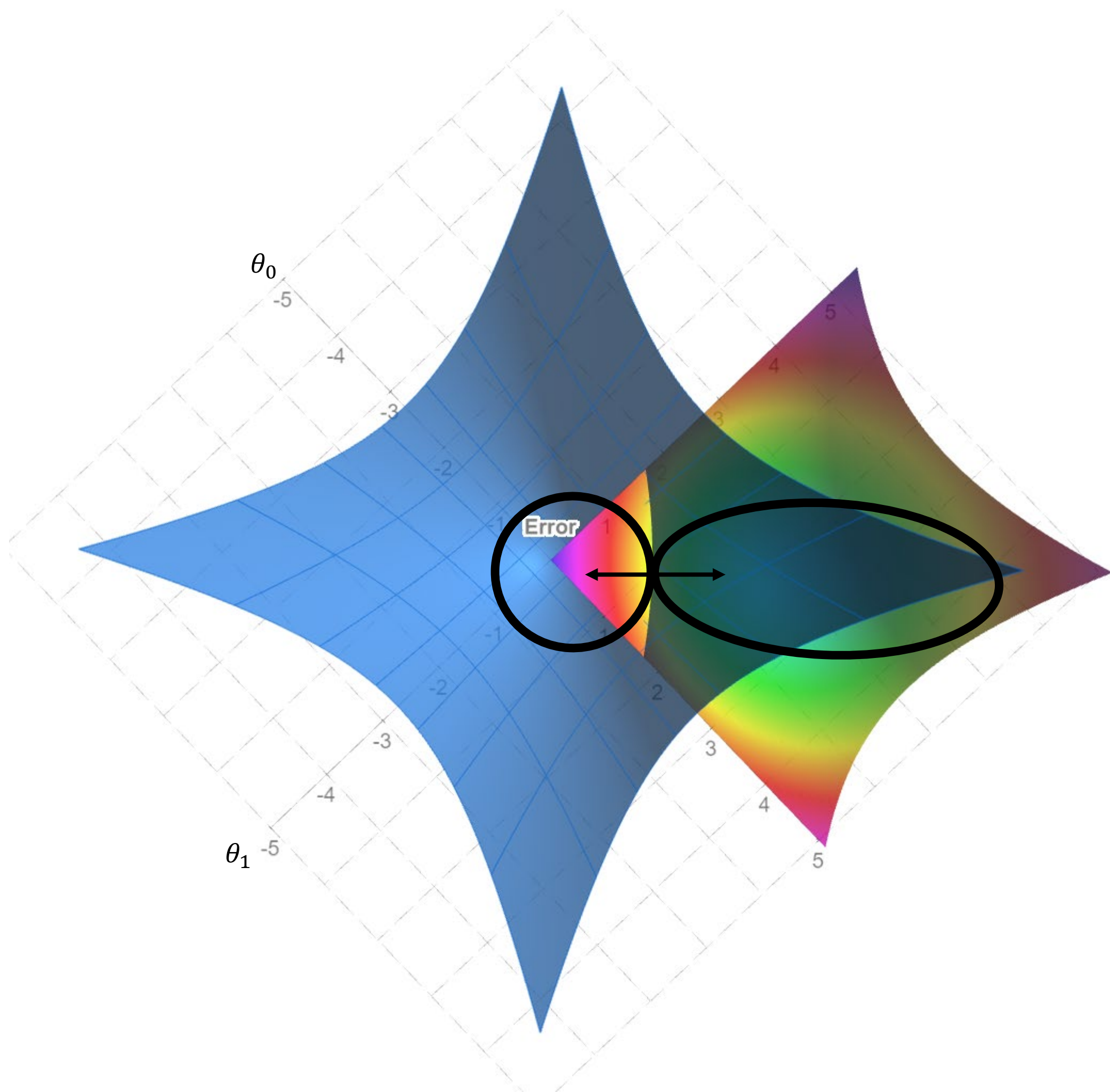
$$\nabla L(\theta, \lambda) = 0 \qquad \nabla [E(\theta) + \lambda \theta^T \theta] = 0$$

$$\nabla [E(\theta)] + \lambda \nabla [\theta^T \theta] = 0$$

How to enforce the gradient of Lagrange function to be zero







# Let's calculate the gradients

Gradient of constraint  $g(\theta)$

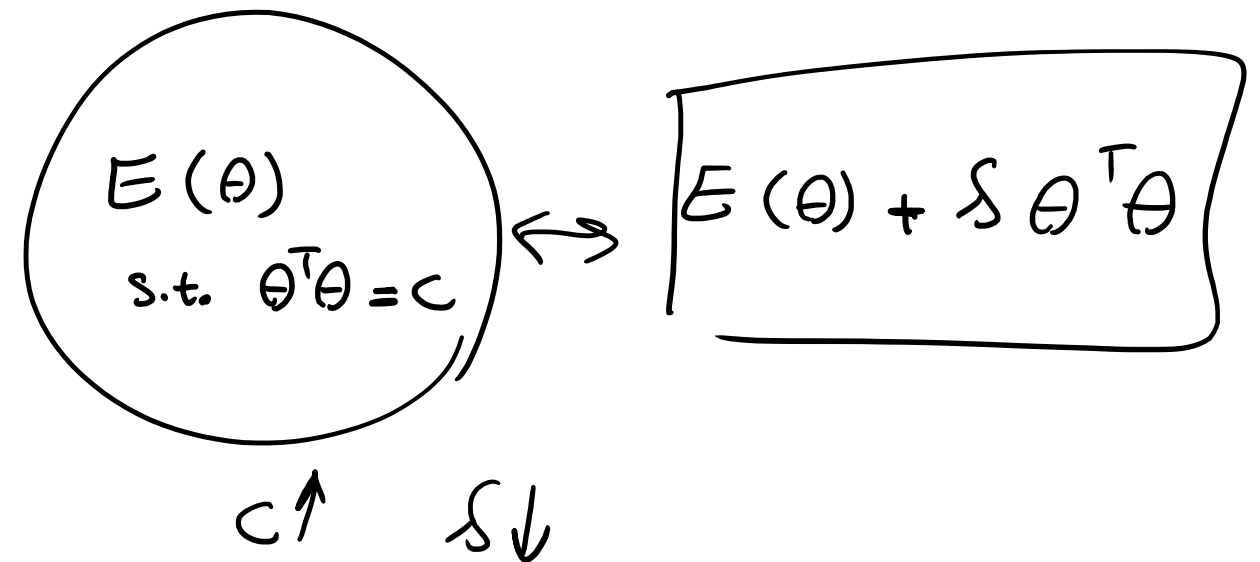
$$\nabla[\theta^T \theta] = 2\theta$$

$$\nabla E(\theta) \sim -\nabla[\theta^T \theta]$$

$$\nabla[E(\theta)] + \lambda \nabla[\theta^T \theta] = 0$$

$$\nabla[E(\theta)] = -\lambda \nabla[\theta^T \theta]$$

$$\nabla E(\theta) = -2\lambda\theta$$



$$\nabla E(\theta) + 2\lambda\theta = 0$$

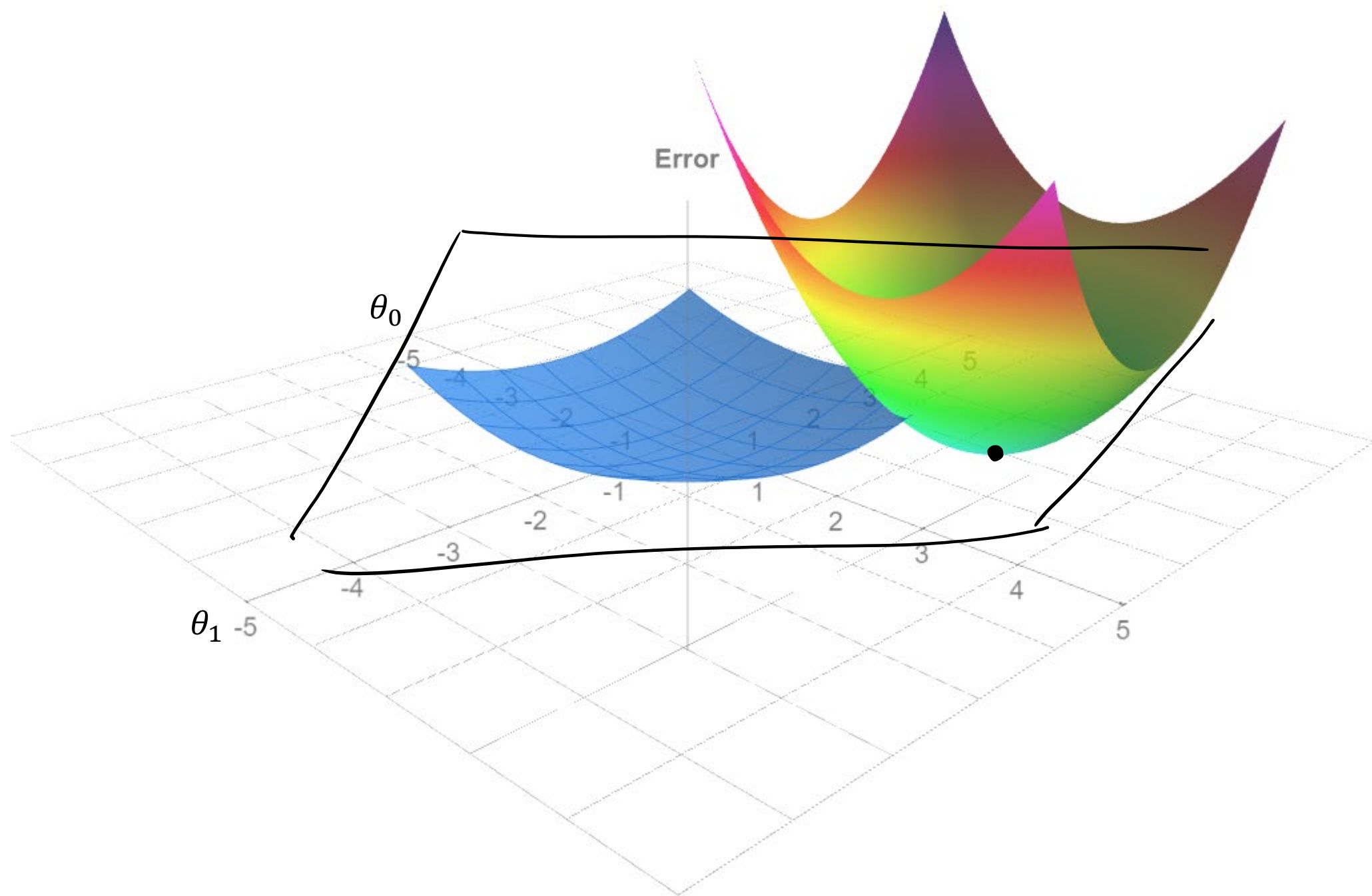
Let's do integration

A diagram showing the minimization of the Lagrangian function. A rounded rectangle contains the expression  $E(\theta) + \lambda \theta^T \theta$ . The word "minimize" is written above the rectangle. The term  $\lambda \theta^T \theta$  is circled, and an arrow points from the word "Regularization term" to this circle.



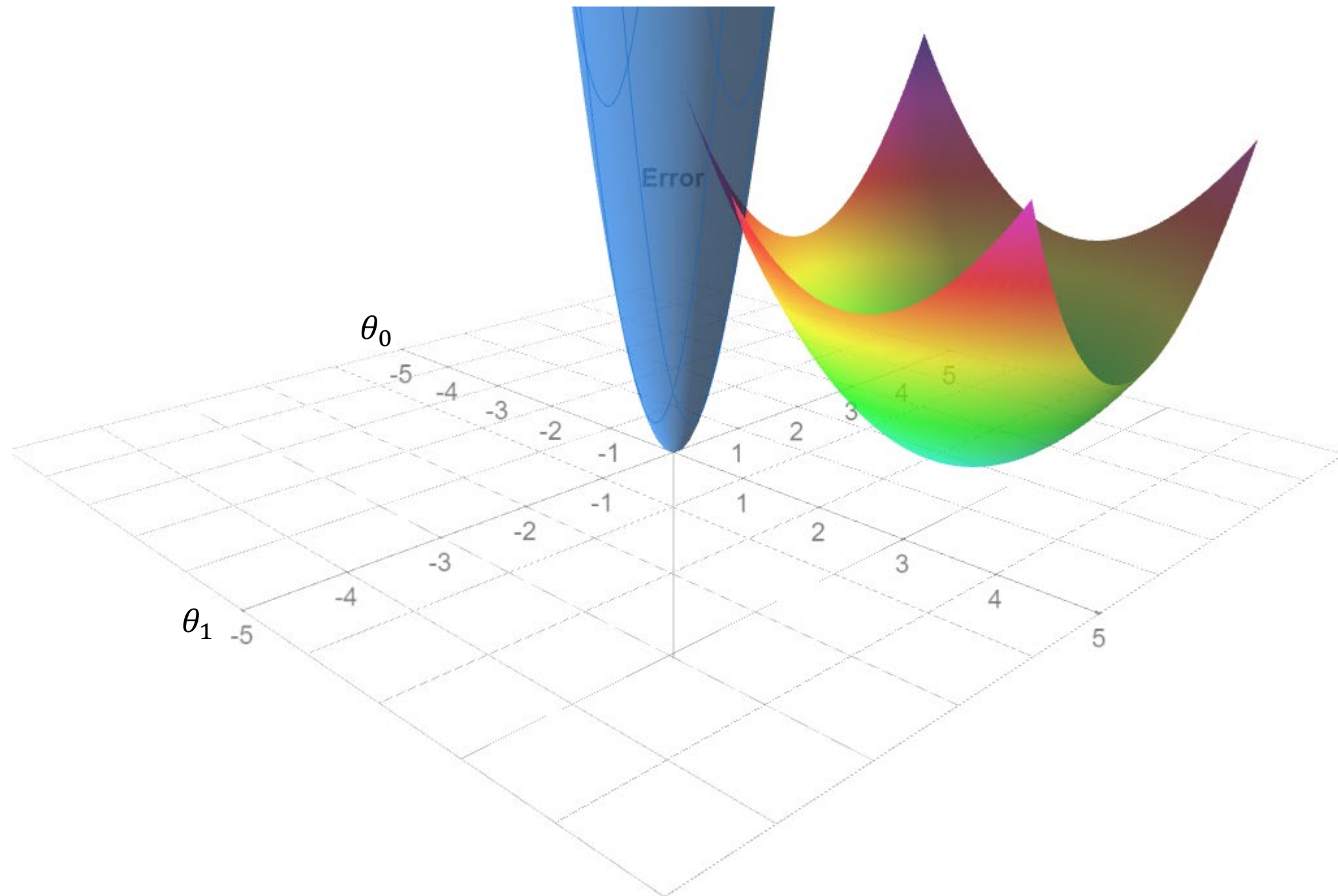
# The effect of low Lambda

$$E(\theta) + \frac{\lambda}{N} \theta^T \theta$$



# The effect of high Lambda


$$E(\theta) + \frac{\lambda}{N} \theta^T \theta$$



# Regularized Learning

Minimize  $E(\theta) + \lambda \theta^T \theta$

Now we know Why this term  
leads to the regularization of  
parameters




Regularized Error

$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^n (y^i - z_i \theta)^2 + \frac{\lambda}{2N} \|\theta\|_2^2$$

L2 Regularization term

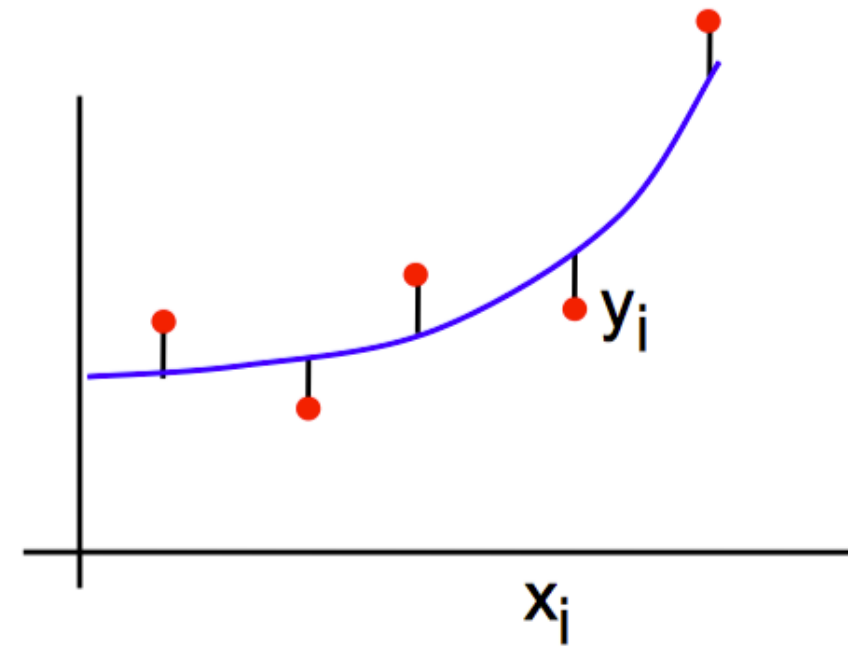


# Outline

- Overfitting and regularized learning
- Ridge regression 
- Lasso regression
- Determining regularization strength

# Ridge Regression

$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^n (y^i - z_i \theta)^2 + \lambda \|\theta\|_2^2$$



$$\theta_0 + \theta_1 z_1 + \theta_2 z_2 + \cdots + \theta_d z_d + \epsilon = \mathbf{z}\boldsymbol{\theta}$$

General form

$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^n (y^i - z_i \theta)^2 + \lambda \|\theta\|_2^2$$

Matrix form

$$\tilde{E}(\theta) = \frac{1}{N} (y - z\theta)^T (y - z\theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

$$\frac{\partial \tilde{E}(\theta)}{\partial \theta} = -z^T (y - z\theta) + \lambda \theta$$

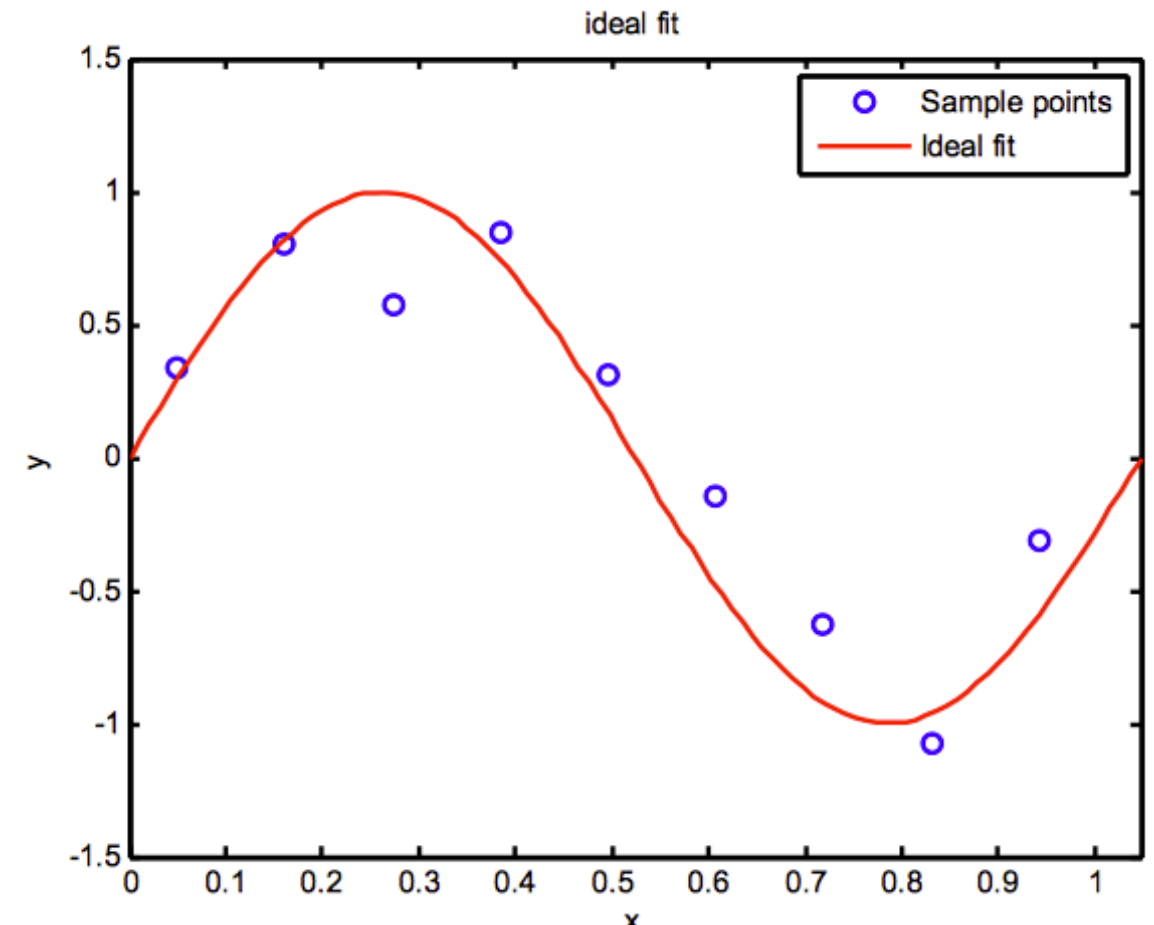
$$(z^T z + \lambda I) \theta = z^T y$$

$$\theta_{\text{reg}} = (z^T z + \lambda I)^{-1} z^T y$$

$$\theta_{\text{overfitted}} = (z^T z)^{-1} z^T y$$

# Ridge Regression Example

- The red curve is the true function (which is not a polynomial)
- The data points are samples from the curve with added noise in y.
- There is a choice in both the degree,  $D$ , of the basis functions used, and in the strength of the regularization



$$f(x, \theta) = z\theta$$

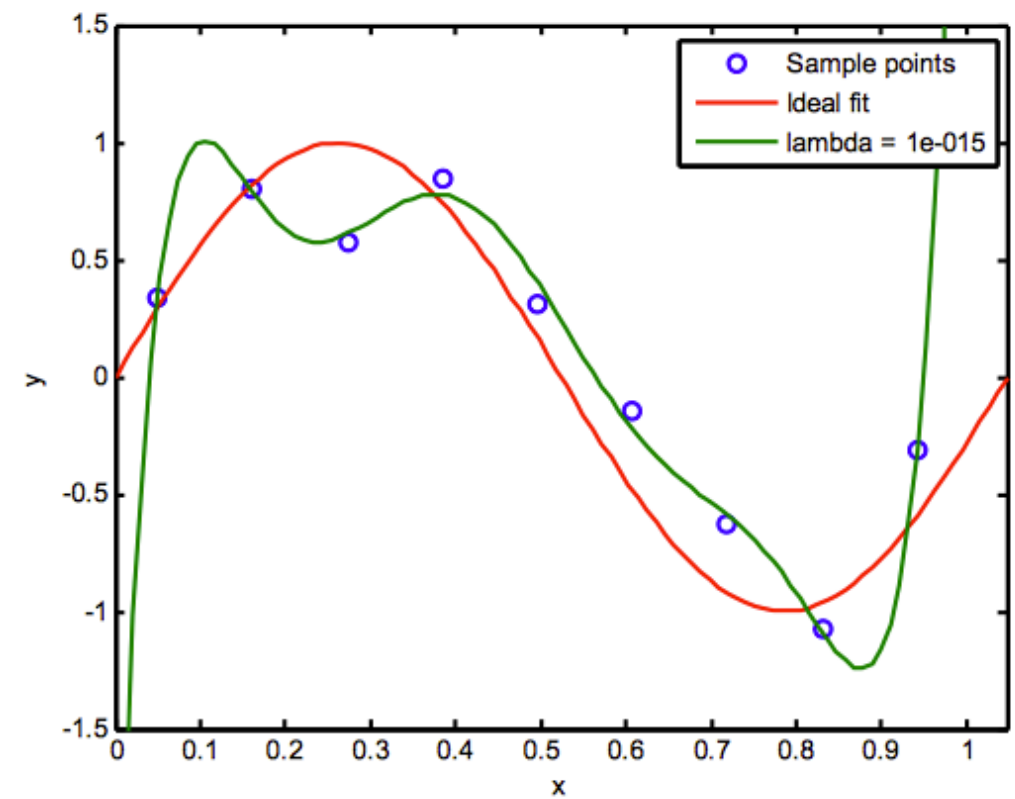
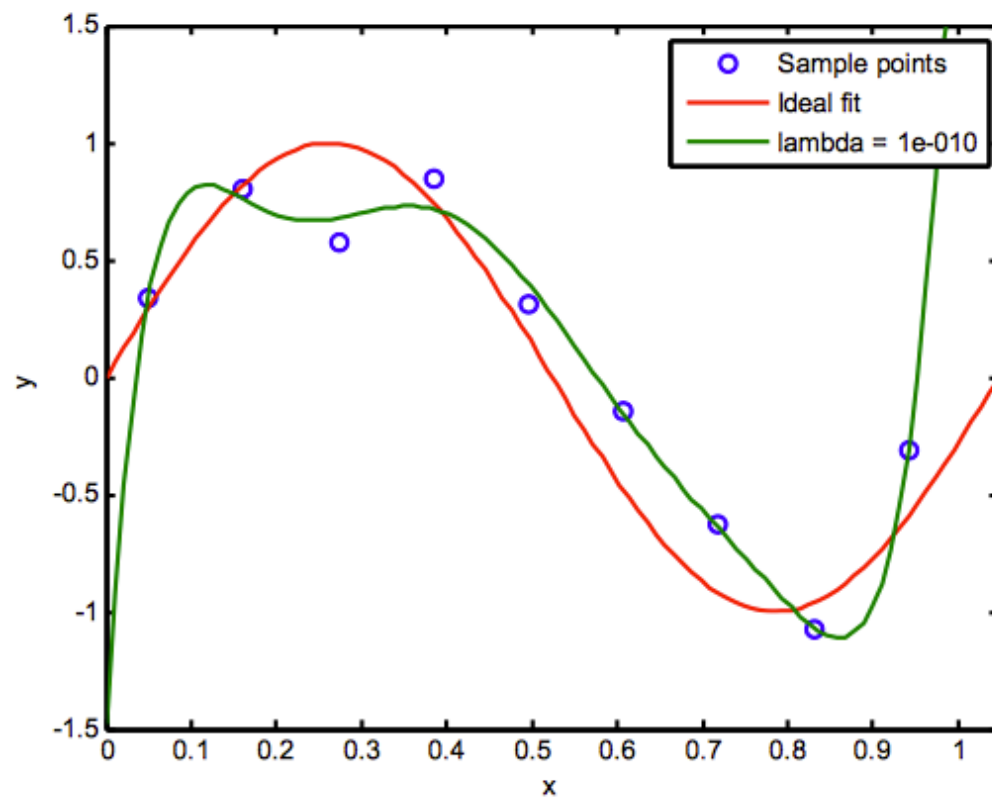
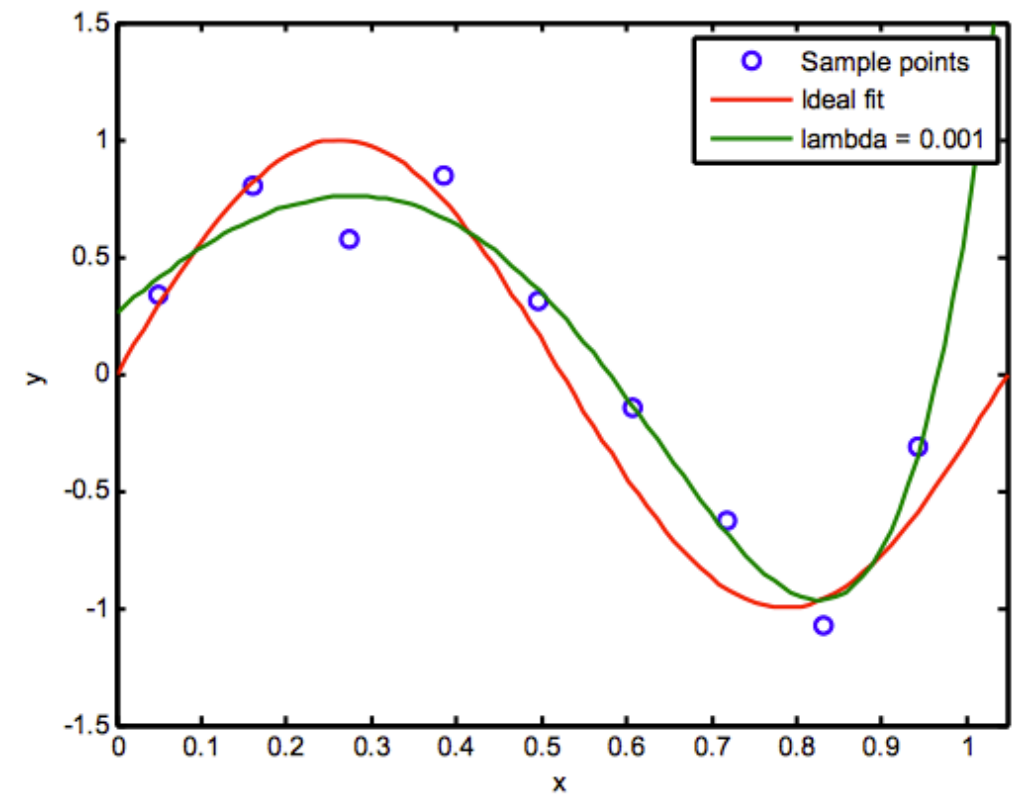
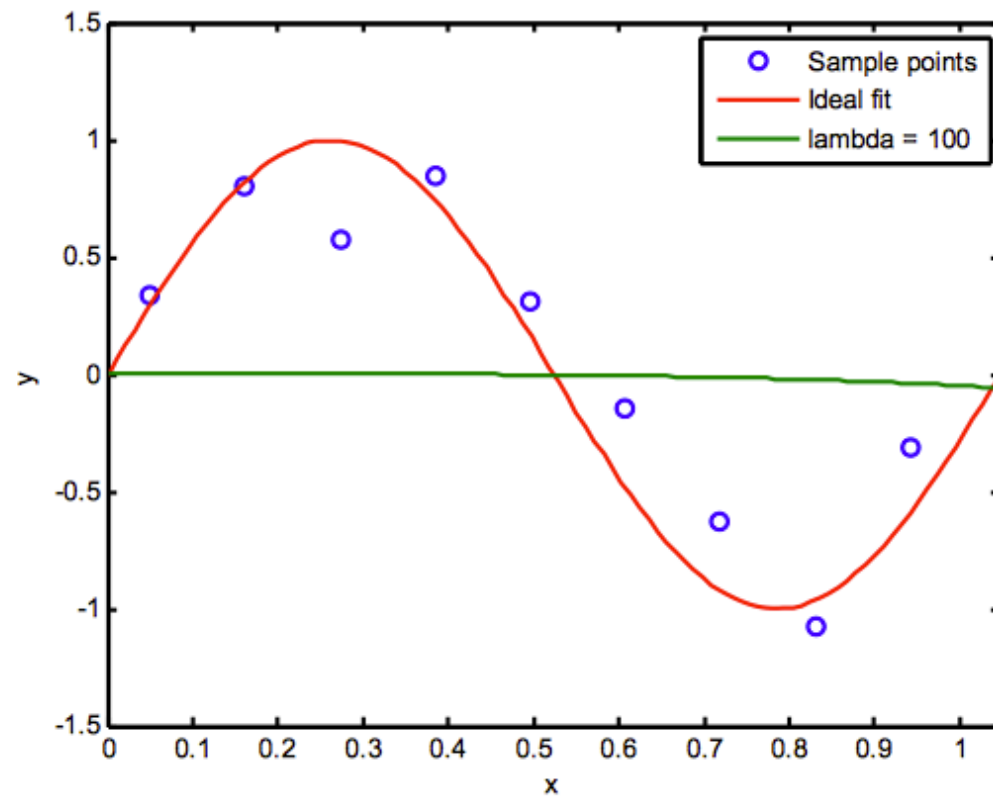
$$z: x \rightarrow z$$

$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^n (y^i - z_i \theta)^2 + \lambda \|\theta\|_2^2$$

$\theta \in \mathbb{R}^{D+1}$

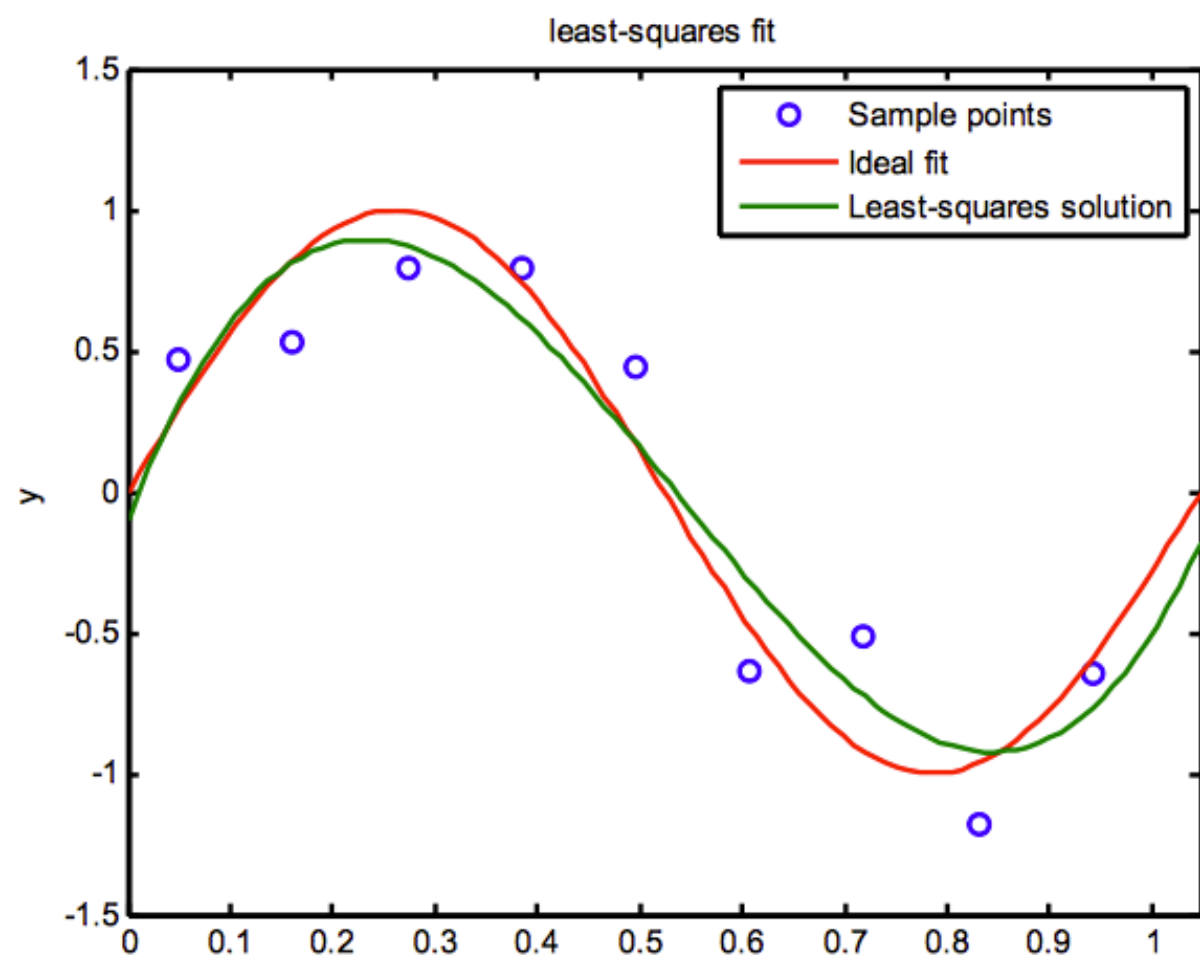
$\frac{\lambda}{2N} \|\theta\|_2^2$

$N = 9$  samples,  $D = 7$

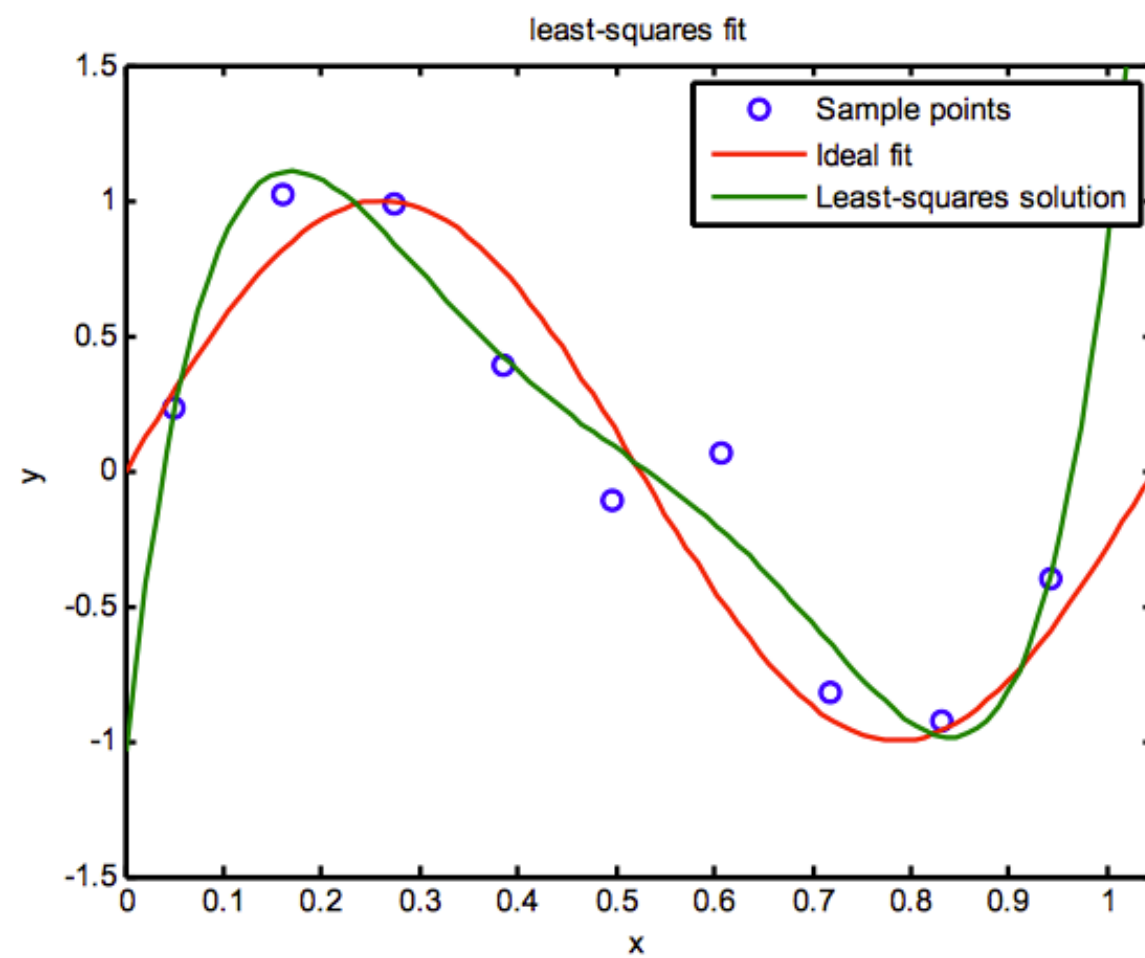





$D = 3$



$D = 5$



# Outline

- Overfitting and regularized learning
- Ridge regression
- Lasso regression 
- Determining regularization strength

# Regularized Regression

$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^n (y^i - z_i \theta)^2 + \lambda \|\theta\|_2^2$$

Squared loss\Error

$$\frac{1}{N} \sum_{i=1}^n (y^i - z_i \theta)^2$$

L2 Regularizer

$$\lambda \|\theta\|_2^2$$

Now let's look at another regularization choice.

# The Lasso Regularization (L1 norm) and sparsity

Lasso = **L**east Absolute **S**hrinkage and **S**election **O**perator

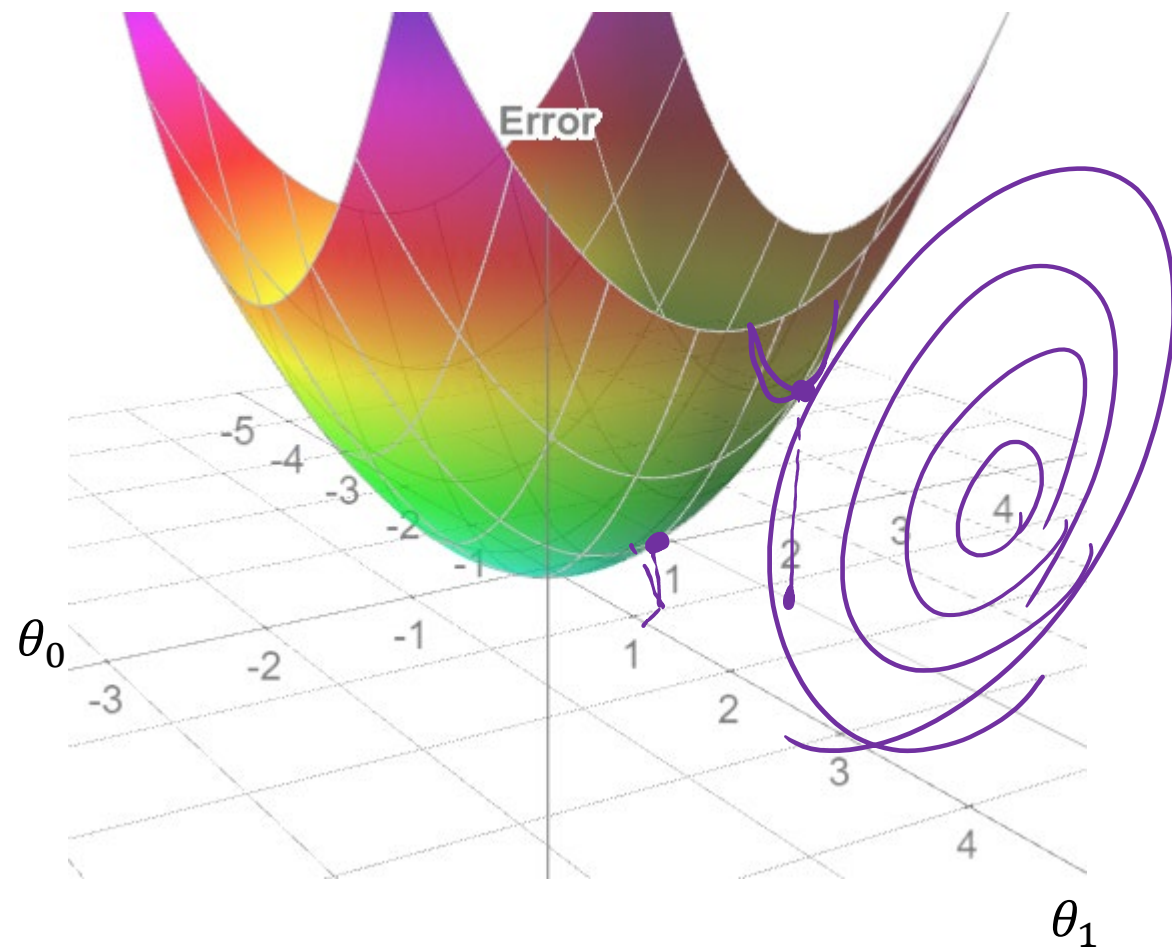
$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^n (y^i - z_i \theta)^2 + \lambda \|\theta\|_1$$

$\|\theta\|_2^2$    
 $\|\theta\|_1$  

L1 norm induces sparsity. This means that some of the weights become zero, and the feature contribution will be completely removed. L1 Regularizer could be used for feature selection

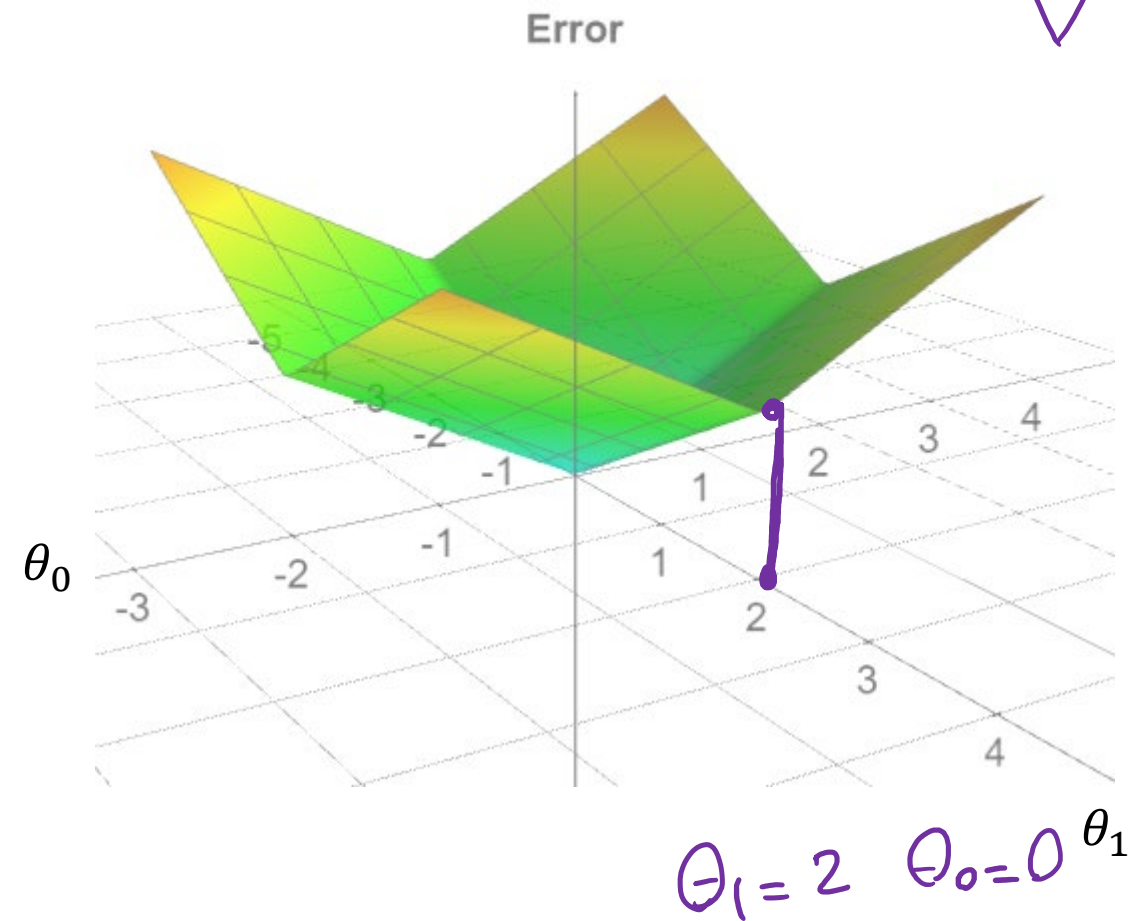
# Ridge Regularizer

$$g(\theta) = \theta_0^2 + \theta_1^2 = \theta^T \theta$$



# Lasso Regularizer

$$g(\theta) = \theta_0 + \theta_1 = \theta$$

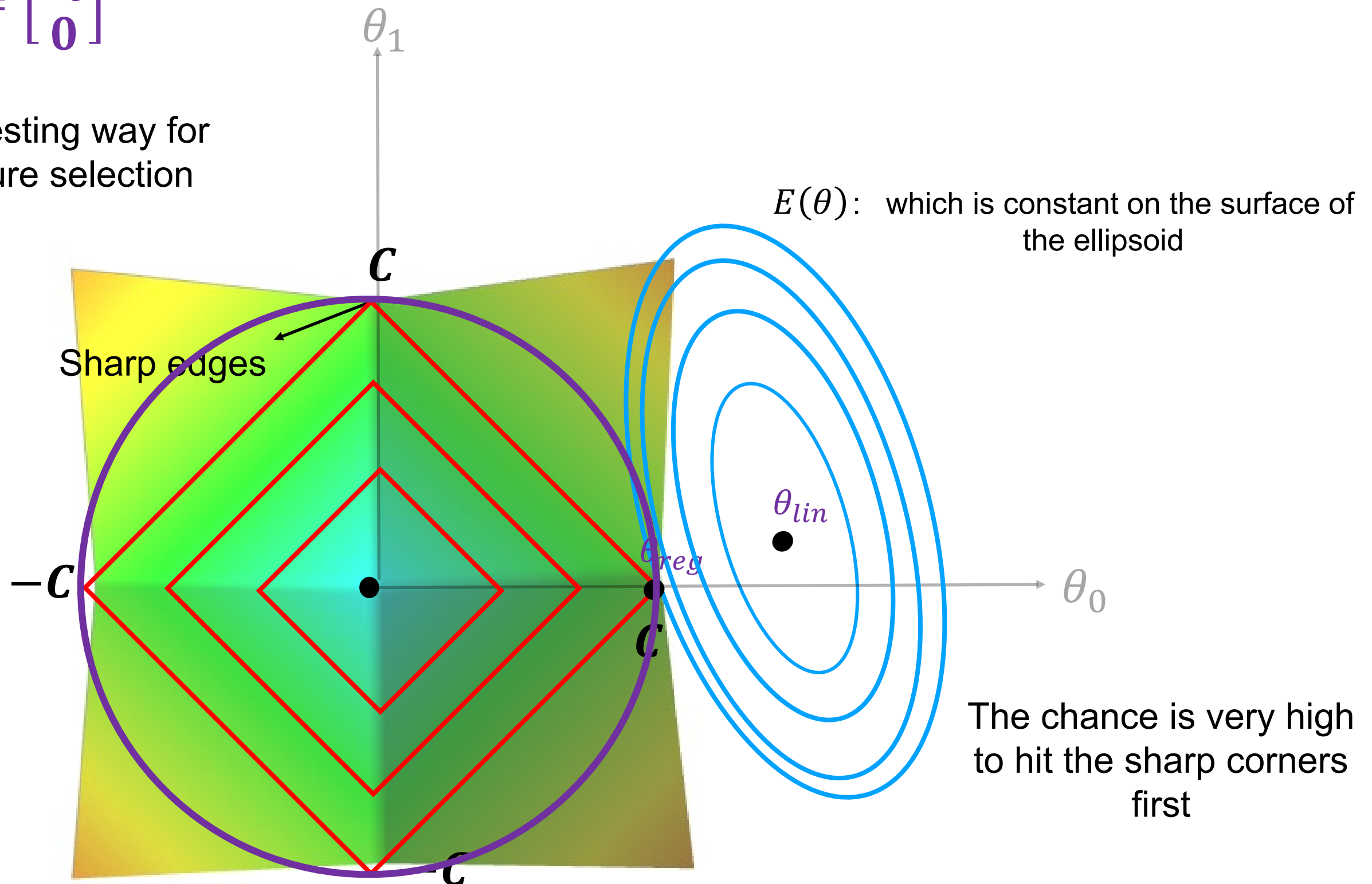


Let's say we have two parameters ( $\theta_0$  and  $\theta_1$ )

$$\text{Min } E(\theta) = \frac{1}{N} (\mathbf{z}\mathbf{w} - y)^T (\mathbf{z}\theta - y) + \lambda \|\theta\|_1$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

Interesting way for  
feature selection



# Ridge versus Lasso

## Ridge

$$\tilde{E}(\theta) = \frac{1}{N} (y - z\theta)^T (y - z\theta) + \lambda \|\theta\|_2^2$$

It is a convex model

Both mean squared error and L2 regularizer are differentiable.

We can get a closed form solution

## Lasso

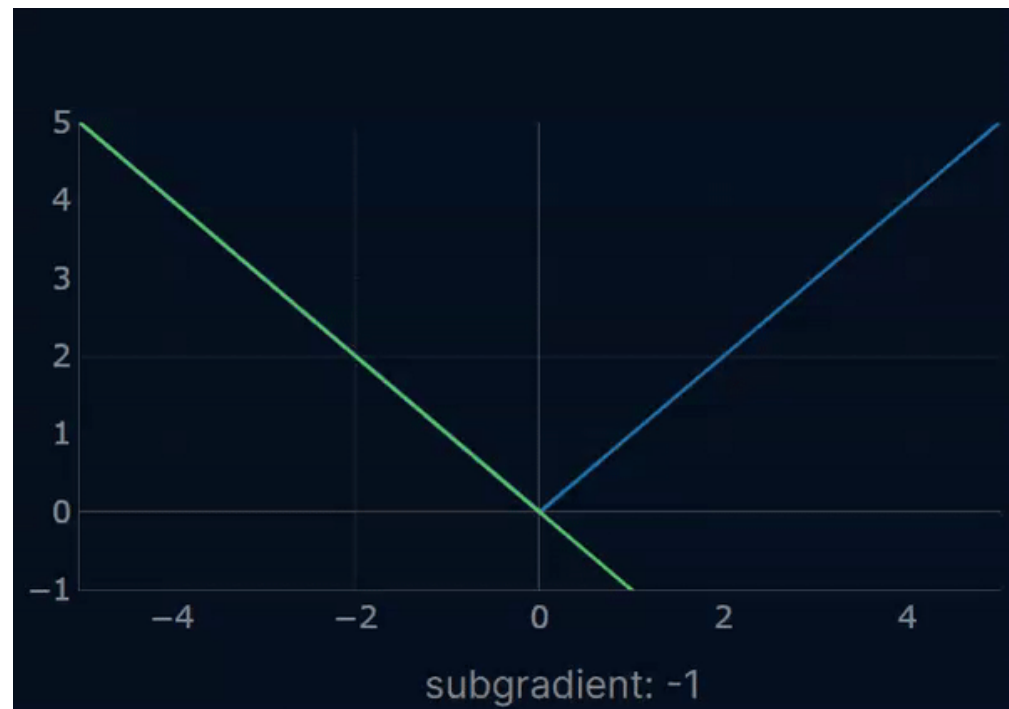
$$\tilde{E}(\theta) = \frac{1}{N} (y - z\theta)^T (y - z\theta) + \lambda \|\theta\|_1$$

It is a convex model

L1 regularizer is NOT differentiable.

We can **NOT** get a closed form solution

# Sub-gradient Descend in Lasso

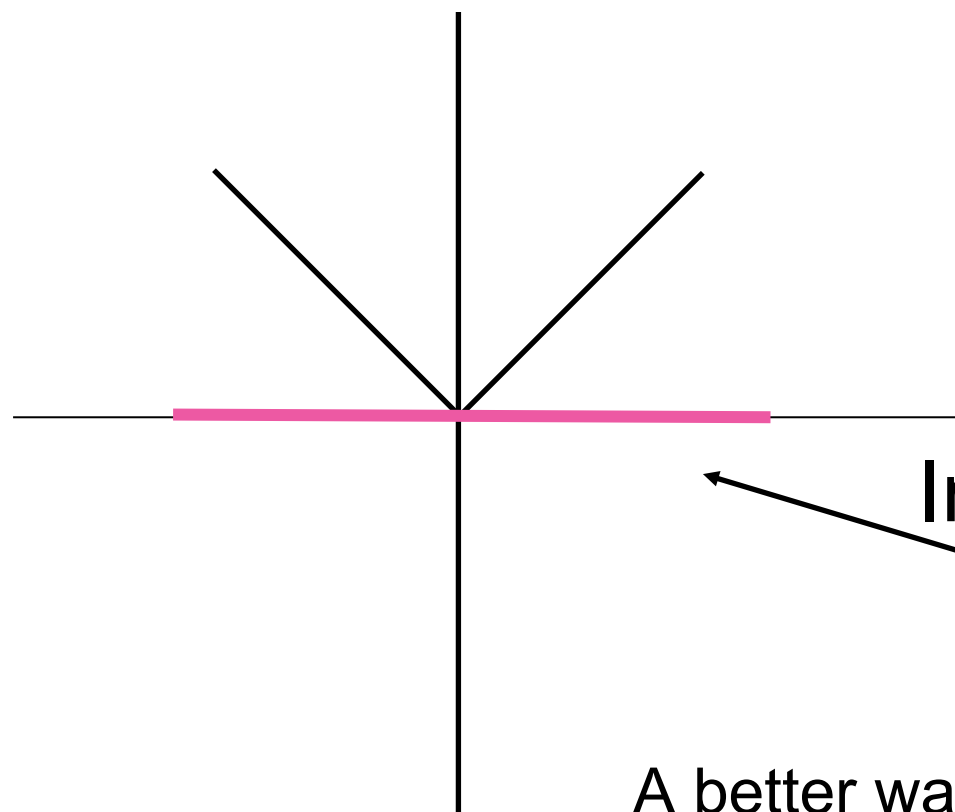


$$\tilde{E}(\theta) = \frac{1}{N} (y - z\theta)^T (y - z\theta) + \lambda \|\theta\|_1$$

$$\frac{\partial \tilde{E}(\theta)}{\partial \theta} = -z^T (y - z\theta) + \frac{\partial (\lambda \|\theta\|_1)}{\partial \theta}$$

Using Sub-gradient

$$\frac{\partial \tilde{E}(\theta)}{\partial \theta} = -z^T (y - z\theta) + \lambda \text{sign}(\theta)$$




In *sign* function, we use this sub-gradient line as our under-estimator (below our function)

A better way: Proximal gradient descent with soft-thresholding



# Outline

- Overfitting and regularized learning
- Ridge regression
- Lasso regression
- Determining regularization strength 

# Leave-One-Out Cross Validation

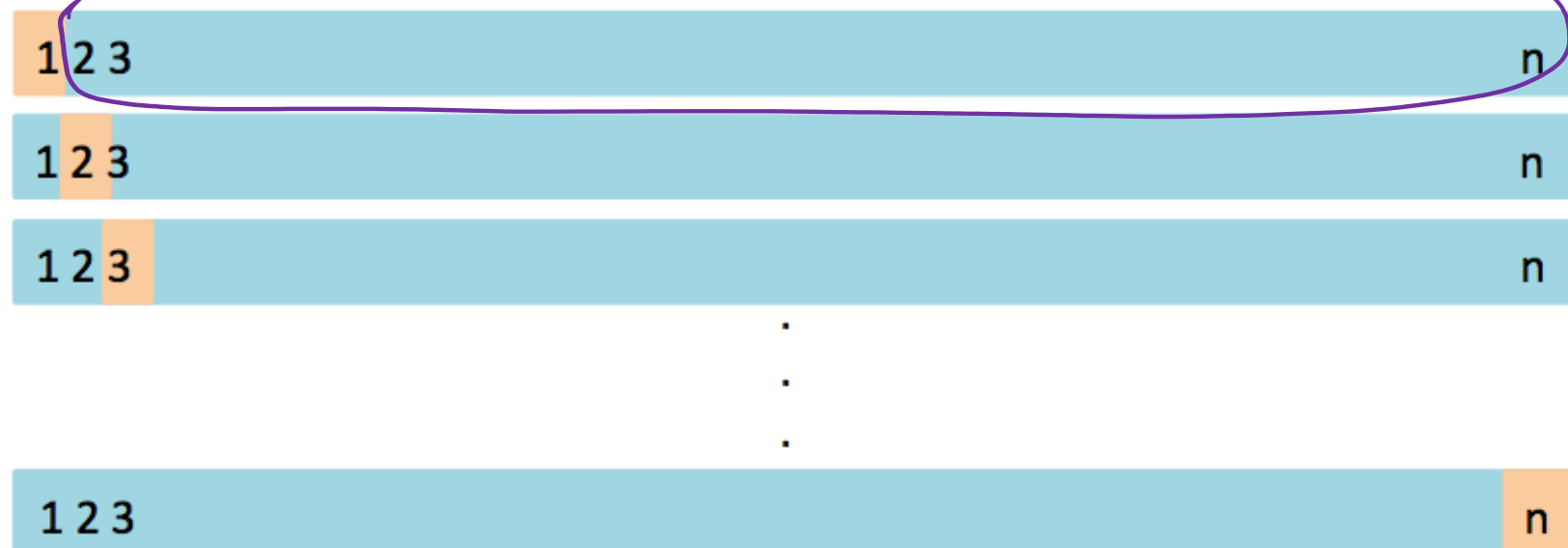
LOOCV

For every  $i = 1, \dots, n$ :

- ▶ train the model on every point except  $i$ ,
- ▶ compute the test error on the held out point.

Average the test errors.

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i^{(-i)})^2$$



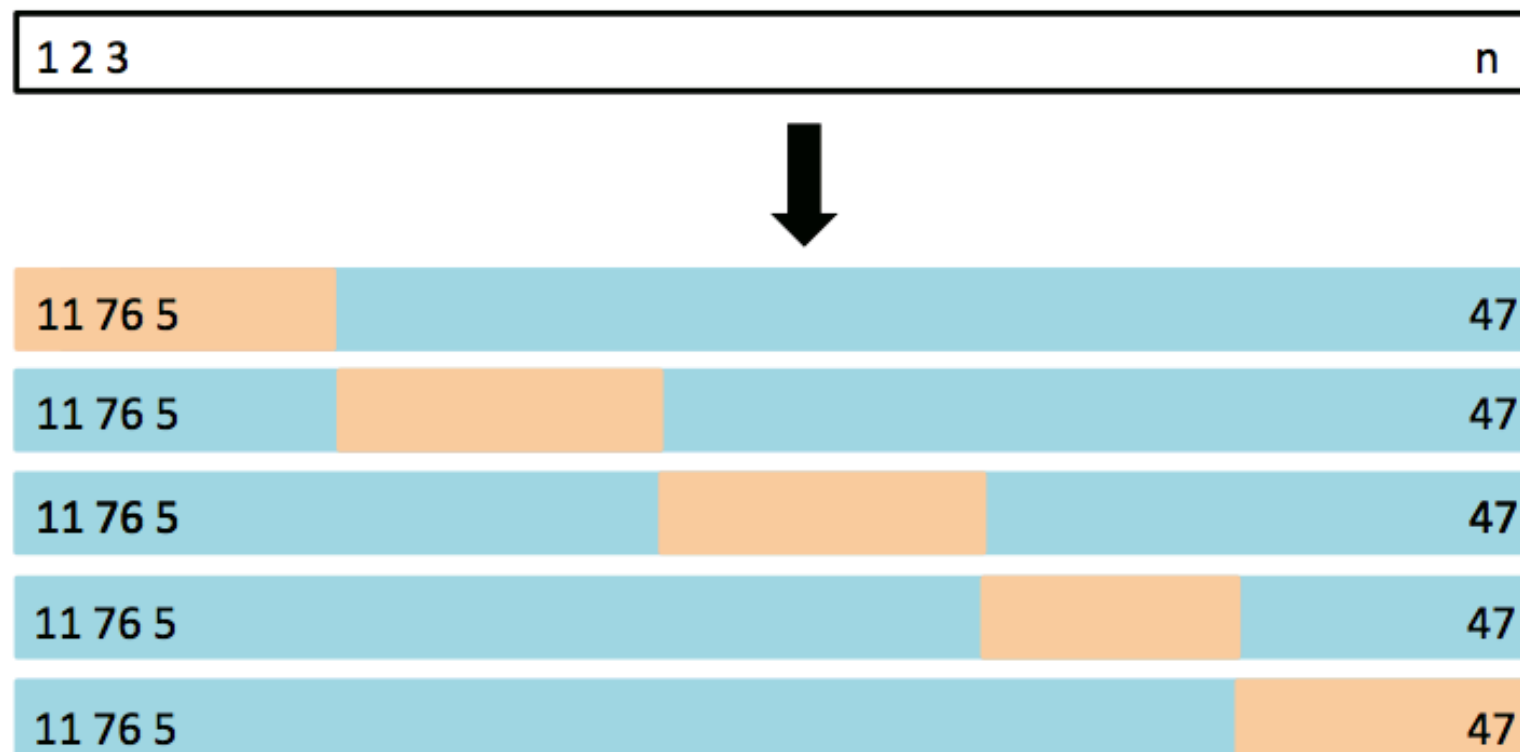
# K-Fold Cross Validation

Split the data into  $k$  subsets or *folds*.

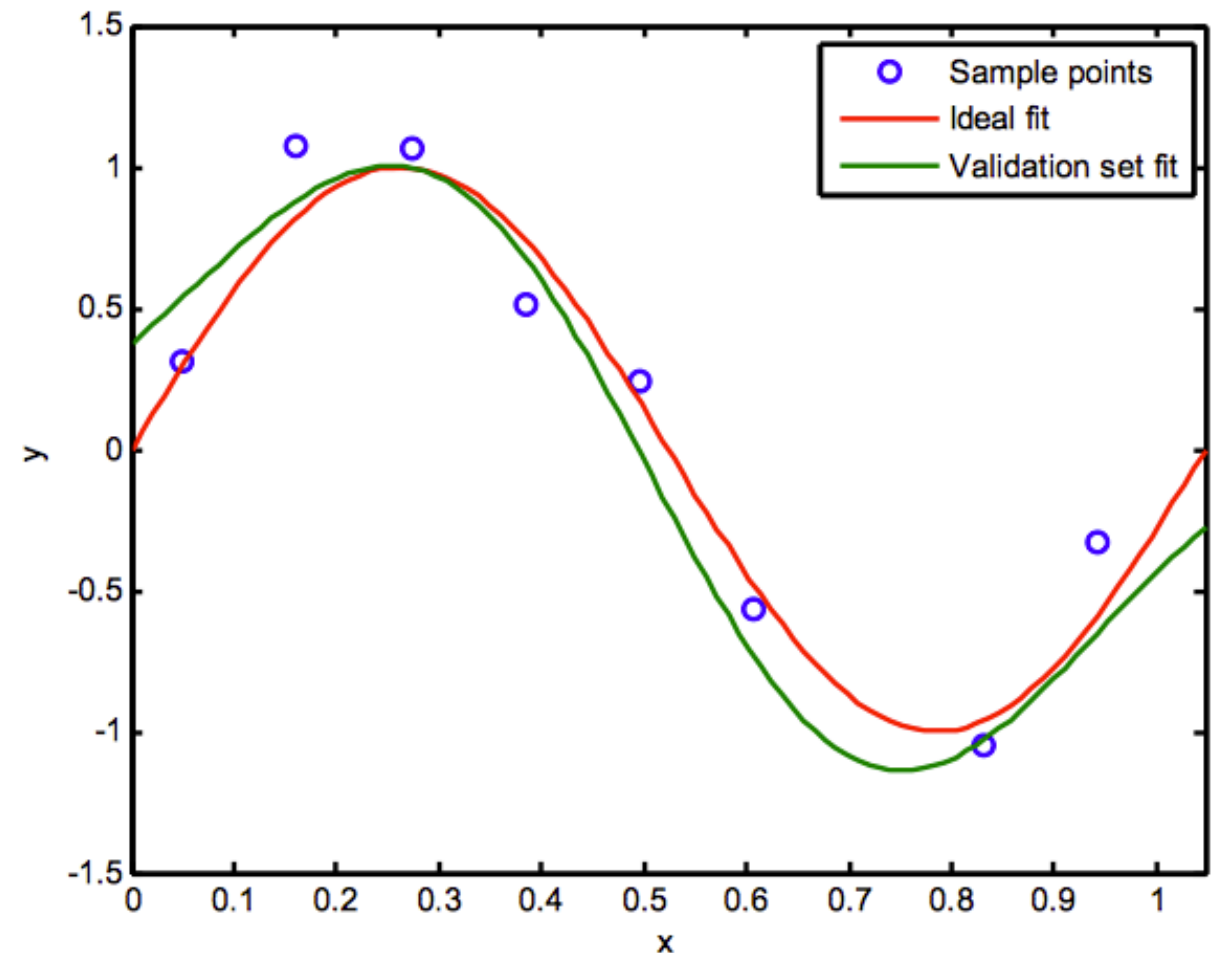
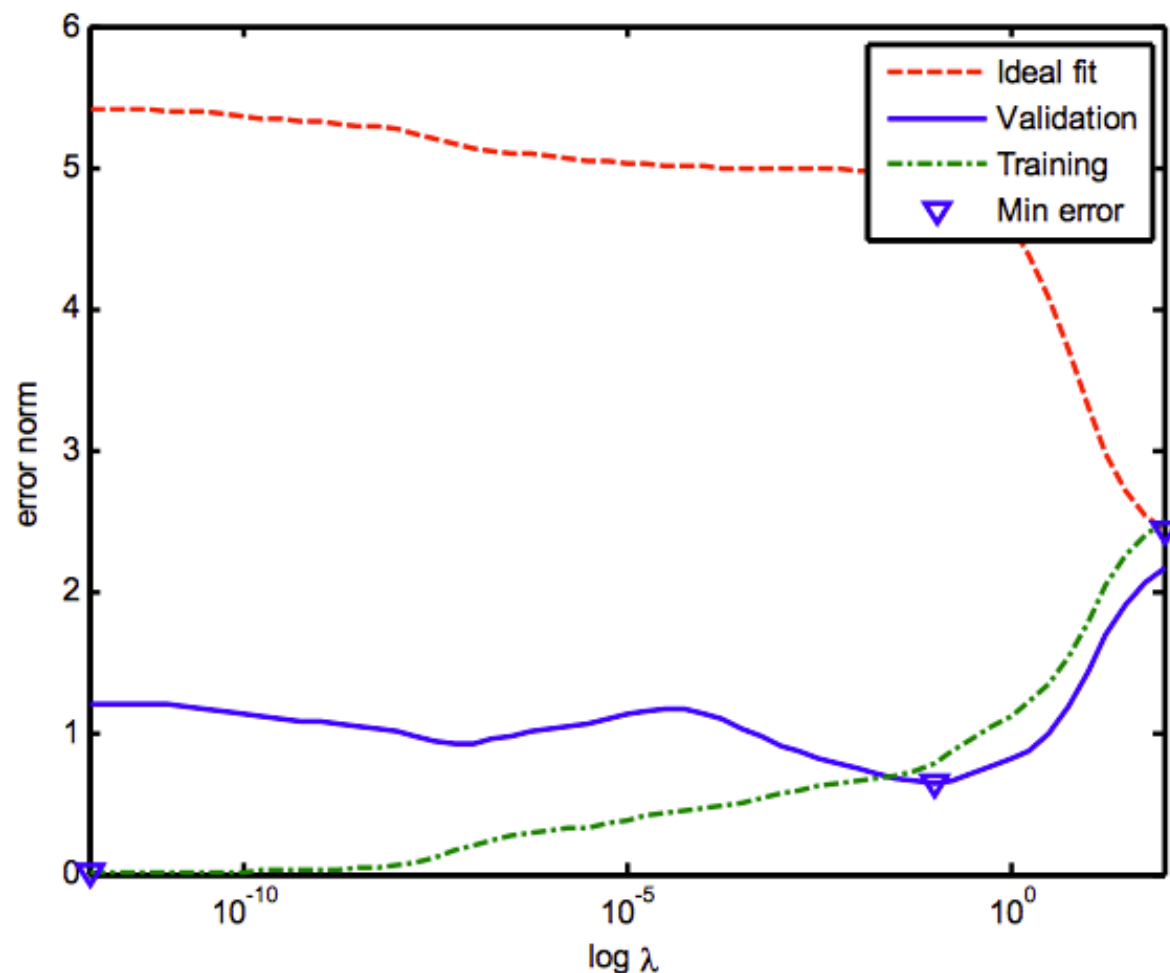
For every  $i = 1, \dots, k$ :

- ▶ train the model on every fold except the  $i$ th fold,
- ▶ compute the test error on the  $i$ th fold.

Average the test errors.



# Choosing $\lambda$ Using Validation Dataset



Pick up the lambda with the lowest mean value of rmse calculated by Cross Validation approach

# Take-Home Messages

- What is overfitting
- What is regularization
- How does Ridge regression work
- Sparsity properties of Lasso regression
- How to choose the regularization coefficient  $\lambda$