


Announcements

- HW 3 Ongoing
- Project Midpoint Report Due Today
- Anonymous Feedback for completion due this Sunday
- Peer Evaluation Released Today. Due 4/3
- No Quiz this Week

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**

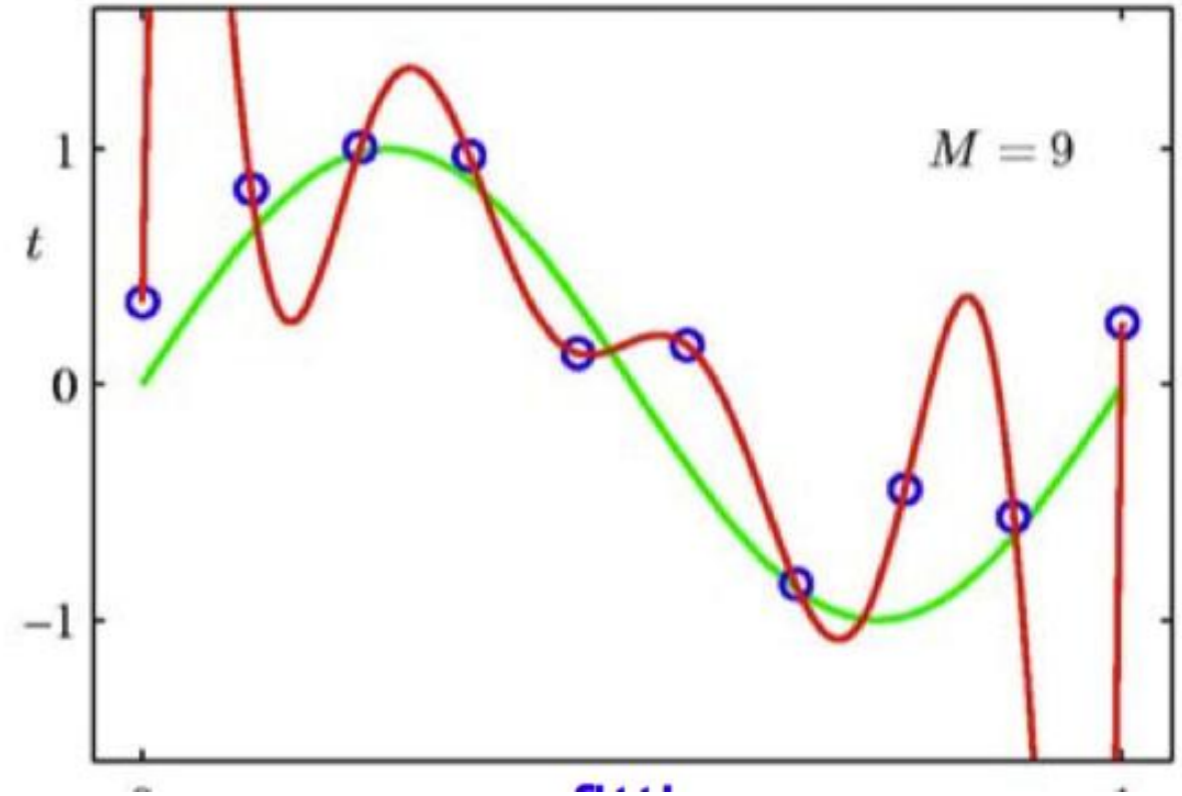
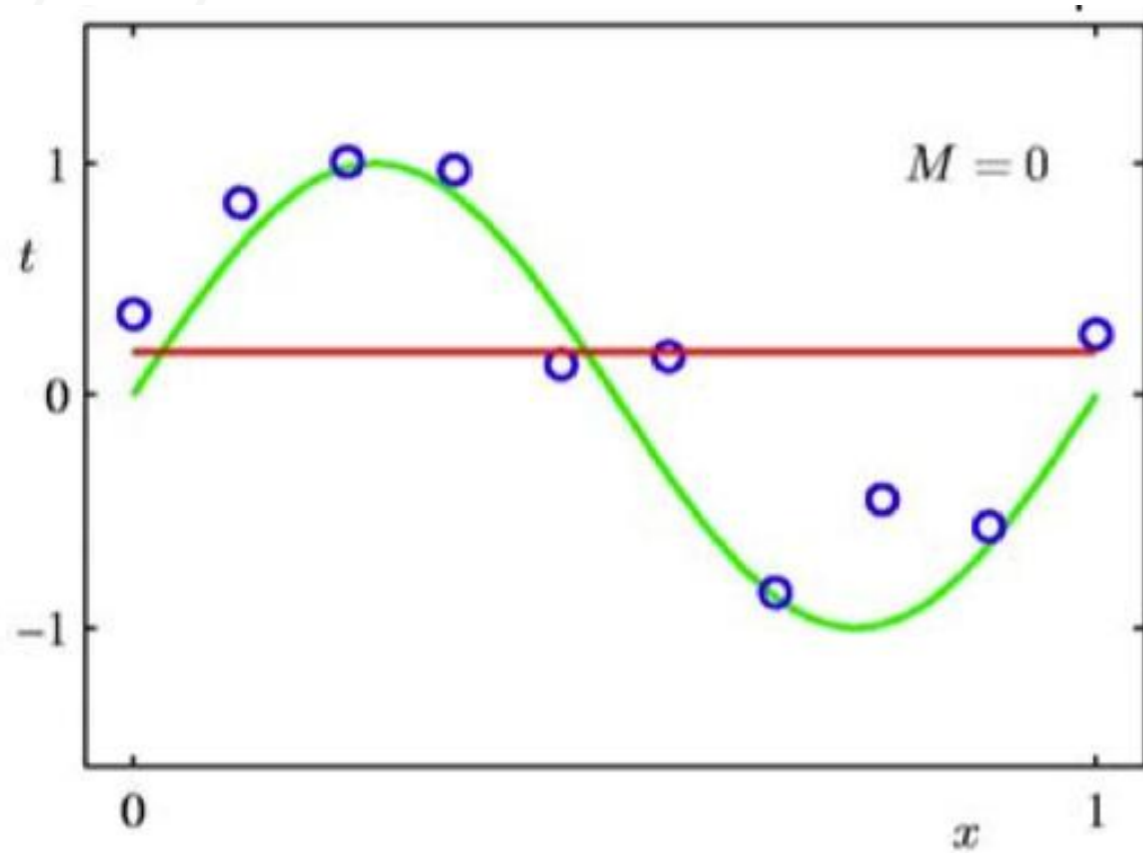
Regularized Linear Regression

Dr. Nimisha Roy
Georgia Tech

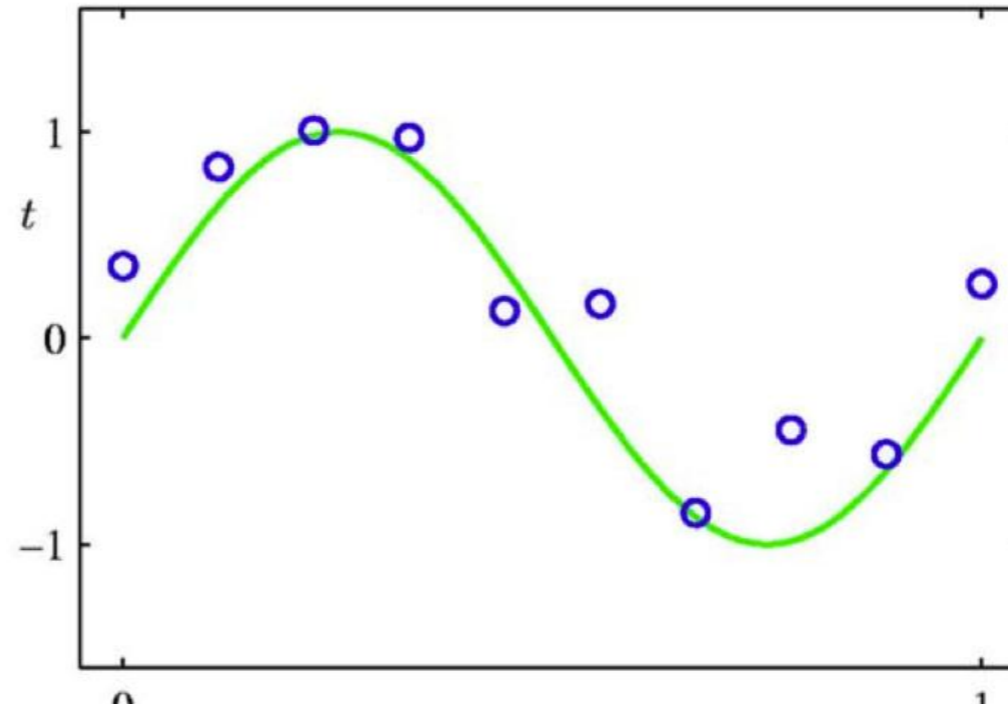
Outline

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

Recap: Overfitting vs Underfitting



Regression: Why is overfitting more common?



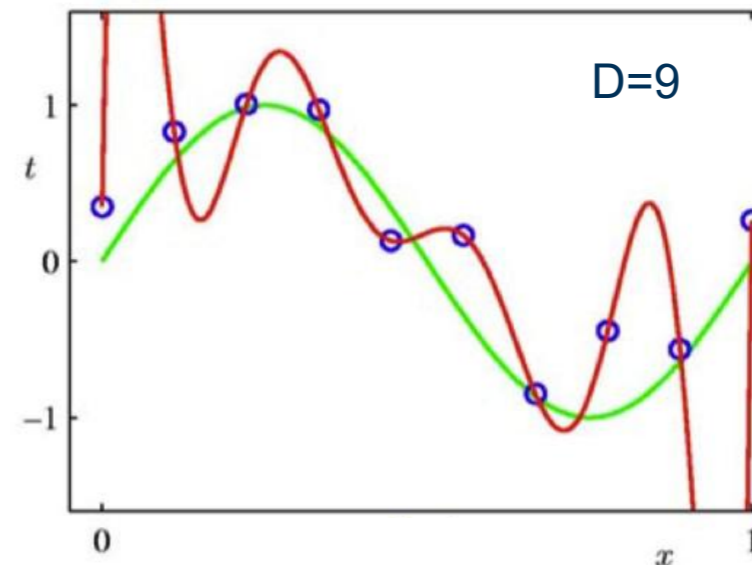
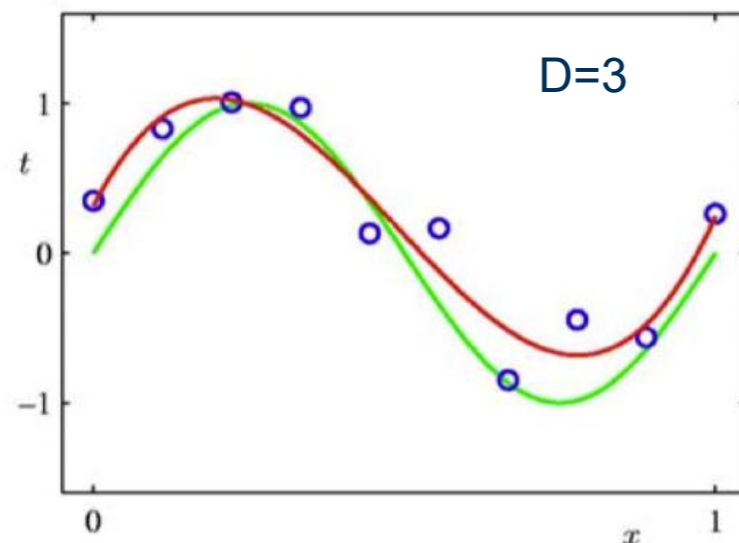
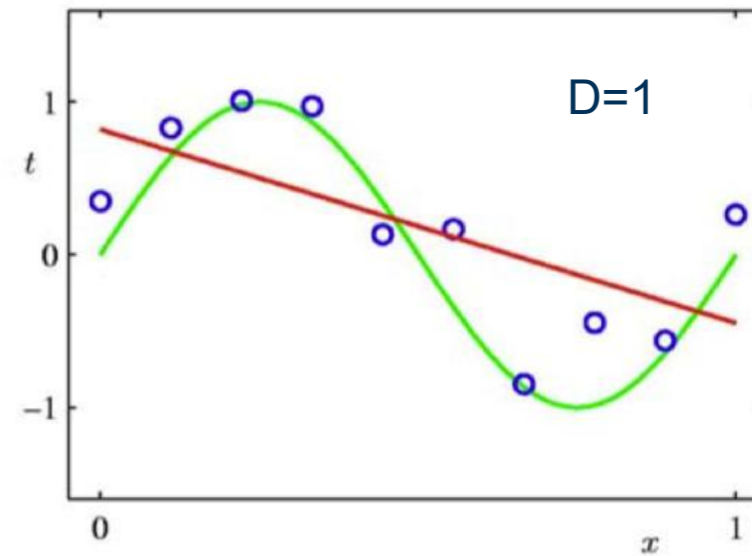
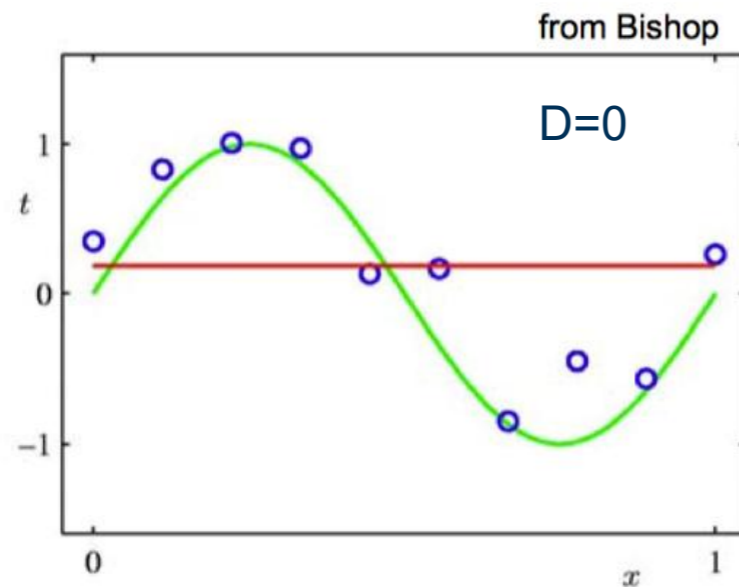
- Want to fit a polynomial regression model

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

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

$$y = z\theta$$

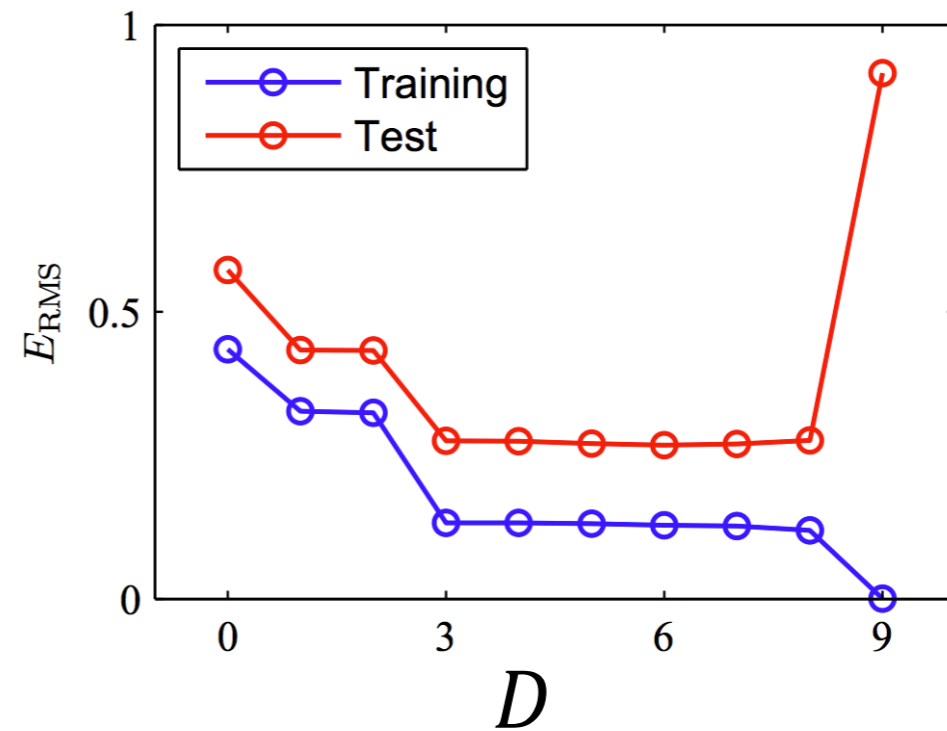
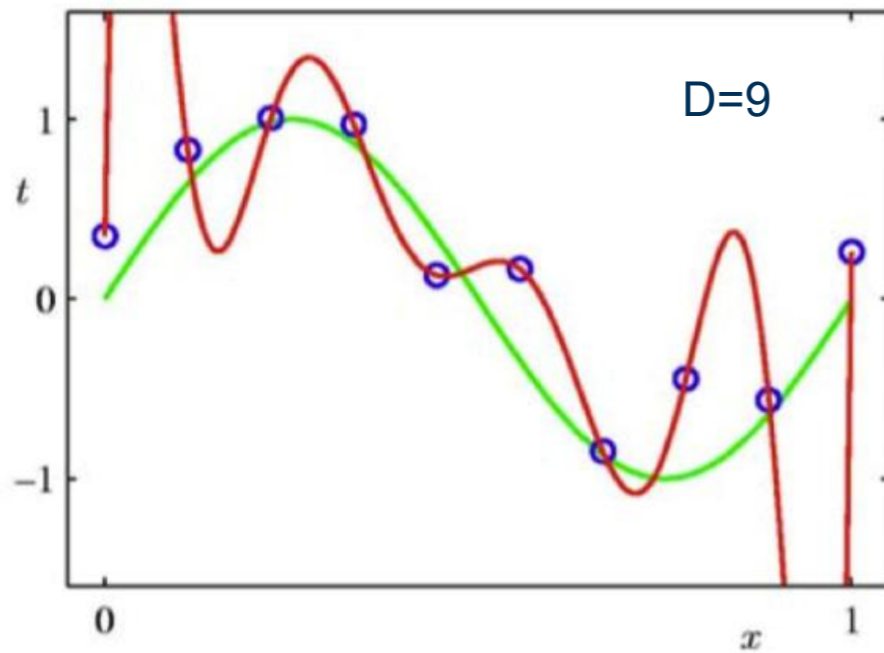
Which One is Better?



- 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

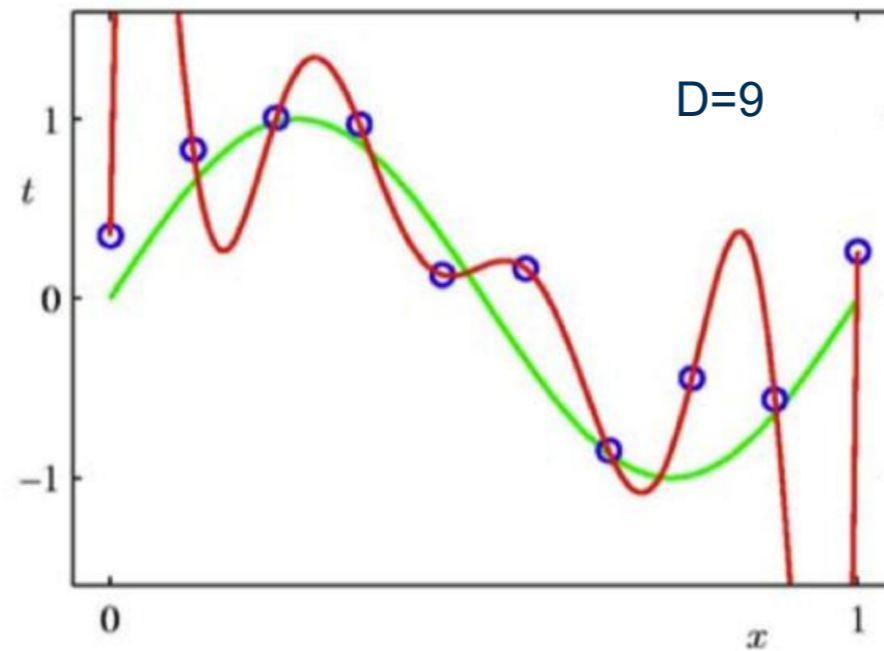


- The training error is very low, but the error on test set is large.
- The model captures not only patterns but also noisy nuisances in the training data.

Root cause of overfitting?

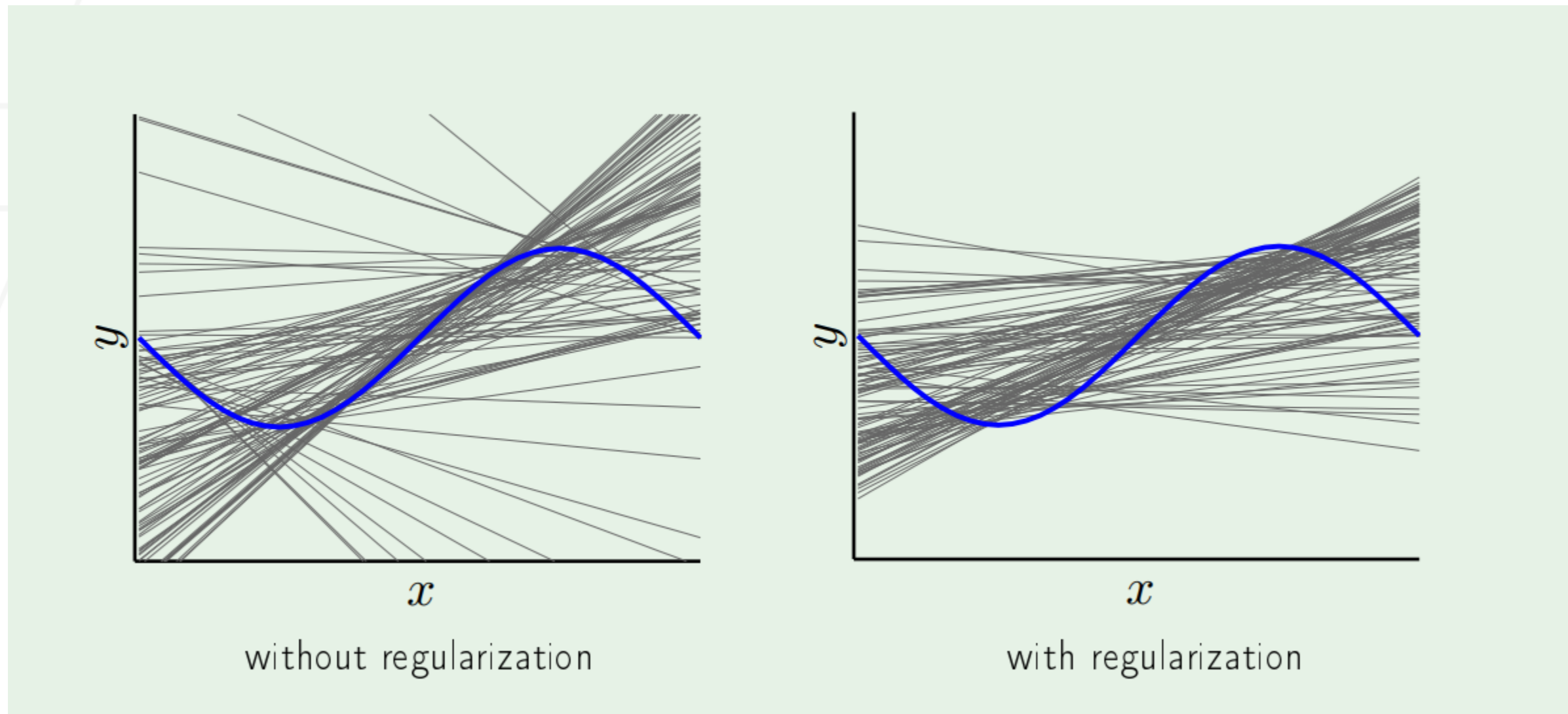
- I want to convert 1 feature in x space to d features in z space

The Overfitting Problem



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

Regularization (smart way to cure overfitting disease)



Put a brake on fitting

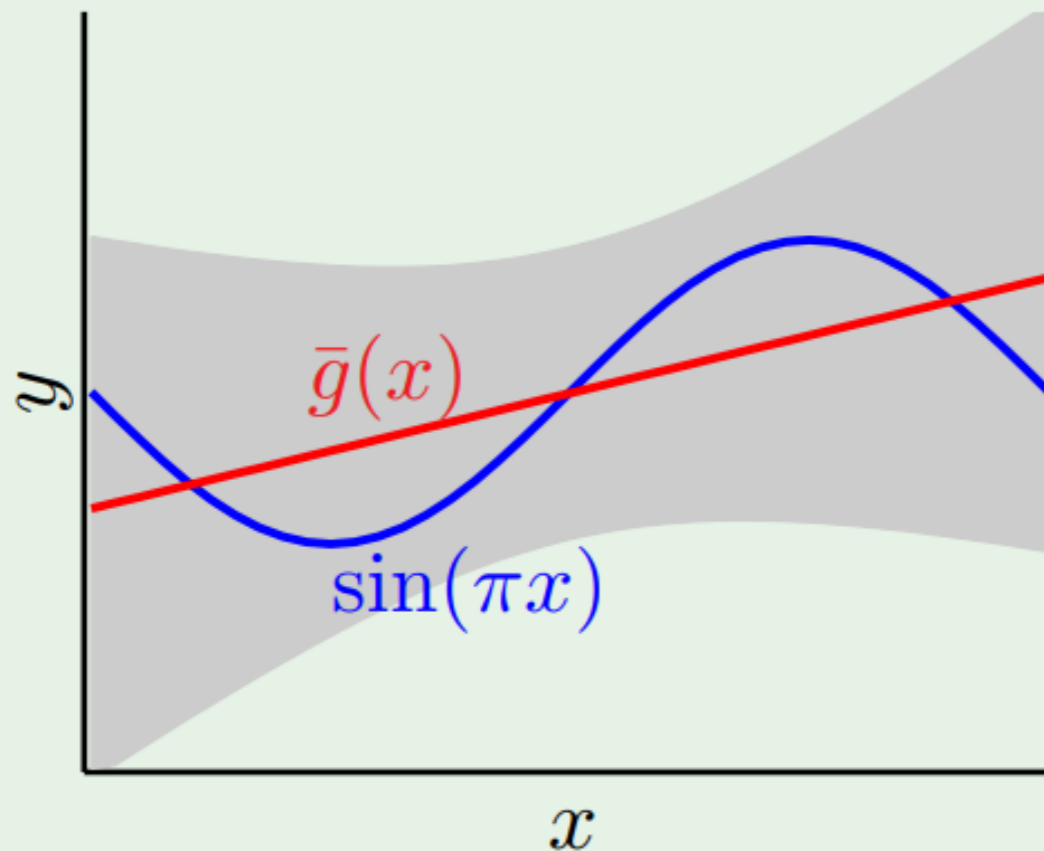


Fit a linear line on sinusoidal with just two data points

Who is the winner?

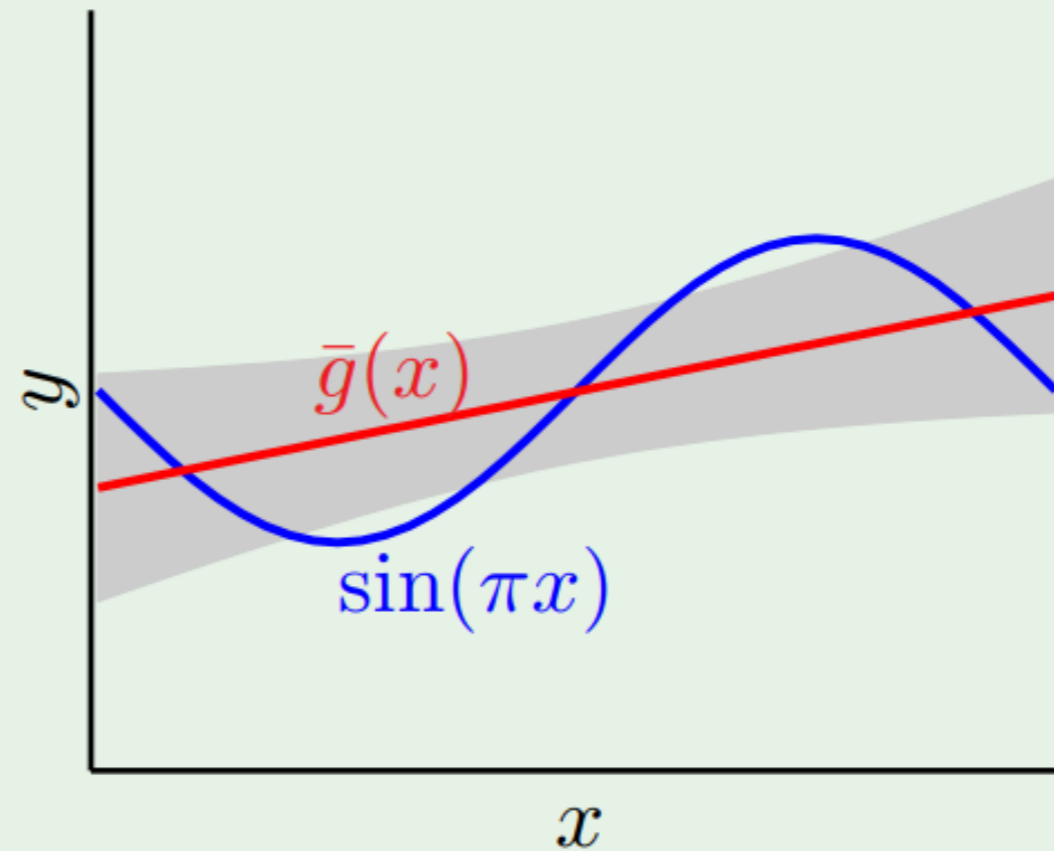
$\bar{g}(x)$: average over all lines

without regularization



bias=0.21; var=1.69

with regularization





bias=0.23; var=0.33

So, what is regularization in general?

Regularization adds a penalty term to the error or objective function:

- The penalty discourages large weights → the model becomes *simpler*.
- A simpler model can't bend perfectly to fit all training data → **bias increases slightly**.
- But since it's not so sensitive to the quirks of any one dataset → **variance decreases significantly**.
- That tradeoff leads to better **generalization**.

Regularization	Bias	Variance	Total Error	Stability
 Without	Low	High	High	Overfits
 With	Slightly higher	Much lower	Lower	Generalizes better

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 (θ)

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

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

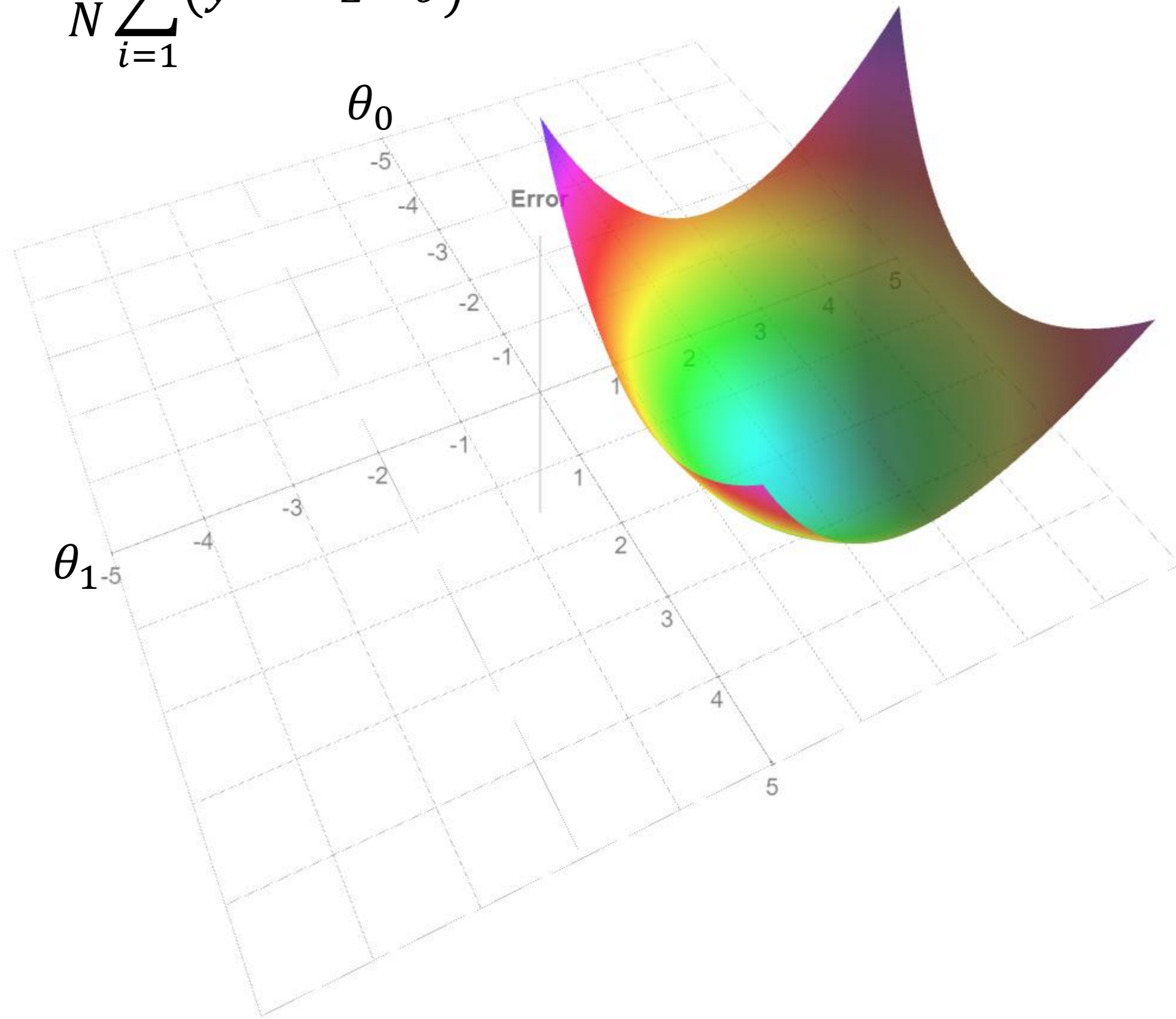
subject to

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

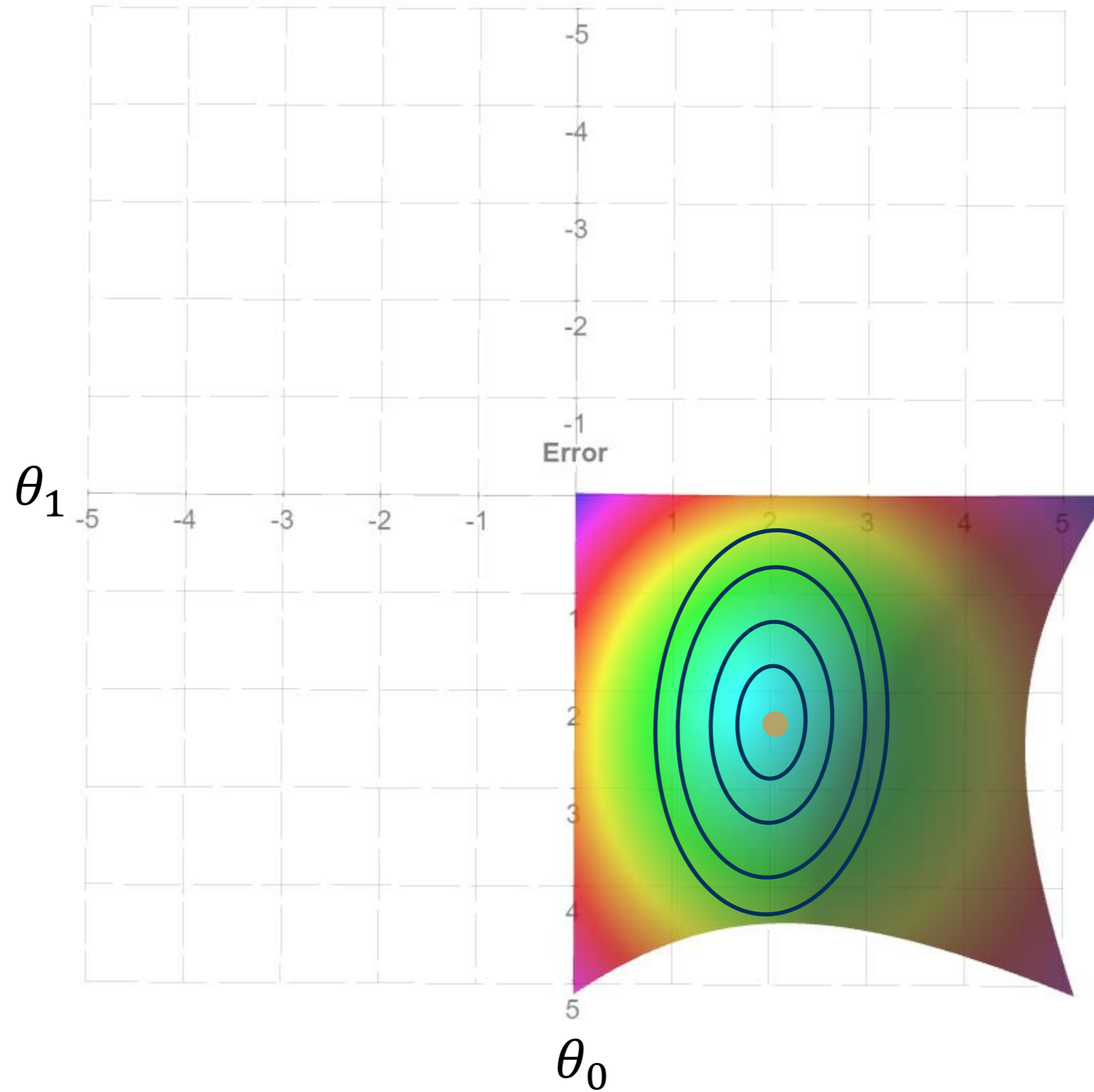


$$y = \theta_0 + \theta_1 z_1 + \theta_2 z_2 + 0 + \cdots + 0$$

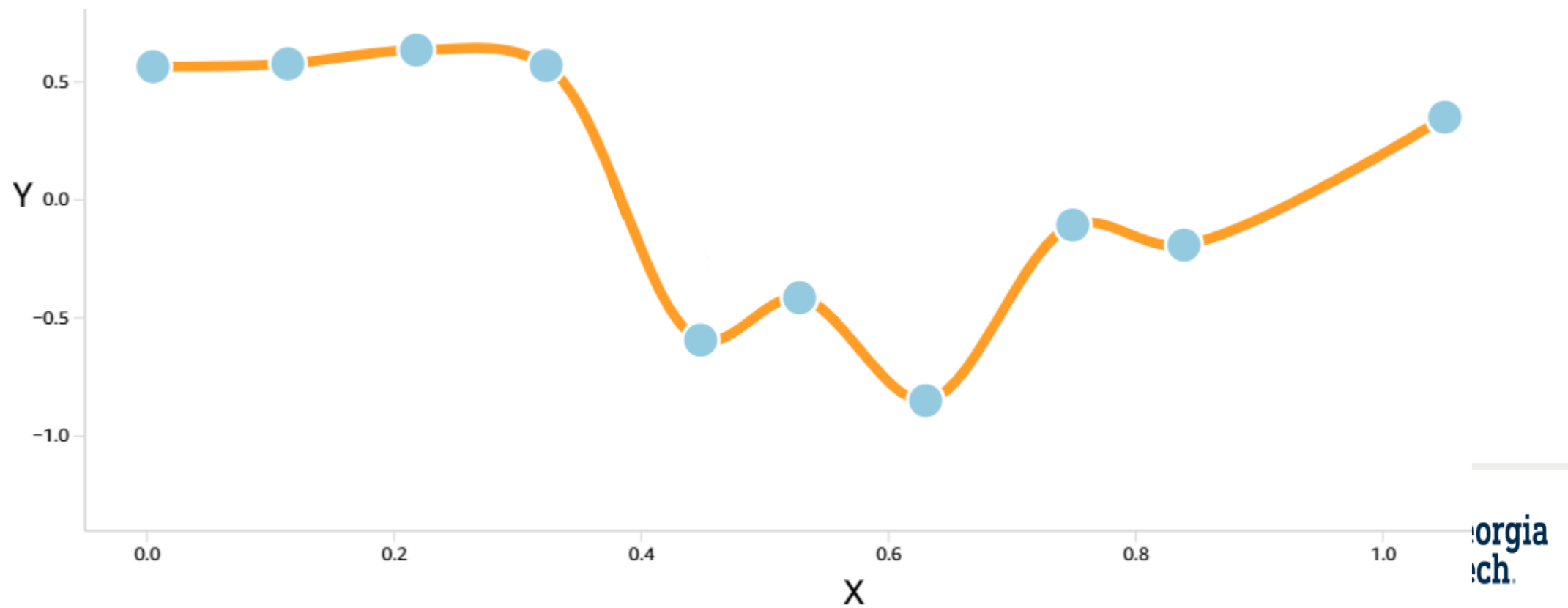
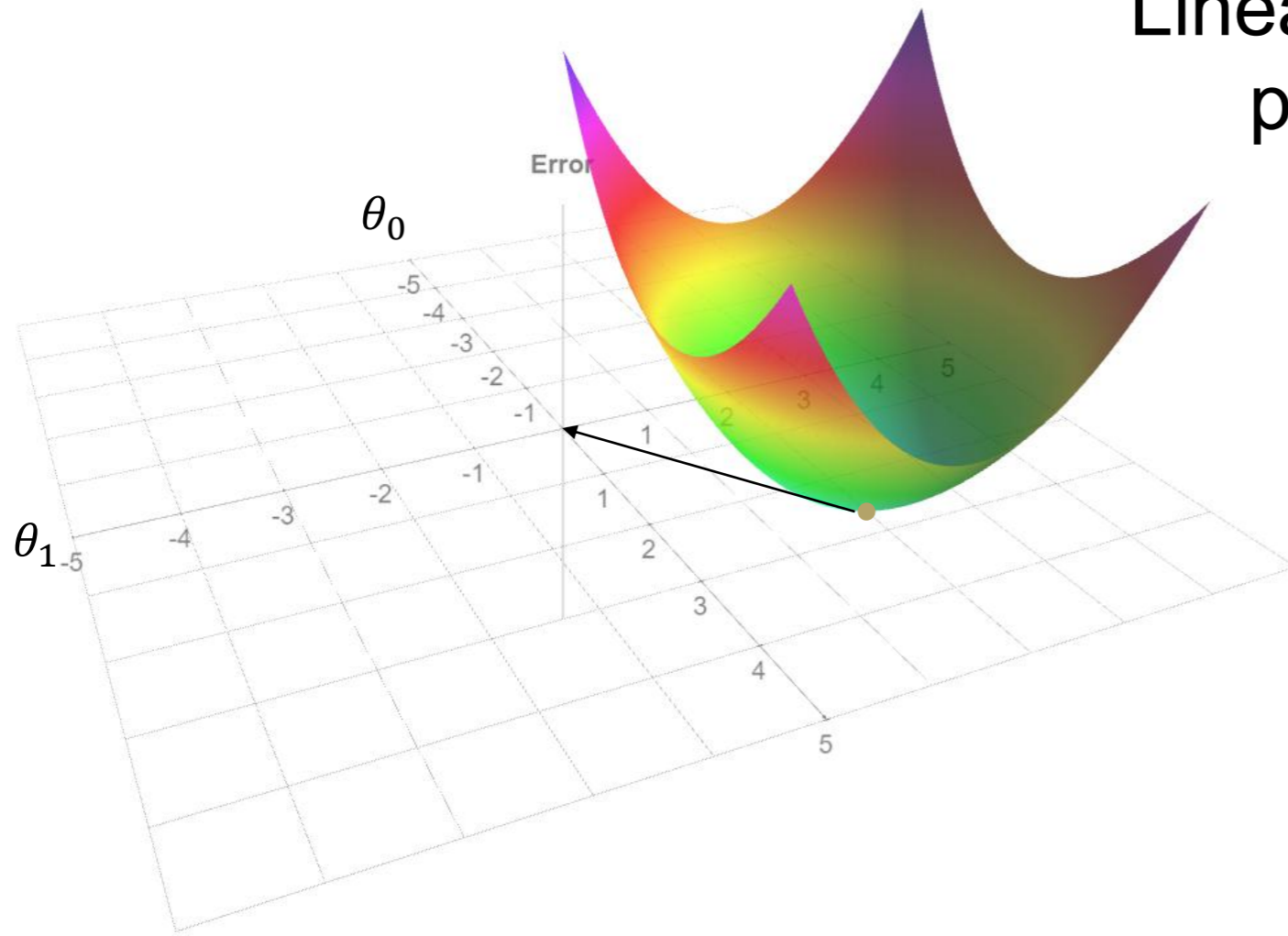
$$E(\theta) = \frac{1}{N} \sum_{i=1}^n (y^{i} - z^{i}\theta)^2$$

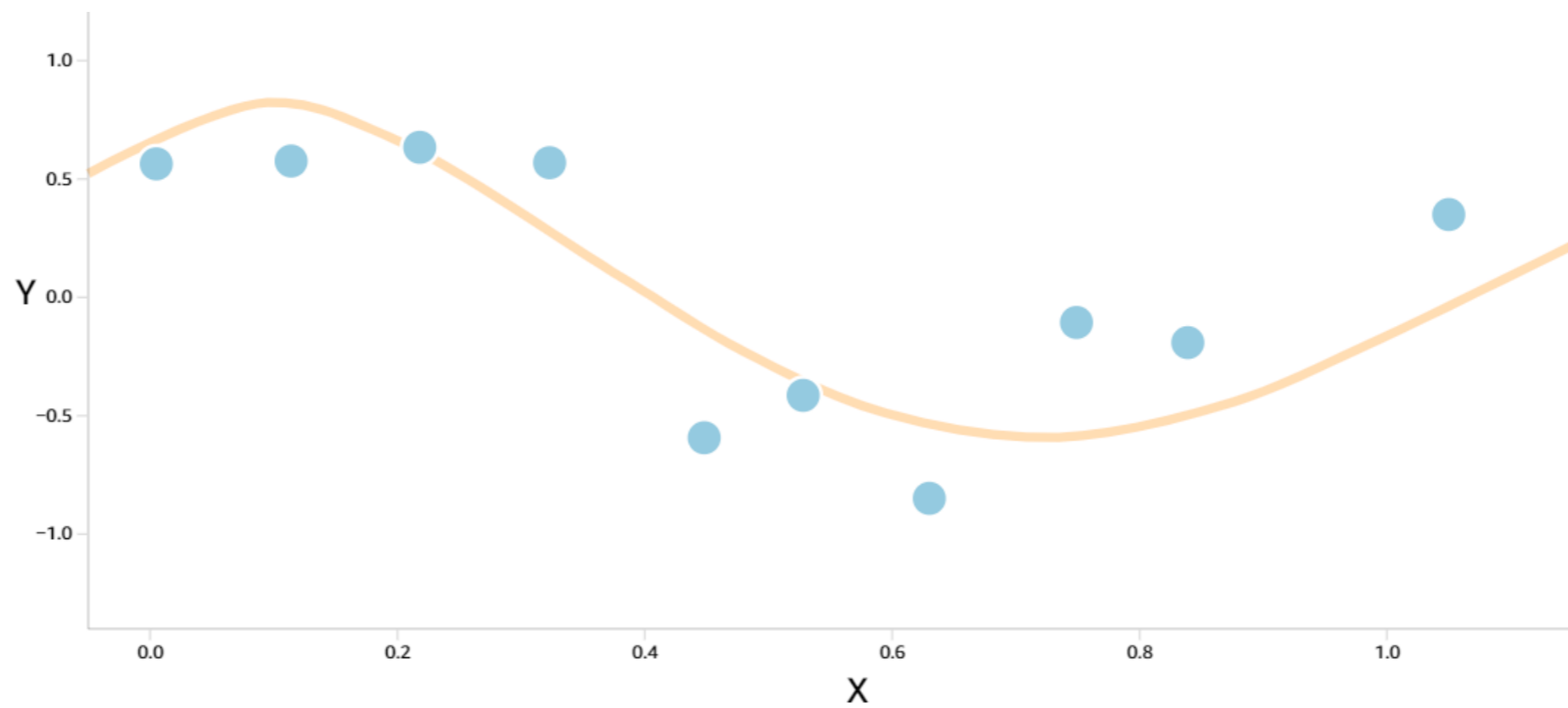
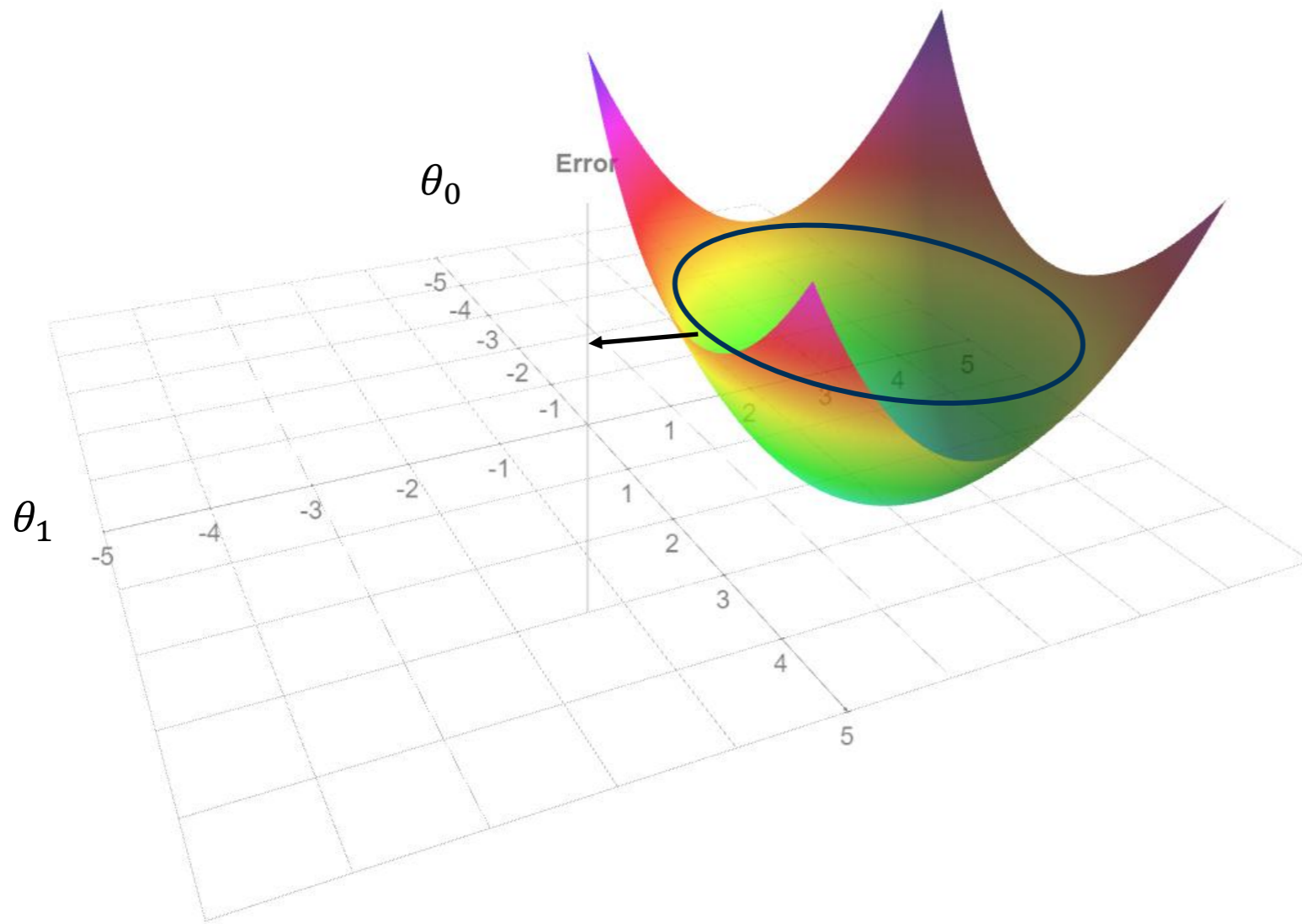


Project the same graph on x-y using contour plot



Linear regression with a very high polynomial degree solution

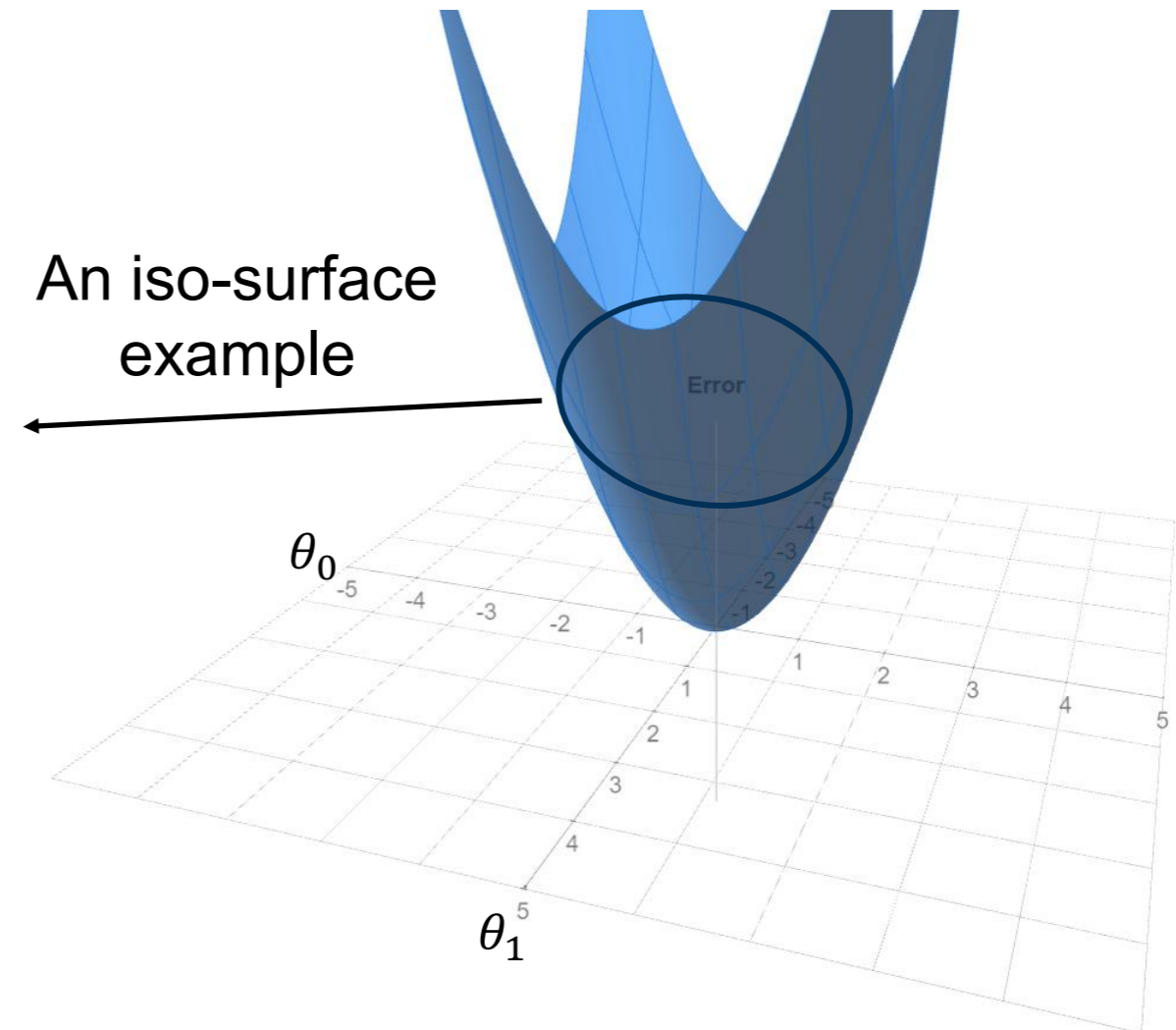




How can we get an optimal solution with a positive error for a model that overfits?

We need to introduce a constraint

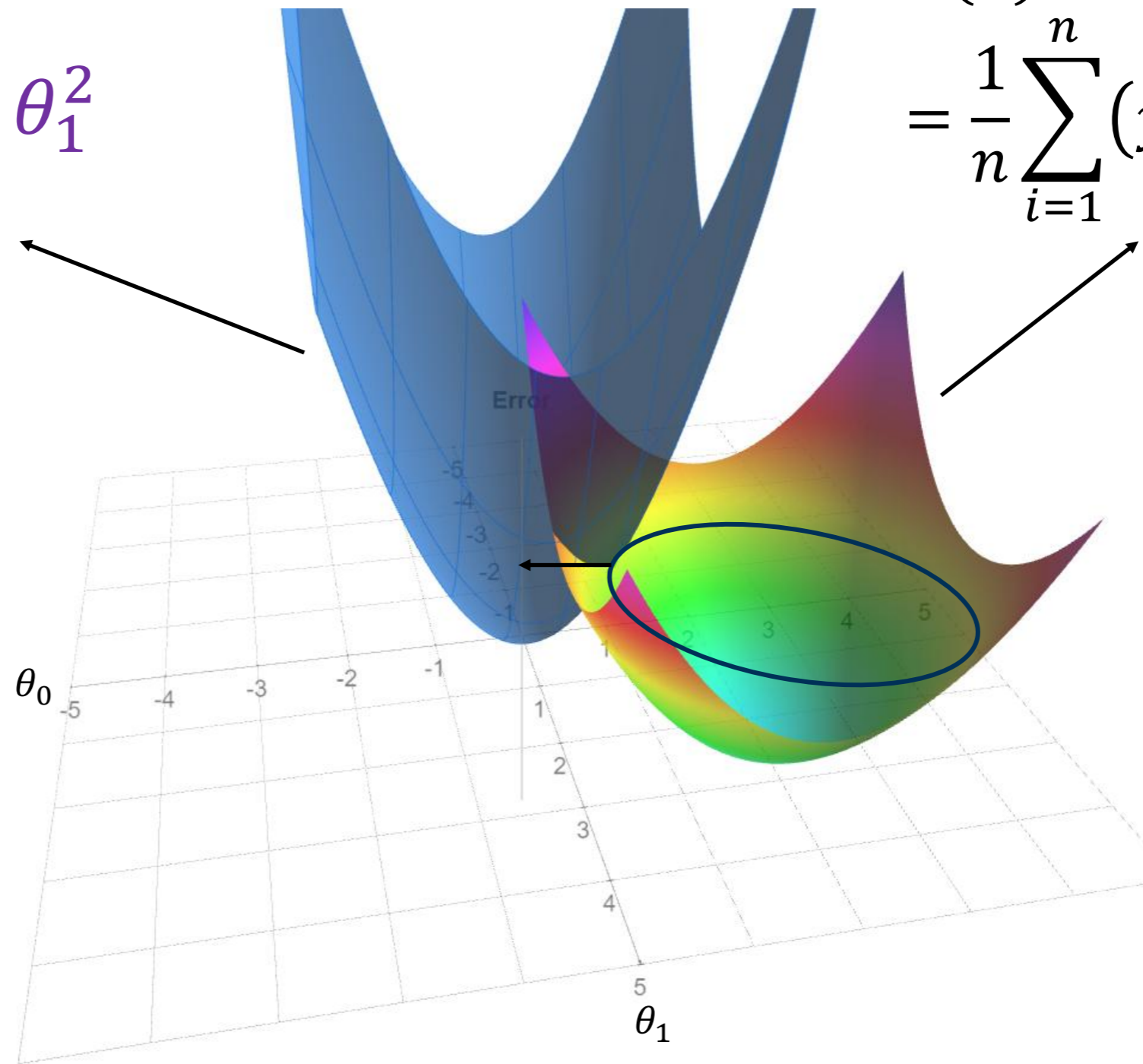
$$\begin{aligned}g(\theta) &= \theta_0^2 + \theta_1^2 \\ &= \theta^T \theta = C\end{aligned}$$



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

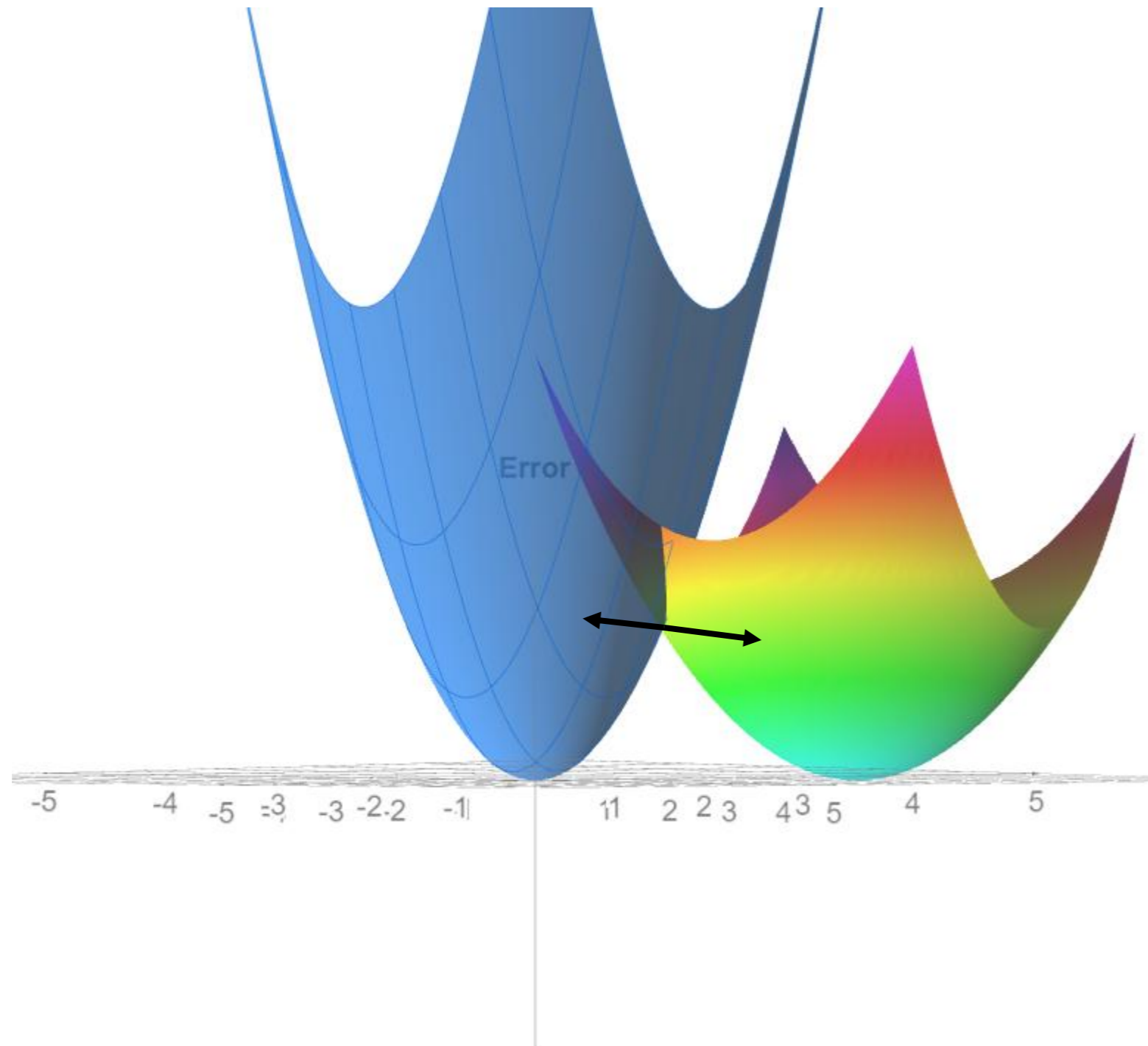
$$L(\theta, \lambda) = E(\theta) + \lambda g(\theta)$$

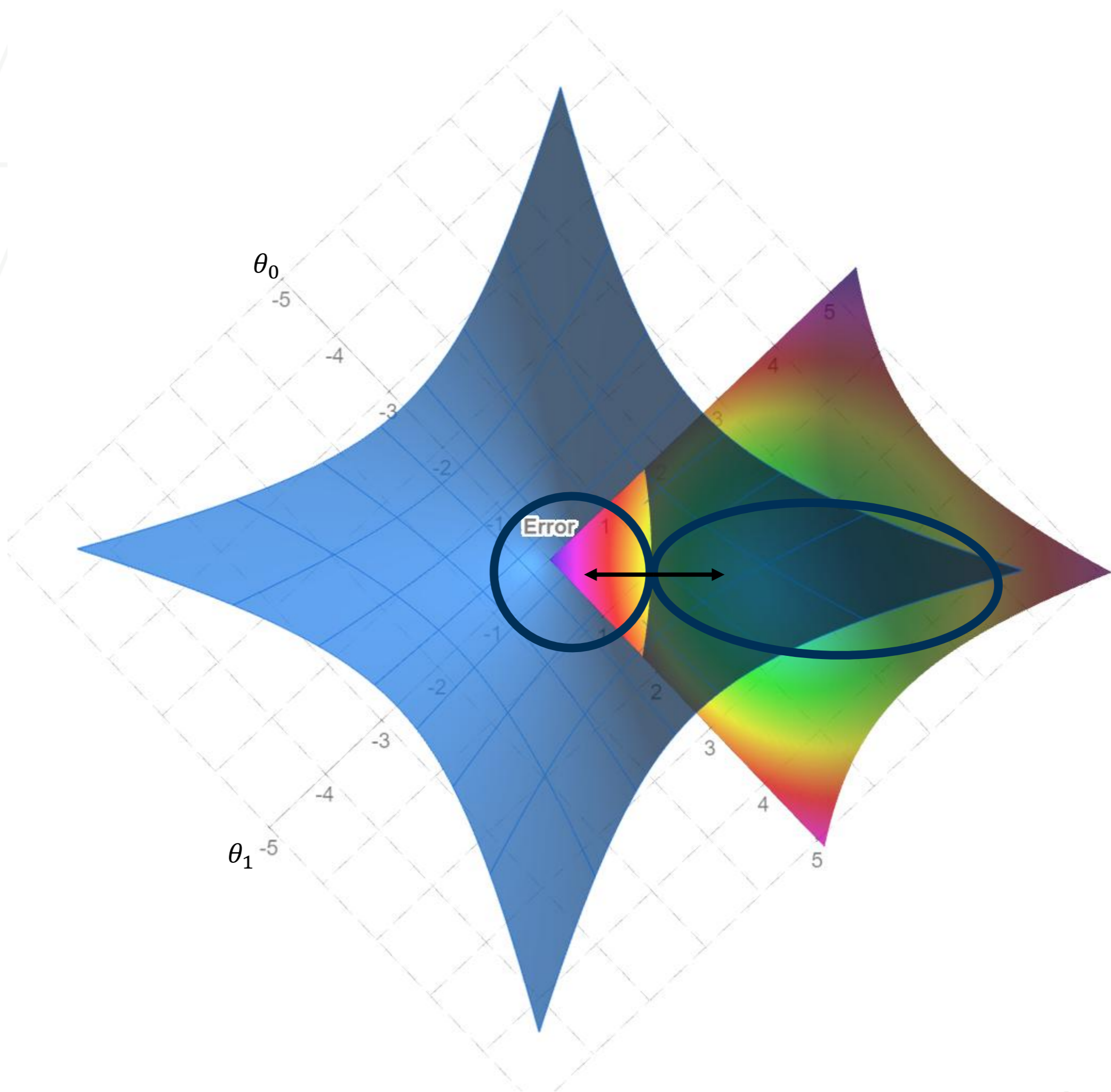
$$L(\theta, \lambda) = E(\theta) + \lambda \theta^T \theta - \lambda C$$

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

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

How to enforce the gradient of Lagrange function to be zero







Let's calculate the gradients

Gradient of constraint $g(\theta)$ $\nabla[\theta^T \theta] = 2\theta$

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

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

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

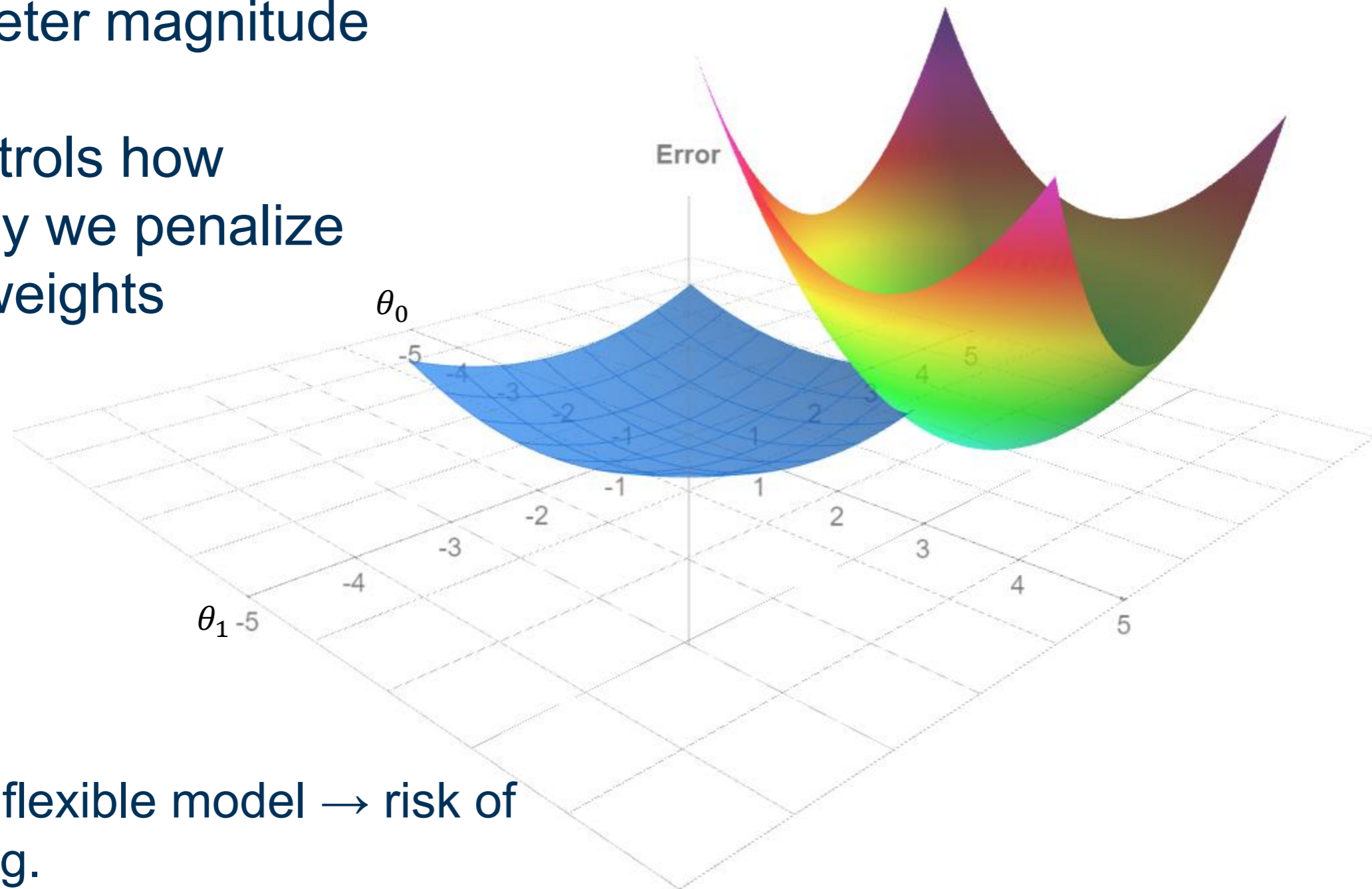
$$\nabla E(\theta) + 2\lambda\theta = 0 \quad \left. \vphantom{\nabla E(\theta) + 2\lambda\theta = 0} \right\} \text{Let's do integration} \quad E(\theta) + \lambda\theta^T \theta$$

The effect of low Lambda

$\frac{\lambda}{N} \theta^T \theta$: The **regularization penalty** on parameter magnitude

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

λ : Controls how strongly we penalize large weights



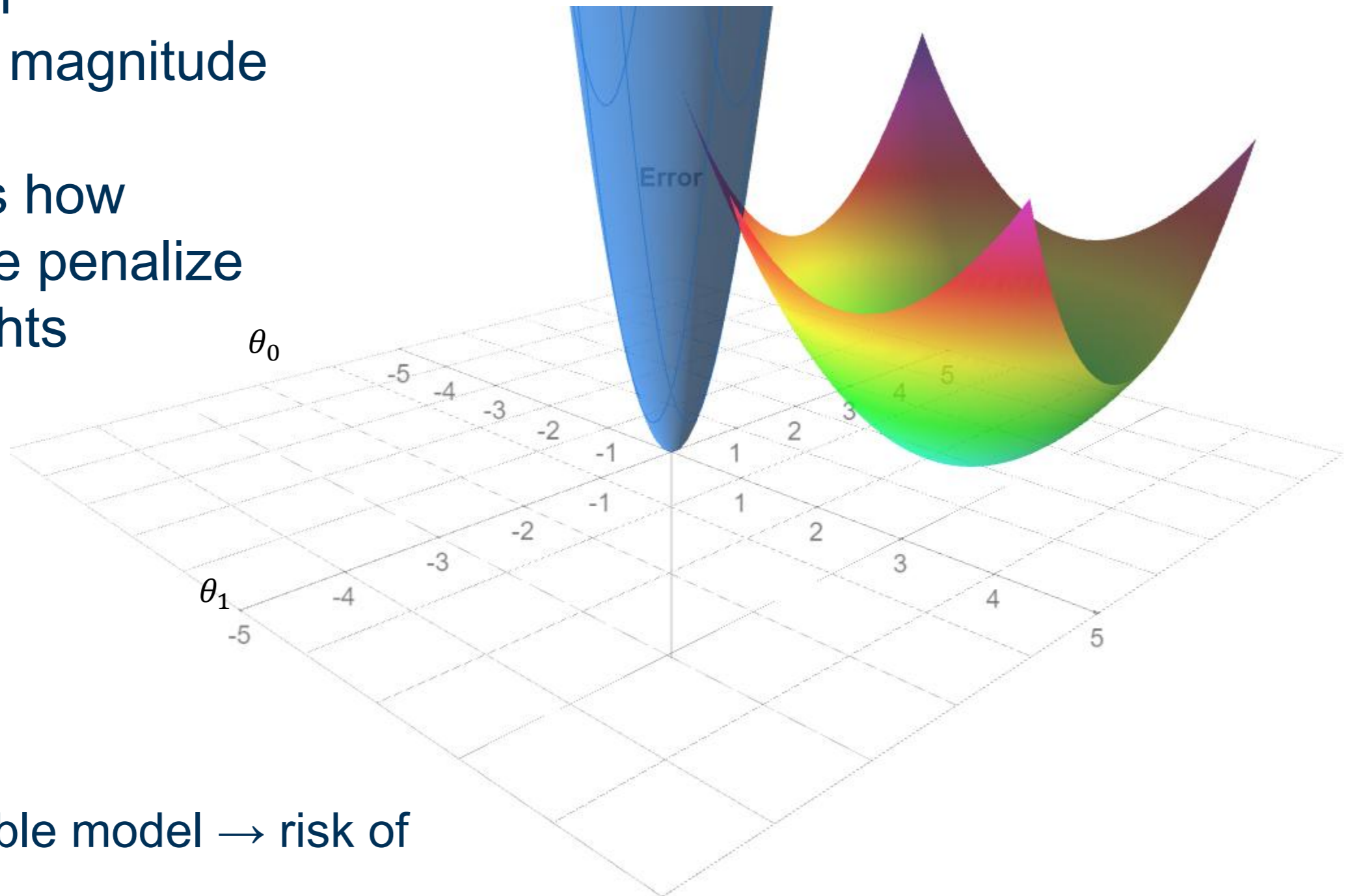
low $\lambda \rightarrow$ flexible model \rightarrow risk of overfitting.

The effect of high Lambda

$\frac{\lambda}{N} \theta^T \theta$: The **regularization penalty** on parameter magnitude

λ : Controls how strongly we penalize large weights

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



high $\lambda \rightarrow$ stable model \rightarrow risk of underfitting.

Regularized Learning

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

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

Regularized Error

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

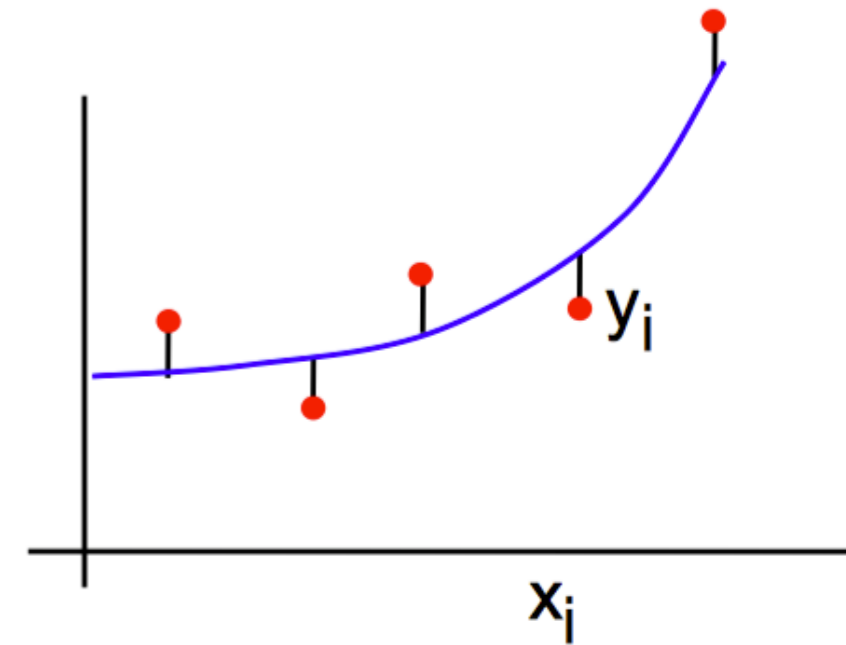
L2 Regularization term

Outline

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

Ridge Regression

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



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

General form

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

Matrix form

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

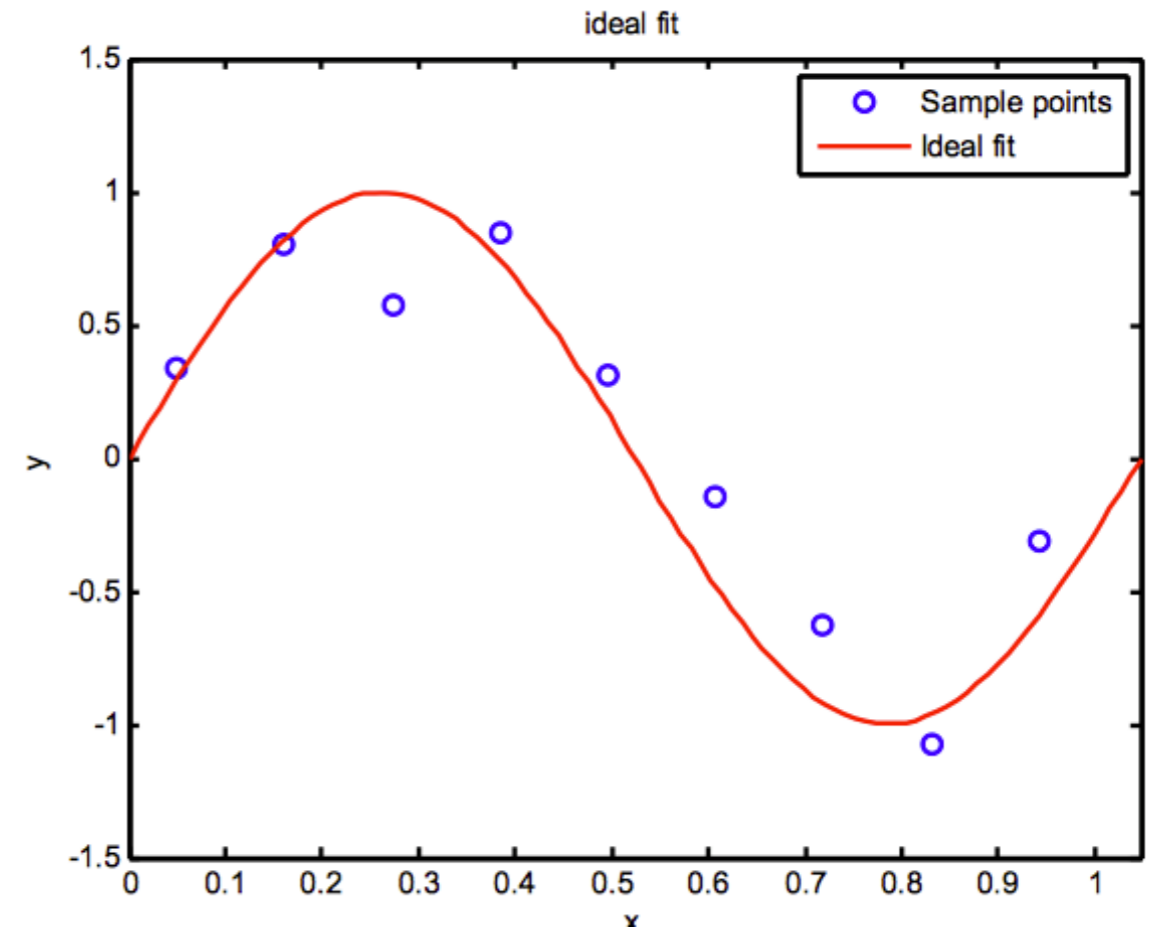
$$\frac{\partial \tilde{E}(\theta)}{\partial \theta} = -z^T (y - z\theta) + \lambda \theta$$

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

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

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

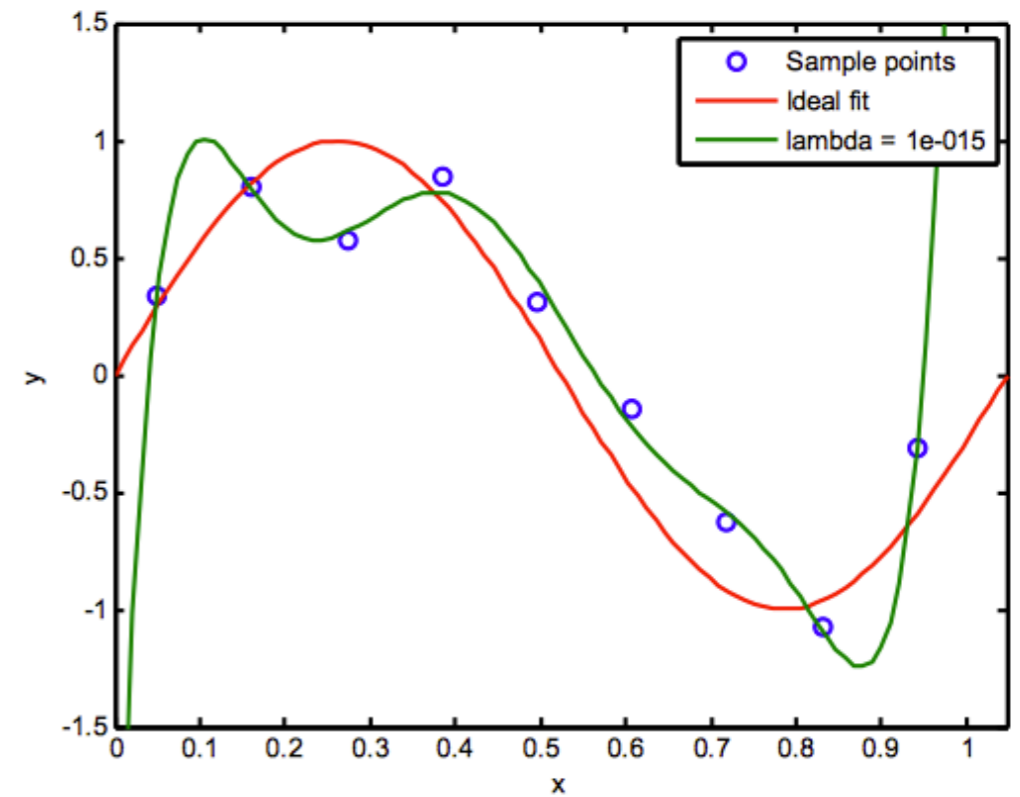
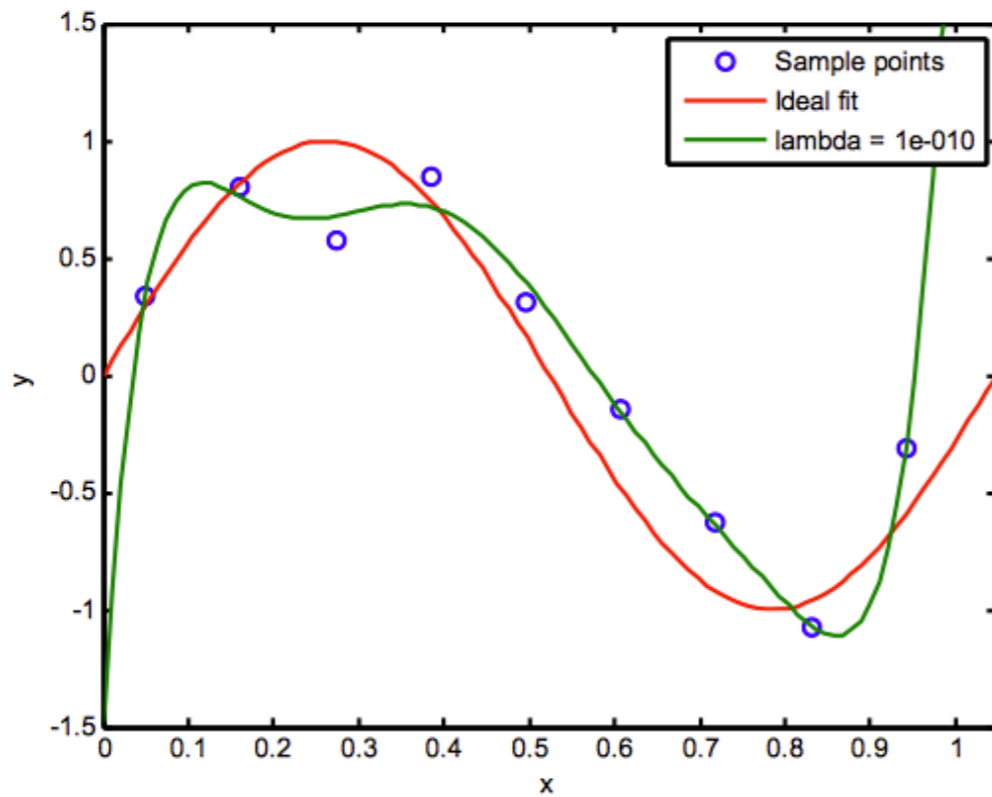
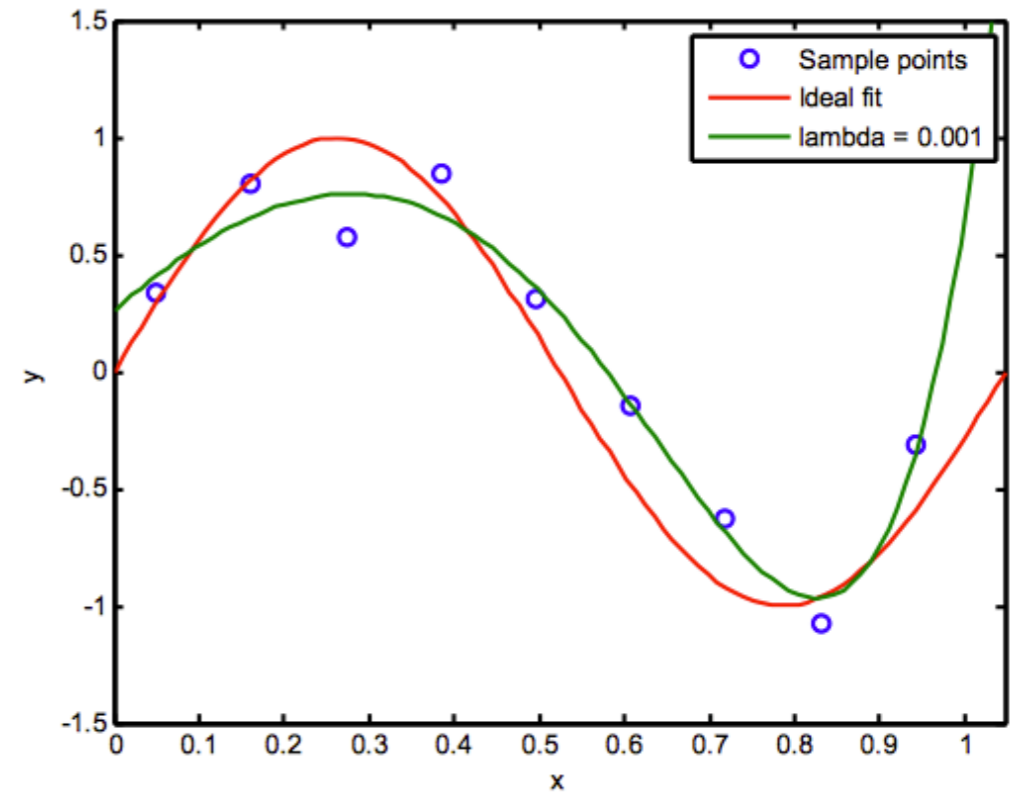
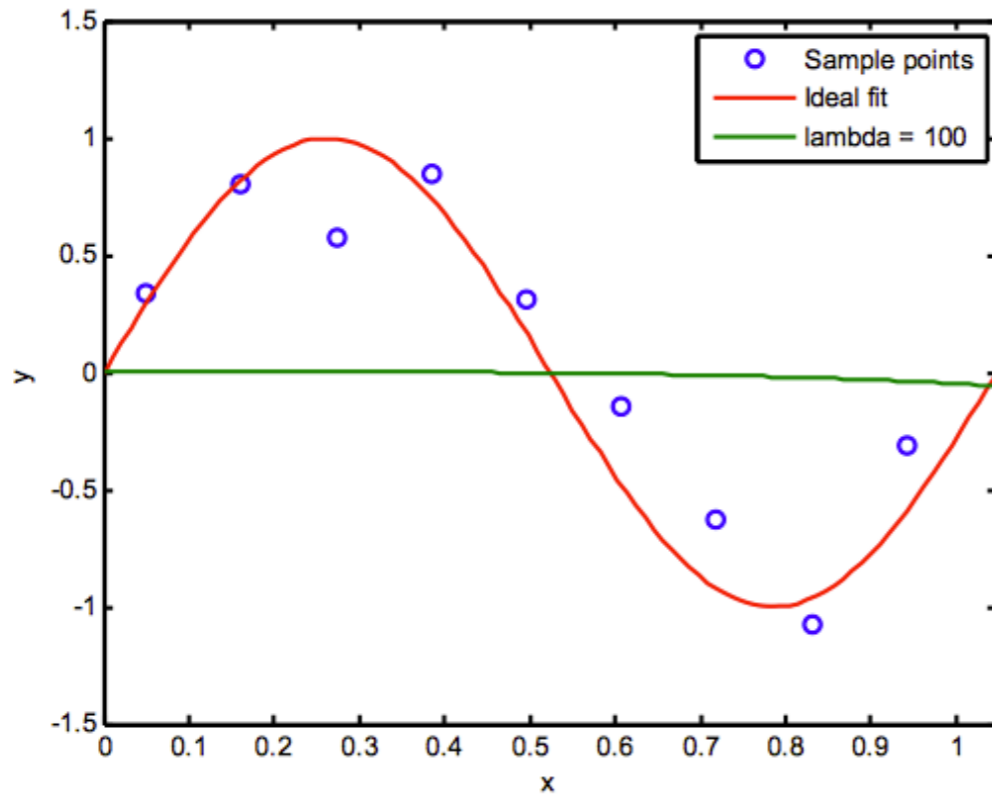


$$f(x, \theta) = z\theta$$

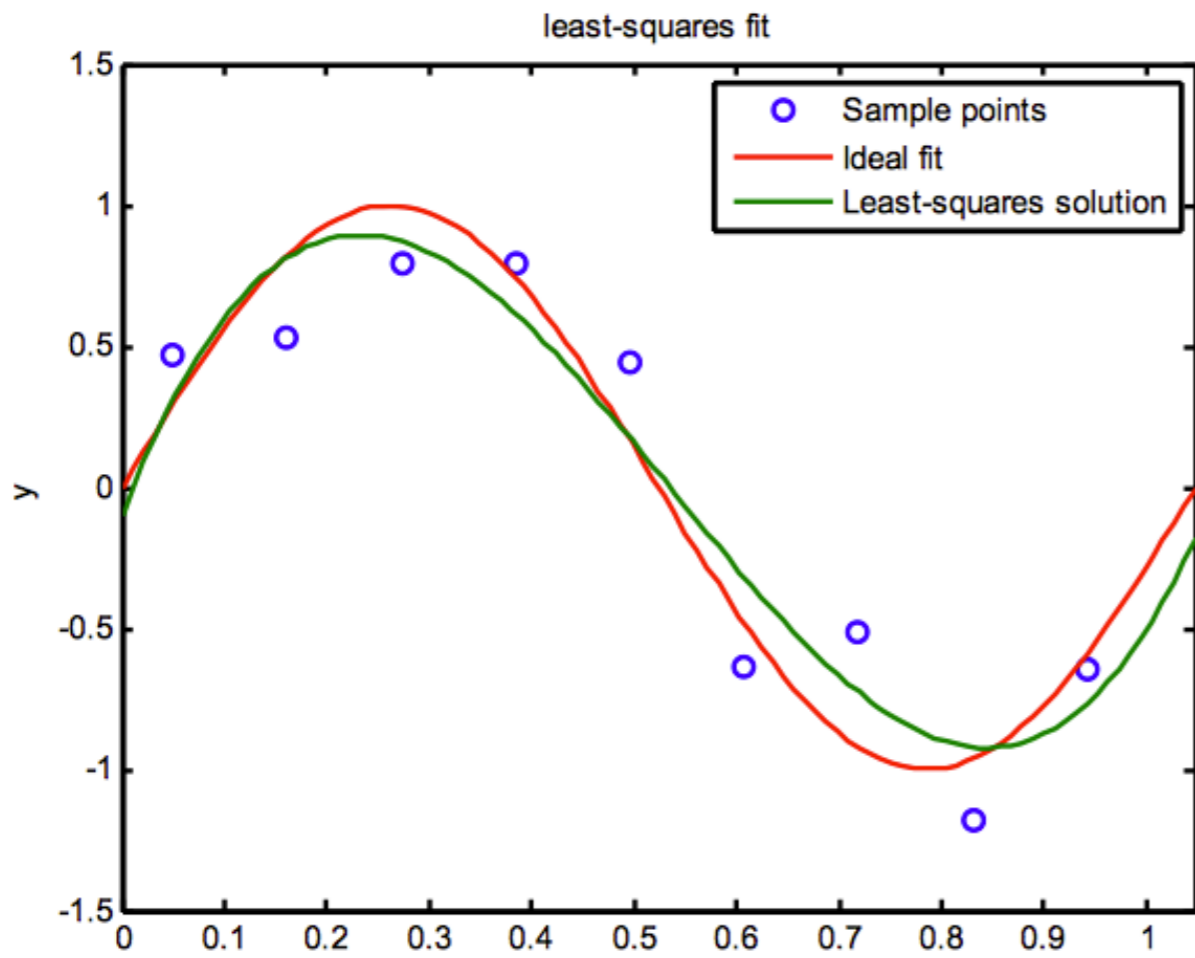
$$z: x \rightarrow z$$

$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^n (y^{\{i\}} - z^{\{i\}}\theta)^2 + \lambda \|\theta\|_2^2 \quad \theta \in \mathbb{R}^{D+1}$$

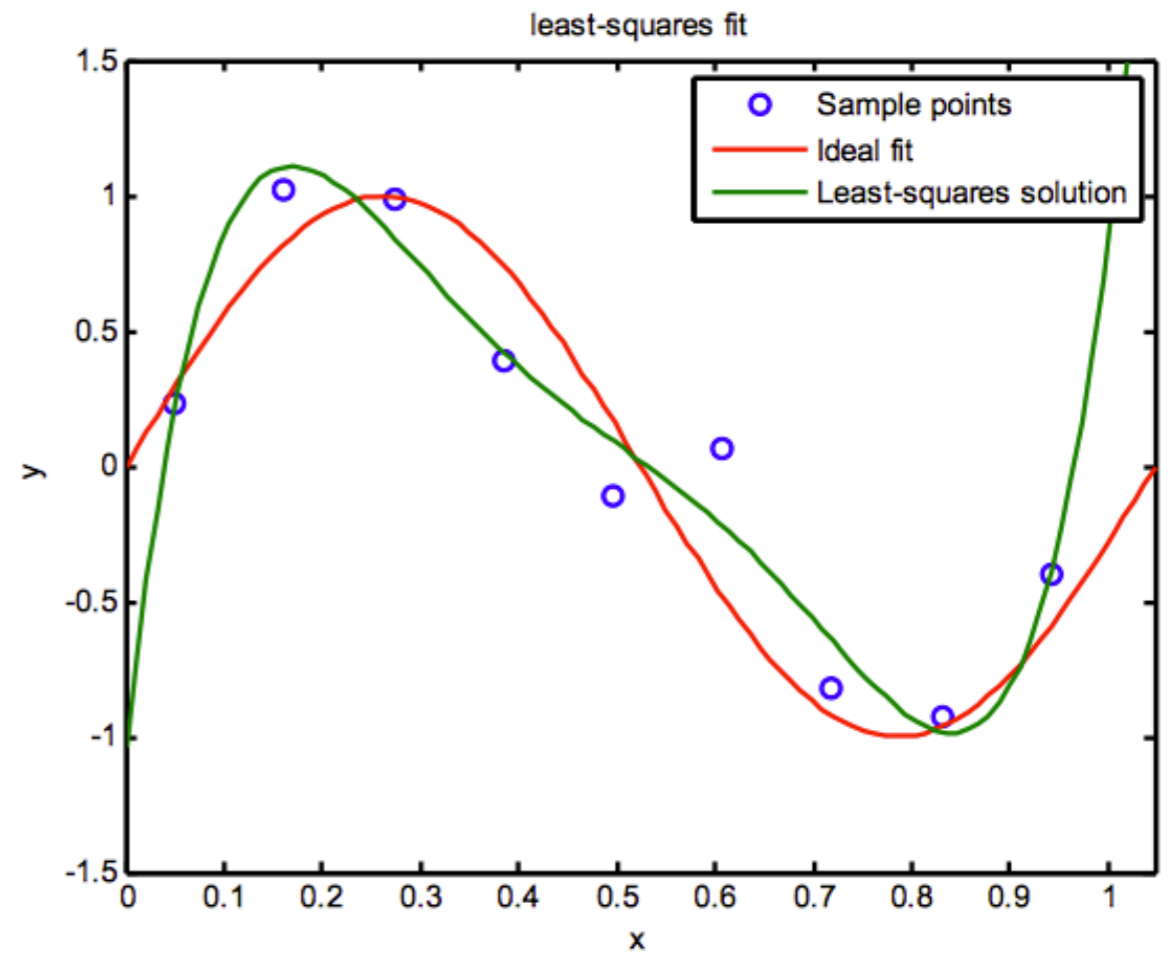
N = 9 samples, D = 7



$D = 3$



$D = 5$



Ridge Regression In Action: Predicting App Engagement

Goal:

Predict user engagement using features like:

- time on app, clicks, scroll depth
- engineered features: time^2 , time^3 (like x, x^2, x^3)

Without Regularization (Overfitting)

- Model learns large weights:
- $\theta_1 = 500, \theta_2 = -480, \theta_3 = 300$
- Small input change \rightarrow huge prediction change (user scrolls slightly more \rightarrow prediction jumps wildly)
- Model is **unstable / wiggly**
- Fits noise in training data

With Ridge Regression

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

- Weights shrink:
 $\theta_1 = 80, \theta_2 = 60, \theta_3 = 40$
- Predictions become **smooth and stable**
- Less sensitive to noise
- Better generalization

Key Takeaway

Ridge regression controls **large θ values**, preventing extreme reactions to inputs and reducing overfitting.

Outline

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

Recap

- Why do we need regularization?
- What is the root cause?
- How do we regularize in ridge regression?
- Why is the penalty term chosen as $\theta^T \theta$?

Recap

- New regularized objective function?
- Why was lambda chosen as a hyperparameter to optimize the objective function?
- Small vs high values of lambda

Recap

- Other forms of regularization covered in previous classes?

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.

Ridge versus Lasso

Ridge

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

It is a convex model

Both mean squared error and L2 regularizer are differentiable.

We can get a closed form solution

Lasso

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

It is a convex model

L1 regularizer is NOT differentiable.

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

The Lasso Regularization (L1 norm) and sparsity

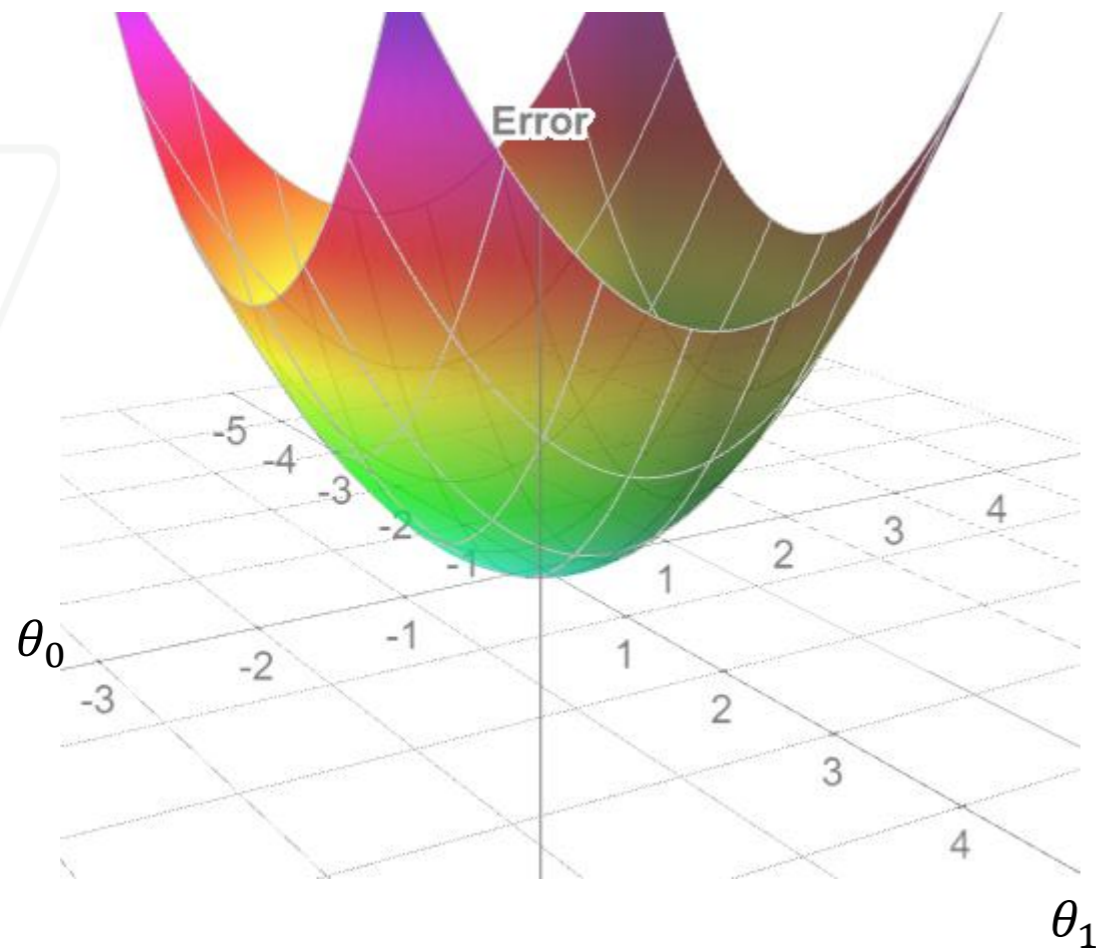
Lasso = Least Absolute Shrinkage and Selection Operator

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

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

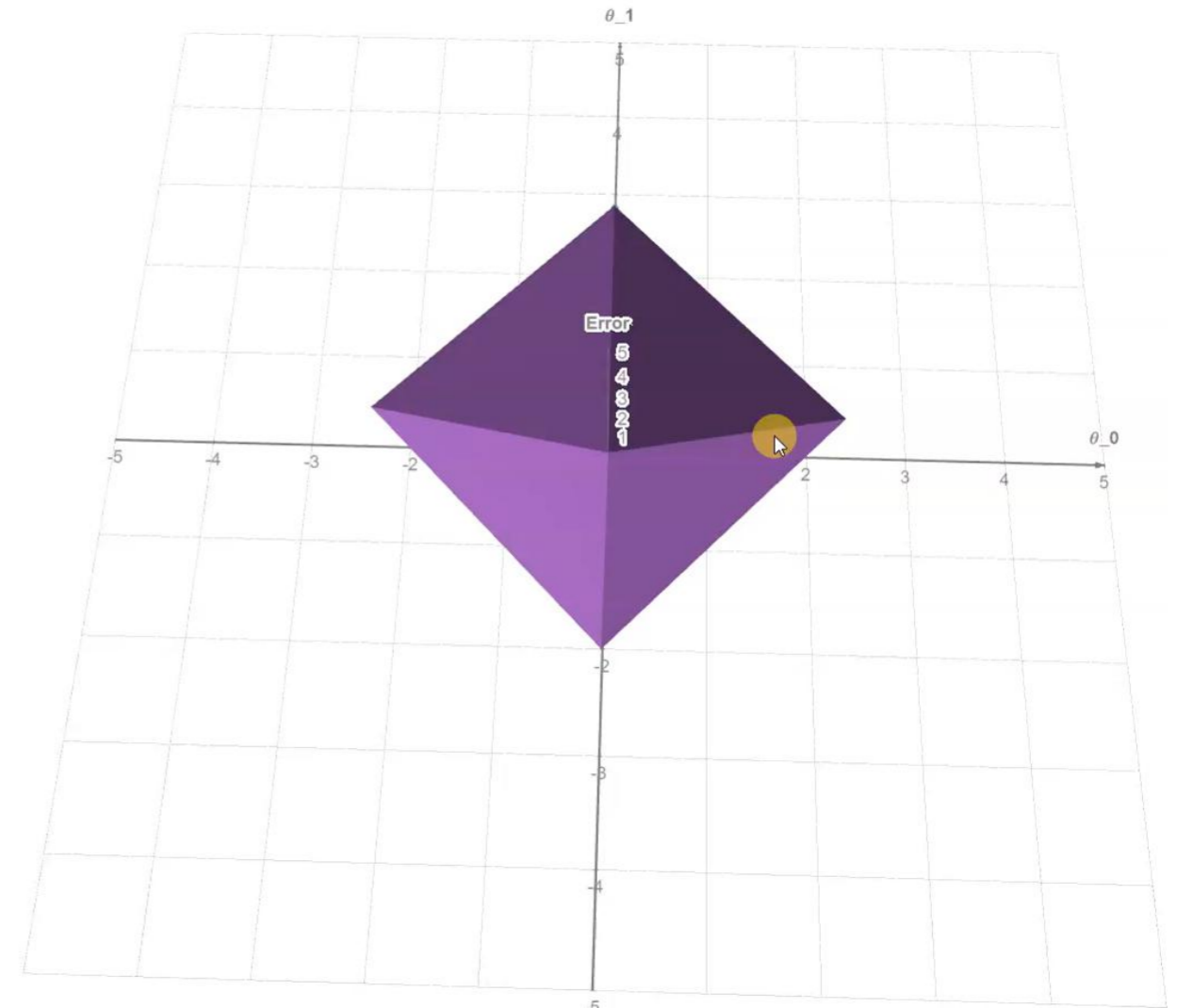
Ridge Regularizer

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



Lasso Regularizer

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



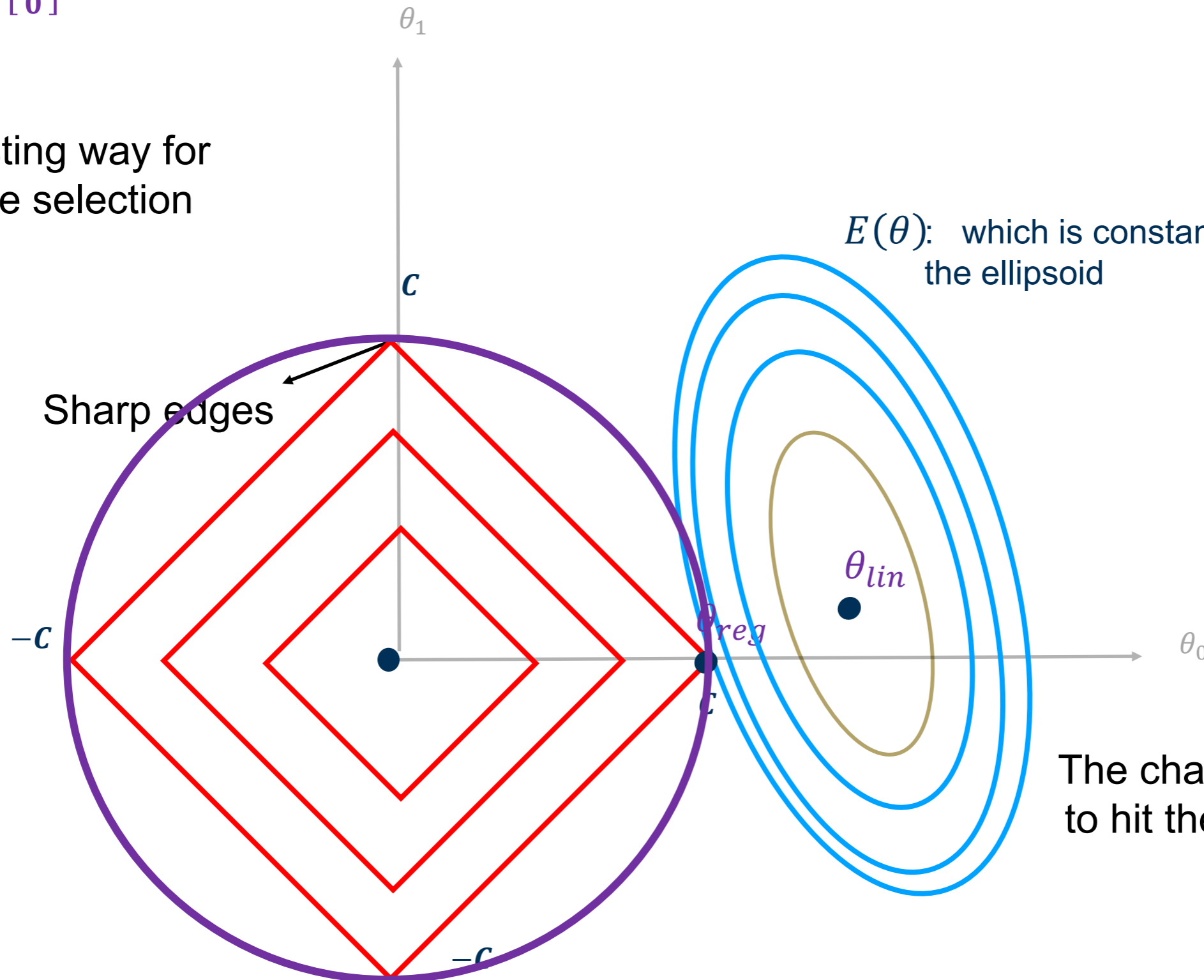
[Animation](#)

Let's say we have two parameters (θ_0 and θ_1)

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

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

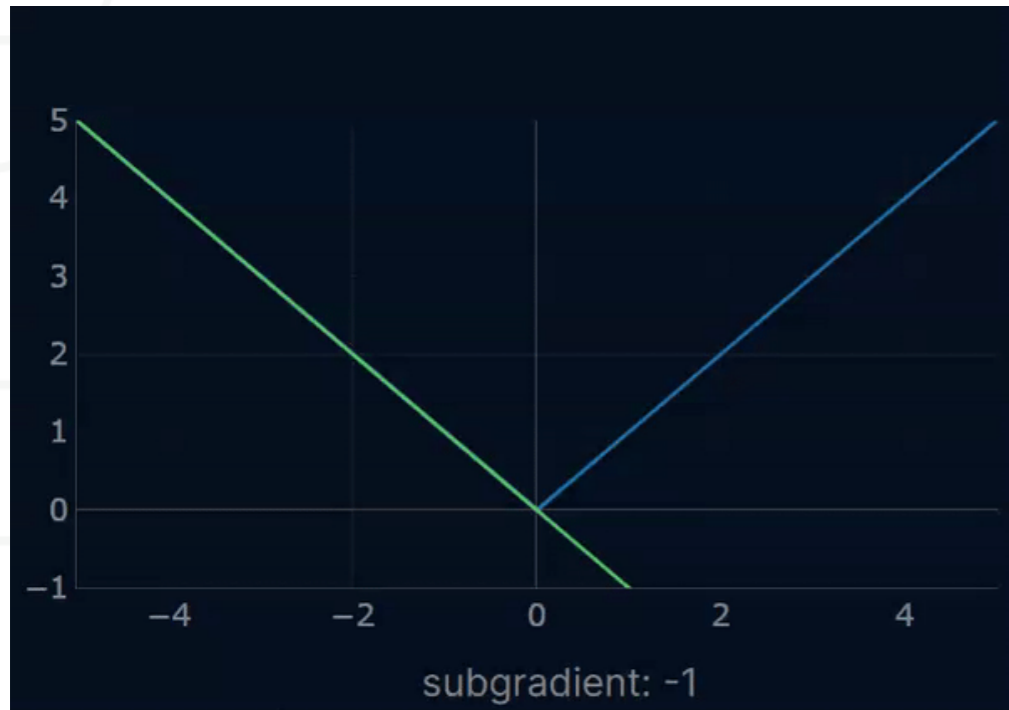
Interesting way for feature selection



$E(\theta)$: which is constant on the surface of the ellipsoid

The chance is very high to hit the sharp corners first

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)

Regularization	Bias	Variance	Total Error	Stability
 Without	Low	High	High	Overfits
 With	Slightly higher	Much lower	Lower	Generalizes better

Outline

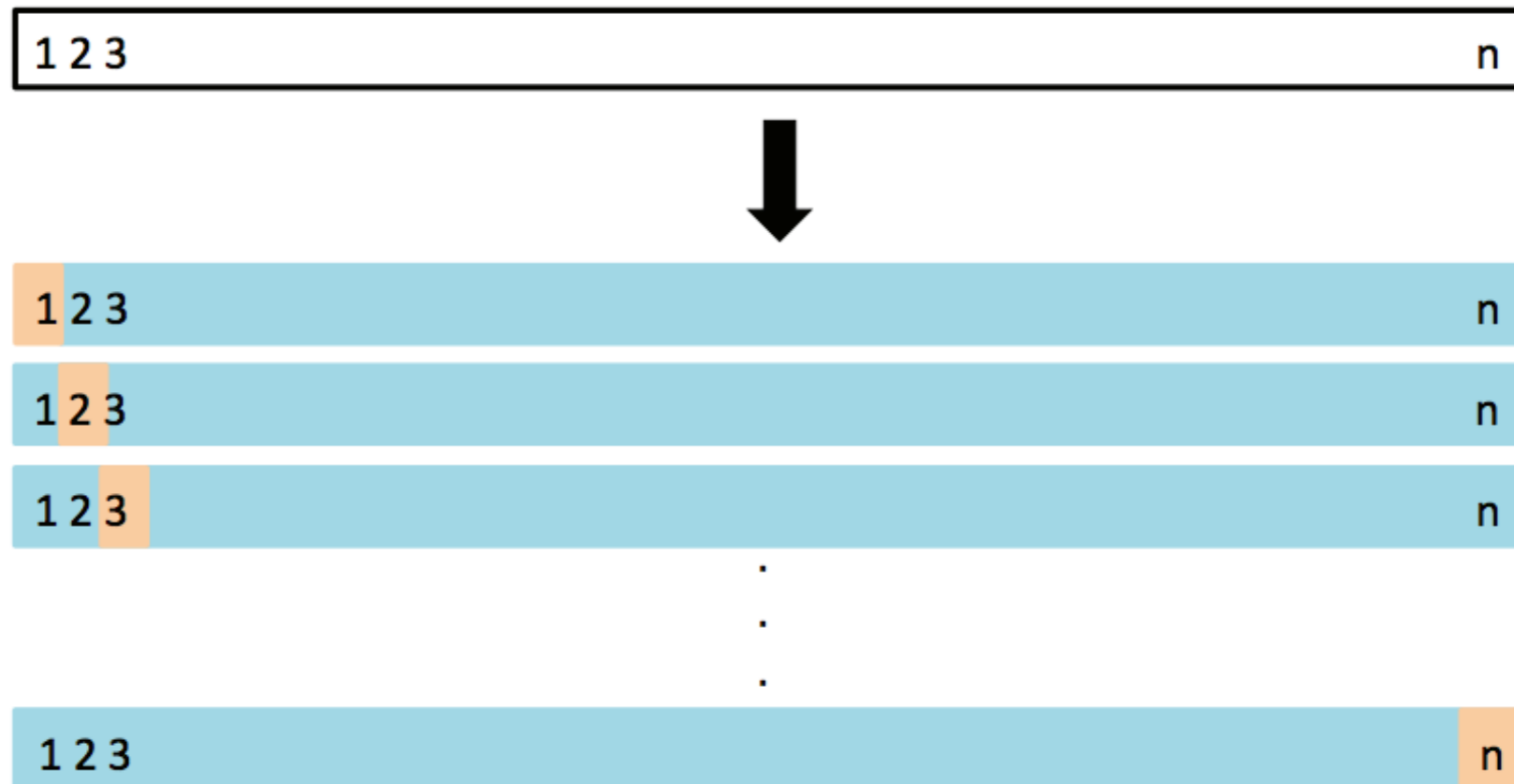
- Overfitting and regularized learning
- Ridge regression
- Lasso regression
- Determining regularization strength ←

Leave-One-Out Cross Validation

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



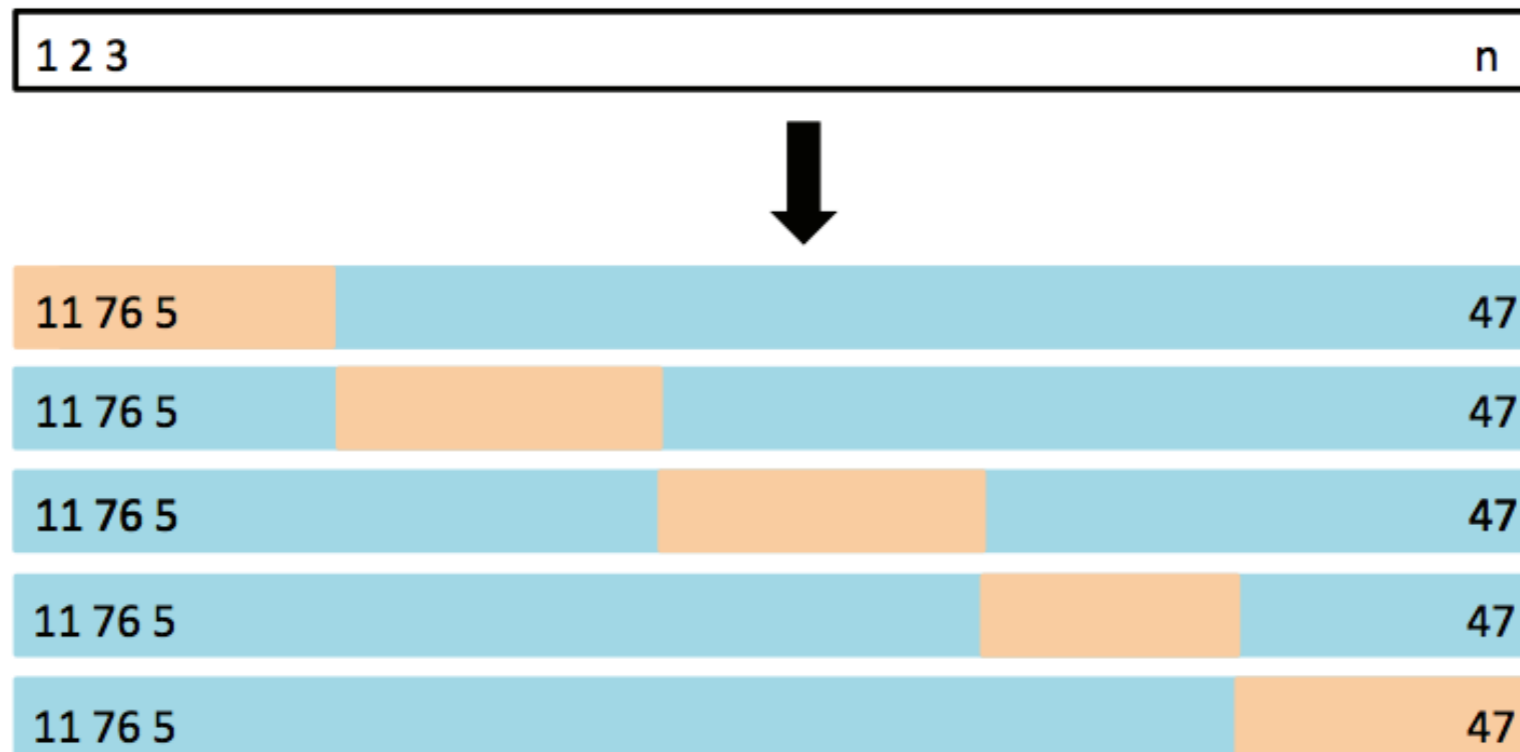
K-Fold Cross Validation

Split the data into k subsets or *folds*.

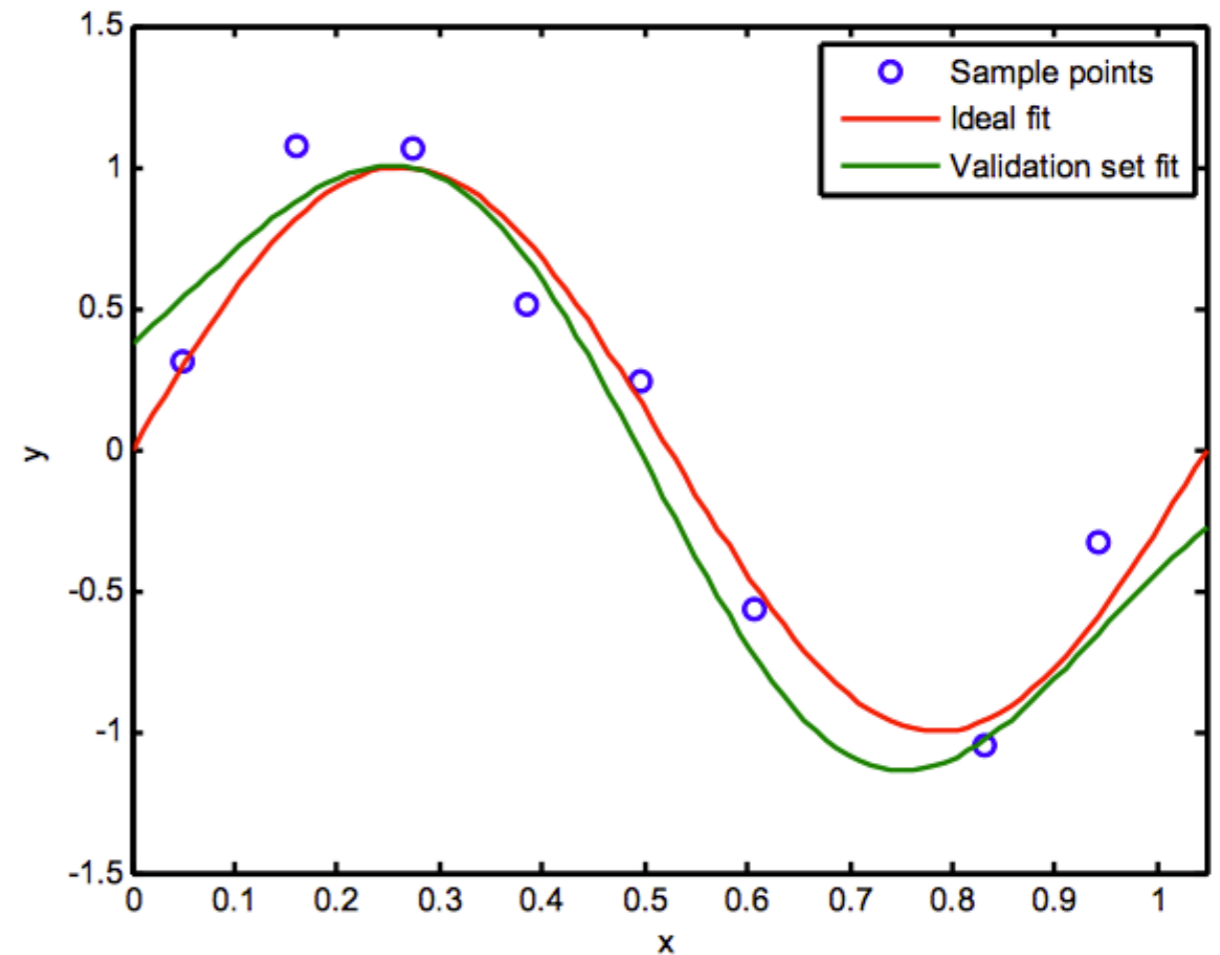
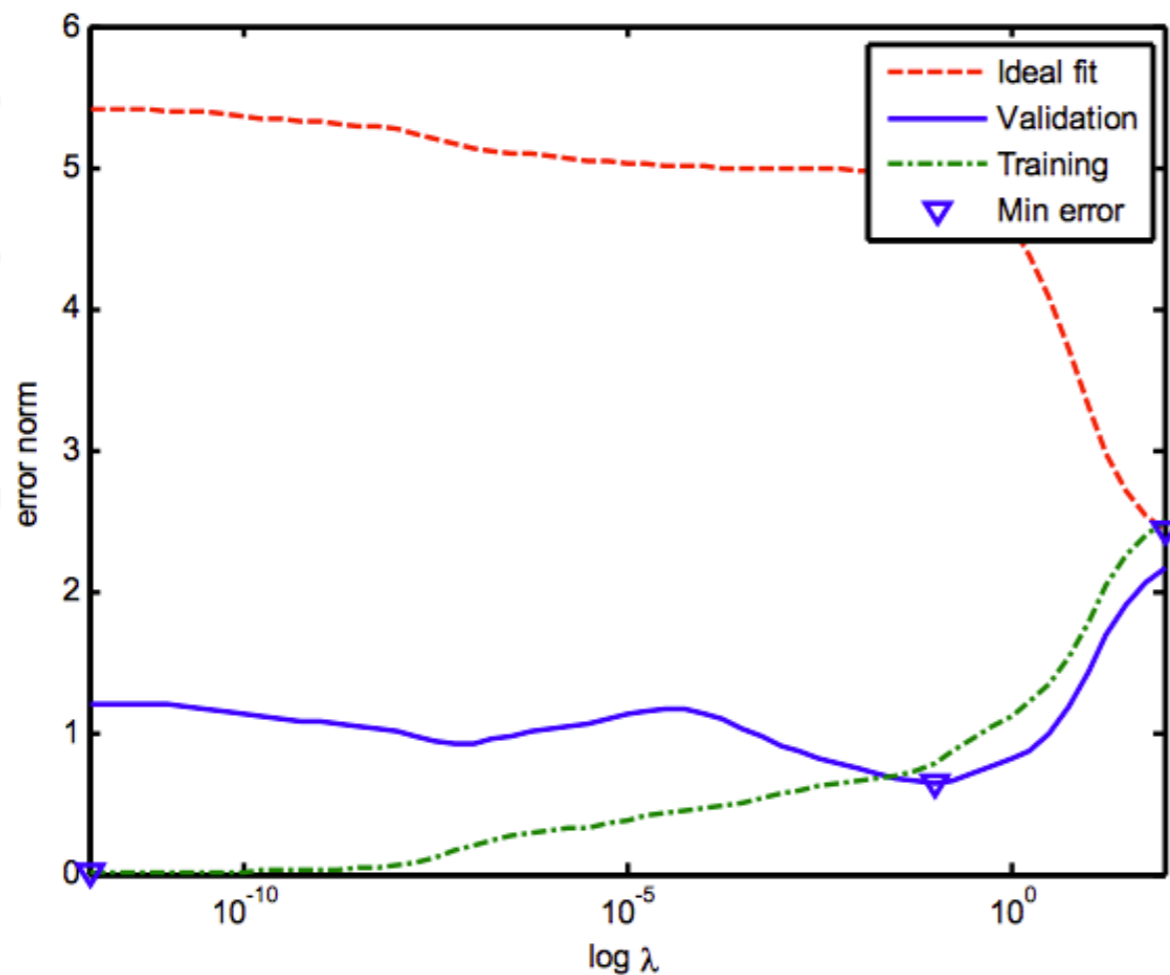
For every $i = 1, \dots, k$:

- ▶ 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 λ Using Validation Dataset



Pick up the lambda with the lowest mean value of mse 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 λ

Quick Knowledge Check

1. Which of the following is the *main reason* for overfitting in regression? A. Too few features. B. Too small a training dataset. C. Large model weights. D. High bias
2. When λ (lambda) is **very small**, what happens to the model? A. High bias, low variance. B. Low bias, high variance. C. Both bias and variance decrease. D. Both bias and variance increase
3. Which regularization technique results in **sparser** models? A. Lasso (L1). B. Ridge (L2). C. Both equally. D. Neither
4. In **k-fold cross-validation**, what does each fold provide? A. A new model parameter. B. A separate test set. C. A different λ value. D. A validation error estimate
5. Fill in the blank for ridge regression:

$$J(\theta; \lambda) = \frac{1}{2n} \|y - X\theta\|^2 + \frac{\lambda}{2} \text{---}$$

- A. $\|\theta\|$ B. $X^T \theta$ C. $\|\theta\|^2$ D. $|\theta|$