

# Regularized Linear Regression

Nakul Gopalan  
Georgia Tech

# Recap

- Linear regression:
- $Y = \theta X$  *→ Matrix*
- MSE

$$\begin{matrix} x_1 & \dots & x_n & \underline{1} \\ \theta_1 & \dots & \theta_n & \theta_{n+1} \end{matrix} \quad \text{constant}$$

$$\hat{y} = \hat{\theta} \cdot x$$

↑

$$\hat{y} \approx y$$

# Polynomial regression

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \dots$$

$$y = \theta_0 z_0 + \theta_1 z_1 + \dots$$

$$z_n = x^n$$

# Polynomial regression when order not known


$$g \propto x^{-2}$$

$x \rightarrow$  distance  
between objects

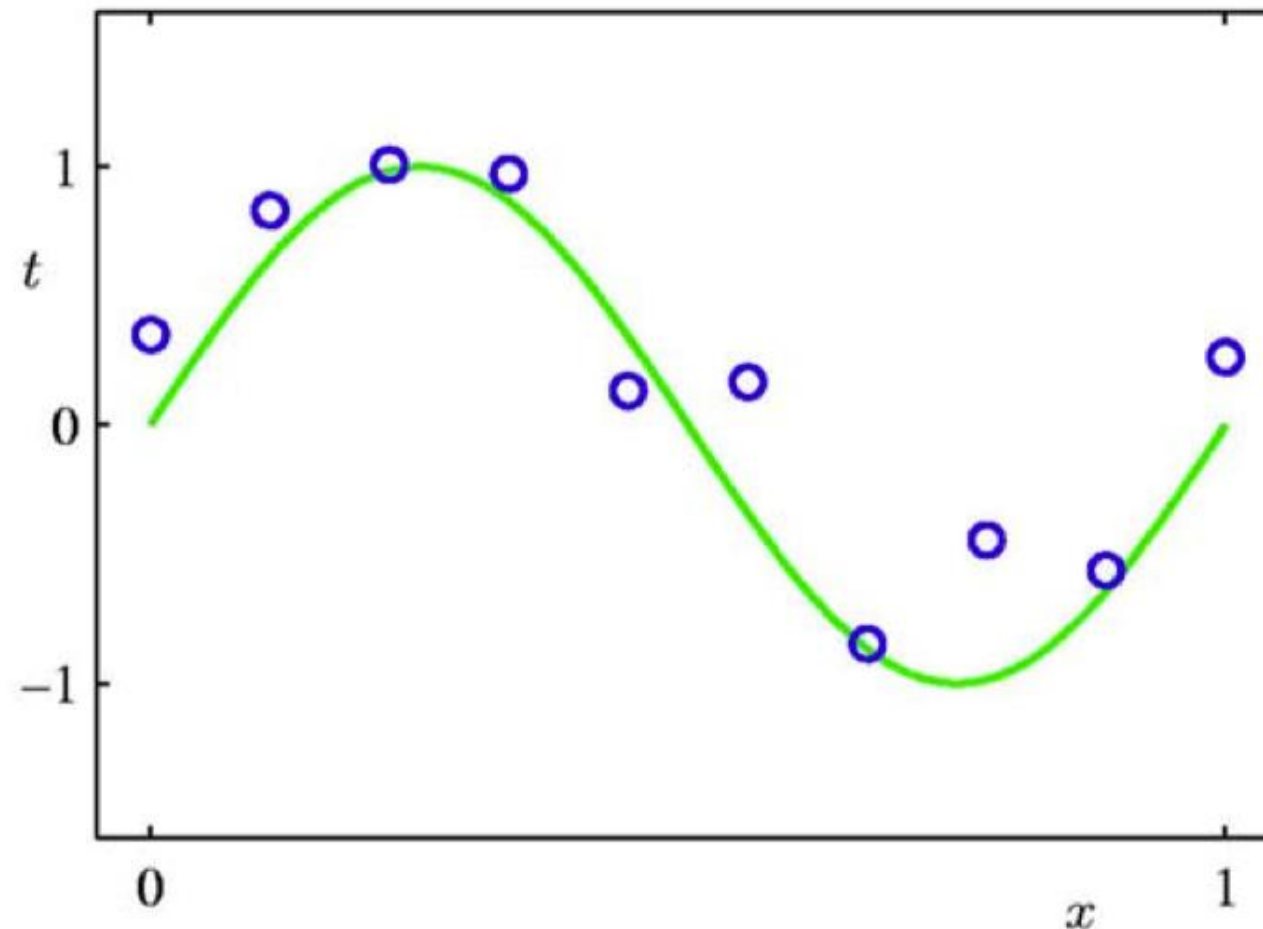
$$g = \underset{\substack{\uparrow \\ \text{Fit}}}{\theta_{-100}} x^{-100} + \theta_{-99} x^{-99} \dots + \theta_{100} x^{100}$$

$g \rightarrow$  gravitational  
force b/w them

# Outline

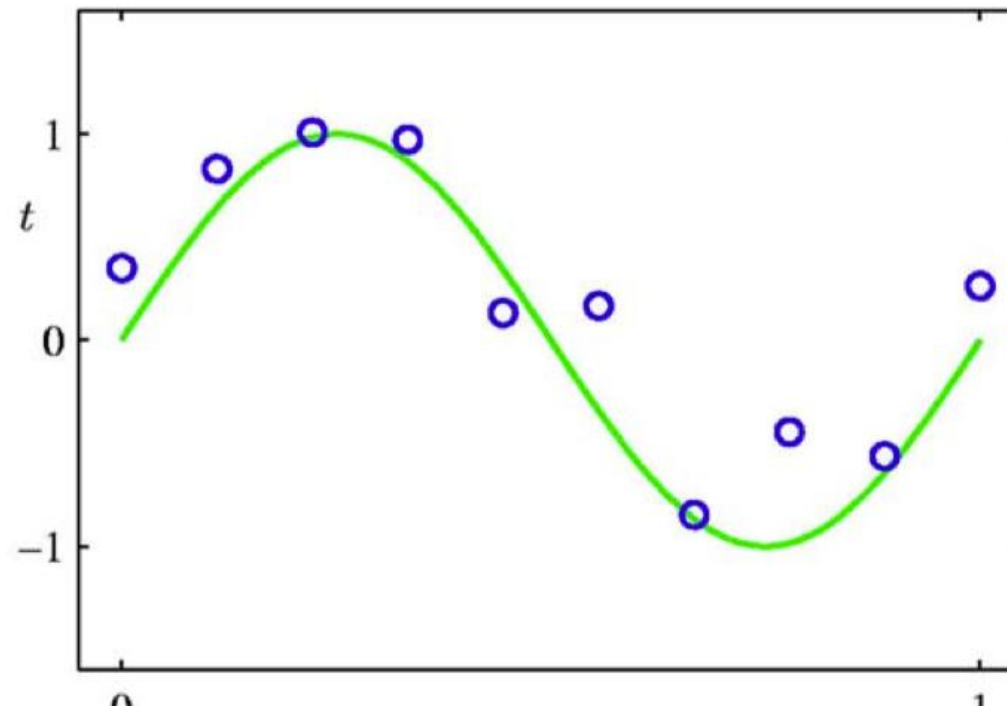
- Overfitting and regularized learning 
- Ridge regression
- Lasso regression
- Determining regularization strength

# 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

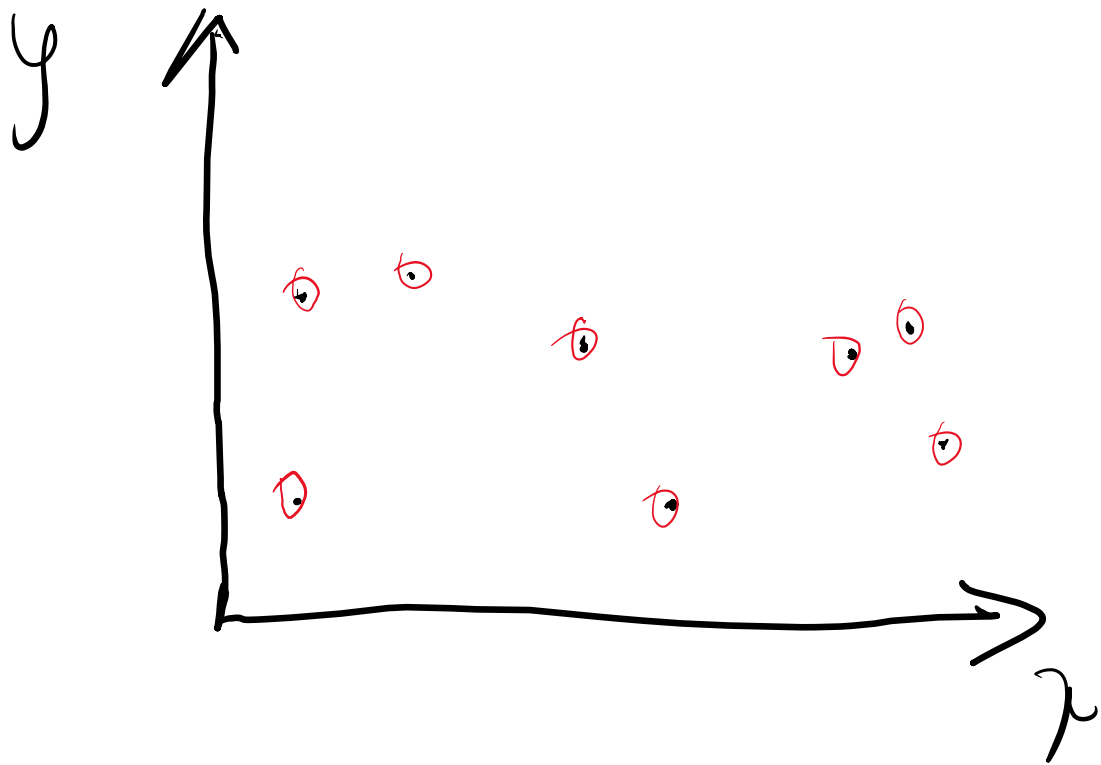
$$\underline{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_d x^d + \underline{\epsilon}$$

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

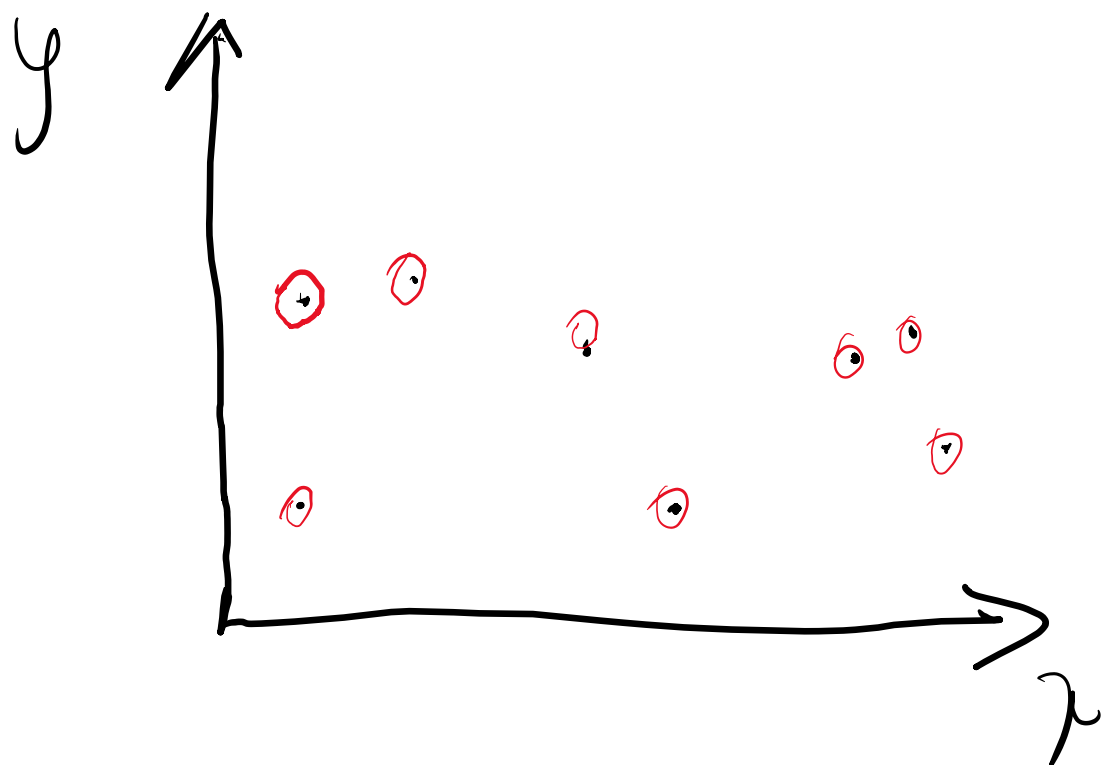
linear ↗

$$\underline{y} = \underline{z\theta}$$

Real  
Regression  
Problem!!!

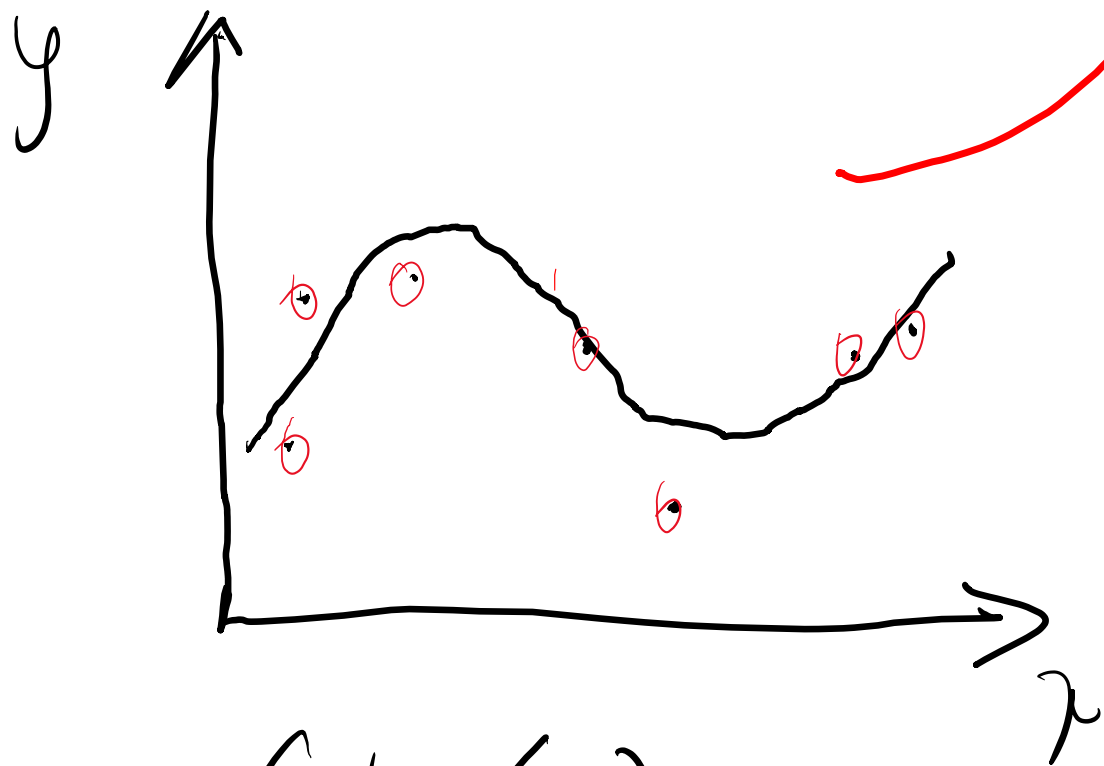
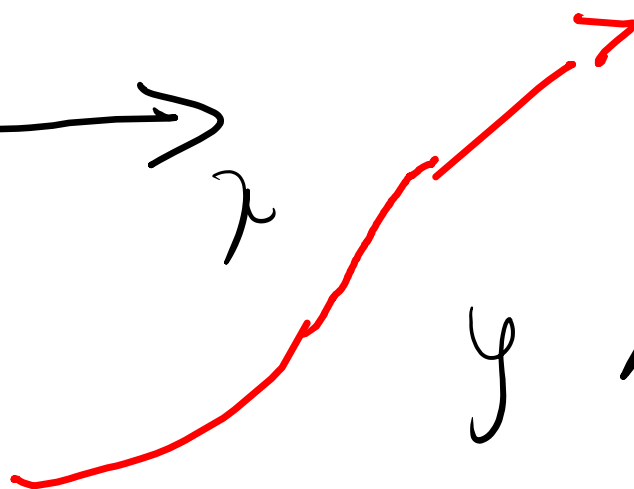




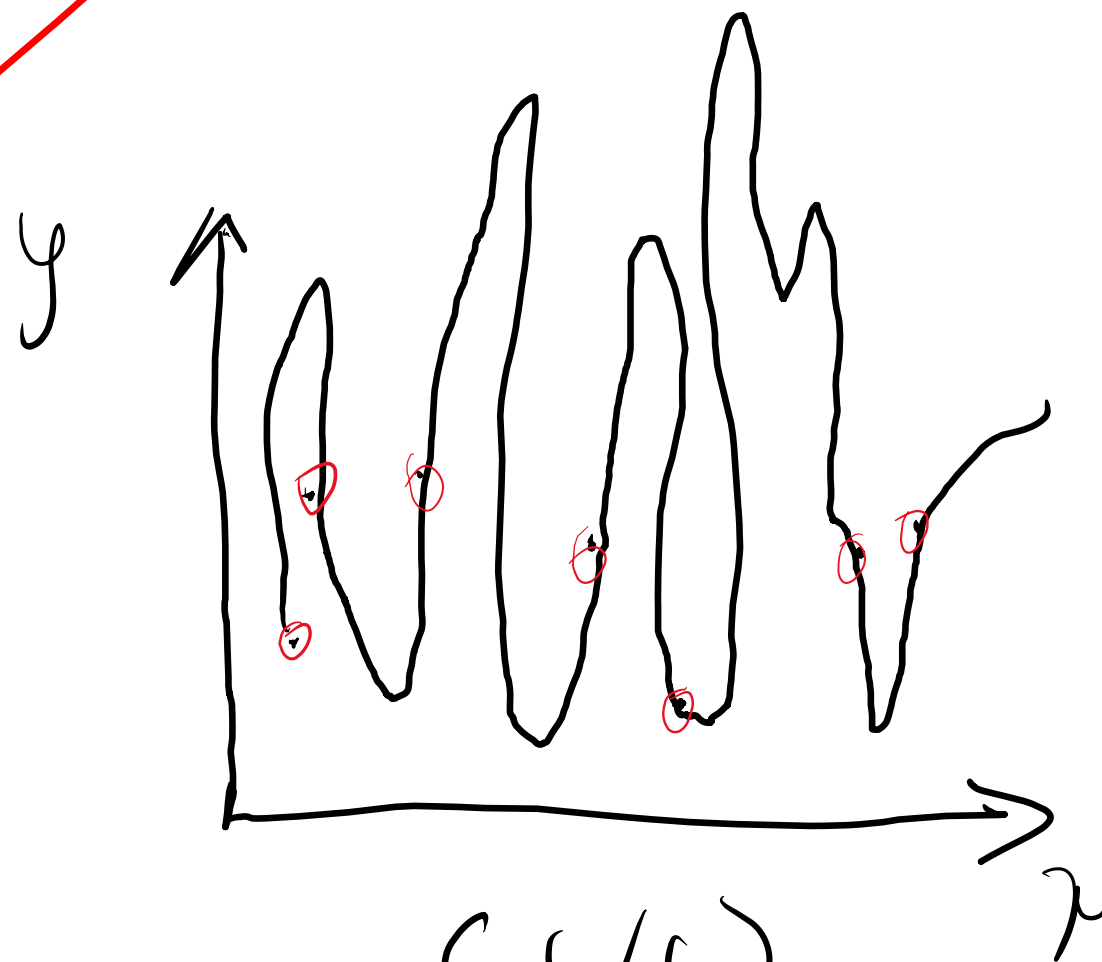


generalization

because no overfitting

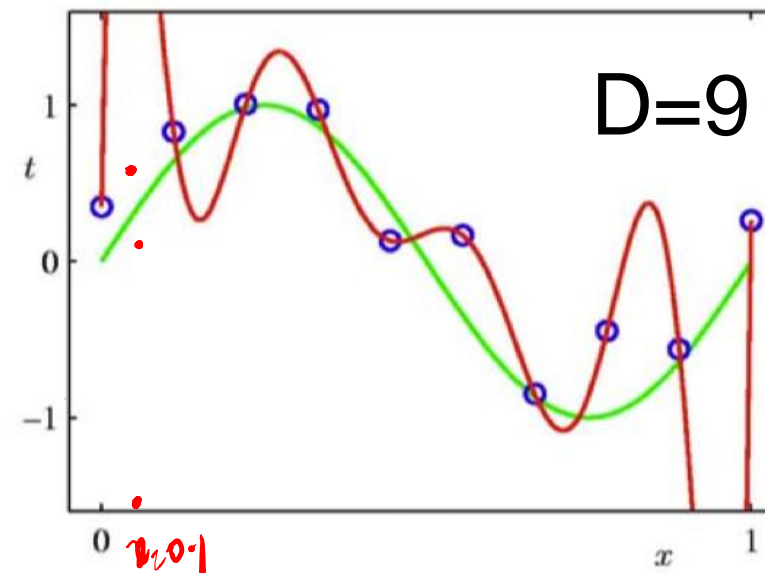
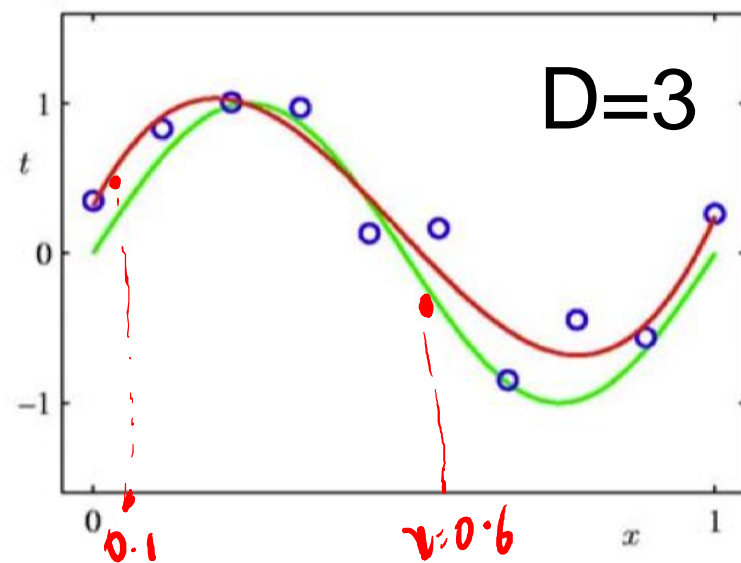
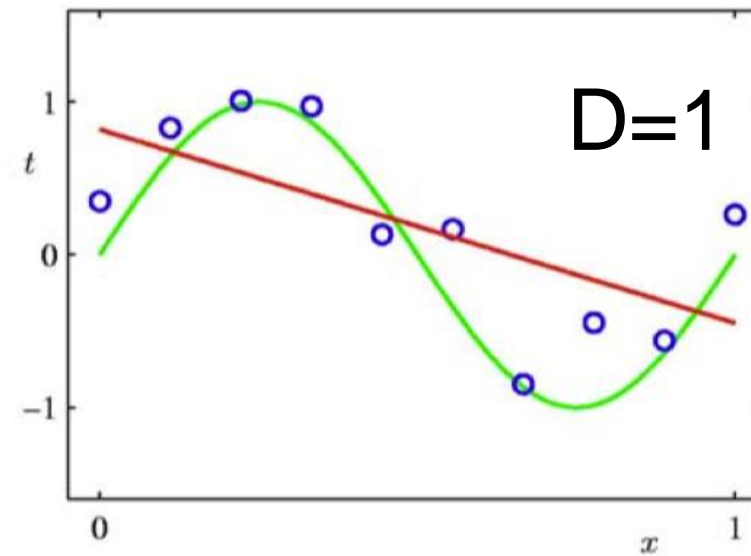
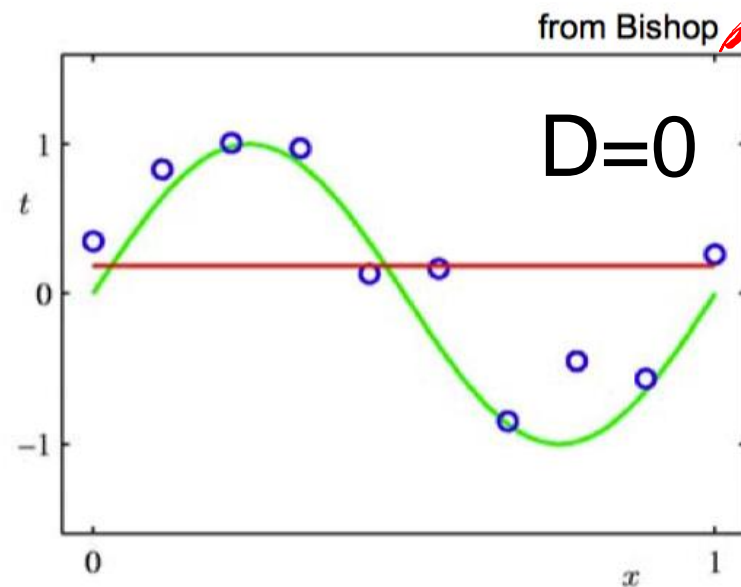


Sol. (a)



Sol. (b)

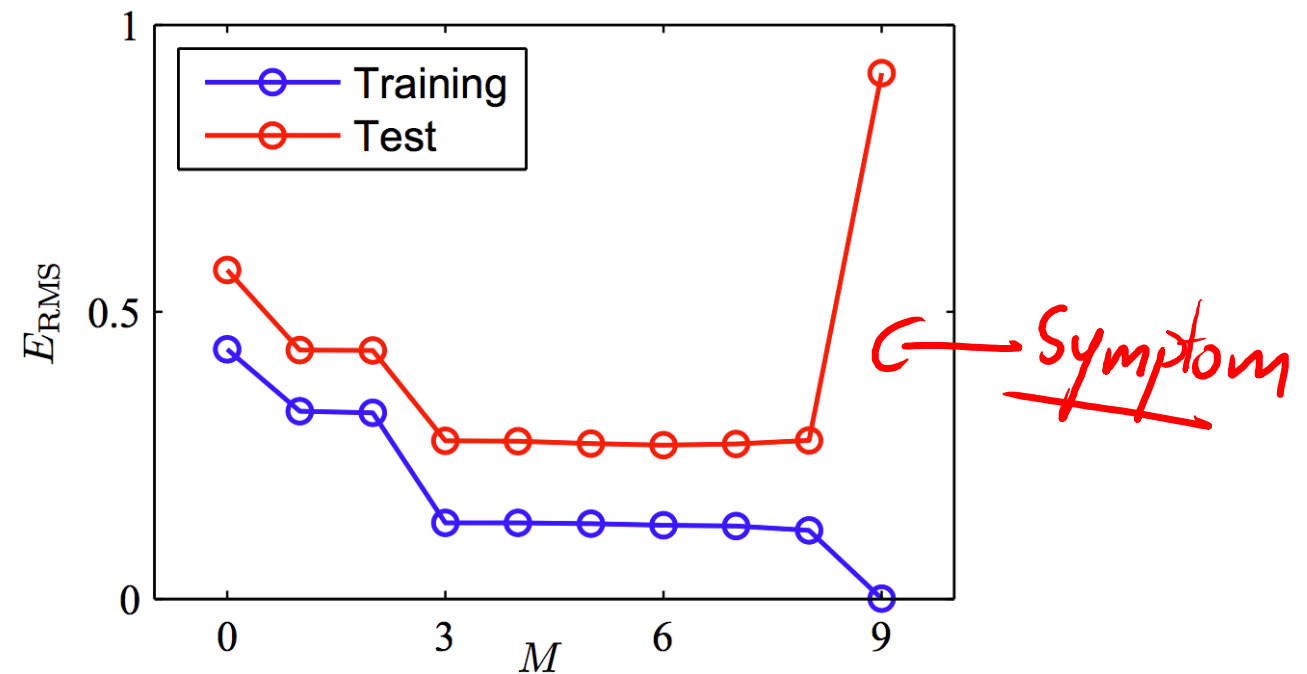
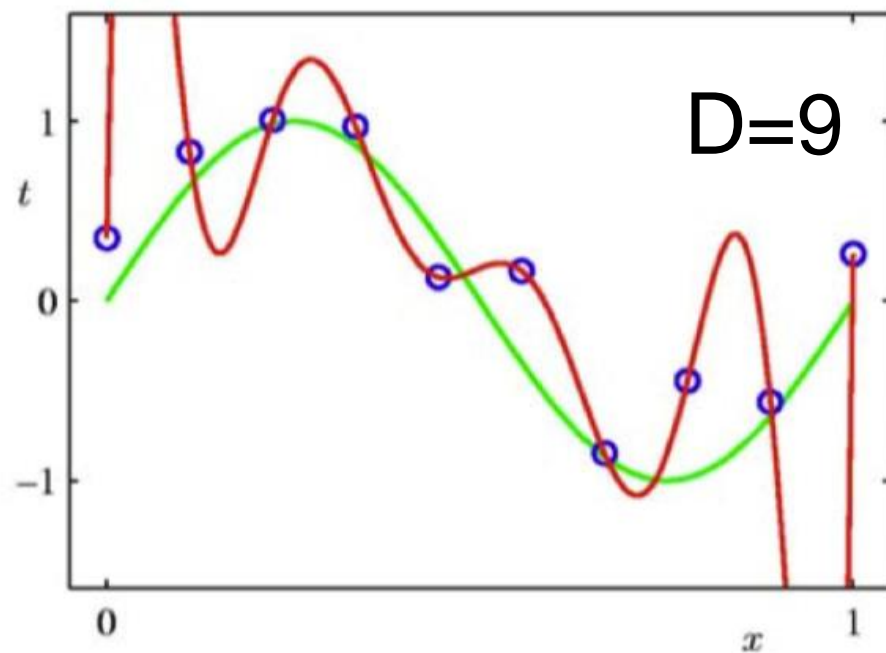
# Which One is Better?



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

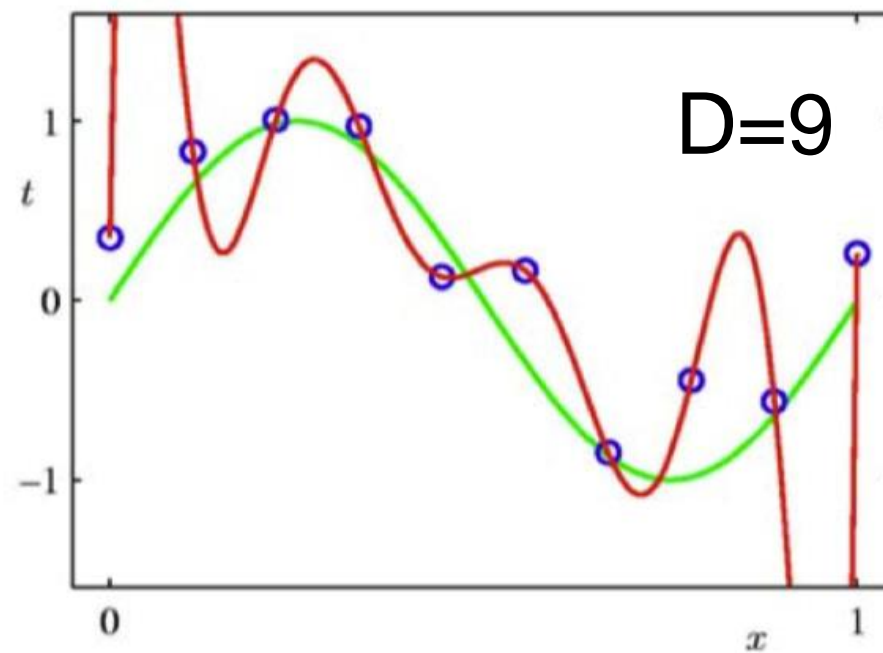
No, this can lead to overfitting!

# 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.

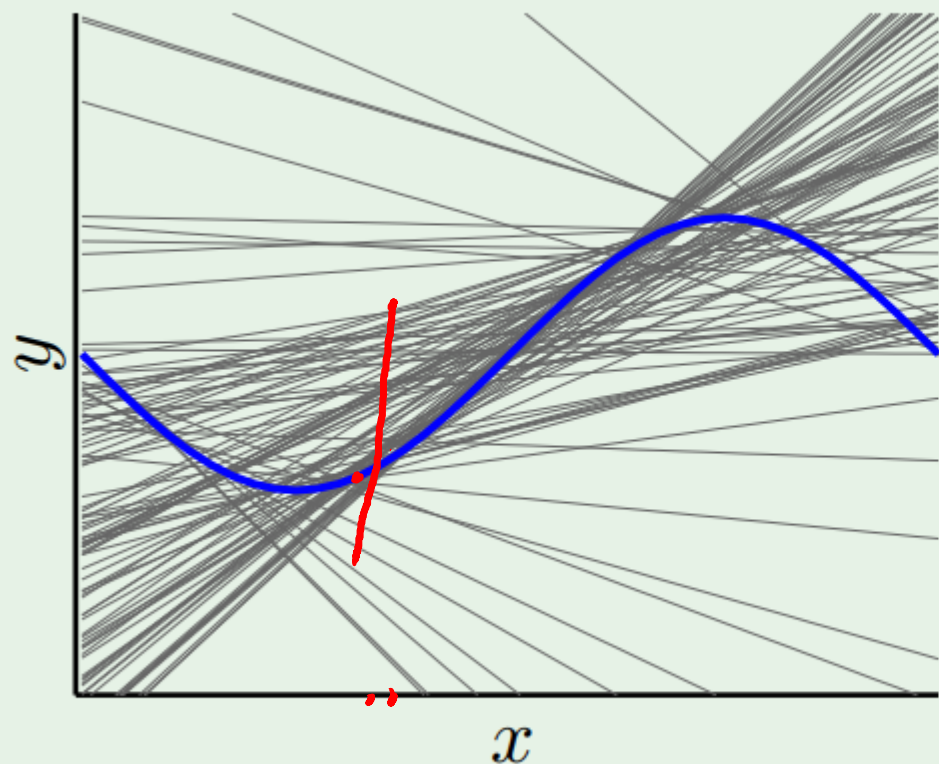
# The Overfitting Problem



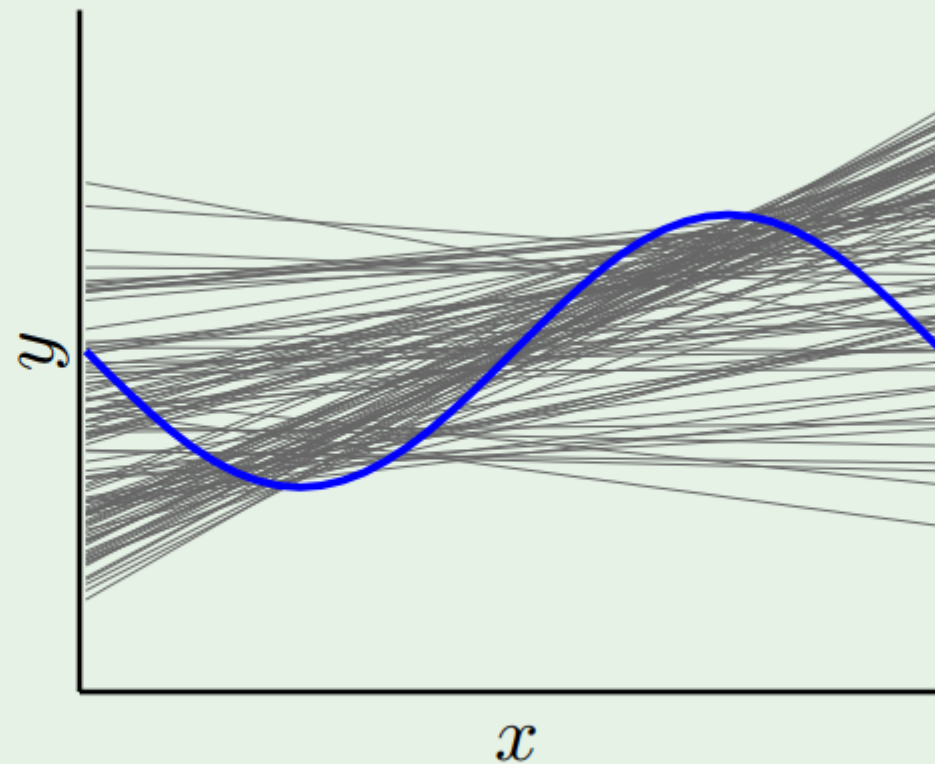
- In regression, overfitting is often associated with large Weights (**severe oscillation**)
- How can we address overfitting?

# Regularization

(smart way to cure overfitting disease )



without regularization



with regularization

$$y(0.005) = 10$$

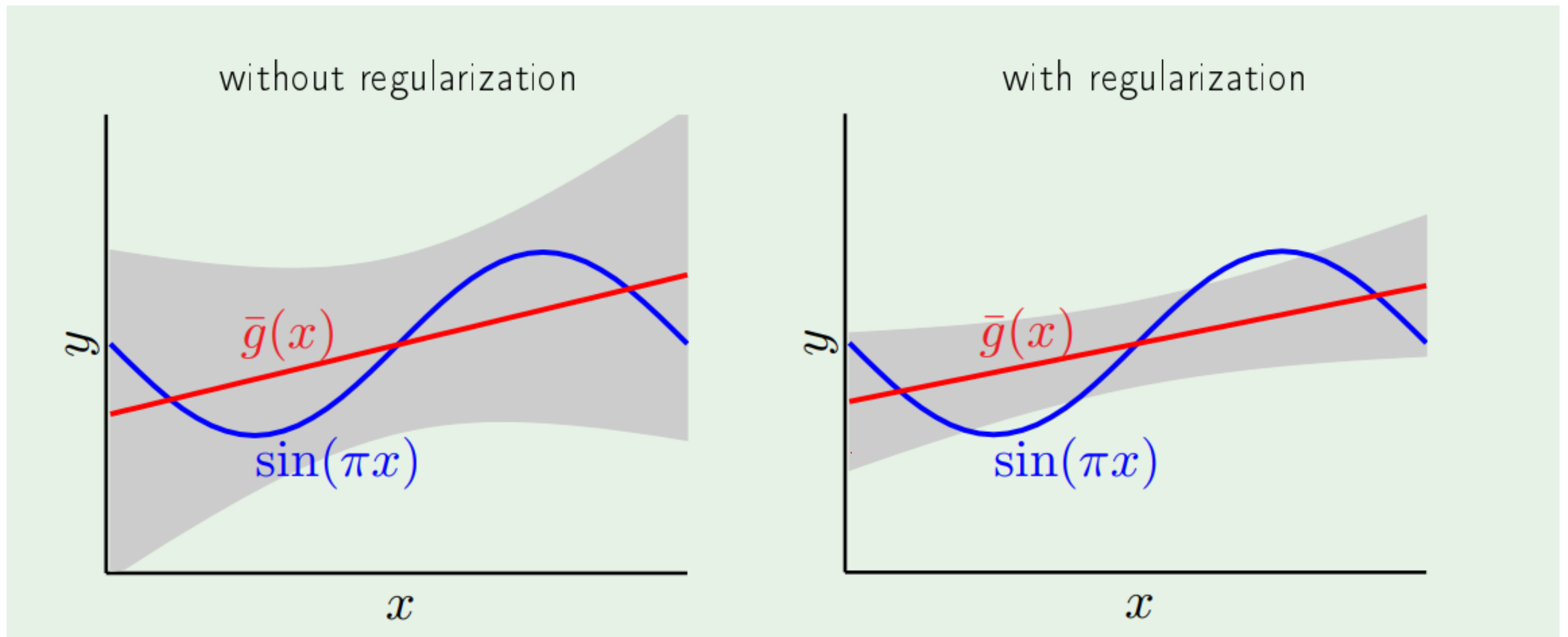
$$y(0.0051) = 10,000$$

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



bias=0.21; var=1.69

bias=0.23; var=0.33

# Regularized Learning

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

Why this term leads to regularization of parameters

- Cost function – squared loss:

$$\underline{\tilde{E}}(\underline{\theta}) = \frac{1}{N} \sum_{i=1}^N \underbrace{\{f(x_i, \underline{\theta}) - \overset{\text{target value}}{y_i}\}^2}_{\text{loss function}} + \underbrace{\frac{\lambda}{N} \|\underline{\theta}\|^2}_{\text{regularization}}$$

# 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:

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

$$y = \mathbf{z}\boldsymbol{\theta}$$



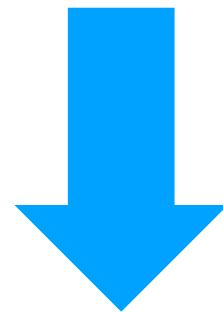
# Regularizing is just constraining the weights ( $\theta$ )

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

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



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

Let's not penalize  $\theta$  in such a harsh way  
let's cut them some slack

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

$$\text{Minimize } \frac{1}{N} (z\theta - y)^T (z\theta - y)$$

$$\text{Subject to } \theta^t \theta \leq c$$

↑  
matrix form

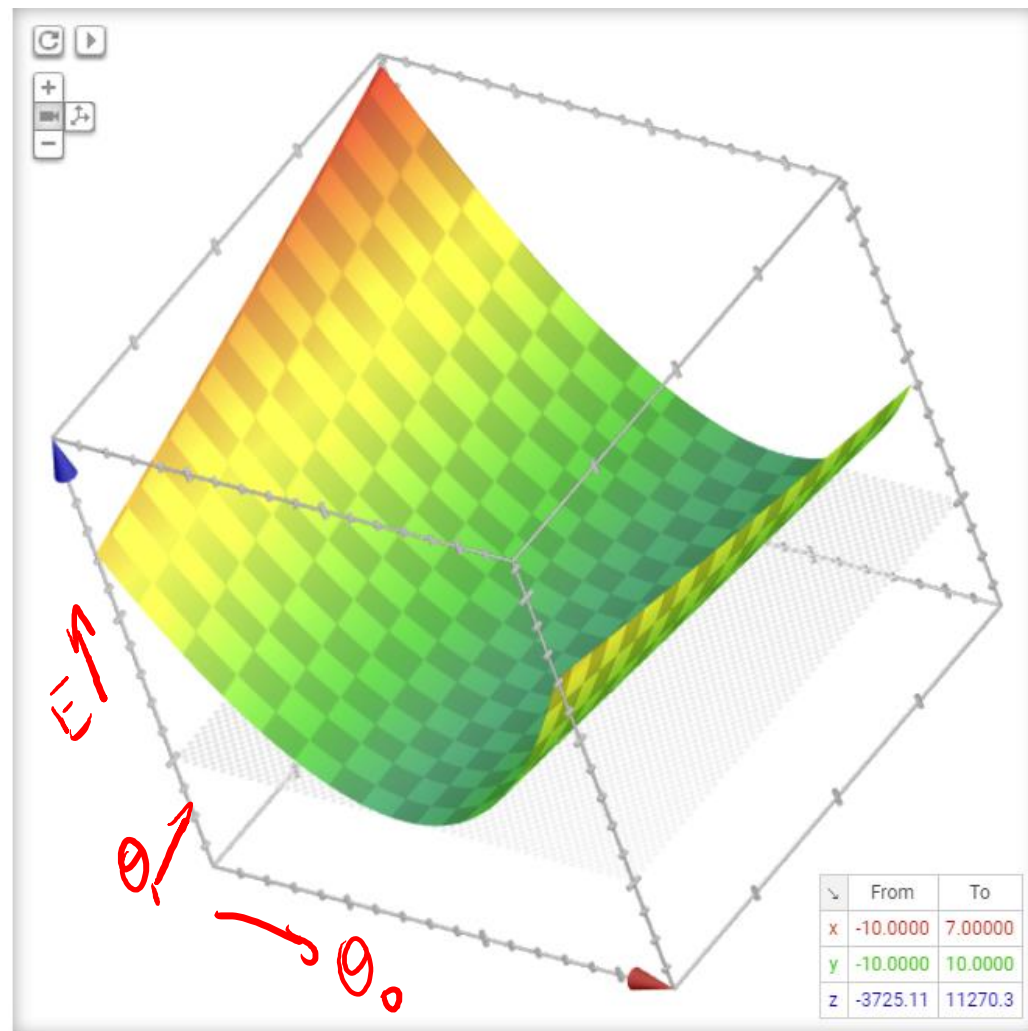
$$\sum \theta_i^2 \leq c$$

↑  
constant

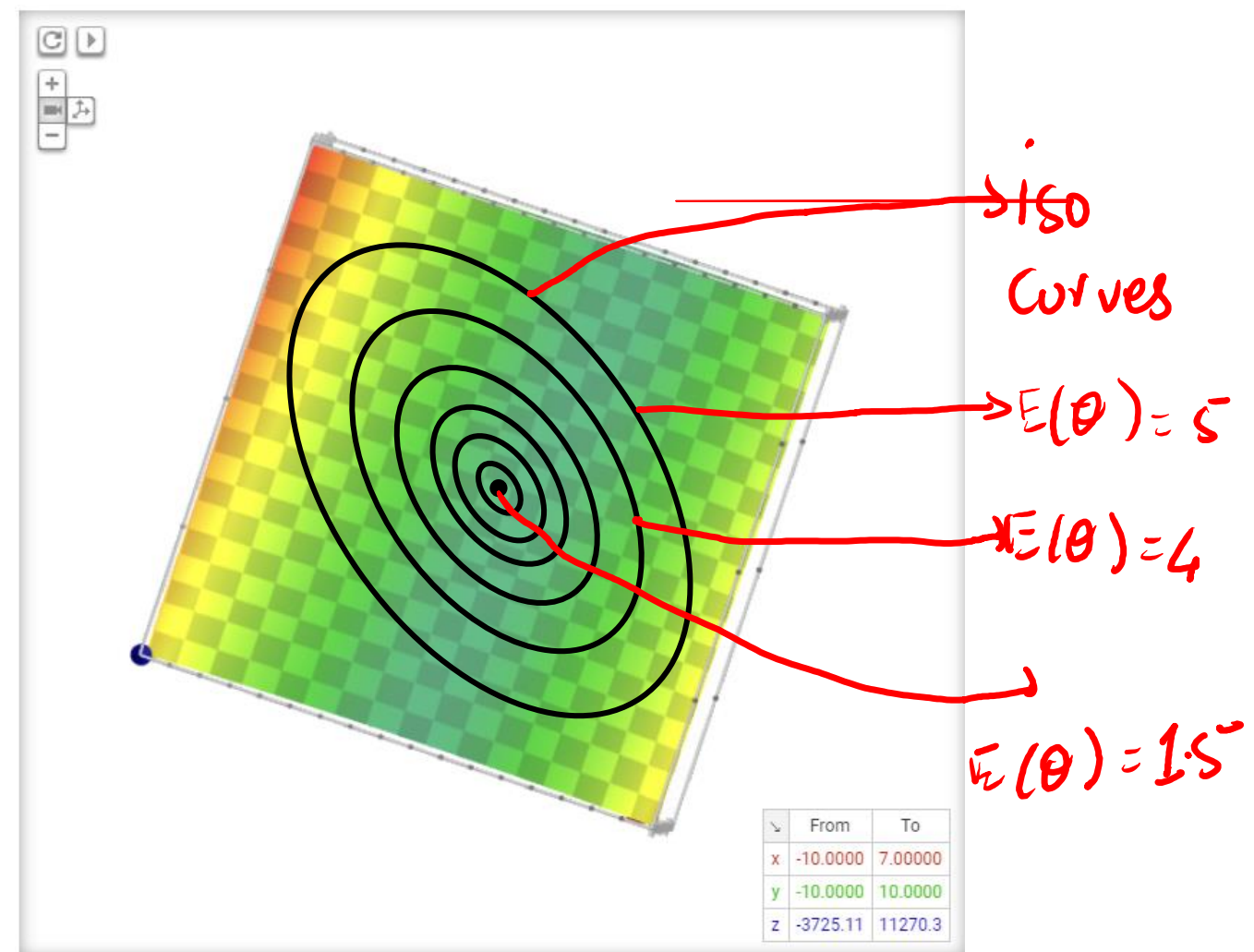
For simplicity let's call  $\theta_{lin}$  as weights' solution for non constrained one  
and  $\theta$  for the constrained model.

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

Possible graph for  $E(\theta)$  for different values of  $\theta_0$  and  $\theta_1$  and given observation data



3D view



Top view

# Gradient of $\theta^T \theta$

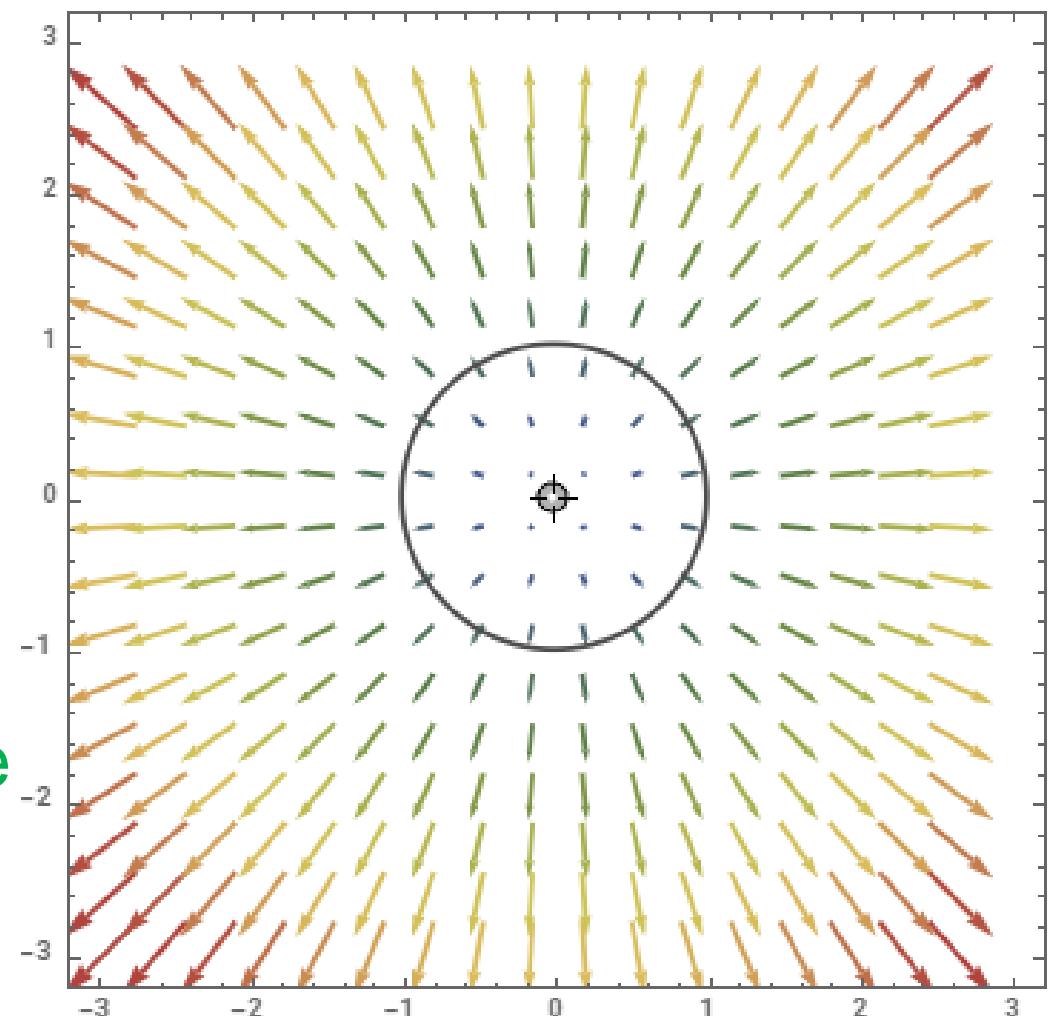
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \Rightarrow \theta^T \theta = \theta_0^2 + \theta_1^2$$

If you imagine standing at a point  $(\theta_0, \theta_1)$ ,  
 $\nabla(\theta^T \theta)$  tells you which direction you should  
travel to increase the value of  $\theta^T \theta$  most rapidly.

$$\nabla(\theta^T \theta) = \begin{bmatrix} \frac{\partial}{\partial(\theta_0)} (\theta^T \theta) \\ \frac{\partial}{\partial(\theta_1)} (\theta^T \theta) \end{bmatrix} = \begin{bmatrix} 2\theta_0 \\ 2\theta_1 \end{bmatrix} \approx \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

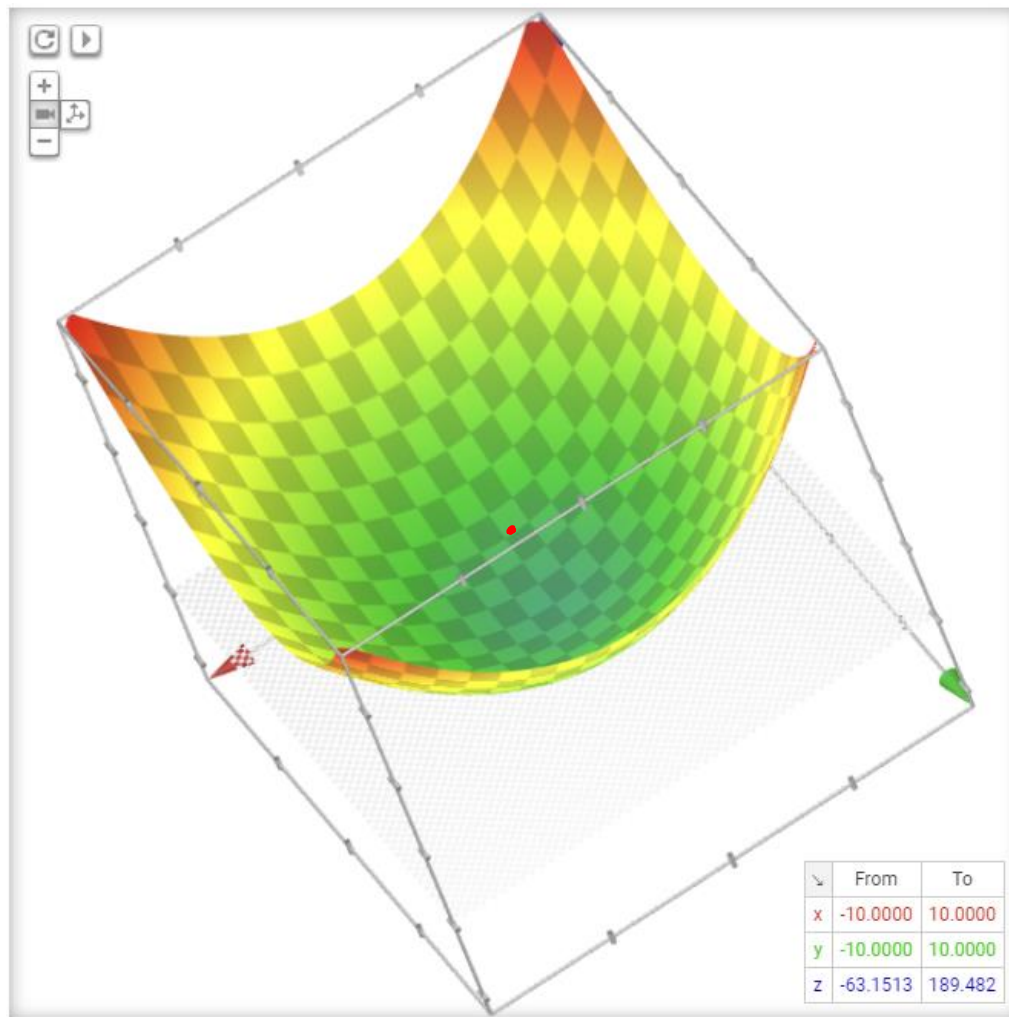
$\nabla(\theta^T \theta)$  is a vector field

any line passing through the center of the  
circle

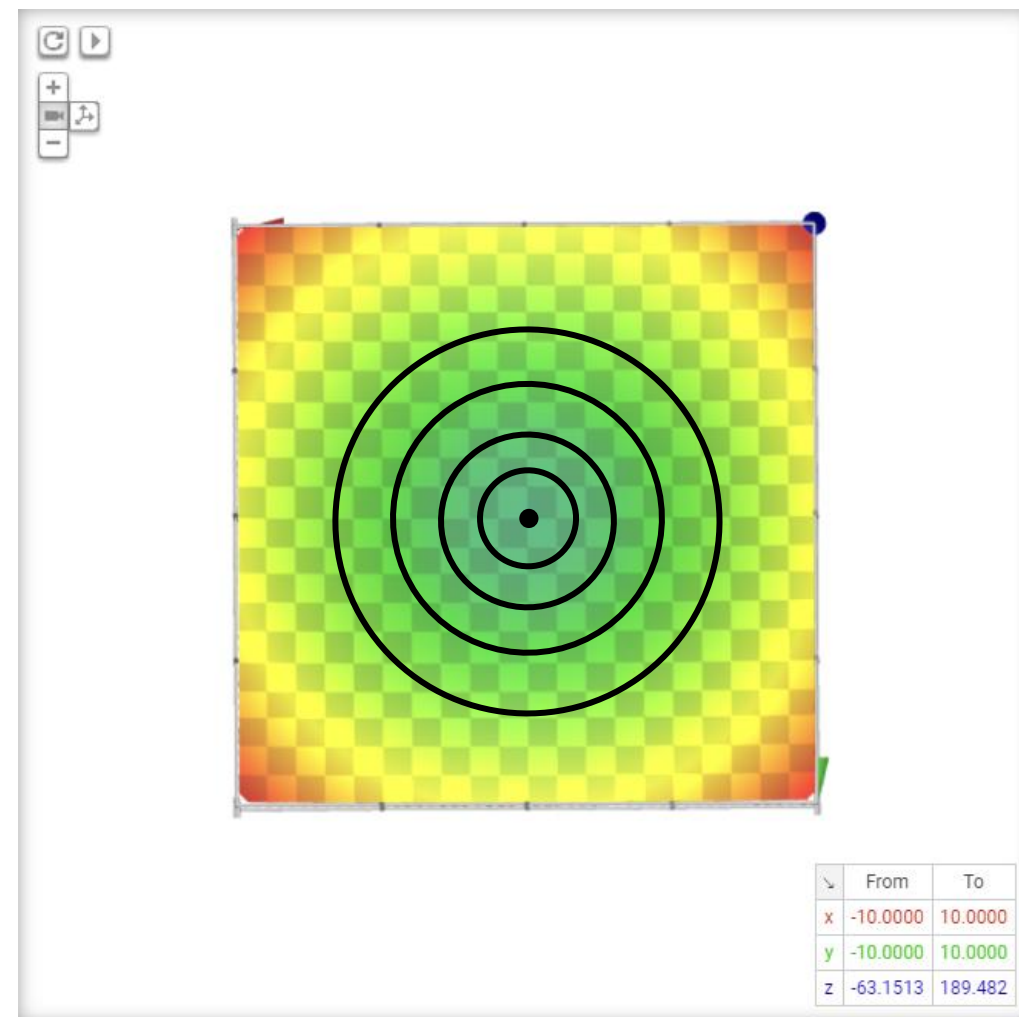


# Plotting the regularization term $\theta^t \theta$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \Rightarrow \theta^t \theta = \underline{\theta_0^2} + \theta_1^2 \rightarrow \text{eq. for a circle}$$



3D view



Top view

$\Sigma$

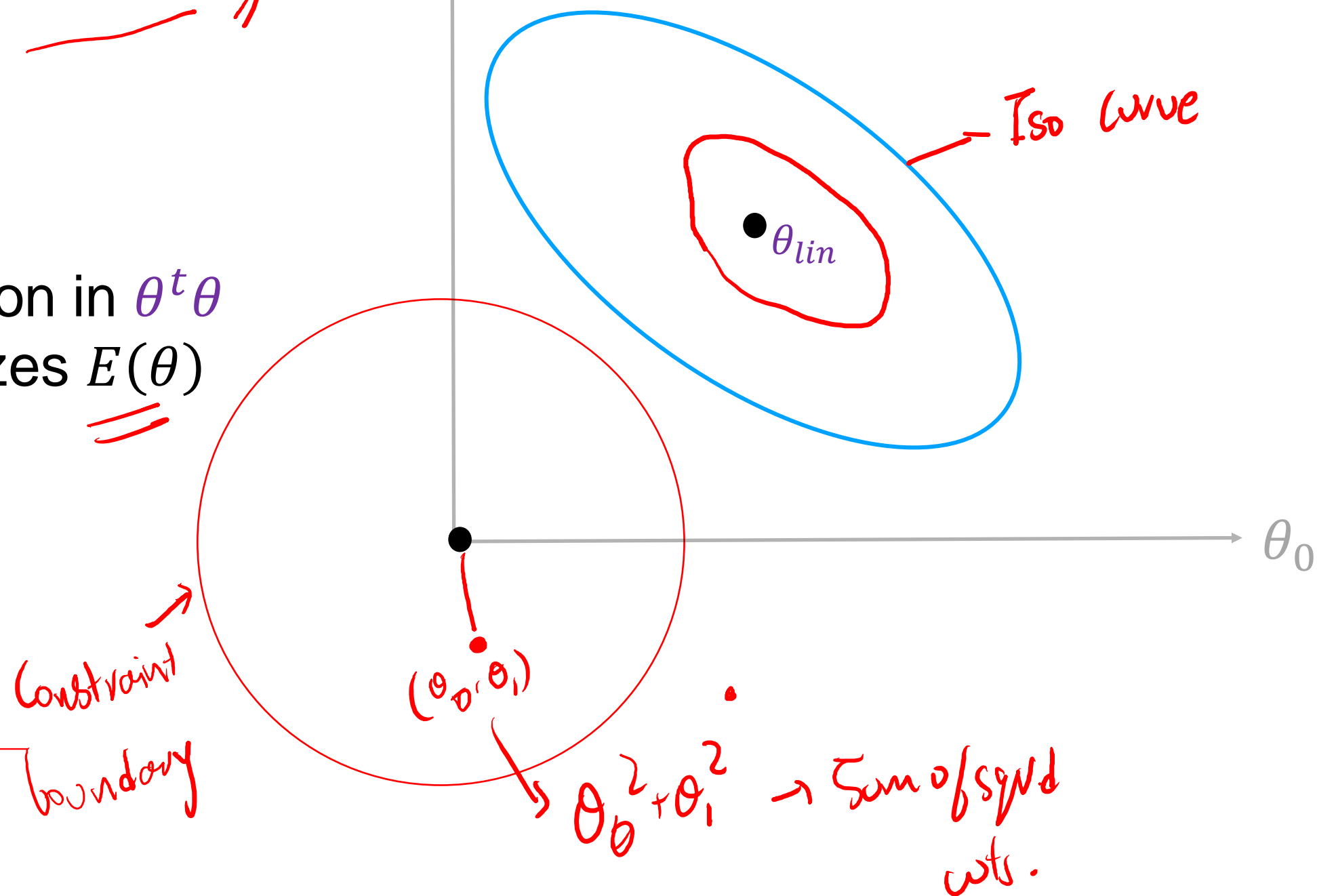
$$E(\theta) = \frac{1}{N} (z\theta - y)^T (z\theta - y)$$

$\theta_{lin}$  is the solution (min absolute)

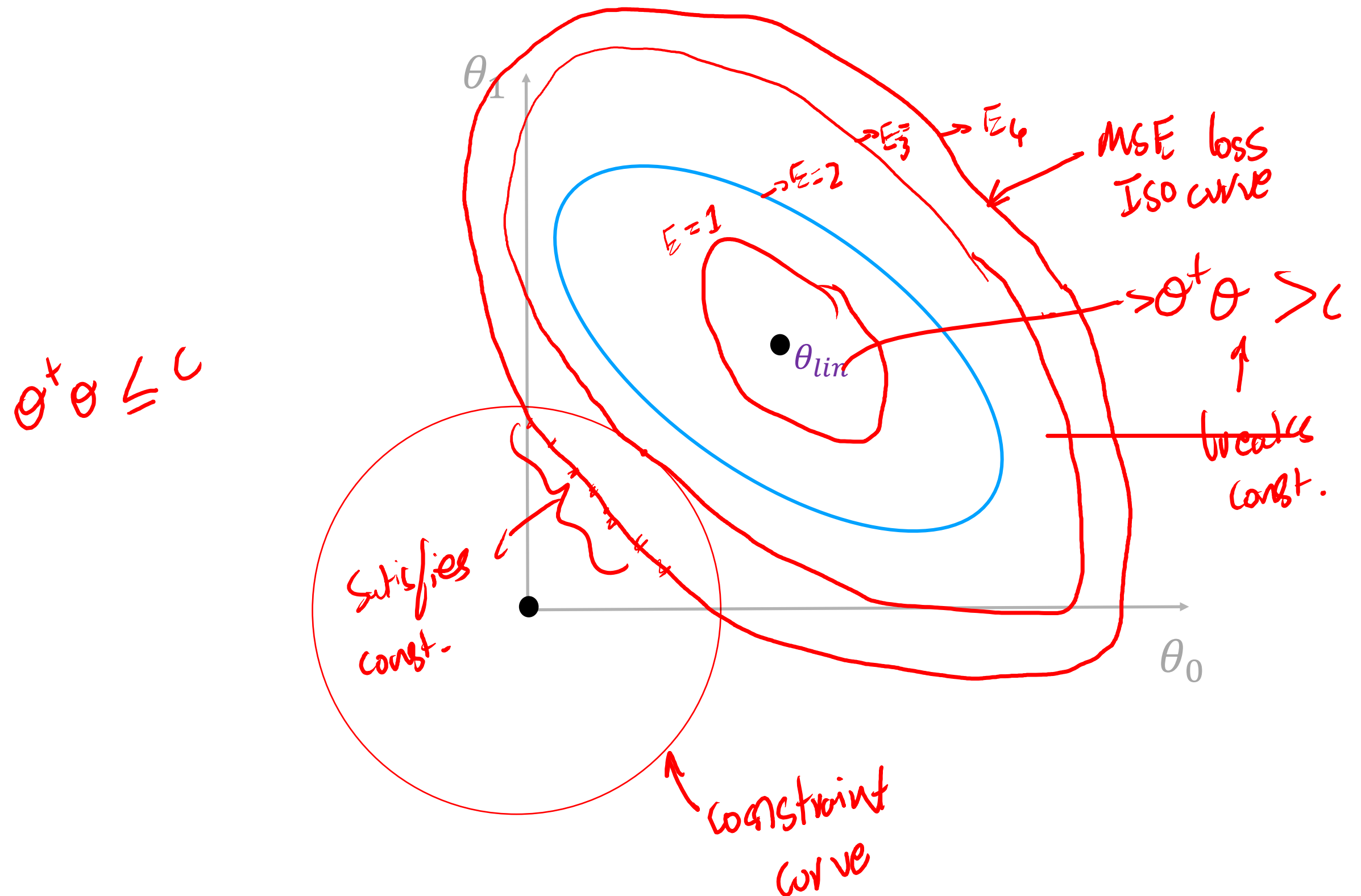
Subject to  $\theta^t \theta \leq C$

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

Find a solution in  $\theta^t \theta$   
that minimizes  $E(\theta)$



# Constraint and Loss



Considering the below  $E(\theta)$  and  $\mathcal{C}$   
what is a  $\theta$  candidate here?

$\nabla E$ : the gradient (rate) in objective function  
which minimizes error (orthogonal to ellipse.  
Changes happen in orthogonal direction)

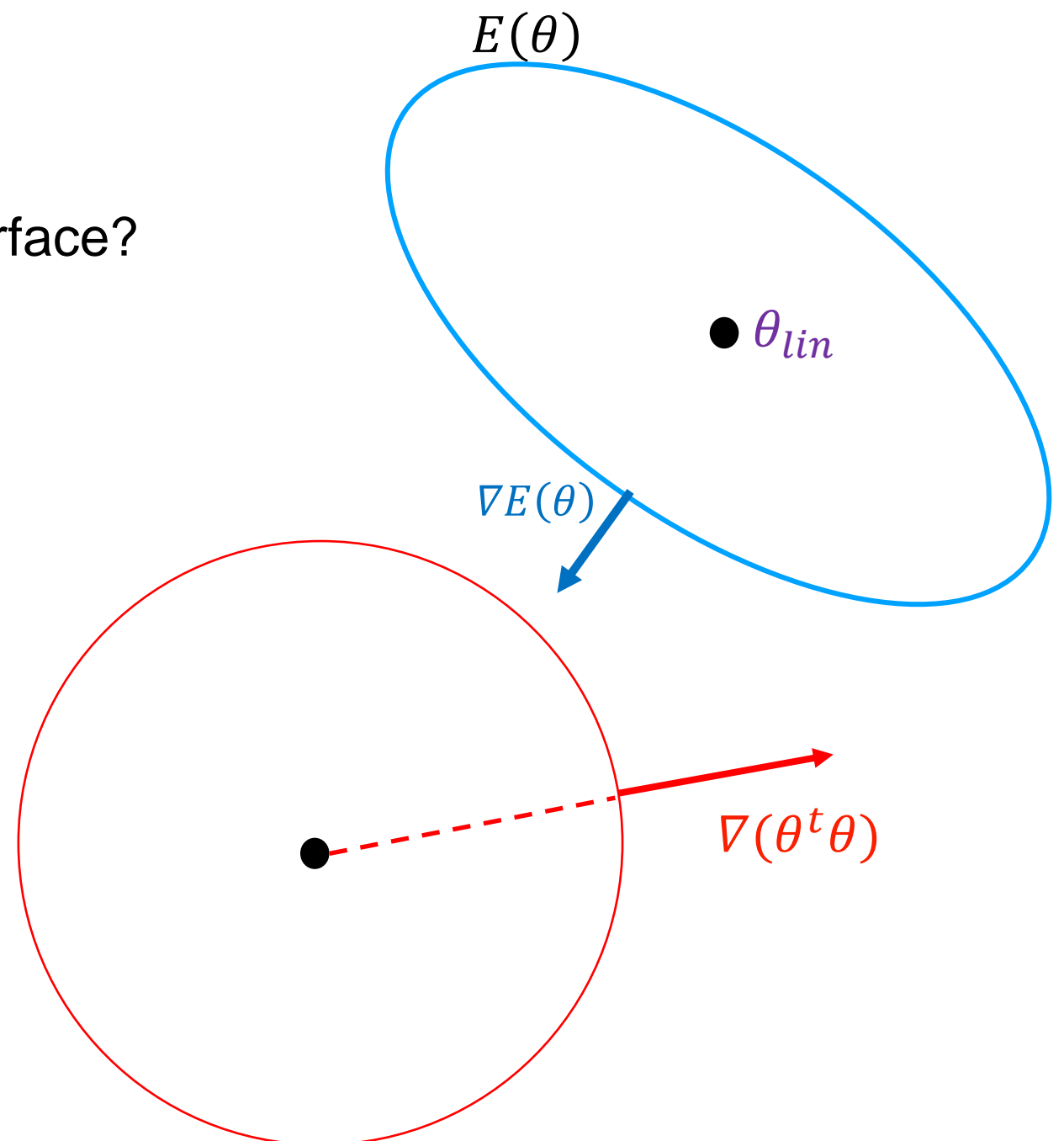
$$\theta^t \theta = \text{Constraint} = \mathcal{C}$$

What is the orthogonal direction on the other surface?

It is just  $\theta$ , a line passing through center of  
the circle

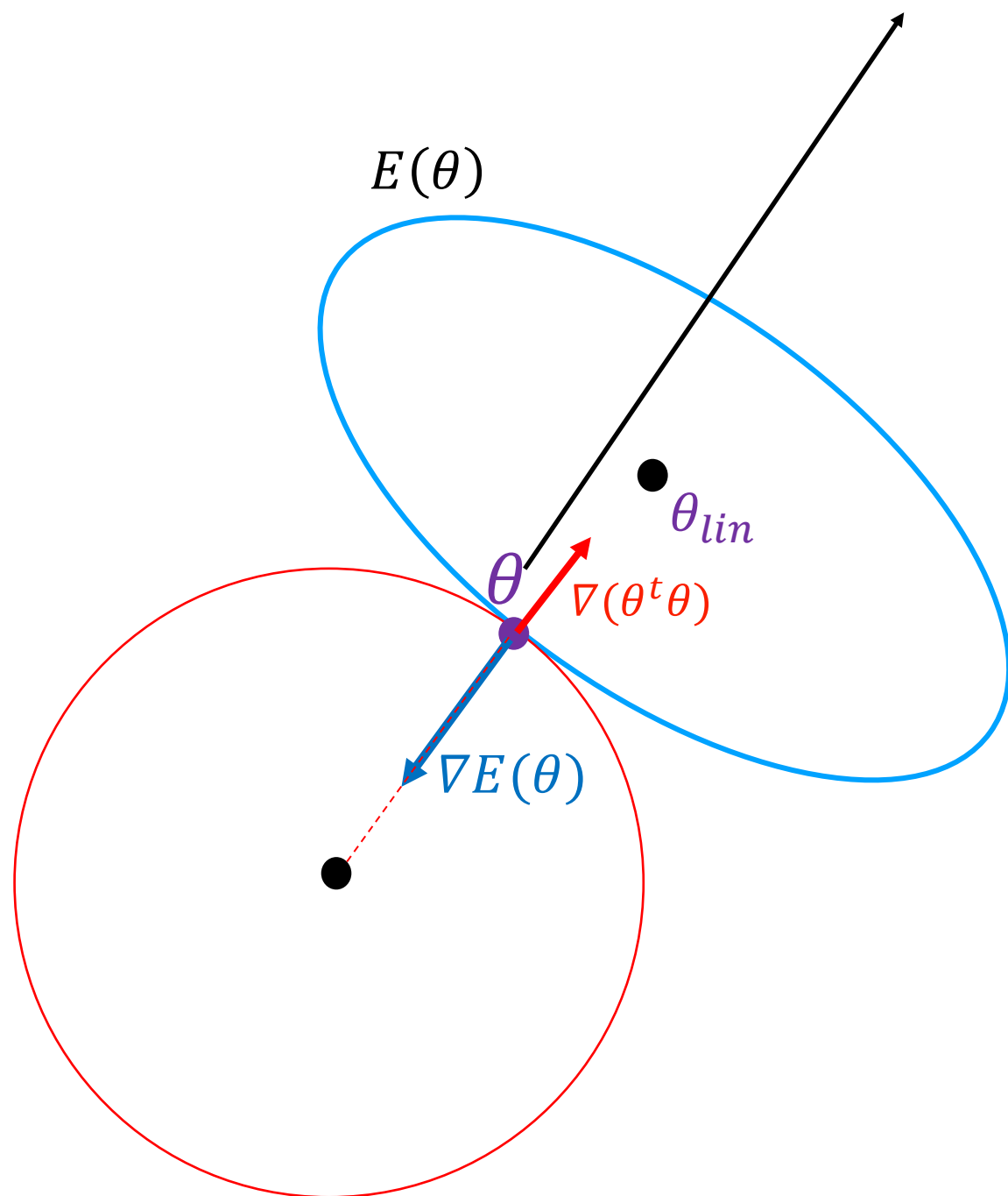
Applying a constrain  $\theta^t \theta$ , where  
the best solution happens?

On the boundary of the circle, as it is the  
closest one to the minimum absolute





Considering the below  $E(\theta)$  and  $C$   
what is the best  $\theta$  solution here?



$$\nabla E(\theta) \propto -\theta \quad \leftarrow \nabla(\theta^T \theta) = 2\theta$$

$$\nabla E(\theta) = -2 \frac{\lambda}{N} \theta$$

$$\nabla E(\theta) + 2 \frac{\lambda}{N} \theta = 0$$


regularization term

Let's do integration

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

$C \uparrow \lambda \downarrow$

# Outline

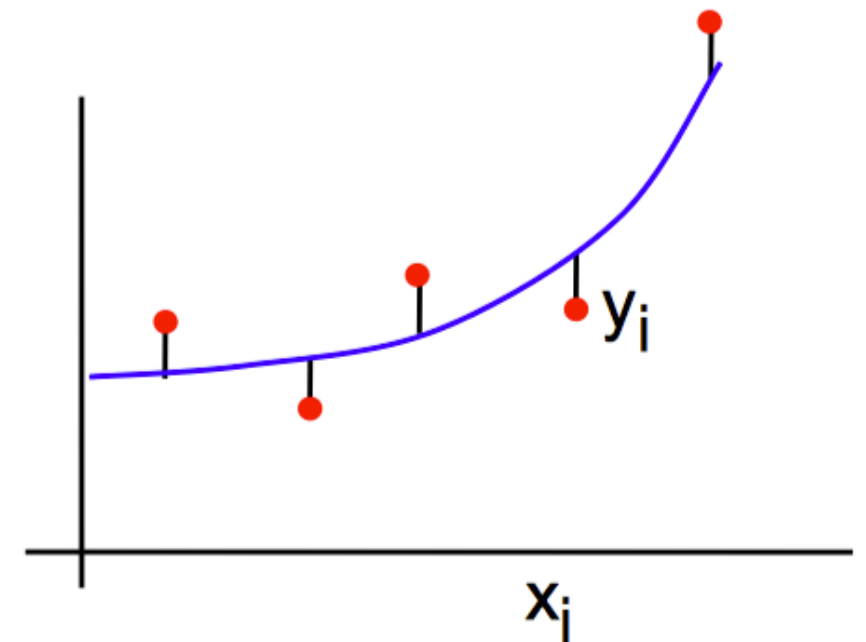
- Overfitting and regularized learning
- Ridge regression 
- Lasso regression
- Determining regularization strength

# Ridge Regression

- Cost function – squared loss:

$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^N \underbrace{\{f(x_i, \theta) - y_i\}^2}_{\text{loss function}} + \underbrace{\frac{\lambda}{N} \|\theta\|^2}_{\text{regularization}}$$

target value
y<sub>i</sub>



- Regression function for x (1D):

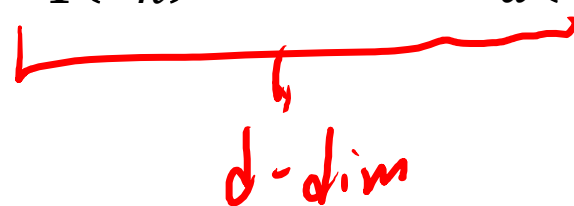
$$f(x, \theta) = \theta_0 + \theta_1 z_1 + \theta_2 z_2 + \cdots + \theta_d z_d + \epsilon = \underline{\underline{z\theta}}$$

$$z_n \sim x^n$$

# Solving for the Weights $\theta$

Notation: write the target and regressed values as  $N$ -vectors


$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \quad \mathbf{f} = \begin{pmatrix} z(x_1)\theta \\ z(x_2)\theta \\ \vdots \\ z(x_n)\theta \end{pmatrix} = \mathbf{z}\theta = \begin{bmatrix} 1 & z_1(x_1) & \dots & z_d(x_1) \\ 1 & z_1(x_2) & \dots & z_d(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_1(x_n) & \dots & z_d(x_n) \end{bmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{pmatrix}$$


  
 $d$ -dim

$\mathbf{z}$  is an  $N \times D$  design matrix

e.g. for polynomial regression with basis functions up to  $x^2$

$$\underline{\underline{\mathbf{z}\theta}} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 \end{bmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{pmatrix}$$


  
 data point

$$\begin{aligned}
\tilde{E}(\theta) &= \frac{1}{N} \sum_{i=1}^N \{f(\underline{x_i}, \underline{\theta}) - y_i\}^2 + \frac{\lambda}{N} \|\theta\|^2 \\
&= \frac{1}{N} \sum_{i=1}^N (y_i - z_i \theta)^2 + \frac{\lambda}{N} \|\theta\|^2 \\
&= \frac{1}{N} (y_i - \underset{\substack{\uparrow \\ \text{matrix form}}}{z} \theta)^2 + \frac{\lambda}{N} \|\theta\|^2
\end{aligned}$$

Now, compute where derivative w.r.t.  $\theta$  is zero for minimum

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

Hence

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

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

no regularization EQ.  
 $\downarrow$   
 $\theta = (z^T z)^{-1} z^T y$

D basis functions, N data points

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

$$\begin{matrix} \text{Dx1} & & \text{DxD} & & \text{DxN} & \text{Nx1} \end{matrix}$$

assume  $N > D$

- This shows that there is a unique solution.

- If  $\lambda = 0$  (no regularization), then

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

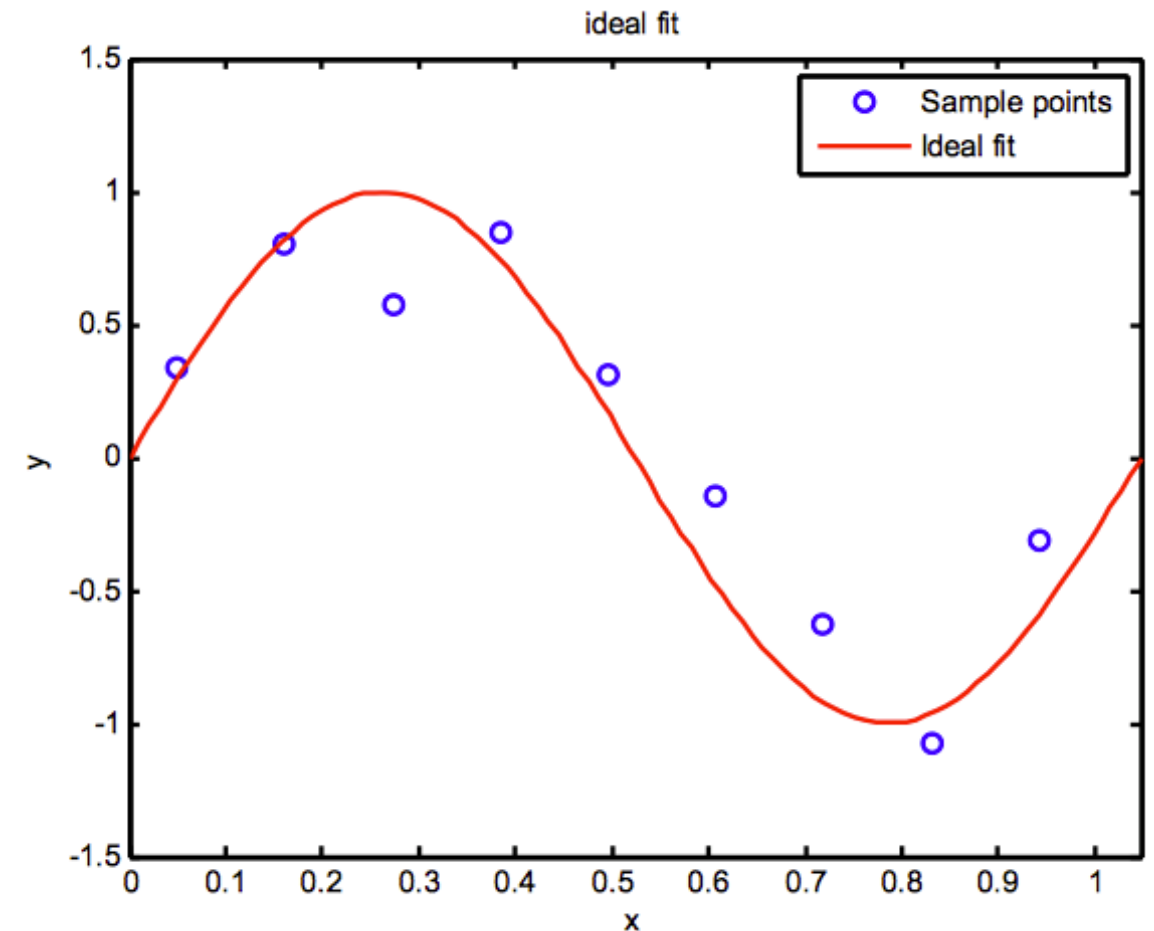
~~old unreg. sol.~~

where  $Z^+$  is the pseudo-inverse of  $Z$  (pinv in Matlab)

- Adding the term  $\lambda I$  improves the conditioning of the inverse, since if  $Z$  is not full rank, then  $(Z^T Z + \lambda I)$  will be (for sufficiently large  $\lambda$ )
- As  $\lambda \rightarrow \infty$ ,  $\theta \rightarrow \frac{1}{\lambda} Z^T y \rightarrow 0$

# 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



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

$$z: x \rightarrow z$$

$$\mathbb{R} \rightarrow \mathbb{R}^{D+1}$$

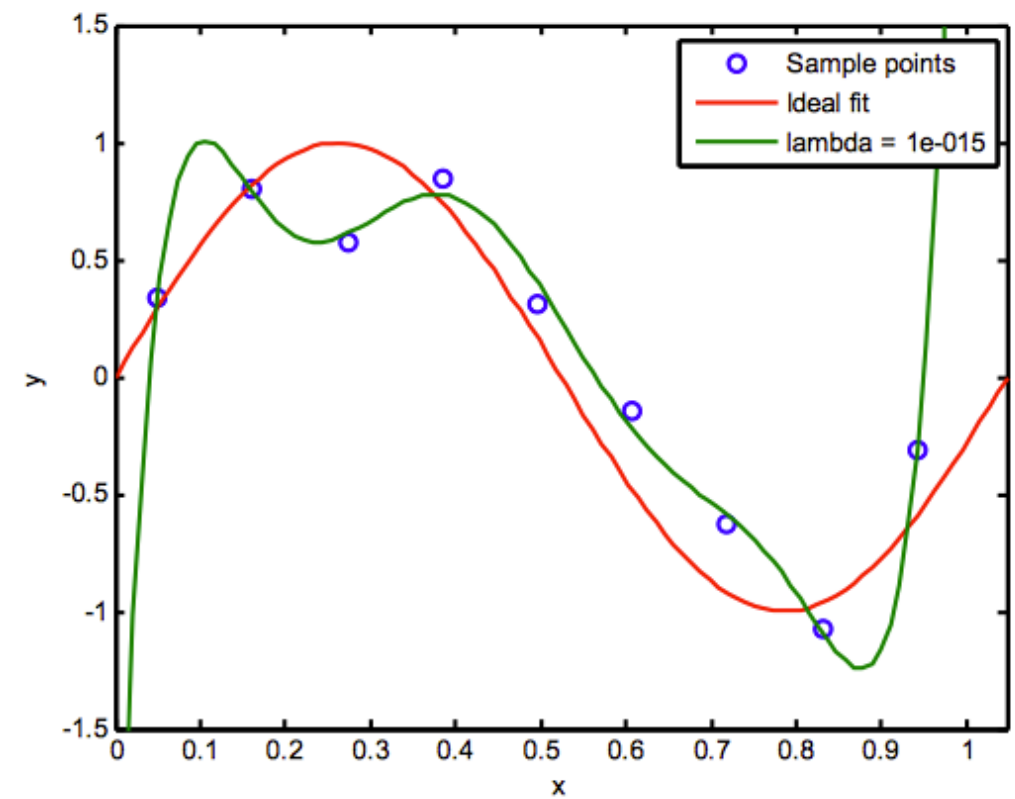
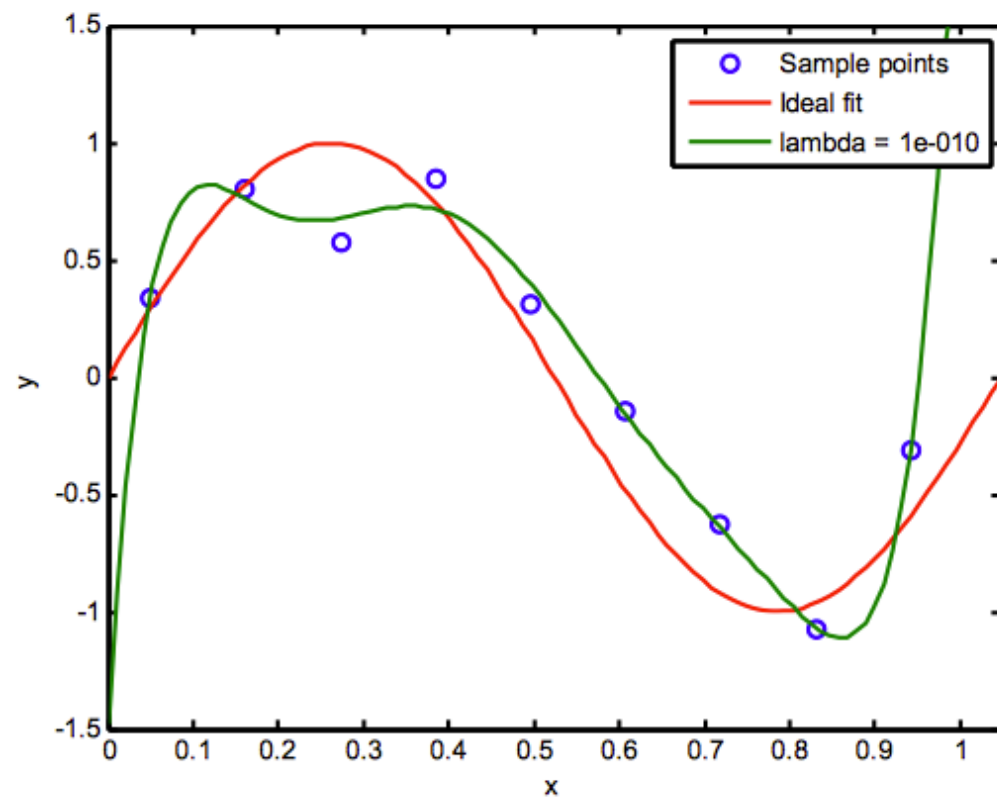
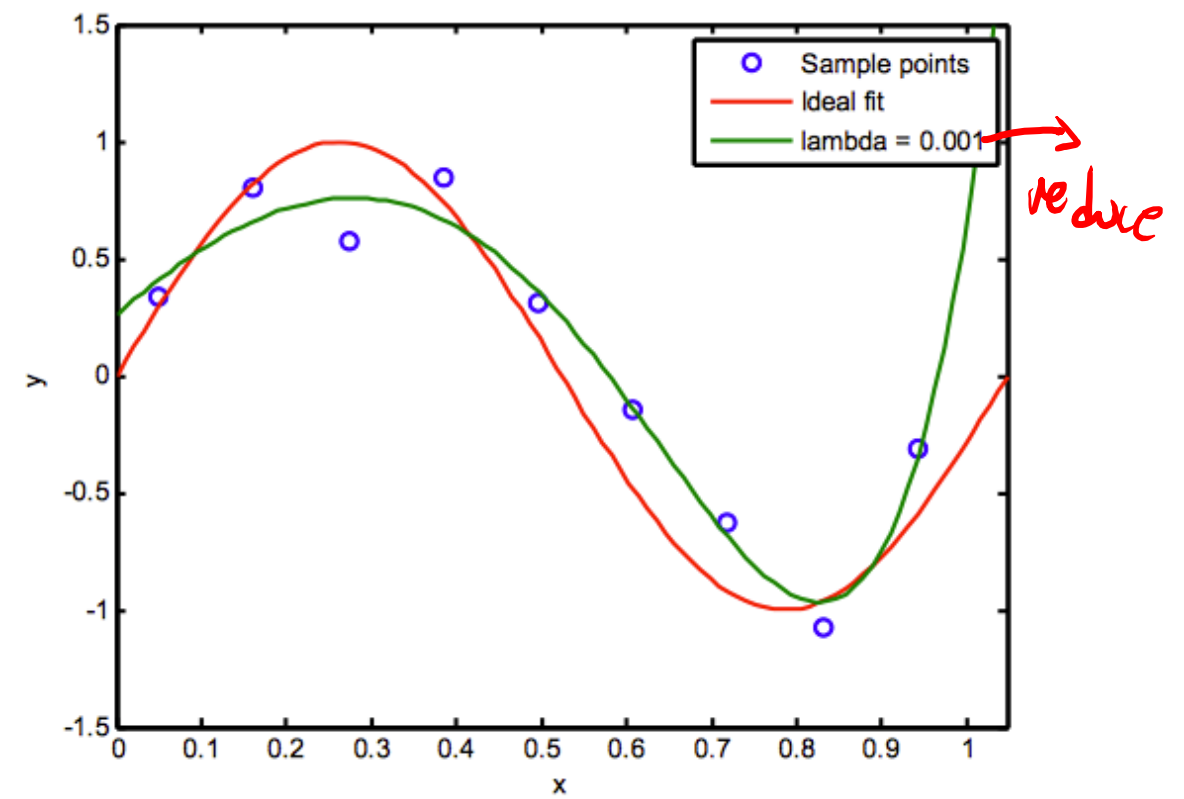
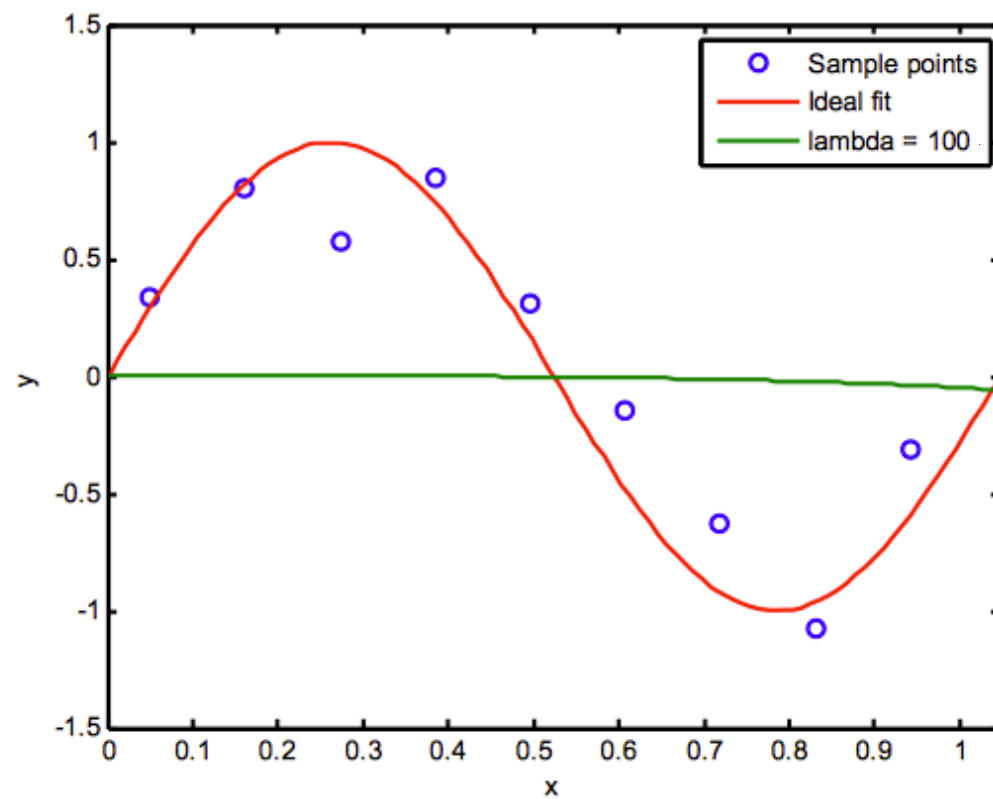
$$\tilde{E}(\theta) = \frac{1}{N} \sum_{i=1}^N \{f(x_i, \theta) - y_i\}^2 + \frac{\lambda}{N} \|\theta\|^2$$

*data points* (handwritten red arrow pointing to the sum index)

*$\sum_{j=1}^{D+1} \theta_j^2$*  (handwritten red note under the regularization term)

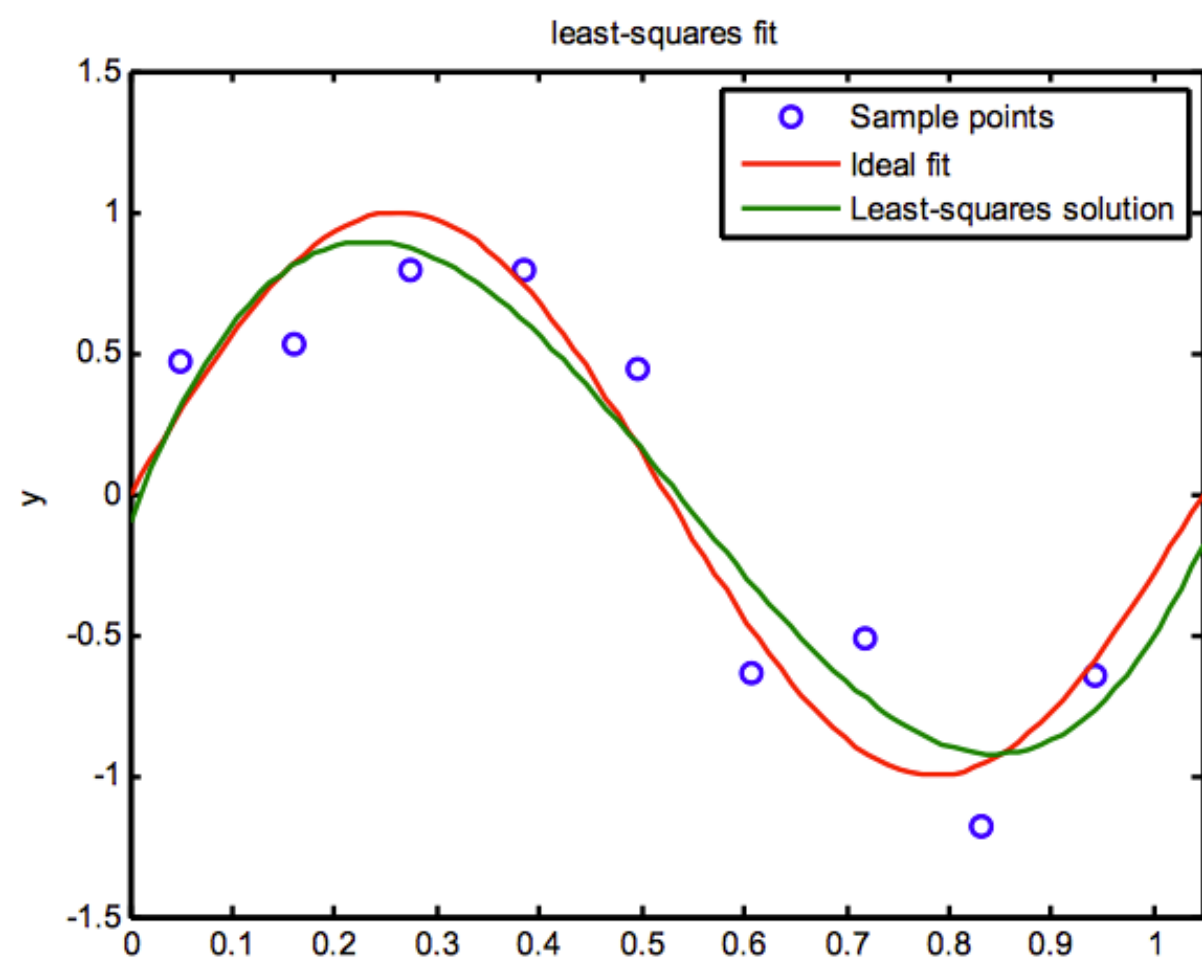
$\theta$  is a  $D+1$  dimensional vector

N = 9 samples, D = 7

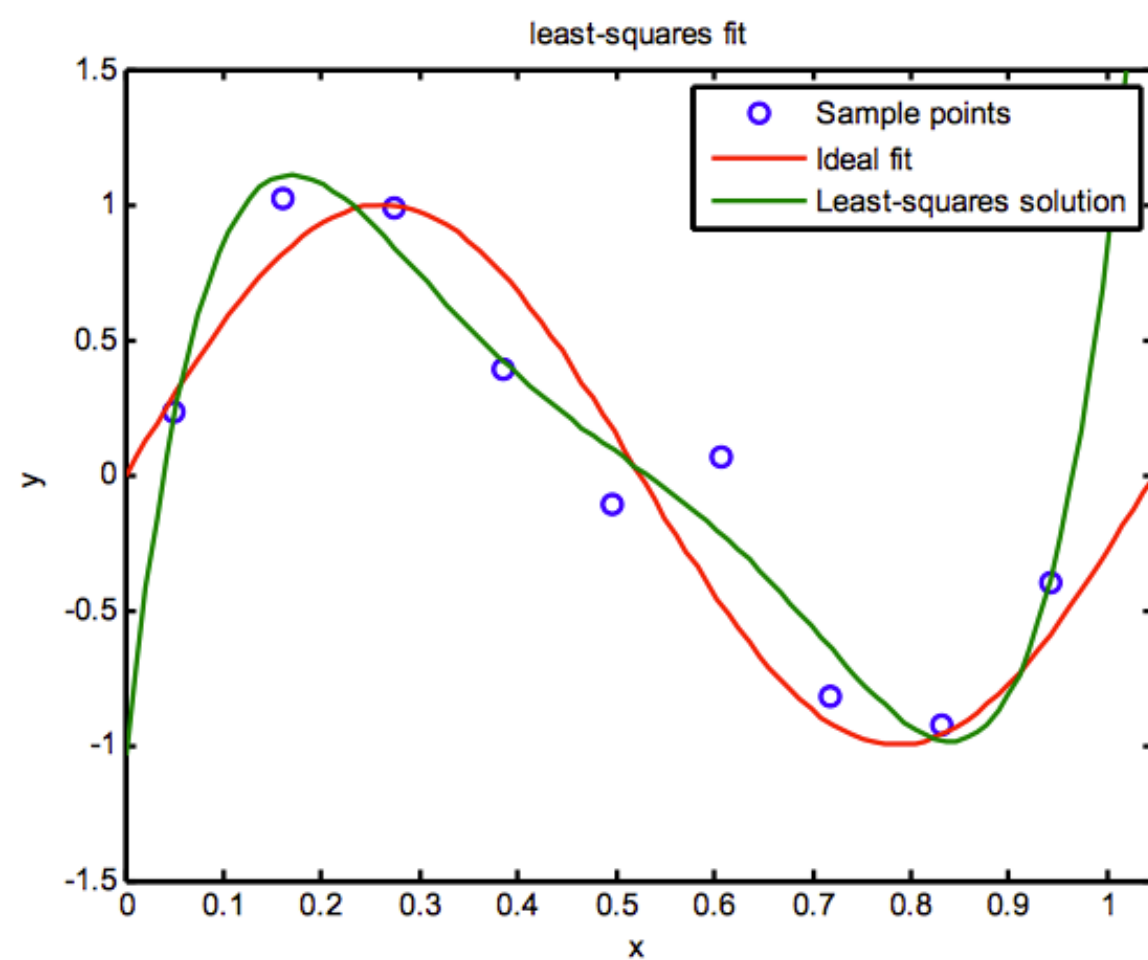





$D = 3$



$D = 5$



# Outline

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

# Regularized Regression

Minimize with respect to  $\theta$

$$\sum_{i=1}^N \underbrace{l(f(\mathbf{x}_i, \theta), y_i)}_{\text{loss function}} + \underbrace{\lambda R(\theta)}_{\text{regularization}}$$

*MSF*  
*regularization const./strength*  
*loss by Regularized*

- There is a choice of both loss functions and regularization
- So far we have seen – “ridge” regression

- squared loss:  $\sum_{i=1}^N (y_i - f(x_i, \theta))^2$

- squared regularizer:  $\lambda \|\theta\|^2$

Now let's look at another regularization choice.

# The Lasso Regularization (norm one)

- LASSO = Least Absolute Shrinkage and Selection

Minimize with respect to  $\theta$

$$\sum_{i=1}^N \underbrace{l(f(\mathbf{x}_i, \theta), y_i)}_{\text{loss function}} + \underbrace{\lambda R(\theta)}_{\text{regularization}}$$

- This is a quadratic optimization problem
- There is a unique solution
- p-Norm definition:  $\|\theta\|_p = \left( \sum_{j=1}^d |\theta_j|^p \right)^{\frac{1}{p}}$

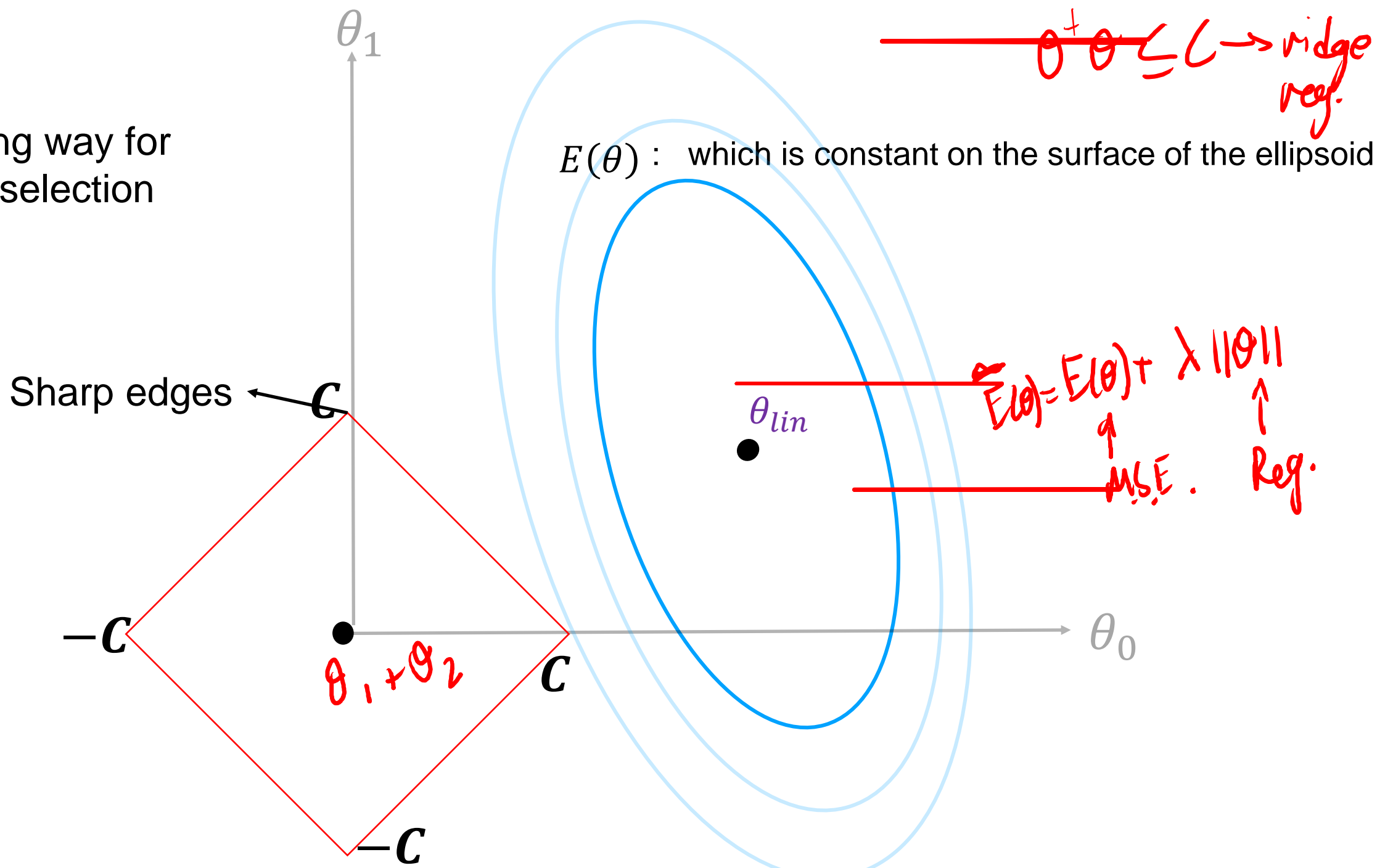
Let's say we have two parameters ( $\theta_0$  and  $\theta_1$ )

$$\text{Minimize } E(\theta) = \frac{1}{N} (z^T \theta - y)^T (z^T \theta - y)$$


Subject to  $\theta \leq C$

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

Interesting way for  
feature selection



# Outline

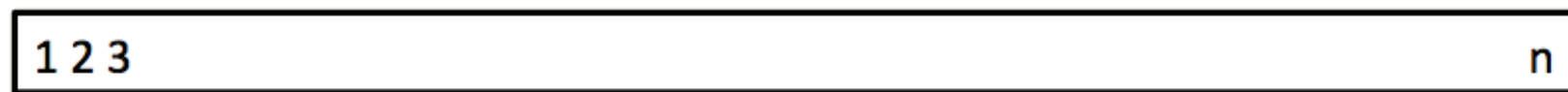
- Overfitting and regularized learning
- Ridge regression
- Lasso regression
- Determining regularization strength 

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



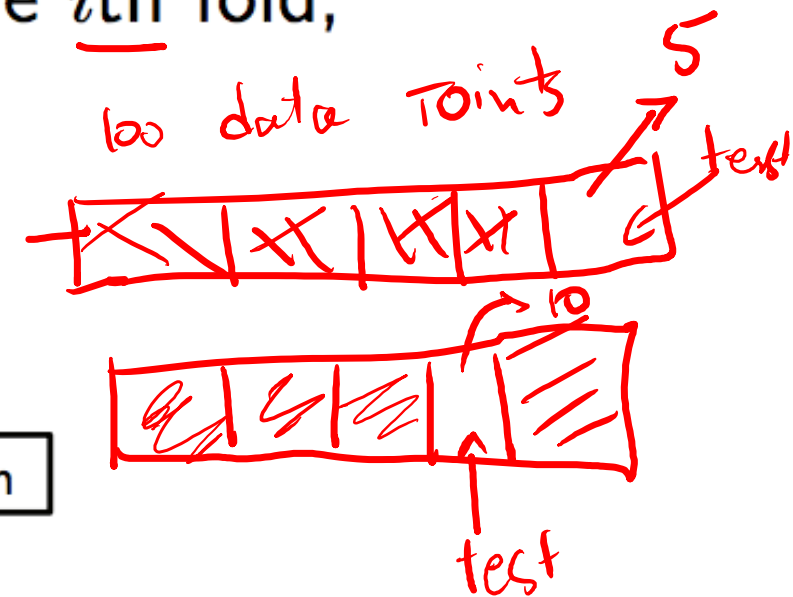
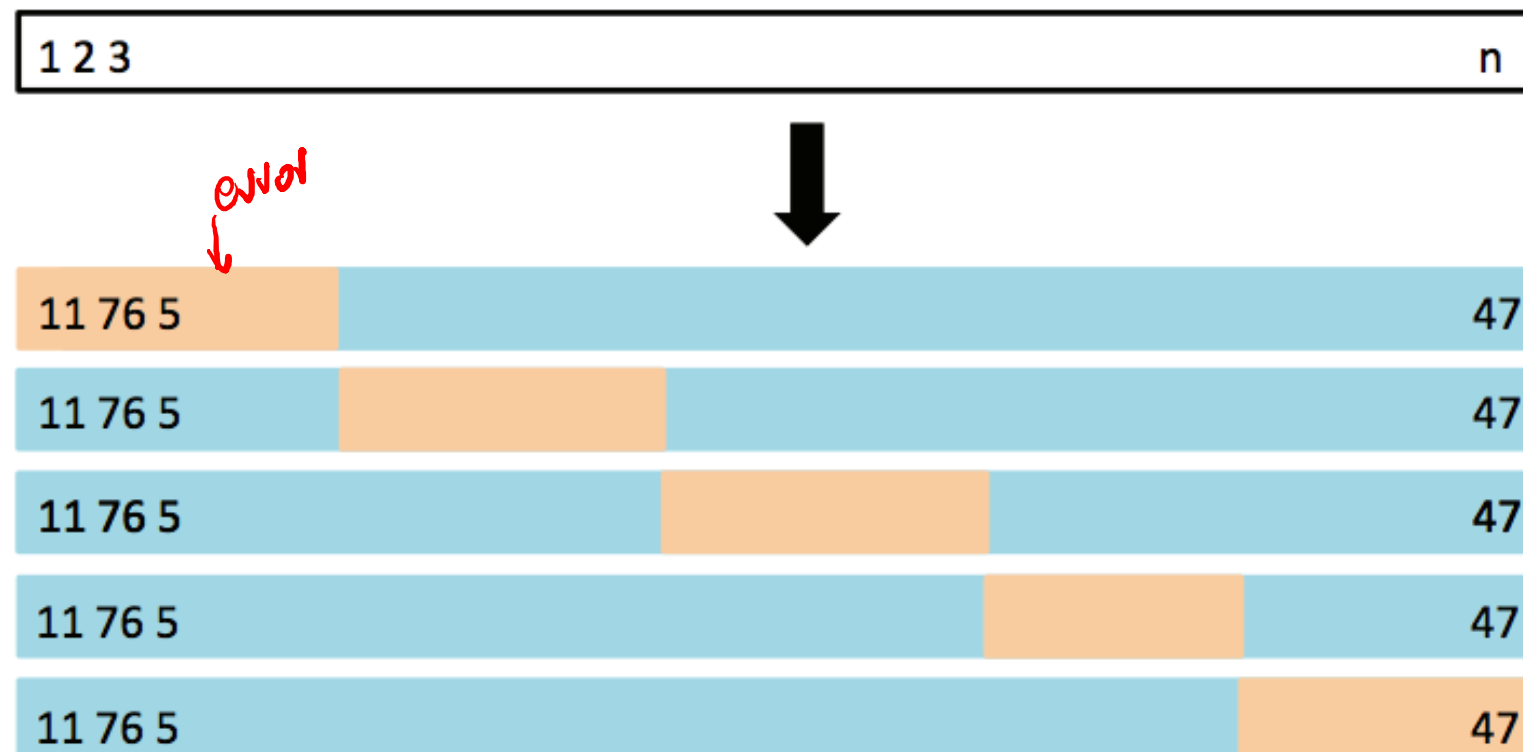
# K-Fold Cross Validation

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.

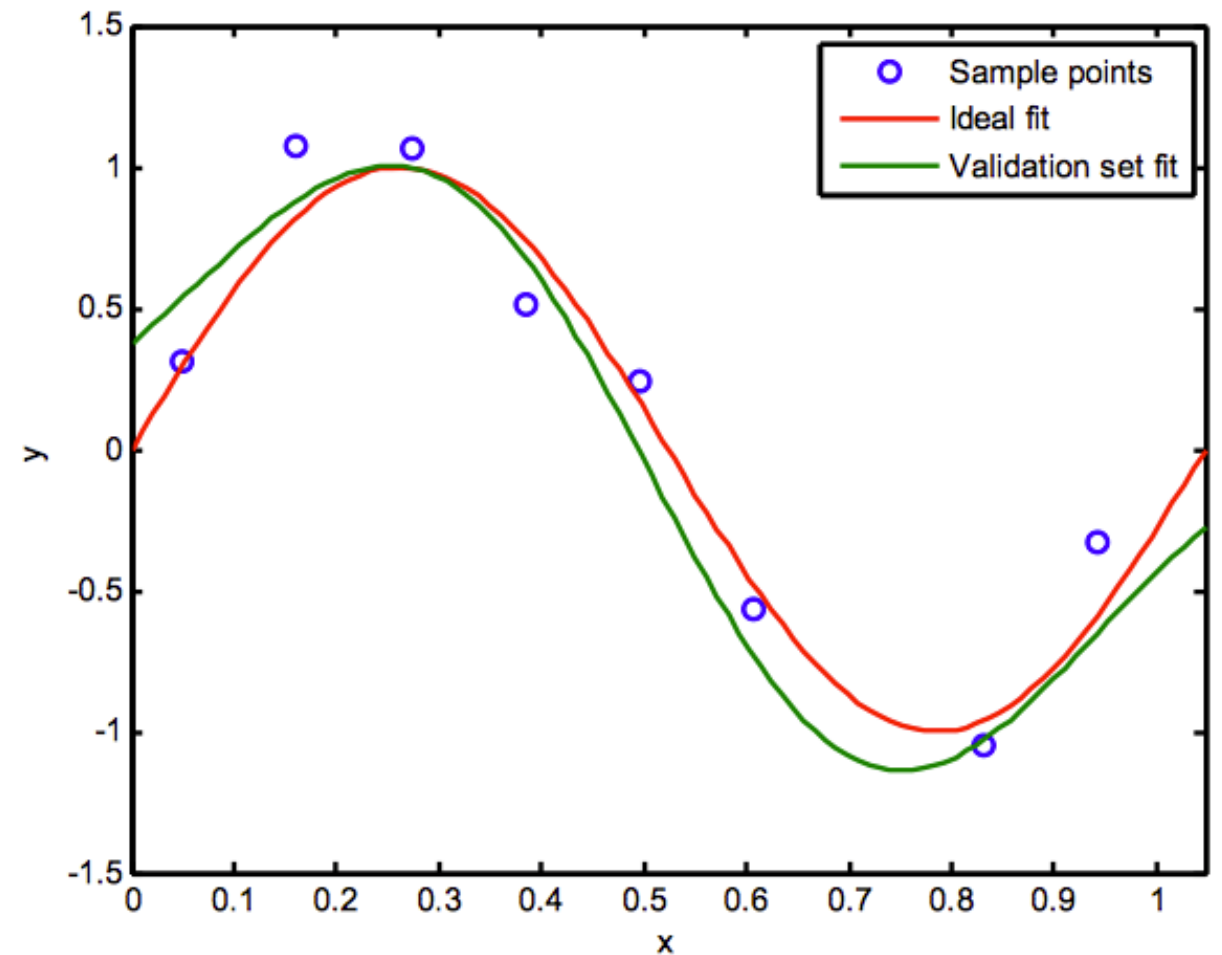
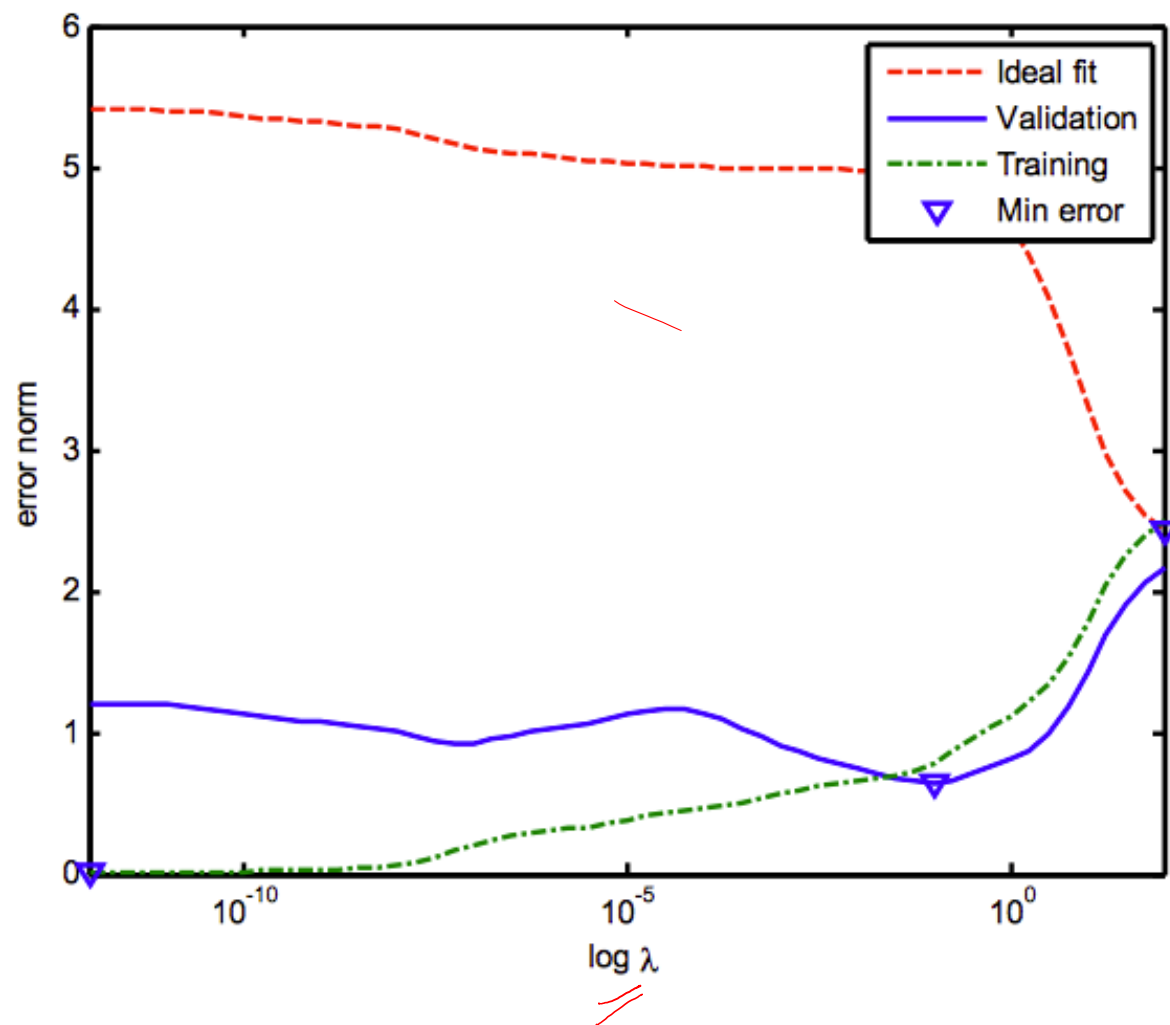


$$CV_{\text{test}} = \sum_{i=1}^k CV_i =$$

error



# 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  $\lambda$