

# **Artificial Intelligence and Machine Learning**

## **Linear Regression**

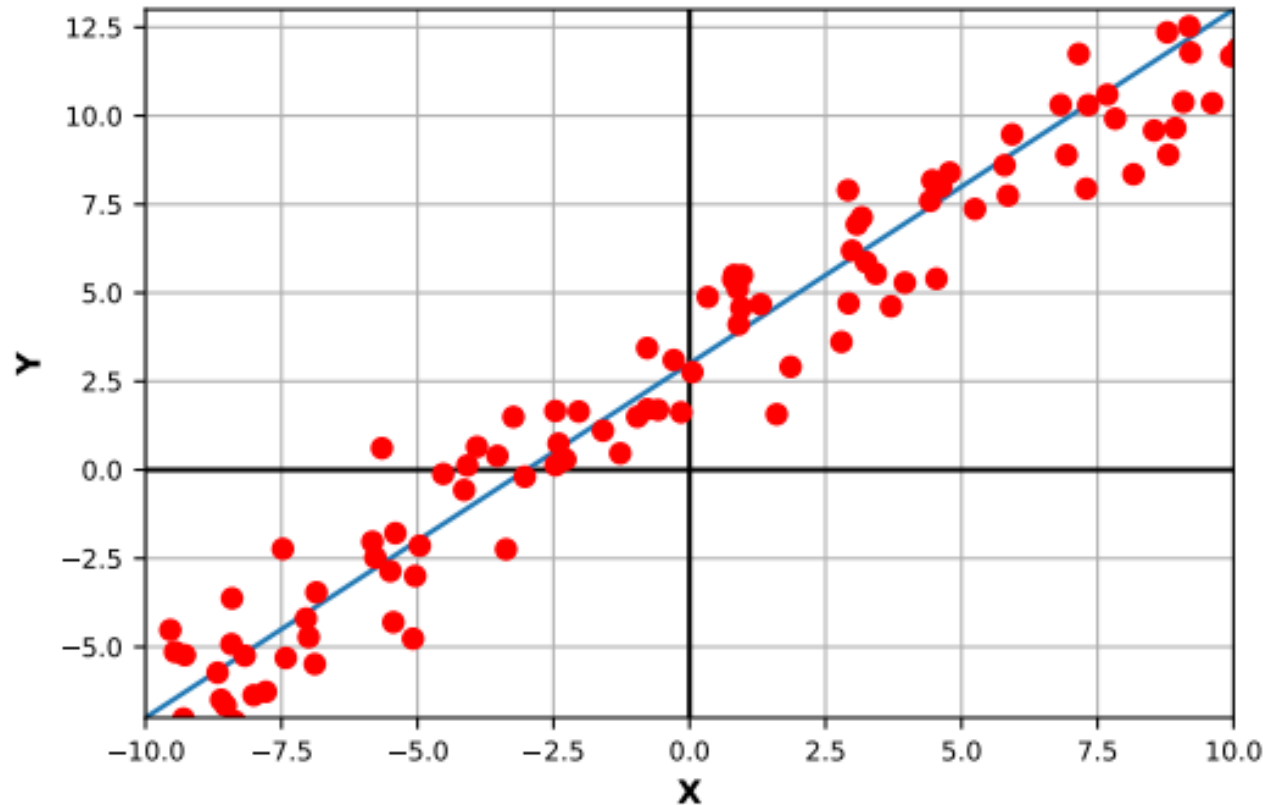
# Lecture 1: Outline

- Linear Regression
- Optimization
- Applications

# Motivation

- Linear Regression is “still” one of the more widely used ML/DL Algorithms
- Easy to understand and implement
- Efficient to Solve
- We will use Linear Regression to Understand the concepts of:
  - Data
  - Models
  - Loss
  - Optimization

# Simple Linear Regression



Model (*Linear*)

$$Y = mX + b$$

Y: Response Variable

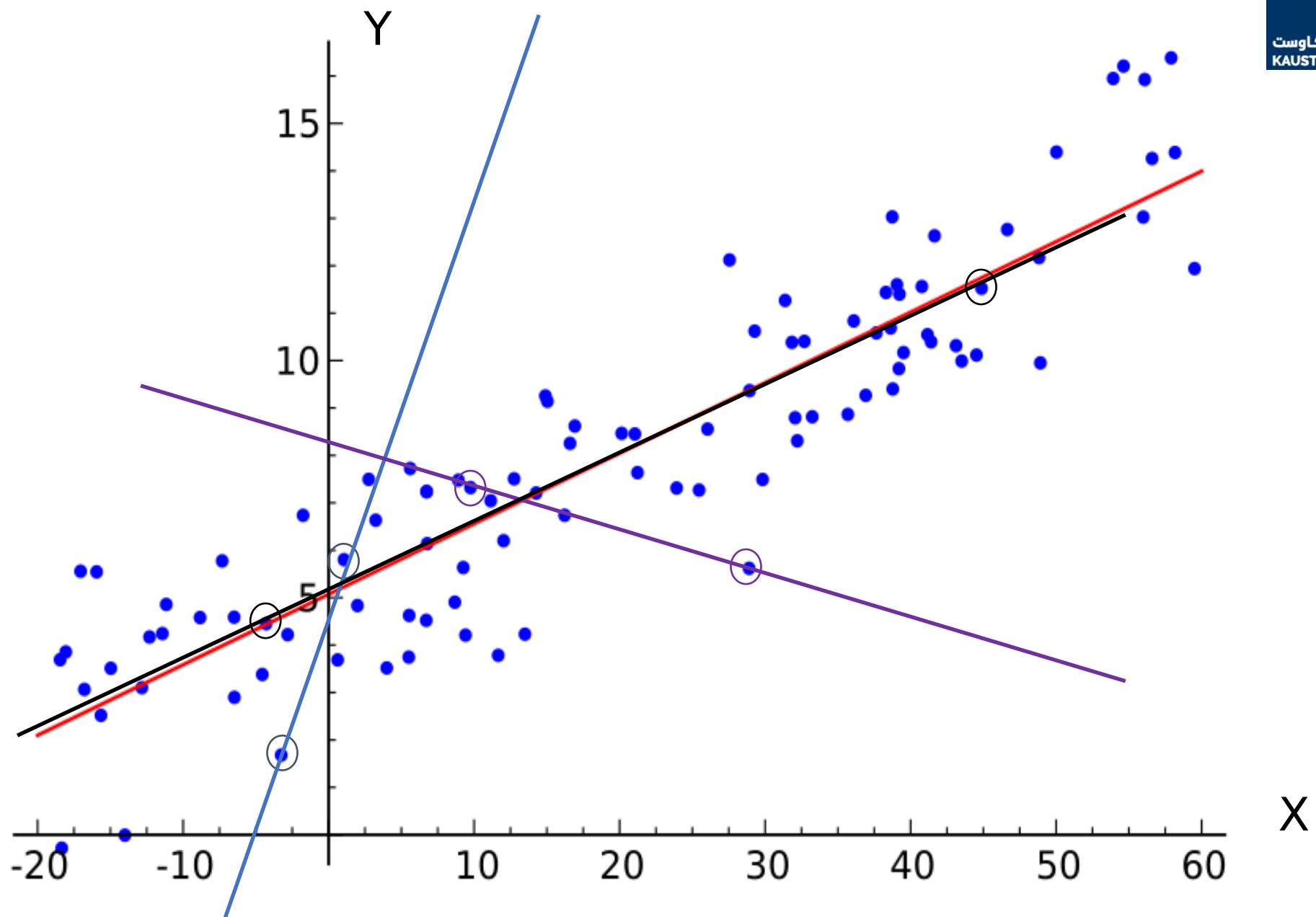
X: Covariate / Ind.,  
var/Regressors

m: slope

b: bias  $\theta = \{m, b\}$

# Simple Linear Regression

- Input: data (),
  - Goal: learn values of variable ()
- 
- Question: How many points in a plane do we need to fit a line through it?



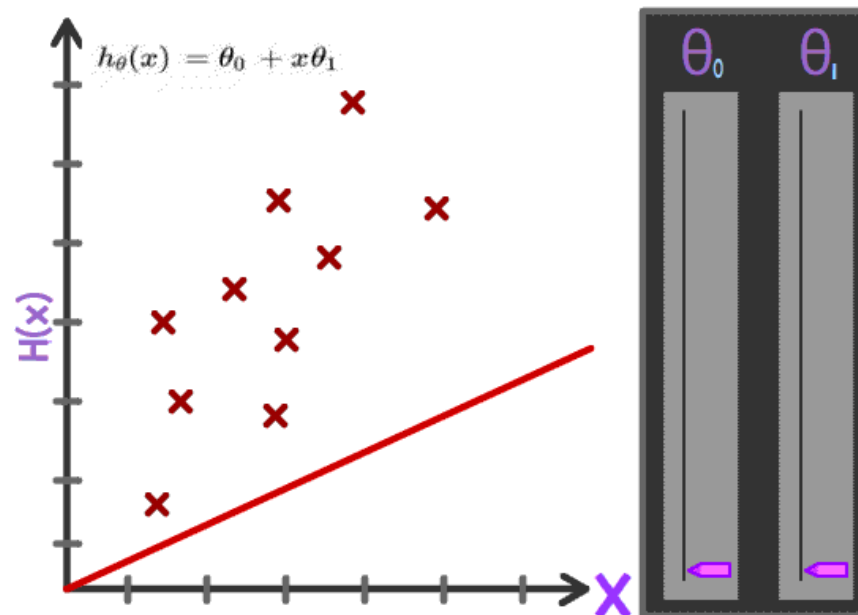
# Notation

- Some clarification about the notation we will use for this course
- $i$  is the index of the data,  $j$  is the feature number, and  $p$  is the power.

# Solution Strategy for Solving the Problem



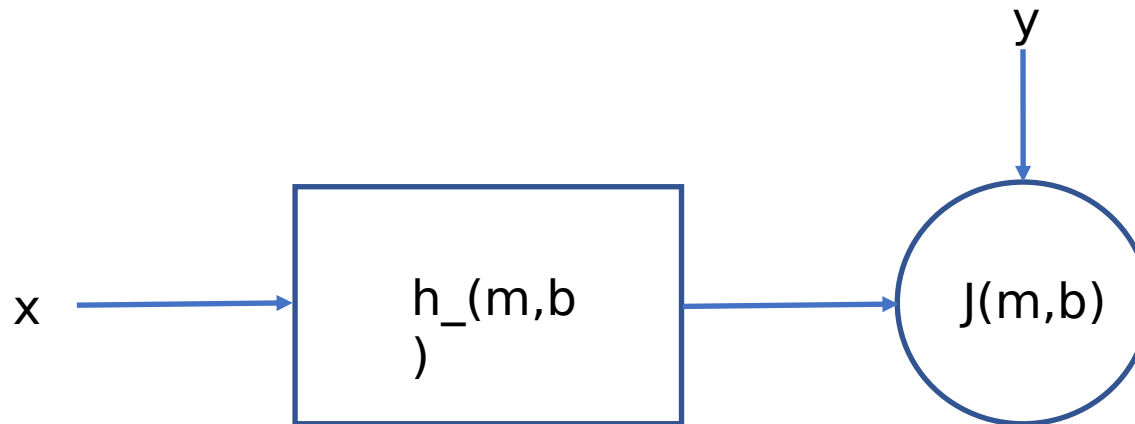
- We want a line which is in some sense the “average line” that represents the data.
- Any ideas as to how we can do it?





# Cost Function

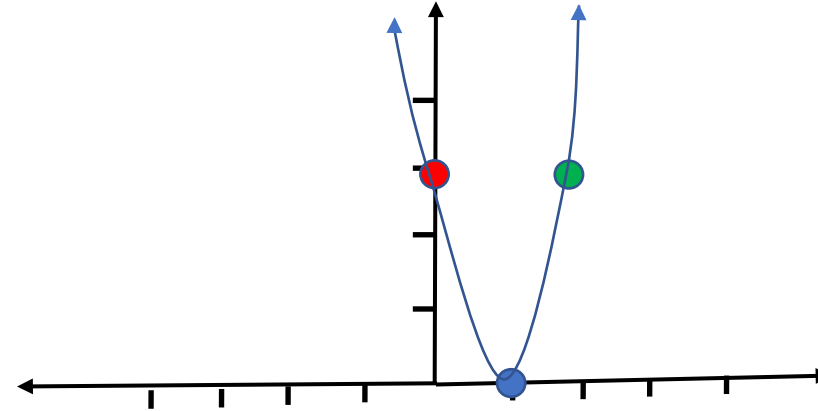
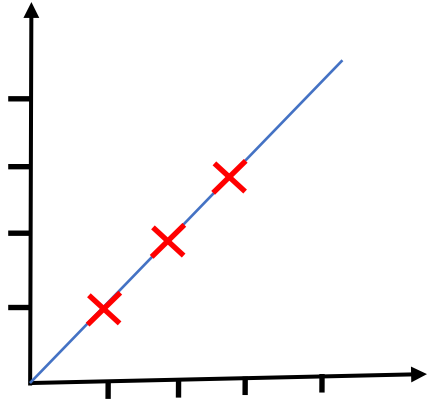
- We want to minimize the discrepancy between our model hypothesis and the observed label.



# Intuition of Cost Function

$$h(x) = mx$$

$$J(m) = \sum_{i=1}^3 (y_i - mx_i)^2$$



Cost  
Function

Hypothesis

$$J(1) = 0$$

$$J(0) = 14$$

$$J(2) = 14$$

# How to find minima of a function (Review):



$$J(m) = \sum_{i=1}^3 (y_i - mx_i)^2 \quad J(m) = \sum_{i=1}^3 (i - mi)^2$$

$$\frac{dJ(m)}{dm} = \frac{d}{dm} \sum_{i=1}^3 (i - mi)^2 \quad \frac{dJ(m)}{dm} = \sum_{i=1}^3 \frac{d}{dm} (i - mi)^2$$

$$\frac{dJ(m)}{dm} = \sum_{i=1}^3 -2i(i - mi) \quad -2 \sum_{i=1}^3 i^2 + 2m \sum_{i=1}^3 i^2 = 0 \quad m = 1$$



# Hypothesis Function with 2 Variables

- Let's setup regression for linear function in two variables:
- The hypothesis function is:

$$\hat{y}_i = mx_i + b$$

- Similar to the previous problem our loss function is:

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

- Let's calculate the partial derivatives of the loss function w.r.t.

# Gradient of the cost function

- We get the following expressions for the gradient of the cost function

$$\frac{\partial J}{\partial m} = \frac{1}{N} \sum_{i=1}^N -2(y_i - \hat{y}_i)x_i$$

$$\frac{\partial J}{\partial b} = \frac{1}{N} \sum_{i=1}^N -2(y_i - \hat{y}_i)$$

# Gradient of the cost function

- Simplifying the above expressions, we get:

$$\frac{\partial J}{\partial m} = \frac{-2}{N} \sum_{i=1}^N y_i x_i + \frac{2m}{N} \sum_{i=1}^N x_i^2 + \frac{2b}{N} \sum_{i=1}^N x_i$$

$$\frac{\partial J}{\partial b} = \frac{-2}{N} \sum_{i=1}^N y_i + \frac{2m}{N} \sum_{i=1}^N x_i + \frac{2b}{N} \sum_{i=1}^N 1$$



# Gradient of the cost function

- Setting the Gradient equal to 0, and solving for m and b, we get

$$\underbrace{\begin{bmatrix} \frac{\sum_i x_i^2}{N} & \frac{\sum_i x_i}{N} \\ \frac{\sum_i x_i}{N} & 1 \end{bmatrix}}_{\text{Matrix A}} \underbrace{\begin{bmatrix} m \\ b \end{bmatrix}}_{\text{Vector X}} = \underbrace{\begin{bmatrix} \frac{\sum_i x_i y_i}{N} \\ \frac{\sum_i y_i}{N} \end{bmatrix}}_{\text{Vector Y}}$$

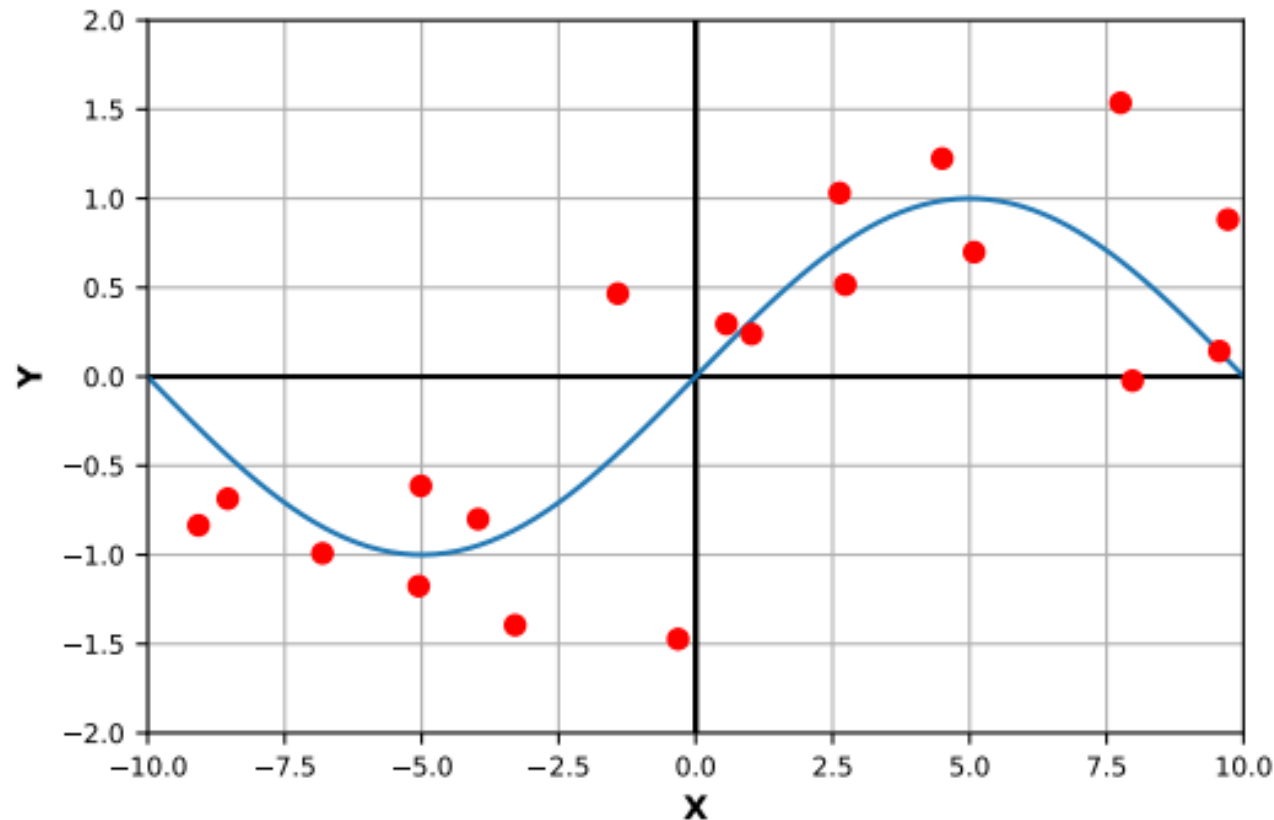
# Issues with the Approach

- Calculating gradients like this can quickly become tedious
- Each term on either side of the expression can be written as a dot product of two vectors (maybe we can calculate it more efficiently)?
- Let's explore if we can do something better through vectorization
- Before, we step into vectorization, let's consider regression for non-linear functions



# Fitting Non-linear Data

- What if  $y$  is a non-linear function of  $x$ , will this approach still work?



# Transforming the Feature Space



- We can transform features  $x_i$

$$x_i = (x_i^1, x_i^2, x_i^3, \dots, x_i^m)$$

- We will apply some non-linear transformation  $\phi$ :

$$\phi : \mathbb{R}^m \rightarrow \mathbb{R}^M$$

- For example, Polynomial transformation:

$$\phi(x_i) = \{1, x_i^1, x_i^{1,[2]}, \dots, x_i^{1,[k]}, x_i^2, x_i^{2,[2]}, \dots, x_i^{2,[k]}, \dots, x_i^m, x_i^{m,[2]}, \dots, x_i^{m,[k]}\}$$

- others: cosine, splines, radial basis functions, etc.
- Expert engineered features (modeling)

# Vectorization

- To truly appreciate the power of vectorization. Let's make the problem a little more complex. The hypothesis function is now
- Where  $w$  are the unknown weights of the data features of the input
- Next, we denote the discrepancy between  $y$  and  $\hat{y}$  as

# Vectorization

- Now let's collect the above equation for all datapoints

$$y_1 = \hat{y}_1 + \epsilon_1$$

$$y_2 = \hat{y}_2 + \epsilon_2$$

.

.

.

$$y_N = \hat{y}_N + \epsilon_N$$



# Vectorization

- Replacing the values of , we get:

$$y_1 = w_0 + w_1x_1^1 + w_2x_1^2 + \dots + w_Mx_1^M + \epsilon_1$$

$$y_2 = w_0 + w_1x_2^1 + w_2x_2^2 + \dots + w_Mx_2^M + \epsilon_2$$

.

.

.

$$y_N = w_0 + w_1x_N^1 + w_2x_N^2 + \dots + w_Mx_N^M + \epsilon_N$$



# Vectorization

- Collecting the equations in matrix form:

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & x_1^1 & x_1^2 & \dots & x_1^M \\ 1 & x_2^1 & x_2^2 & \dots & x_2^M \\ 1 & x_3^1 & x_3^2 & \dots & x_3^M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N^1 & x_N^2 & \dots & x_N^M \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_N \end{bmatrix}$$

# Vectroization

- Notice the rows of the matrix on the right are data samples:

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} \dots & \mathbf{x}_1 & \dots \\ \dots & \mathbf{x}_2 & \dots \\ \dots & \mathbf{x}_3 & \dots \\ \vdots & \vdots & \vdots \\ \dots & \mathbf{x}_N & \dots \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_N \end{bmatrix}$$

# Vectorization

$$\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$$

- Let's formalize some notations:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \dots & \mathbf{x}_1 & \dots \\ \dots & \mathbf{x}_2 & \dots \\ \dots & \mathbf{x}_3 & \dots \\ \vdots & \vdots & \vdots \\ \dots & \mathbf{x}_N & \dots \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{bmatrix} \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_N \end{bmatrix}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon}$$



# Cost function for the Vectorized form



- Notice that we are using the MSE cost function:

$$J = \frac{1}{N} \sum_i (\mathbf{y}_i - \widehat{\mathbf{y}}_i)^2$$

- Using the definition of epsilon we can write the above as:

$$J = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=1}^N (\epsilon_i)^2$$

- Using the definition of dot product the above can be written as:

$$J = \frac{1}{N} \sum_{i=1}^N (\epsilon_i)^2 = \boldsymbol{\epsilon}^T \boldsymbol{\epsilon}$$



# Optimization

- The optimization problem is now:

$$\min_{\theta} \epsilon^T \epsilon$$

$$\min_{\theta} \epsilon^T \epsilon = \min_{\theta} (\mathbf{y} - (\mathbf{X}\theta))^T (\mathbf{y} - (\mathbf{X}\theta))$$

- We will use chain rule to calculate the gradient of the cost function:

$$\frac{\partial}{\partial \theta} J = \frac{dJ}{d\epsilon} \nabla_{\theta} \epsilon$$



# Linear Least Squares

- We get:

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

- Setting it equal to zero we can solve for  $\theta$  :

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

# Probabilistic Interpretation of Linear Regression and MLE

- We can also look at the probabilistic Interpretation of Linear Regression.
- Keeping everything else same as the previous formulation

$$y_i = \mathbf{x}_i^T \boldsymbol{\theta} + \epsilon_i$$

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad y_i \sim \mathcal{N}(\mathbf{x}_i^T \boldsymbol{\theta}, \sigma^2)$$

- Now assume that , then
- We can write the conditional distribution as :

$$\mathbb{P}(y_i | \mathbf{x}_i) \sim \mathcal{N}(\mathbf{x}_i^T \boldsymbol{\theta}, \sigma^2)$$

# Probabilistic Interpretation of LR

- Let's assume that all data point in the dataset are i.i.d. (independent identically distributed). Then we have:

$$\mathbb{P}(\mathcal{D}) = \prod_{i=1}^N \mathbb{P}(\mathbf{x}_i, y_i)$$

- Using Bayes Theorem we can write:

$$\prod_{i=1}^N \mathbb{P}(\mathbf{x}_i, y_i) = \prod_{i=1}^N \mathbb{P}(\mathbf{x}_i) \mathbb{P}(y_i | \mathbf{x}_i)$$



# Maximum Likelihood Estimator

- In simple words, given the Dataset we want to find the values of the unknown parameters which maximize the probability of the Dataset.

- Using the definition of the conditional distribution we have

$$\mathbb{P}(y_i | \mathbf{x}_i) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left( -\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2}{2\sigma^2} \right)$$

• Using the definition we get

$$\prod_{i=1}^N \mathbb{P}(\mathbf{x}_i, y_i) = \prod_{i=1}^N \mathbb{P}(\mathbf{x}_i) \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} \exp \left( -\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2}{2\sigma^2} \right)$$



# Maximum Likelihood Estimator

- Let's try to maximize:

$$\prod_{i=1}^N \mathbb{P}(\mathbf{x}_i, y_i) = \prod_{i=1}^N \mathbb{P}(\mathbf{x}_i) \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} \exp \left( -\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2}{2\sigma^2} \right)$$

- Note that

$$\arg \max_{\boldsymbol{\theta}} \prod_{i=1}^N \mathbb{P}(\mathbf{x}_i, y_i) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^N \exp \left( -\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2}{2\sigma^2} \right)$$

# Maximum Likelihood Estimator

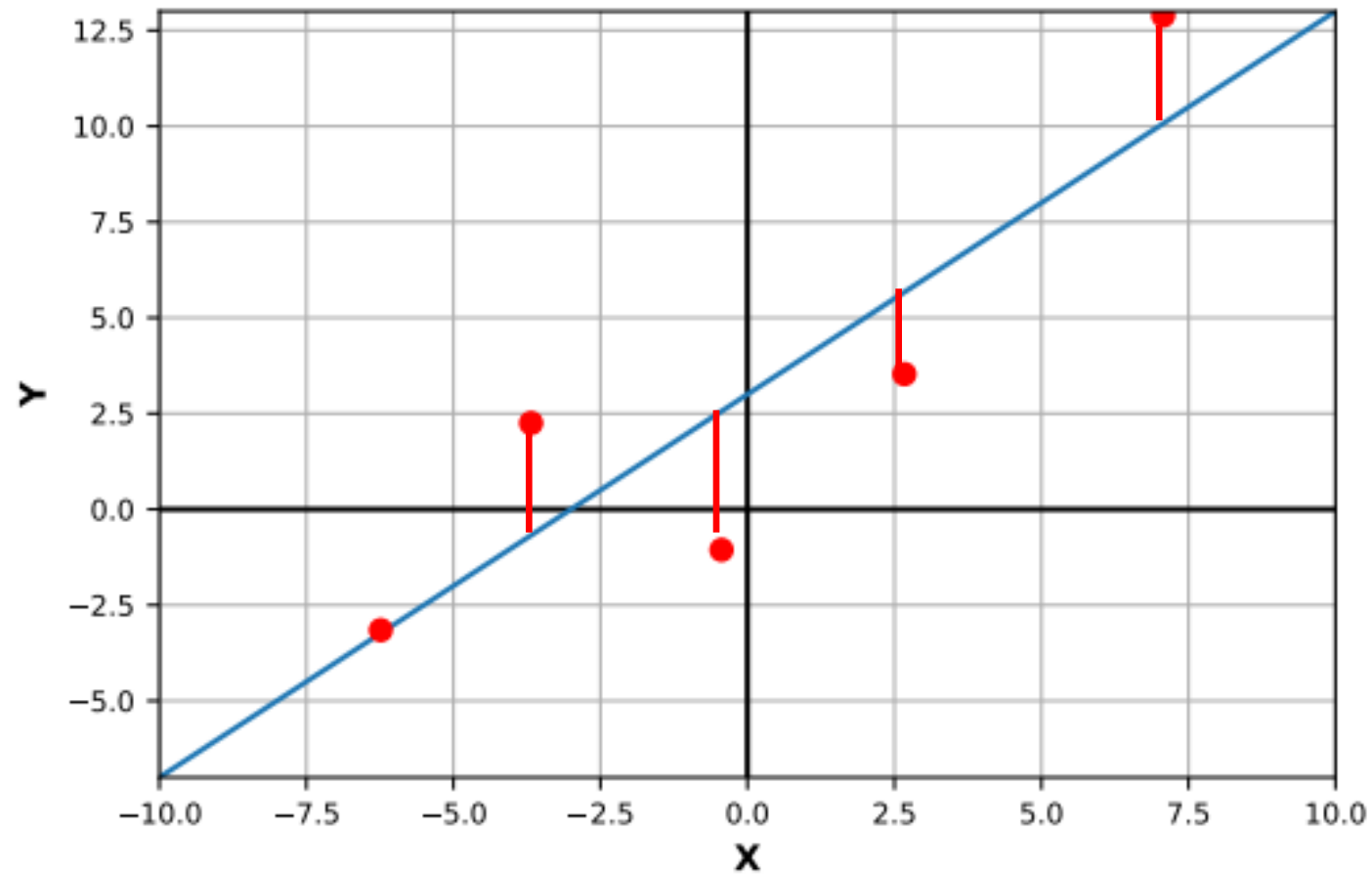
- Furthermore, since the right hand side of the above equation is monotonic in  $\theta$  the arg max will not change if we take log of the expression

$$\arg \max_{\theta} \prod_{i=1}^N \exp \left( - (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2 \right) = \arg \max_{\theta} \sum_{i=1}^N \left( - (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2 \right)$$

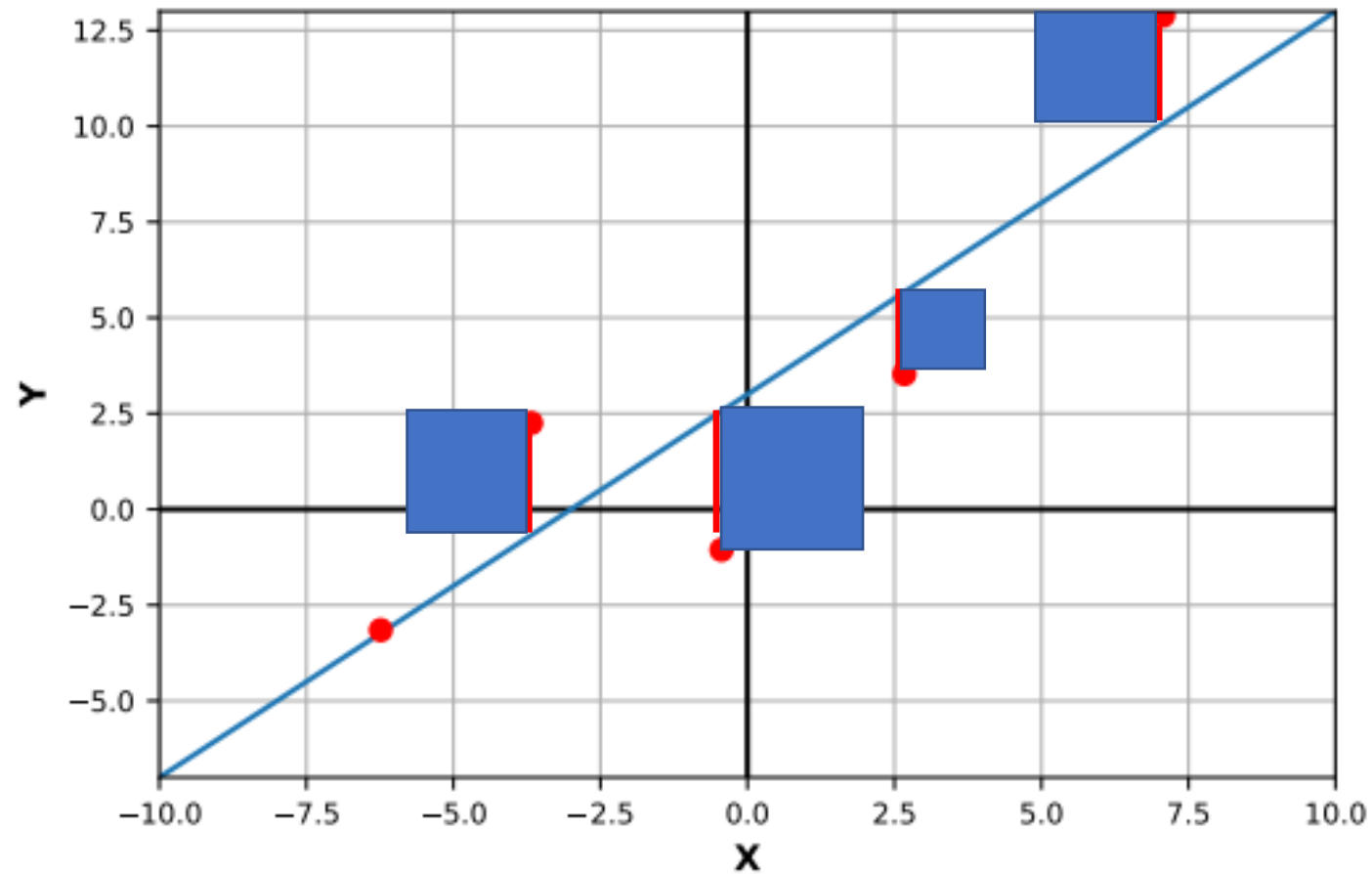
- Notice that the right hand side is minimizing the MSE.
- Hence solution of minimizing the MSE is equivalent to Maximum Likelihood Estimator for linear regression



# Error Visualization

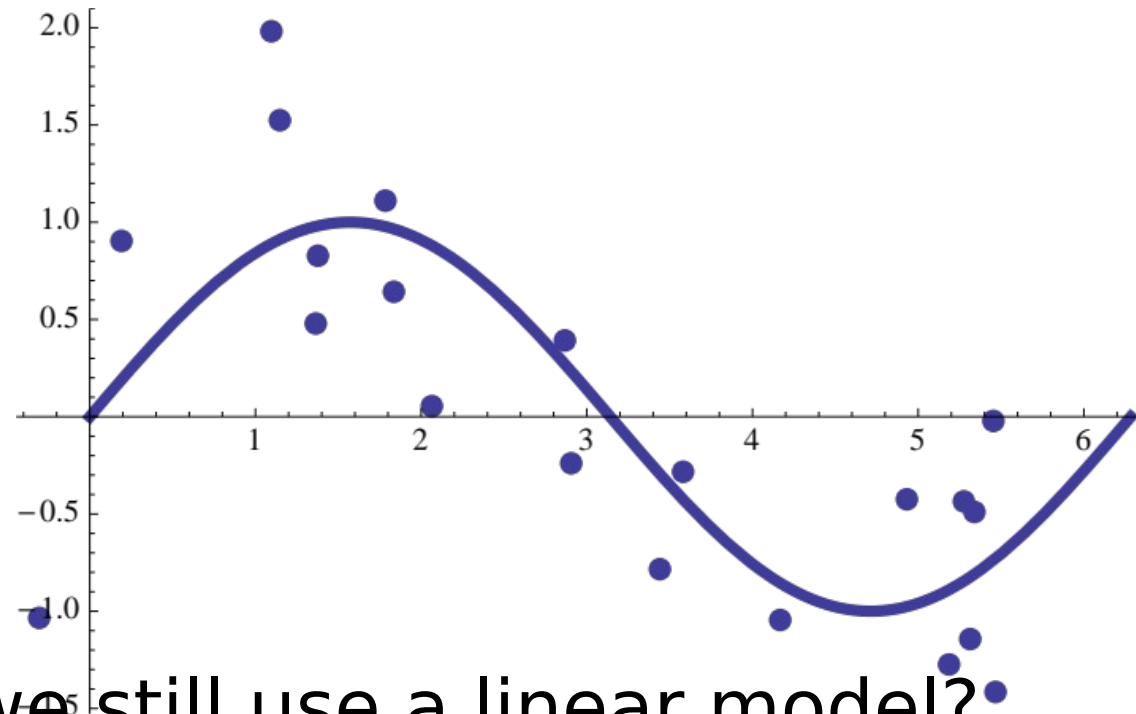


# Error Visualization



# Fitting Non-linear Data

- What if  $Y$  has a non-linear response?



- Can we still use a linear model?

# Transforming the Feature Space

- Transform features  $x_i$

$$x_i = (X_{i,1}, X_{i,2}, \dots, X_{i,p})$$

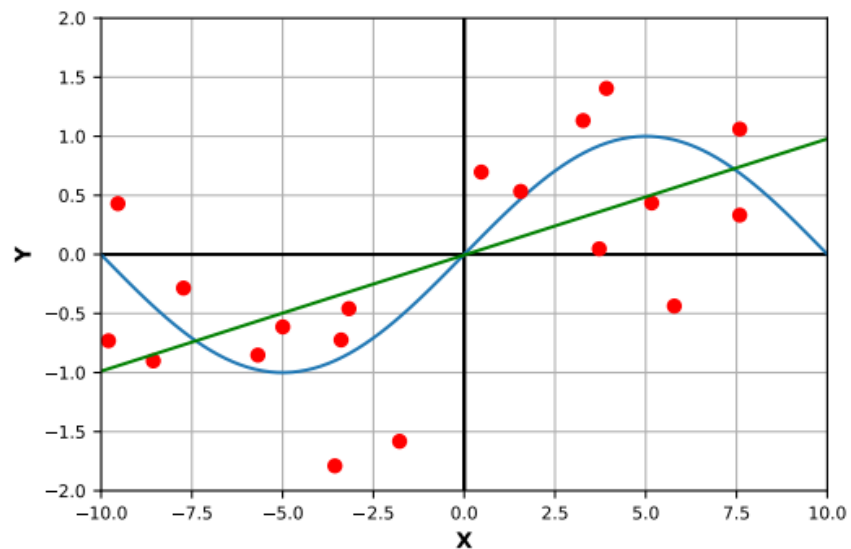
- By applying non-linear transformation  $\phi$ :

- Example:  $\phi : \mathbb{R}^p \rightarrow \mathbb{R}^k$

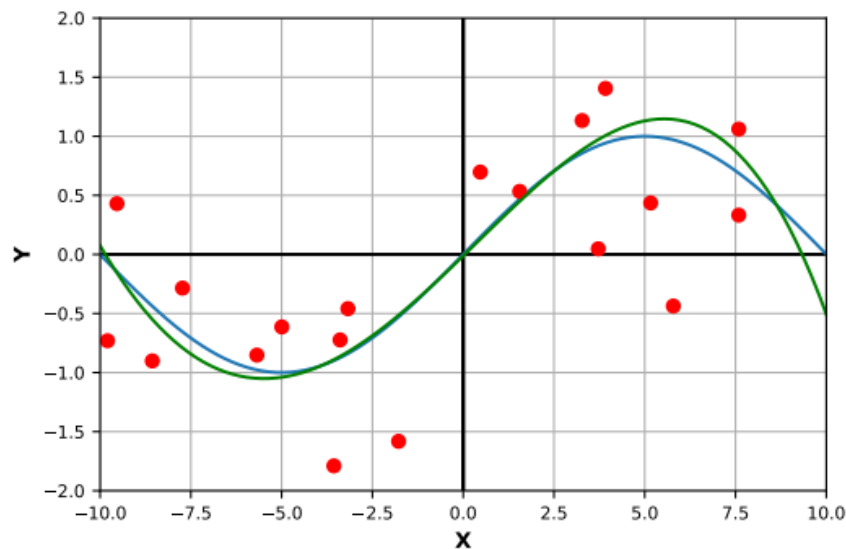
$$\phi(x) = \{1, x, x^2, \dots, x^k\}$$

- others: splines, radial basis functions, ...
- Expert engineered features (modeling)

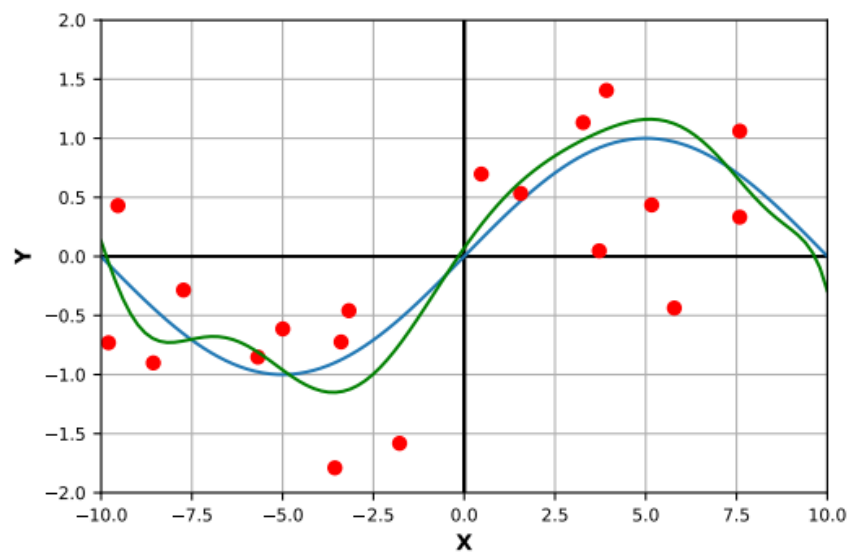
$$\{1, x\}$$



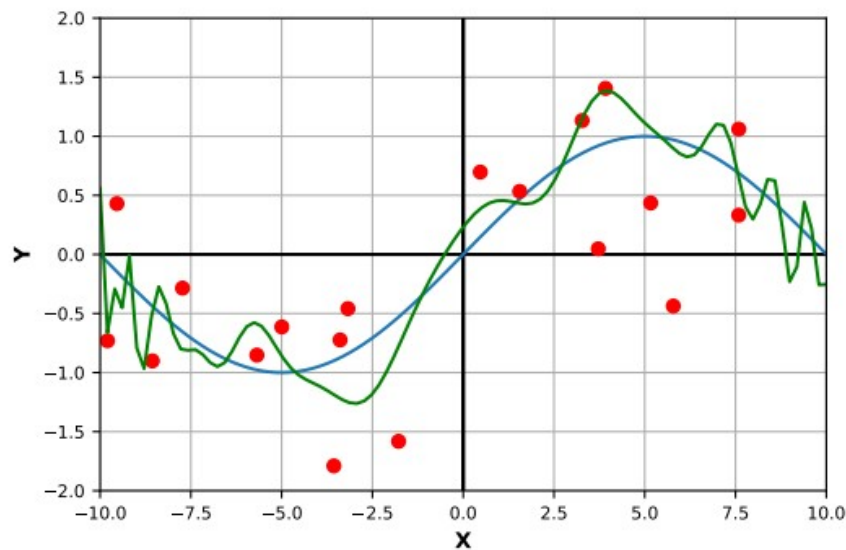
$$\{1, x, x^2, x^3, x^4\}$$



$$\{1, x, x^2, \dots, x^9, x^{10}\}$$



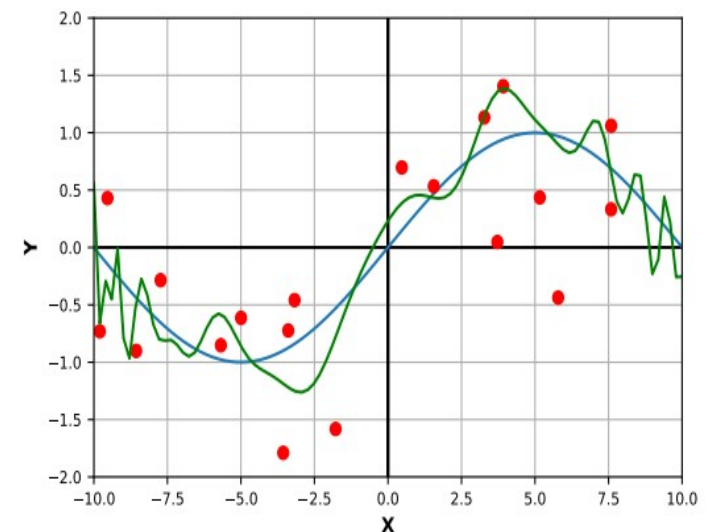
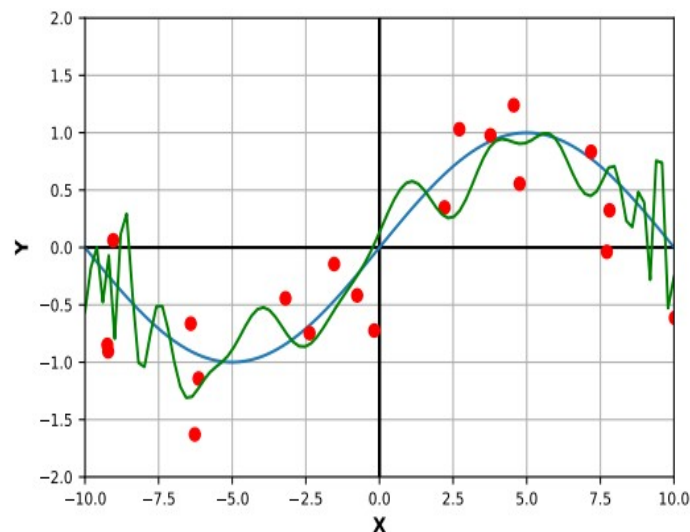
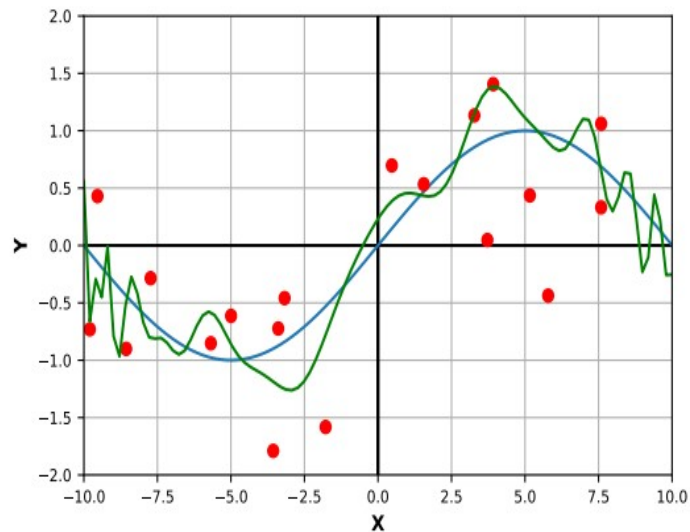
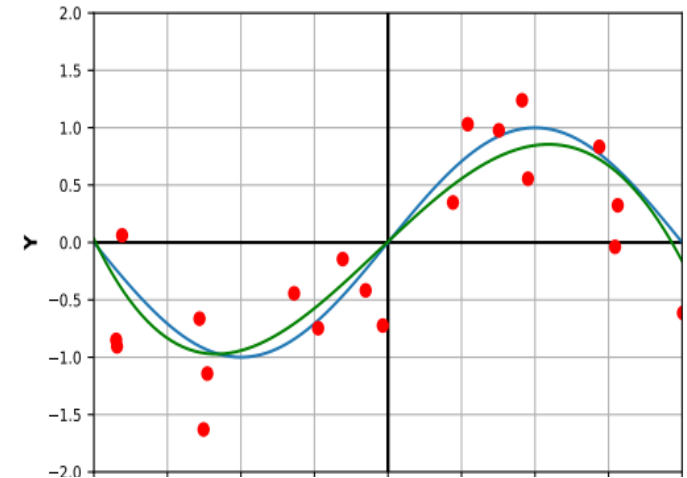
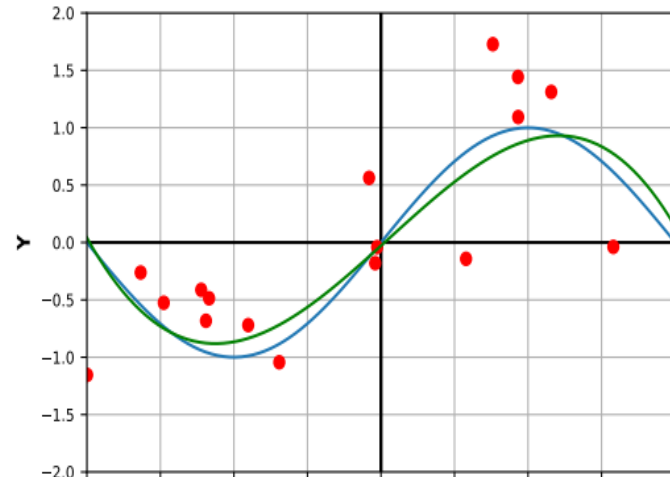
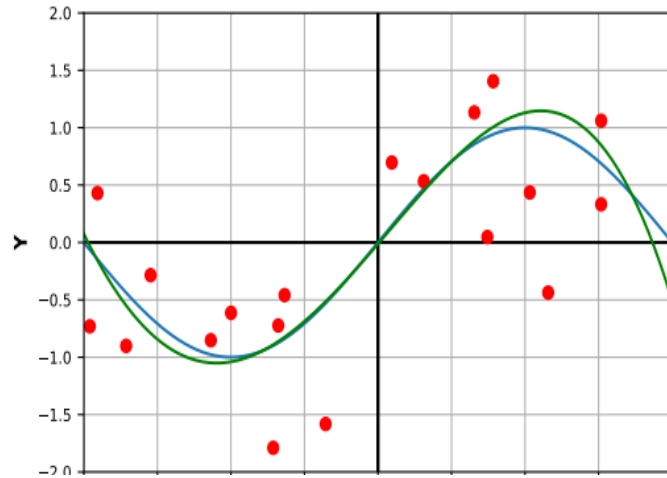
$$\{1, x, x^2, \dots, x^{99}, x^{100}\}$$



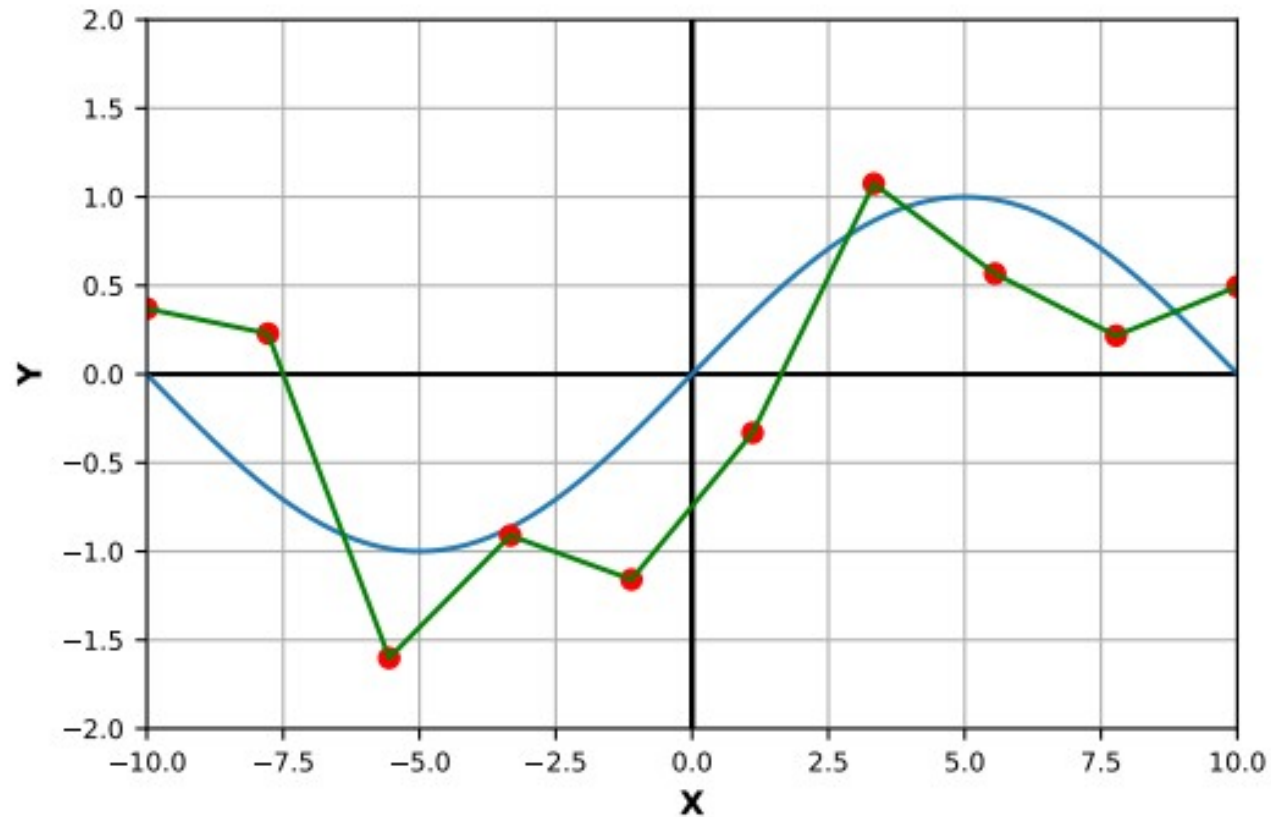
# What is Bias and Variance?



$$\{1, x, x^2, x^3, x^4\}$$



# Real Bad Overfit?



# Bias-Variance Tradeoff



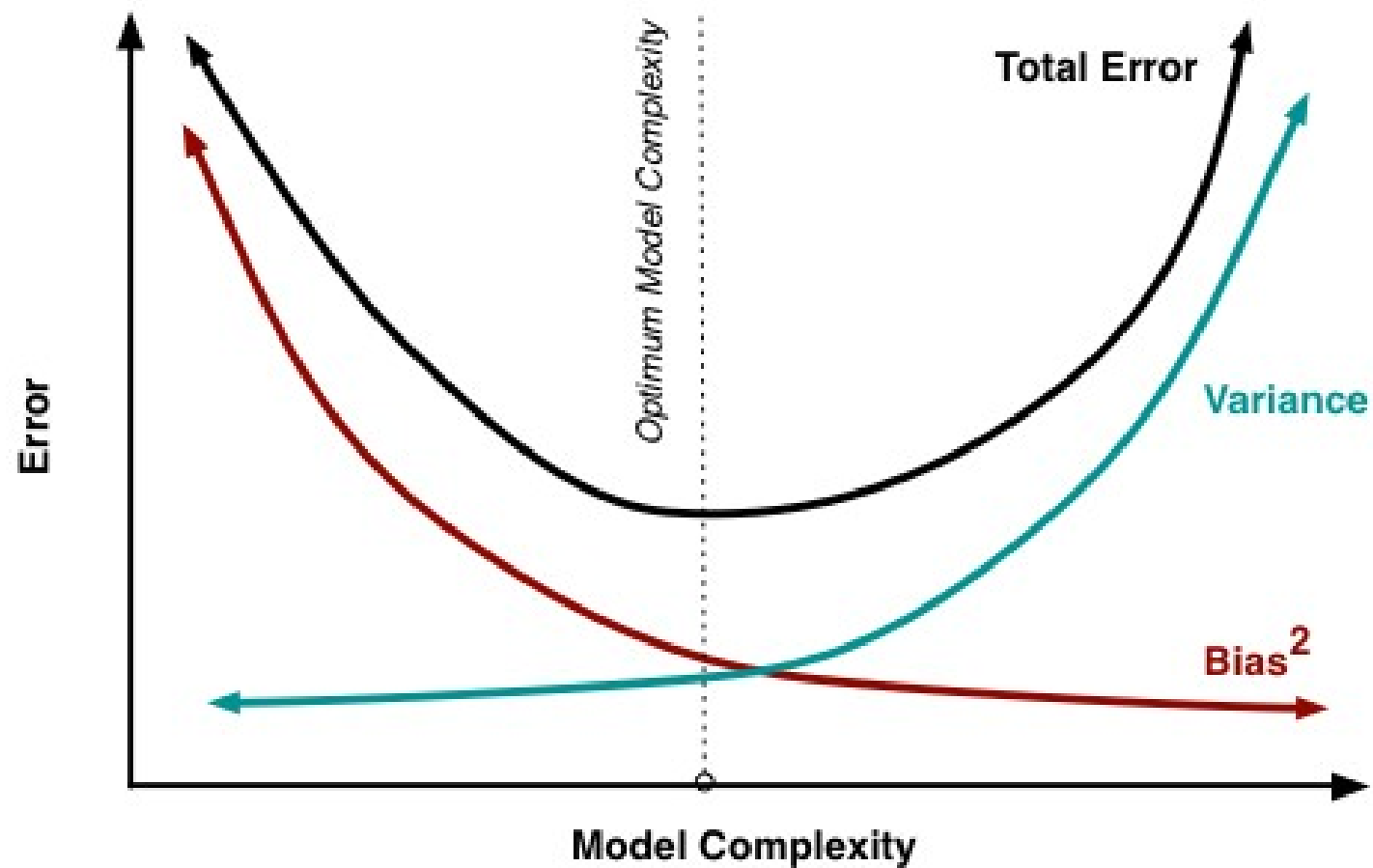
- So far we have minimized the error (loss) with respect to **training data**
  - Low training error does not imply good expected performance: **over-fitting**
- We would like to reason about the **expected loss (Prediction Risk)** over:
  - Training Data:  $\{(y_1, x_1), \dots, (y_n, x_n)\}$
  - Test point:  $(y_*, x_*)$
- We will decompose the expected loss into:

$$\mathbf{E}_{D, (y_*, x_*)} [(y_* - f(x_*|D))^2] = \text{Noise} + \text{Bias}^2 + \text{Variance}$$



# Bias-Variance Trade Off (Intuition)

# Bias Variance Plot



# Bias Variance Application in Training





# Regularization

# Regularization: An Overview



The idea of regularization revolves around modifying the loss function  $L$ ; in particular, we add a regularization term that penalizes some specified properties of the model parameters

$$L_{reg}(\beta) = L(\beta) + \lambda R(\beta),$$

where  $\lambda$  is a scalar that gives the weight (or importance) of the regularization term.

Fitting the model using the modified loss function  $L_{reg}$  would result in model parameters with desirable properties (specified by  $R$ ).

Since we wish to discourage extreme values in model parameter, we need to choose a regularization term that penalizes parameter magnitudes. For our loss function, we will again use MSE.

Together our regularized loss function is:

$$L_{LASSO}(\beta) = \frac{1}{n} \sum_{i=1}^n |y_i - \beta^\top \mathbf{x}_i|^2 + \lambda \sum_{j=1}^J |\beta_j|.$$

Note that  $\sum_{j=1}^J |\beta_j|$  is the  $\ell_1$  norm of the vector  $\beta$

$$\sum_{j=1}^J |\beta_j| = \|\beta\|_1$$

Alternatively, we can choose a regularization term that penalizes the squares of the parameter magnitudes. Then, our regularized loss function is:

$$L_{Ridge}(\beta) = \frac{1}{n} \sum_{i=1}^n |y_i - \beta^\top \mathbf{x}_i|^2 + \lambda \sum_{j=1}^J \beta_j^2.$$

Note that  $\sum_{j=1}^J |\beta_j|^2$  is the square of the  $\ell_2$  norm of the vector  $\beta$

$$\sum_{j=1}^J \beta_j^2 = \|\beta\|_2^2$$



In both ridge and LASSO regression, we see that the larger our choice of the **regularization parameter**  $\lambda$ , the more heavily we penalize large values in  $\beta$ ,

- If  $\lambda$  is close to zero, we recover the MSE, i.e. ridge and LASSO regression is just ordinary regression.
- If  $\lambda$  is sufficiently large, the MSE term in the regularized loss function will be insignificant and the regularization term will force  $\beta_{\text{ridge}}$  and  $\beta_{\text{LASSO}}$  to be close to zero.

To avoid ad-hoc choices, we should select  $\lambda$  using cross-validation.

Solution to ridge regression:

$$\beta = (X^T X + \lambda I)^{-1} X^T Y$$

The solution to the LASSO regression:

LASSO has no conventional analytical solution, as the L1 norm has no derivative at 0. We can, however, use the concept of **subdifferential** or **subgradient** to find a manageable expression. See a-sec2 for details.

The solution of the Ridge/Lasso regression involves three steps:

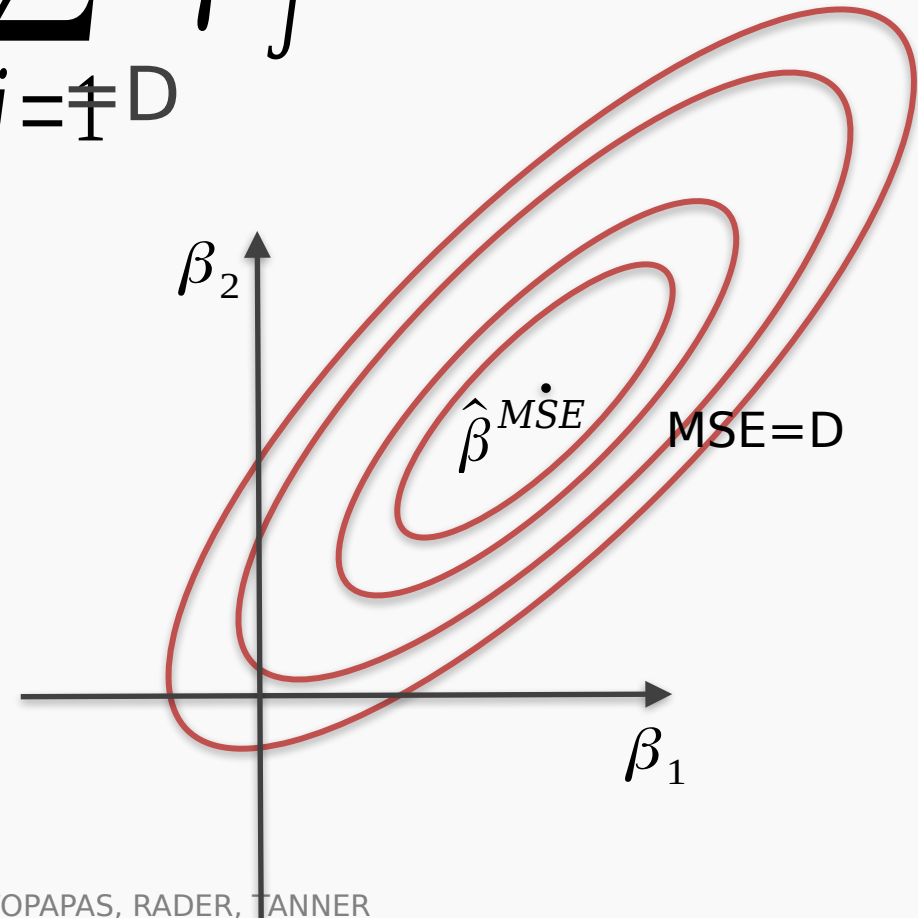
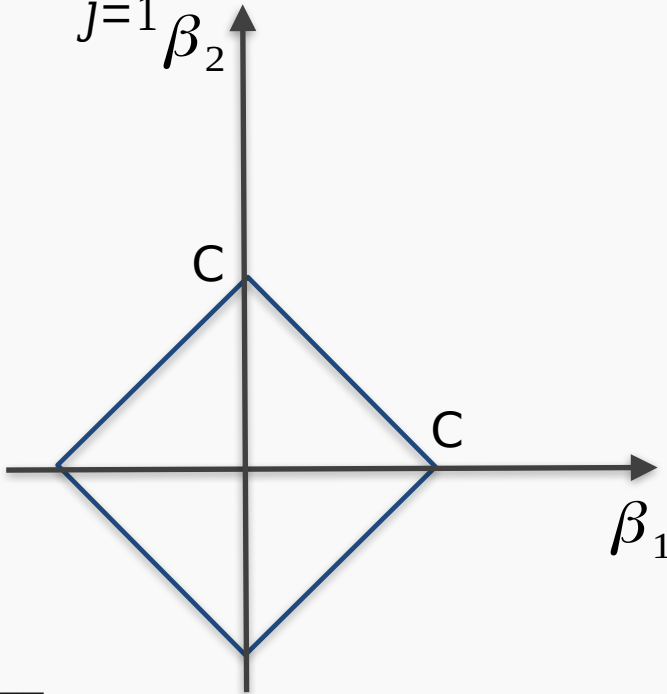
- Select  $\lambda$
- Find the minimum of the ridge/Lasso regression loss function (using the formula for ridge) and record the ***MSE on the validation set***.
- Find the  $\lambda$  that gives the smallest *MSE*

# The Geometry of Regularization (LASSO)

$$L(\beta) = \frac{1}{2} \sum_{i=1}^n |y_i - \beta^T x_i|^2 + \lambda \sum_{j=1}^J |\beta_j|$$

$$\hat{\beta}^{LASSO} = \operatorname{argmin}_{\beta} L_{LASSO}(\beta)$$

$$\lambda \sum_{j=1}^J |\hat{\beta}_j^{LASSO}| = C$$



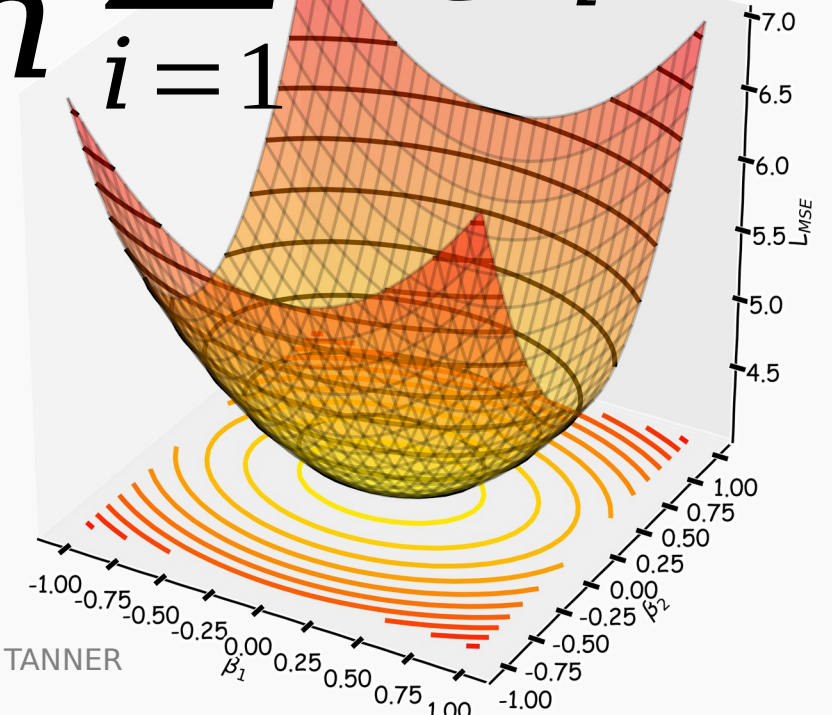
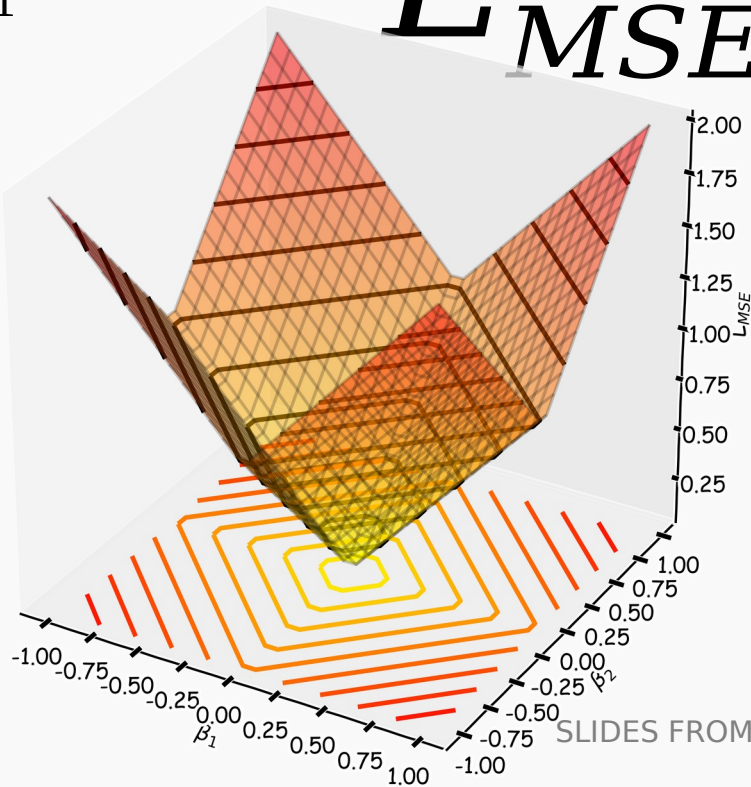
# The Geometry of Regularization (LASSO)

$$L_0(\beta) = \frac{1}{n} \sum_{i=1}^n |y_i - \beta^T x_i|^2 + \lambda \sum_{j=1}^J |\beta_j|$$

$$\hat{\beta}^{\text{LASSO}} = \operatorname{argmin}_{\beta} L_{\text{LASSO}}(\beta)$$

$$L_1 = \lambda \sum_{j=1}^J |\hat{\beta}_j^{\text{LASSO}}|$$

$$L_{\text{MSE}}(\beta) = \frac{1}{n} \sum_{i=1}^n |y_i - \beta^T x_i|^2$$

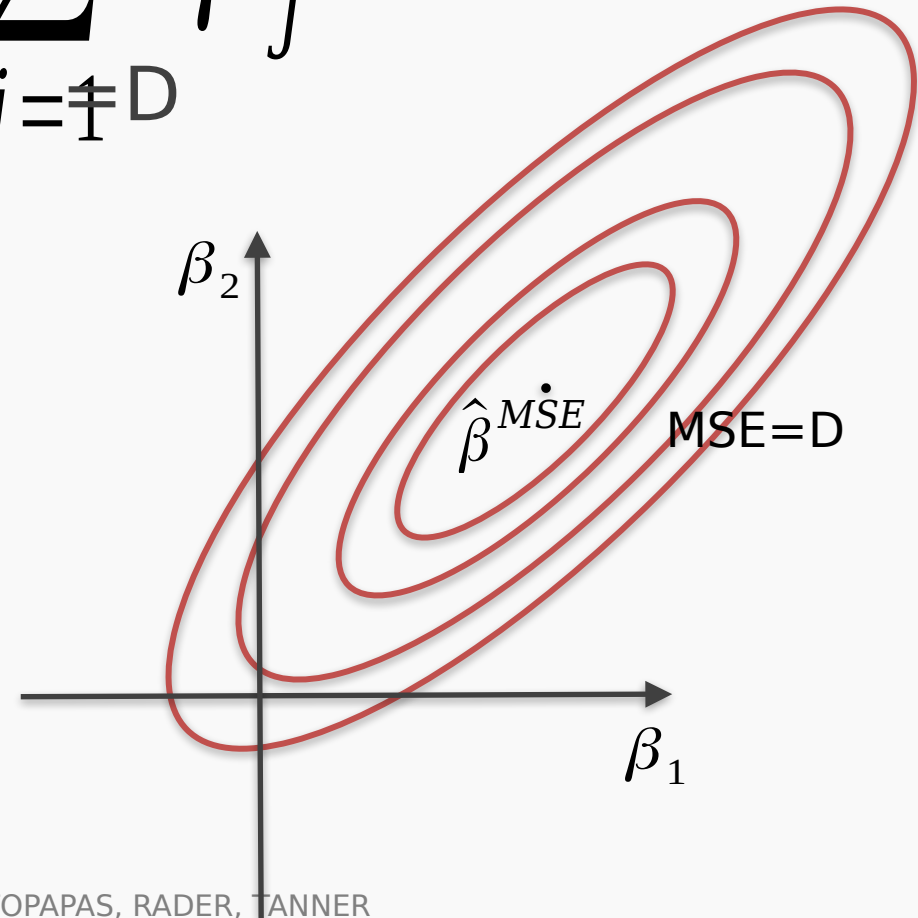
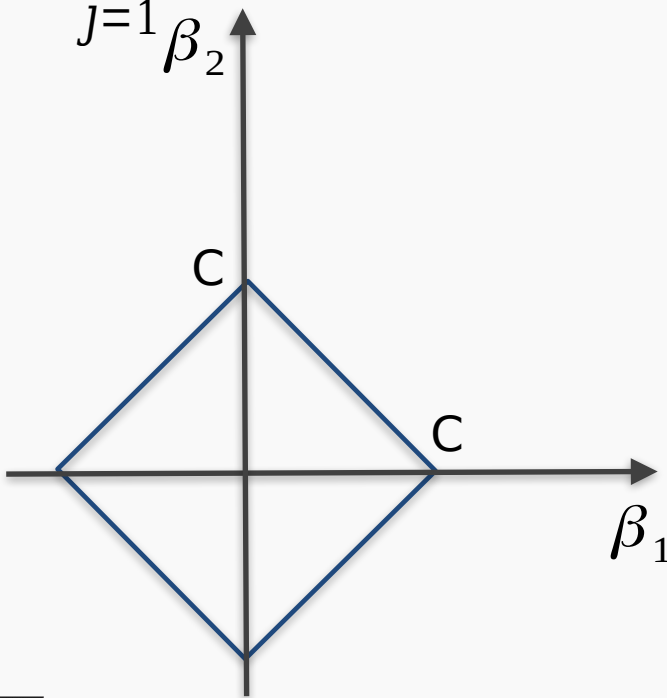


# The Geometry of Regularization (LASSO)

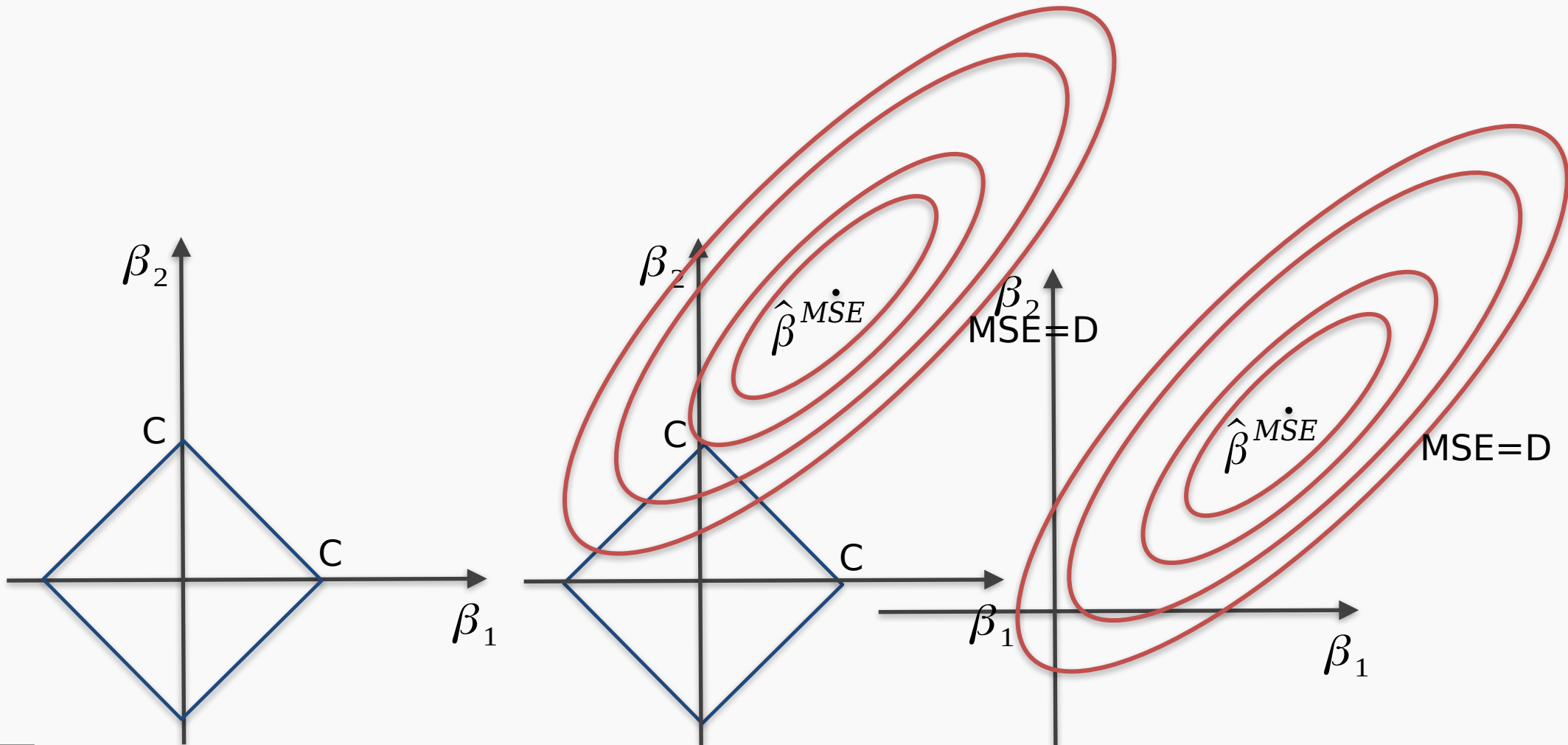
$$L(\beta) = \frac{1}{2} \sum_{i=1}^n |y_i - \beta^T x_i|^2 + \lambda \sum_{j=1}^J |\beta_j|$$

$$\hat{\beta}^{LASSO} = \operatorname{argmin}_{\beta} L_{LASSO}(\beta)$$

$$\lambda \sum_{j=1}^J |\hat{\beta}_j^{LASSO}| = C$$



# The Geometry of Regularization (LASSO)



SLIDES FROM: CS109A, PROTOPAPAS, RADER, TANNER



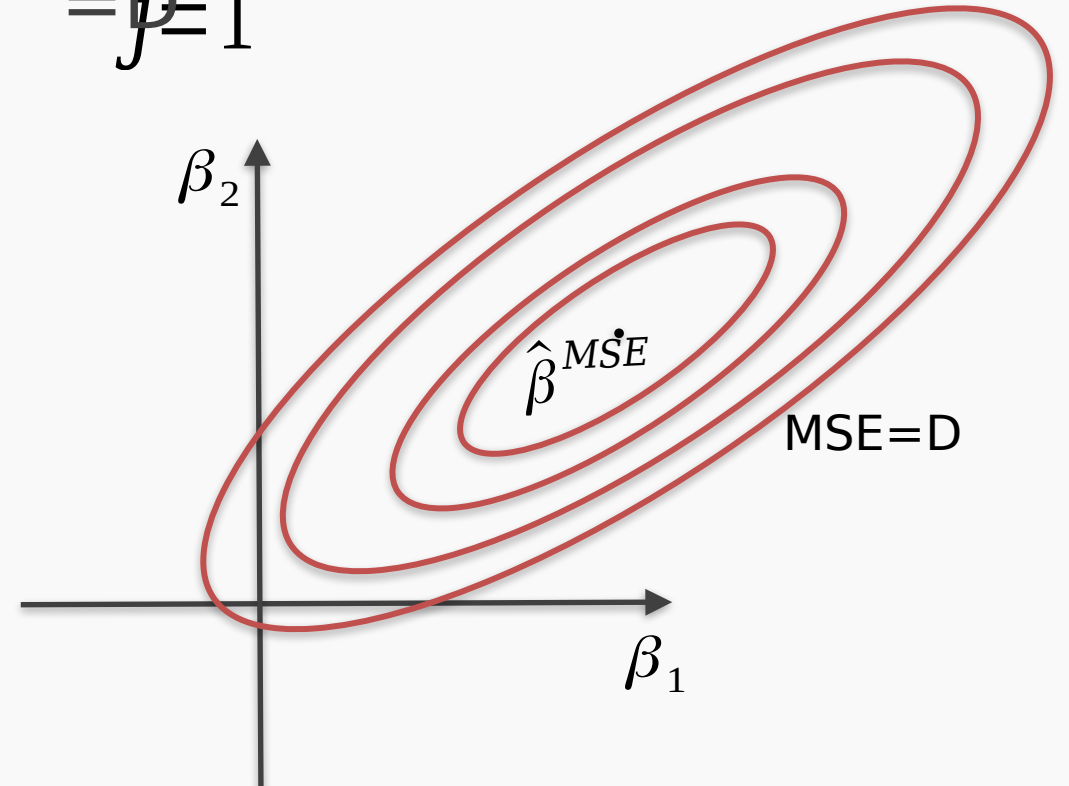
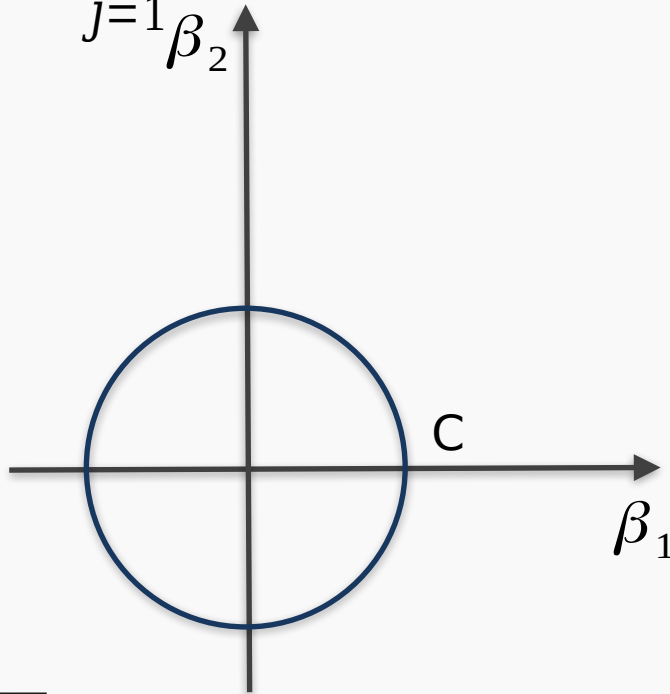
# The Geometry of Regularization (Ridge)



$$L_{\text{Ridge}}(\beta) = \frac{1}{n} \sum_{i=1}^n |y_i - \beta^T x_i|^2 + \lambda \sum_{j=1}^J \beta_j^2$$

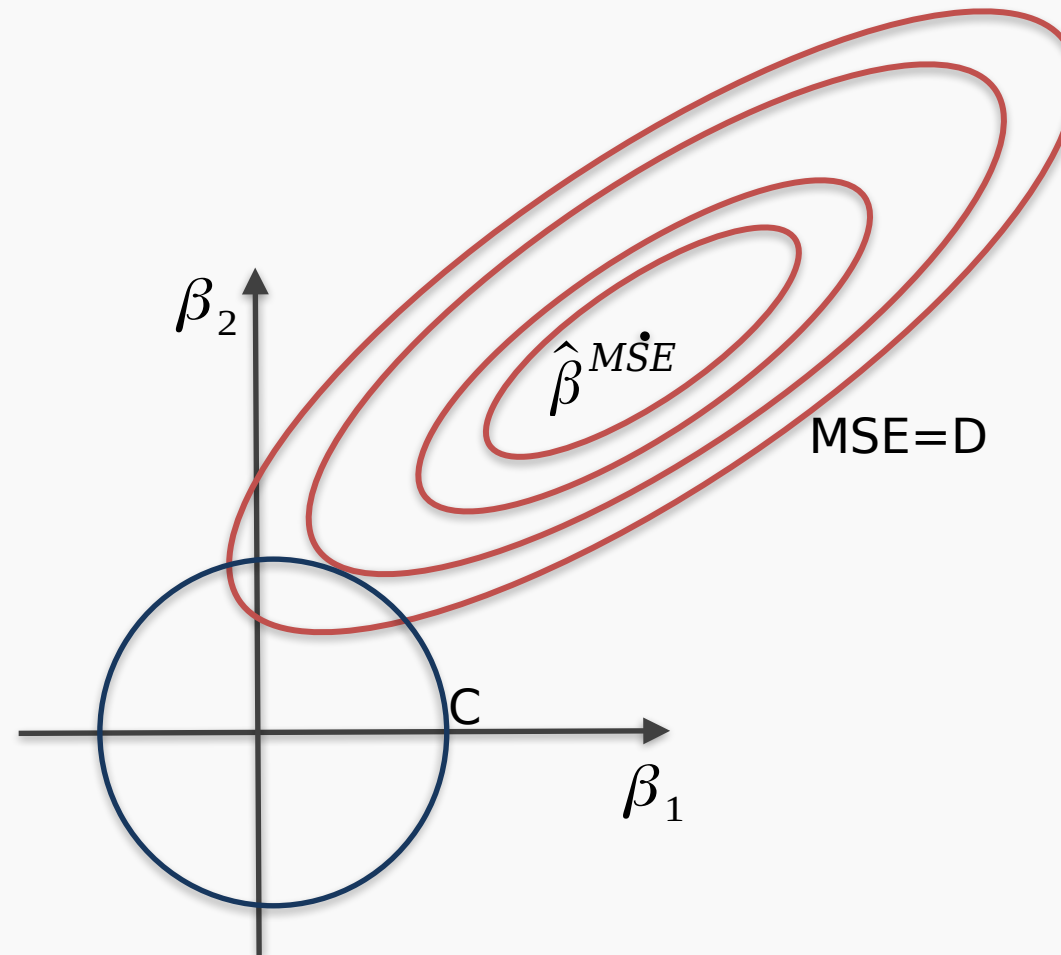
$$\hat{\beta}^{\text{Ridge}} = \underset{\beta}{\text{argmin}} L_{\text{Ridge}}(\beta)$$

$$\lambda \sum_{j=1}^J |\beta_j^{\text{Ridge}}|^2 = C$$



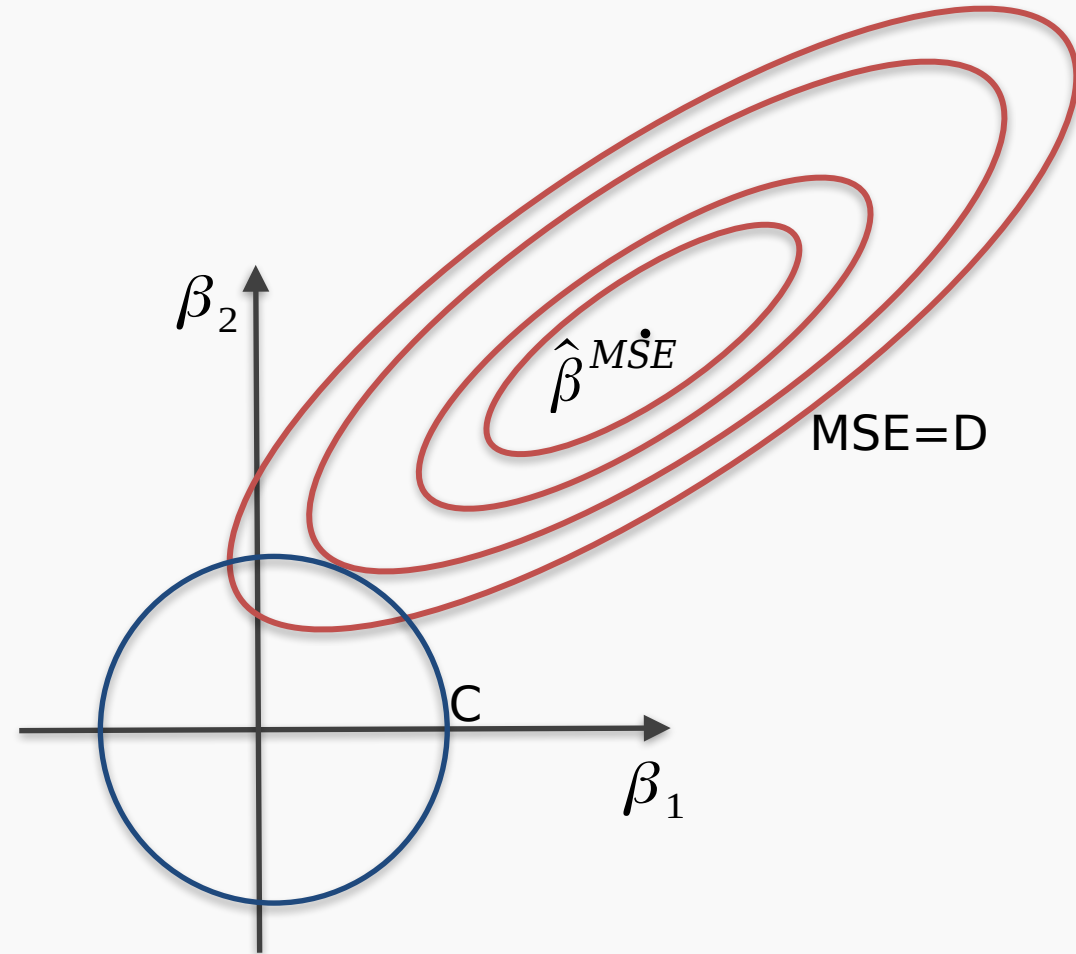
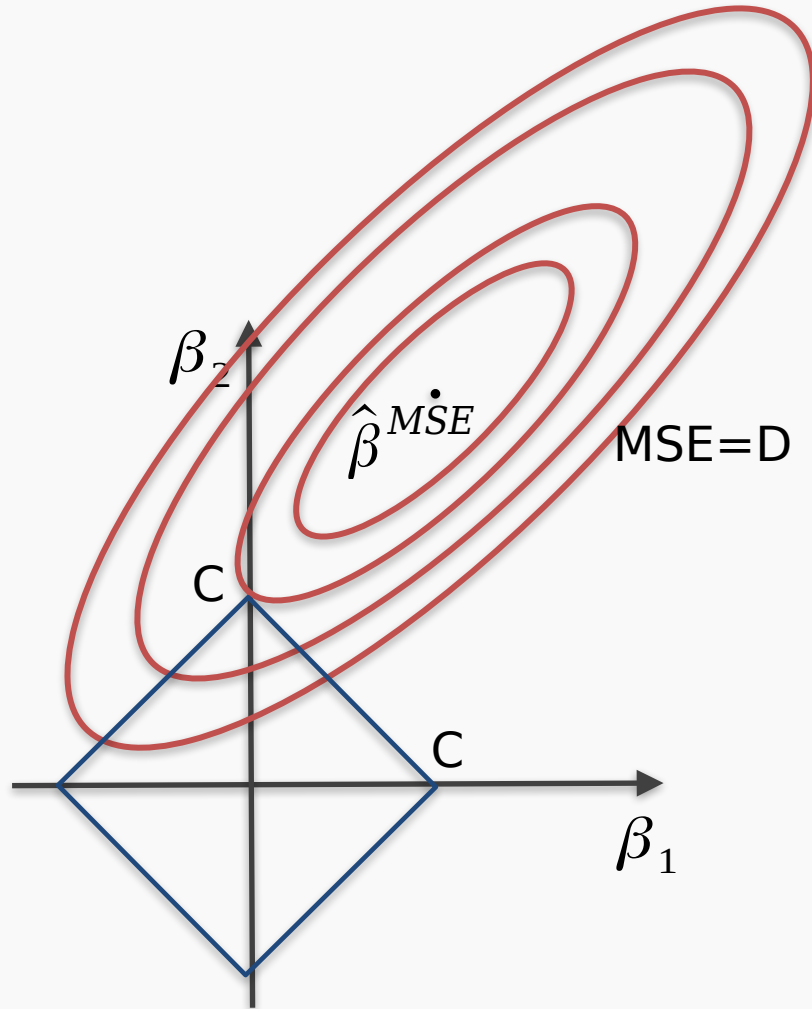


# The Geometry of Regularization (Ridge)



SLIDES FROM: CS109A, PROTOPAPAS, RADER, TANNER

# The Geometry of Regularization



# Examples

```
In [ ]: from sklearn.linear_model import Lasso
```

```
In [22]: lasso_regression = Lasso(alpha=1.0, fit_intercept=True)
lasso_regression.fit(np.vstack((X_train, X_val)), np.hstack((y_train, y_val)))

print('Lasso regression model:\n {} + {}^T . x'.format(lasso_regression.intercept_, lasso_regression.coef_))
```

```
Lasso regression model:
10.424895873901445 + [ 0.24482603  3.48164594  1.84836859 -0.06864603 -0.          -0.
-0.02249766 -0.          0.          0.          0.          0.          ]^T . x
```

```
In [ ]: from sklearn.linear_model import Ridge
```

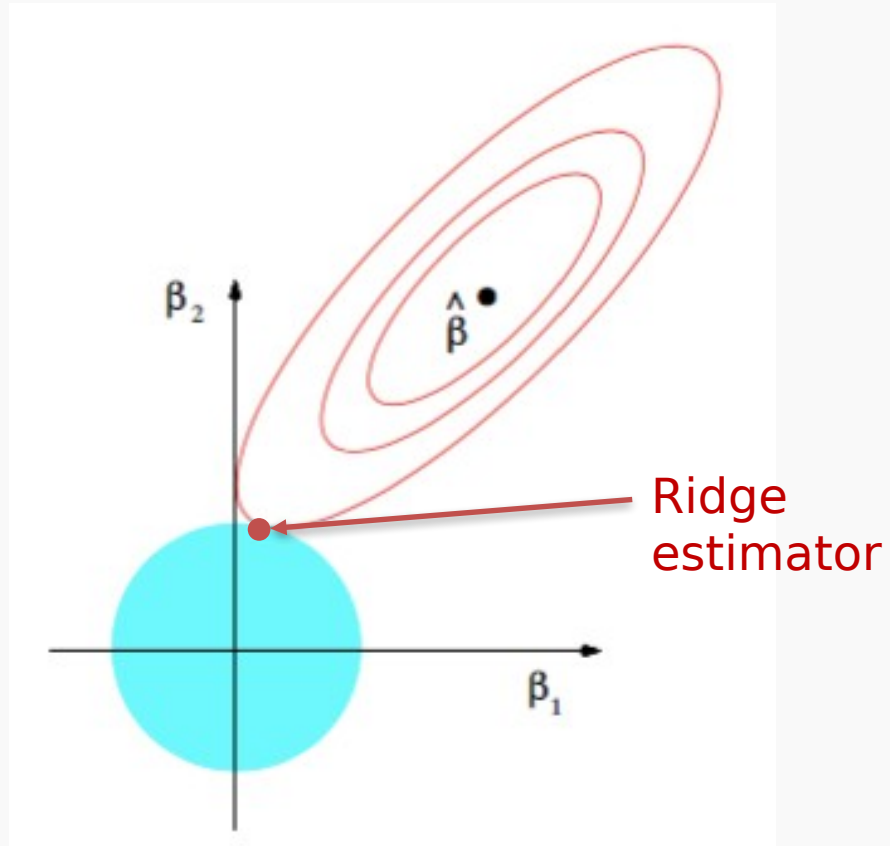
```
In [20]: X_train = train[all_predictors].values
X_val = validation[all_predictors].values
X_test = test[all_predictors].values

ridge_regression = Ridge(alpha=1.0, fit_intercept=True)
ridge_regression.fit(np.vstack((X_train, X_val)), np.hstack((y_train, y_val)))

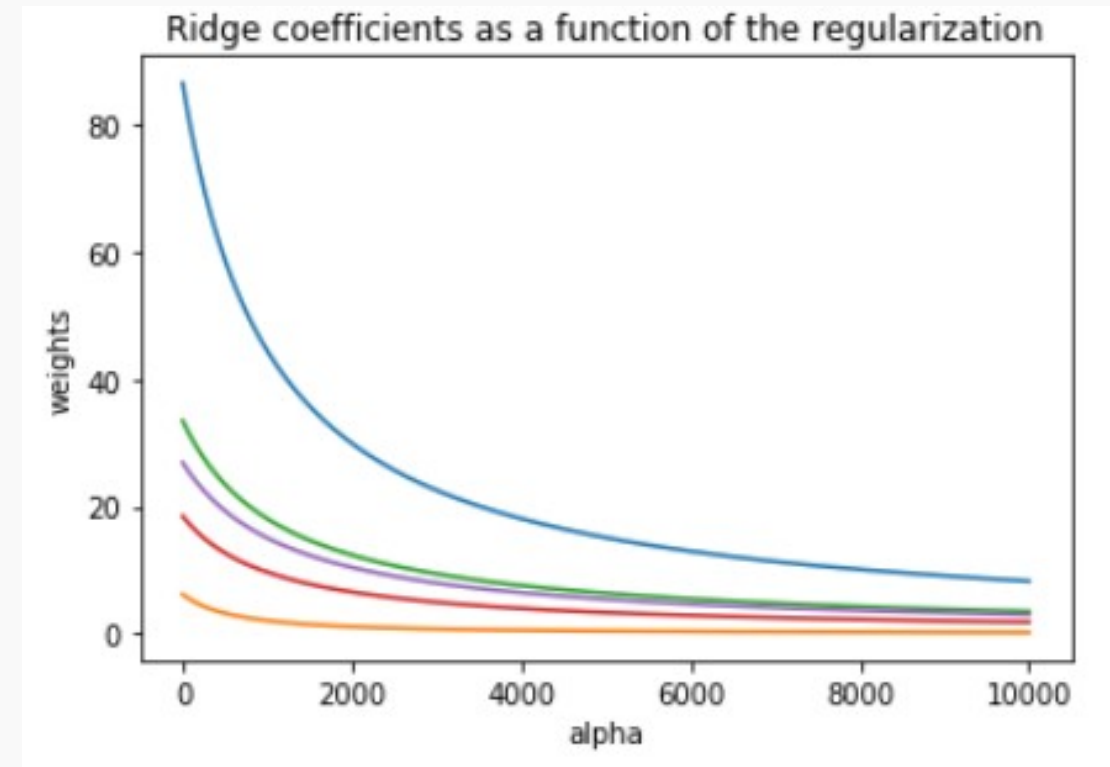
print('Ridge regression model:\n {} + {}^T . x'.format(ridge_regression.intercept_, ridge_regression.coef_))
```

```
Ridge regression model:
-525.7662550875951 + [ 0.24007312  8.42566029  2.04098593 -0.04449172 -0.01227935  0.41902475
-0.50397312 -4.47065168  4.99834262  0.          0.          0.29892679]^T . x
```

# Ridge visualized

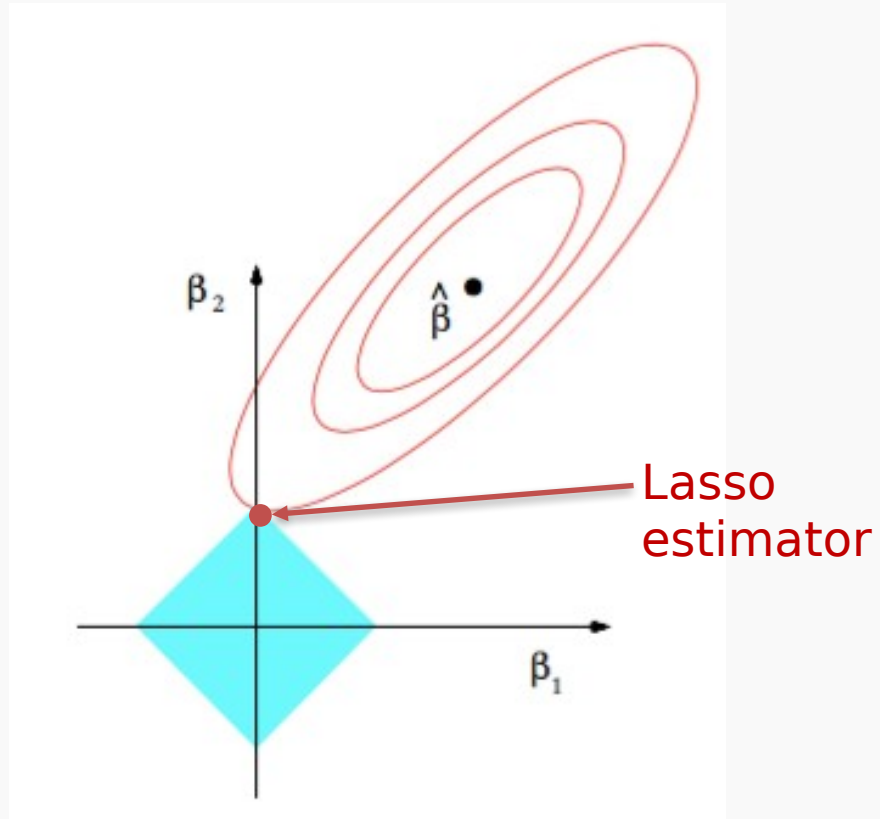


The ridge estimator is where the constraint and the loss intersect.



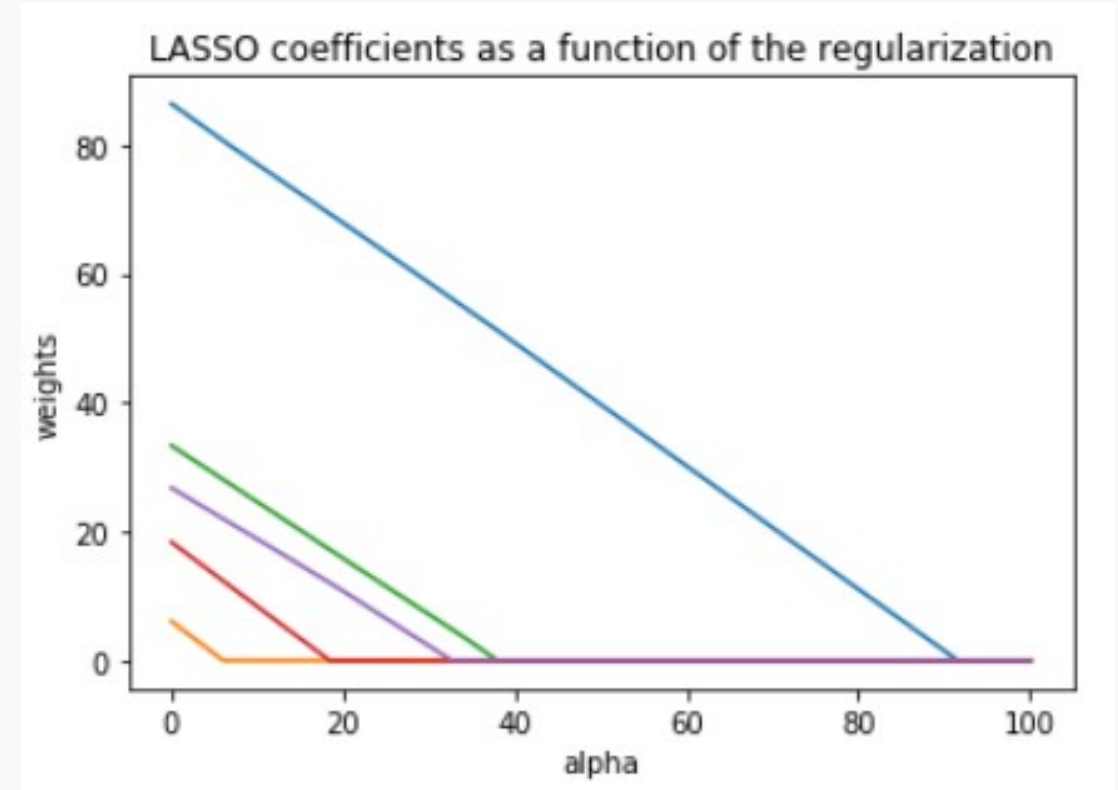
The values of the coefficients decrease as lambda increases, but they are not nullified.

# LASSO visualized



The Lasso estimator tends to zero out parameters as the OLS loss can easily intersect with the constraint on one of the axis.

SLIDES FROM: CS109A, PROTOPAPAS, RADER, TANNER



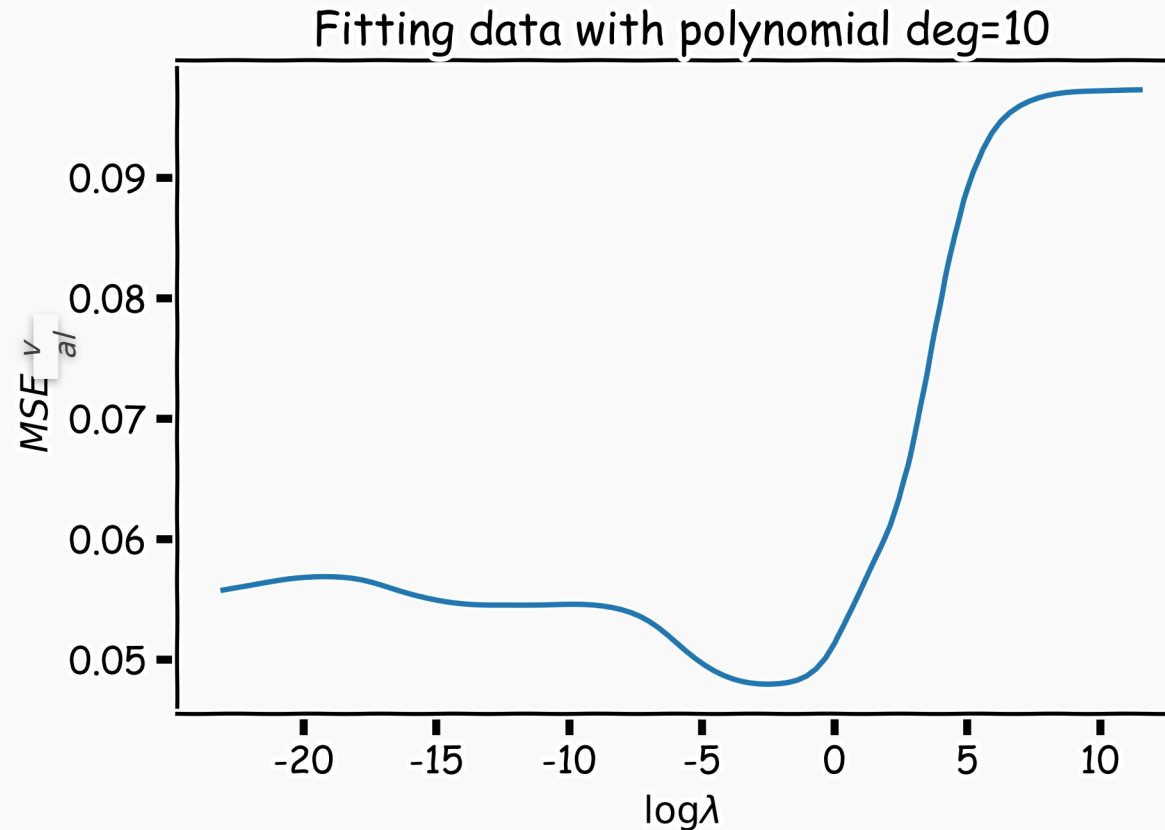
The values of the coefficients decrease as lambda increases, and are nullified fast.

# Ridge regularization with only **validation** : step by step



1. split data into
2. for
  1. determine the  $\lambda$  that minimizes the  $J(\lambda)$ , using the train data.
  2. record using validation data.
3. select the  $\lambda$  that minimizes the loss on the validation data,
4. Refit the model using both train and validation data,  $(X, y)$ , resulting to
5. report MSE or  $R^2$  on  $(X, y)$  given the

# Ridge regularization with **validation** only: step by step



# Lasso regularization with **validation** only: step by step



1. split data into
2. for
  - A. determine the  $\lambda$  that minimizes the  $\text{MSE}_{\text{train}}(\lambda)$ , using the train data. **This is done using a solver.**
  - B. record  $\text{MSE}_{\text{train}}(\lambda)$  using validation data
3. select the  $\lambda$  that minimizes the loss on the validation data,
4. Refit the model using both train and validation data,  $\lambda_{\text{min}}$ , resulting to
5. report MSE or  $R^2$  on  $\text{test}$  given the



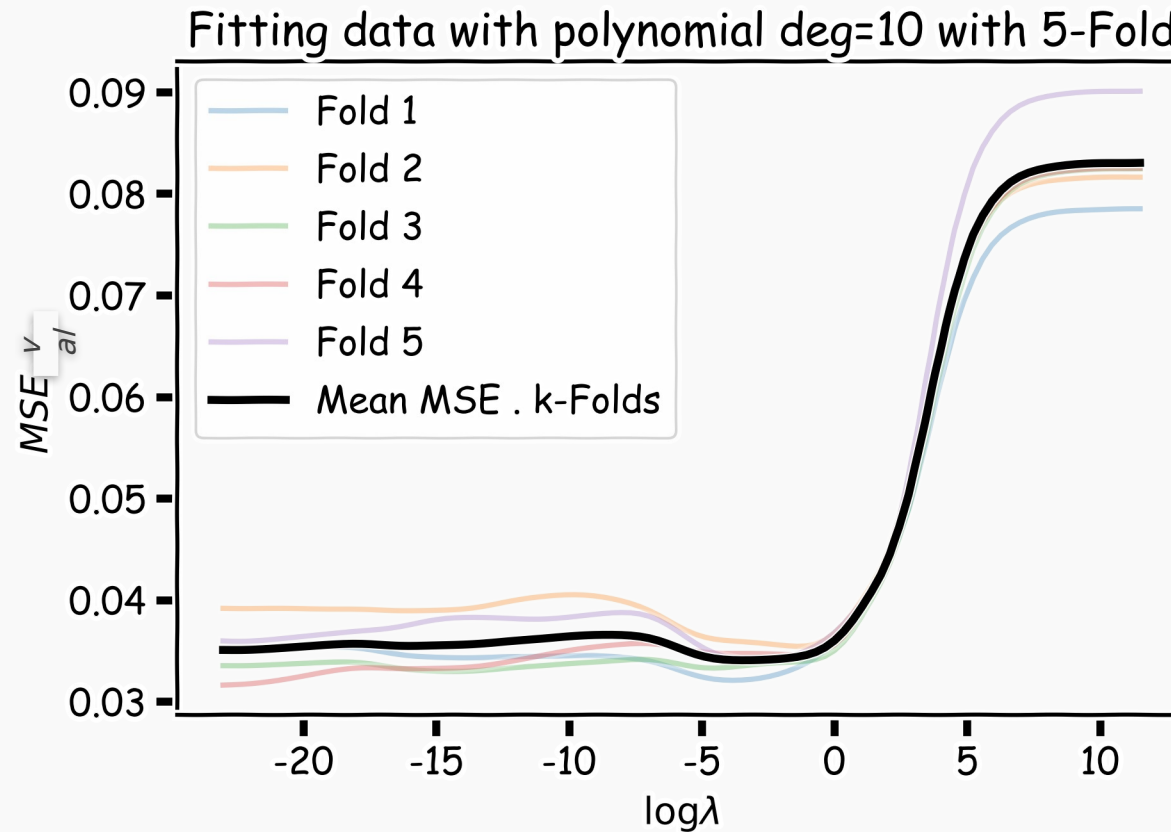
# Ridge regularization with **CV**: step by step

1. remove from data
2. split the rest of data into K folds,
3. for k in
  1. for
    - A. determine the that minimizes the , , using the train data of the fold, .
    - B. record using the validation data of the fold

At this point we have a 2-D matrix, rows are for different k, and columns are for different values.
4. Average the for each , .
5. Find the that minimizes the , resulting to .
6. Refit the model using the full training data, }, resulting to
7. report MSE or  $R^2$  on given the

	...			
	..	...		
	...	..	...	
...	..	...	..	...
	...	...	...	...
E[]	...			

# Ridge regularization with **validation** only: step by step



# Variable Selection as Regularization



Since LASSO regression tend to produce zero estimates for a number of model parameters - we say that LASSO solutions are **sparse** - we consider LASSO to be a method for variable selection.

Many prefer using LASSO for variable selection (as well as for suppressing extreme parameter values) rather than stepwise selection, as LASSO avoids the statistic problems that arises in stepwise selection.

**Question:** What are the pros and cons of the two approaches?

