

# Artificial Intelligence and Machine Learning

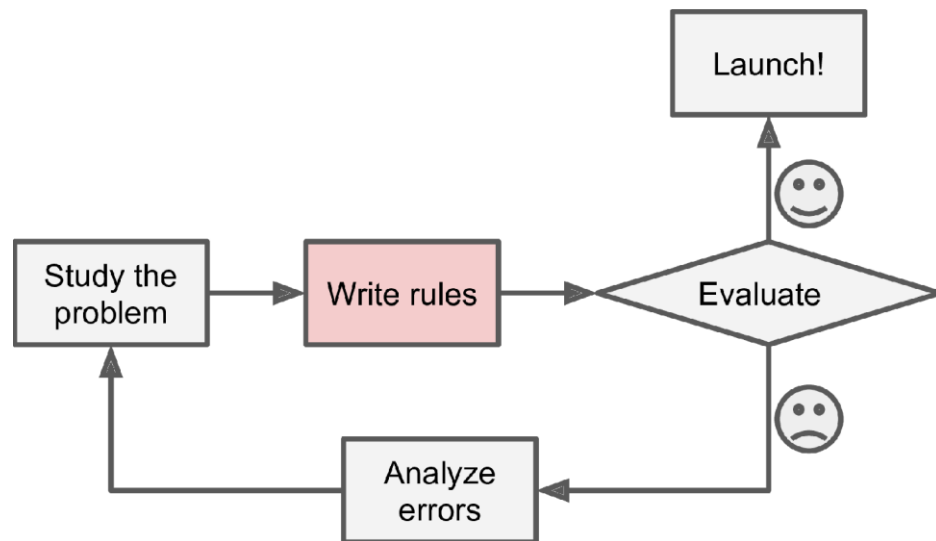
## Linear Regression

# Lecture 1: Outline

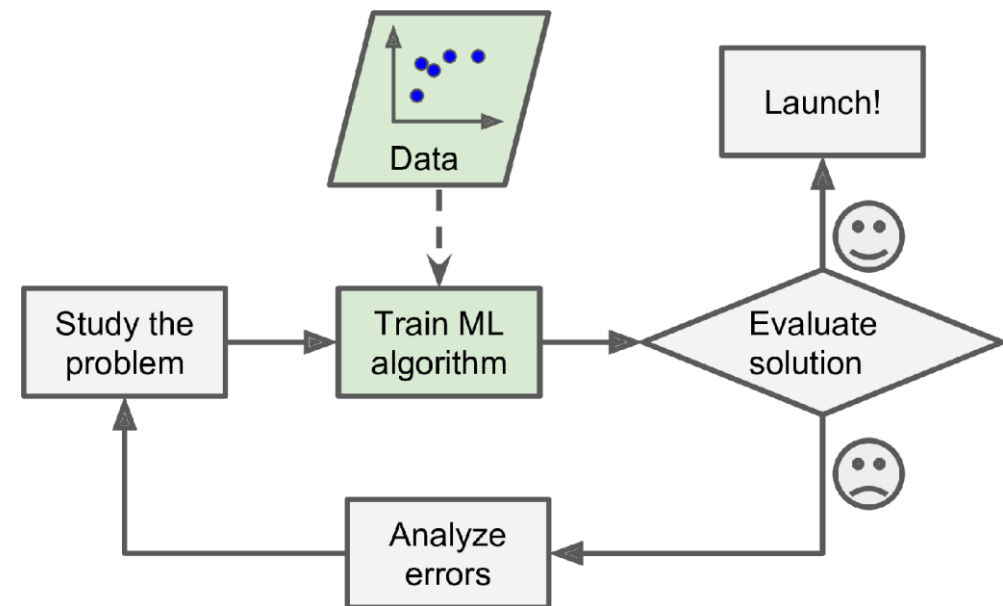
- Introduction to ML
- Linear Regression
- Optimization
- Linear Regression: Probabilistic Interpretation
- Bias-Variance Tradeoff
- Regularization

# Introduction to ML

- **Machine Learning** is the science (and art) of programming computers so they can learn from data.



*The traditional approach*



*The Machine Learning approach*

# Data Types

- **Tabular Data** (e.g., spreadsheets, databases)
  - Note: Columns are called **Features**. Rows are called **Samples**.
- **Time-Series Data** (e.g., stock prices, weather forecasts, IoT sensor data)
- **Text Data** (Natural Language Processing, e.g., emails, social media posts, documents)
- **Images and Videos** (Computer Vision, e.g., medical imaging, surveillance, facial recognition)
- **Audio Data** (Speech Recognition, Music Processing, e.g., voice commands, podcasts, sound classification)

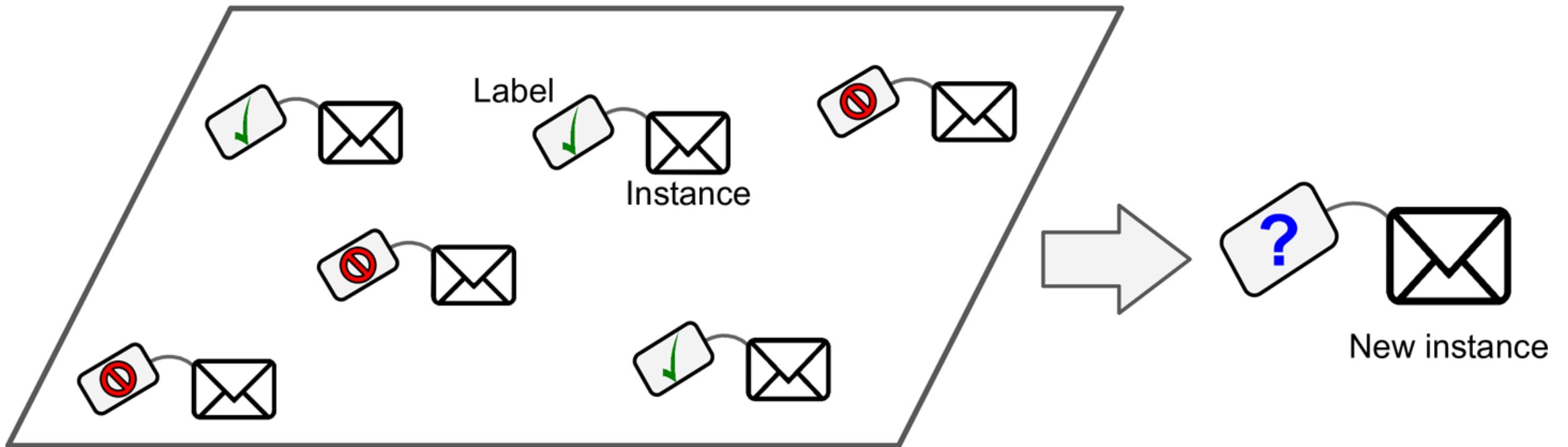
# ML Algorithms Types

- **Supervised:** The algorithm learns from labeled data.
  - **Regression:** Predict continuous value (e.g. house prices).
  - **Classification:** Predict discrete value (e.g. spam/not-spam).
- **Unsupervised:** The algorithm works on unlabeled data. We are interested in things like:
  - **Clustering:** Grouping
  - **Dimensionality Reduction:** Reducing the Dimensions
  - **Anomaly Detection:** Detecting outliers

# ML Algorithms Types cont.

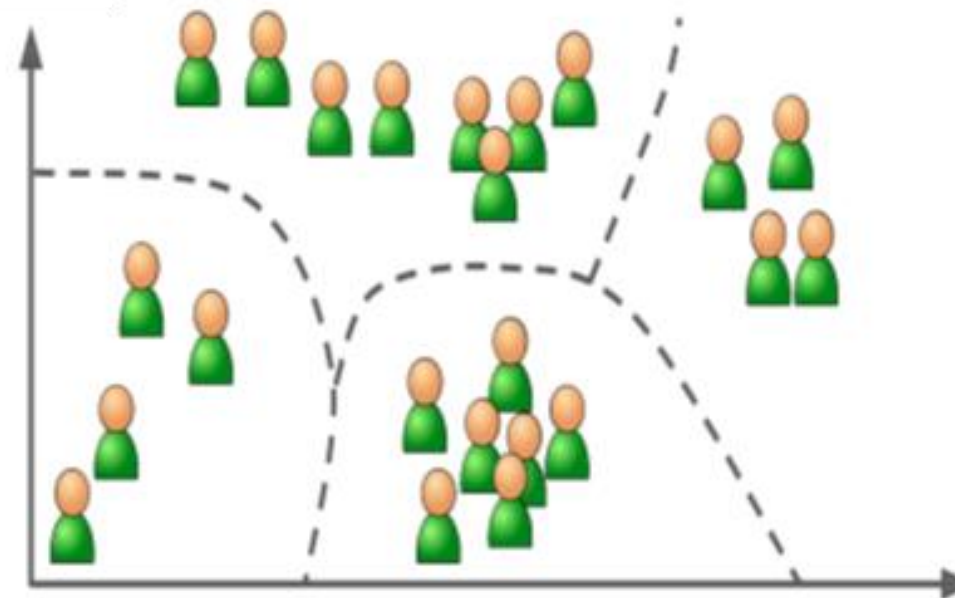
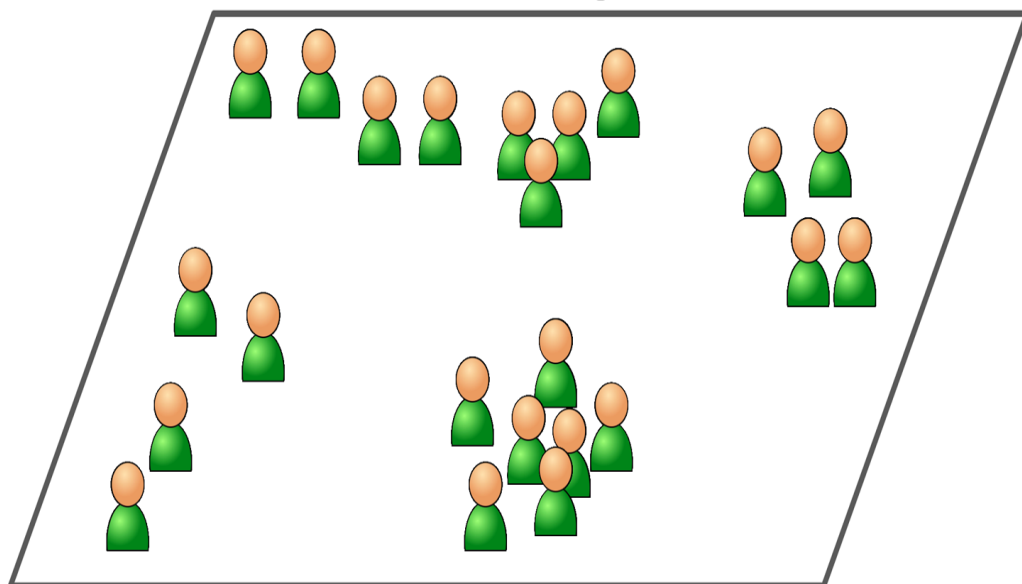
- **Reinforcement Learning:** involves learning to make decisions by interacting with an environment.
  - **Reward Signal:** The agent receives feedback in the form of rewards or penalties, guiding its learning.
  - **Policy:** A strategy the agent learns to decide actions based on the current state.
  - **Value Function:** An estimate of the expected cumulative reward from a state or state-action pair.
  - **Exploration vs Exploitation:** The agent balances exploring new actions to discover rewards and exploiting known actions to maximize them.
  - Really popular in video games and robotics!(Also recently in LLMs, see [RLHF](#))

## Training set



An example of Supervised Learning: Spam Classification

Training set



An example of Unsupervised Learning: Clustering



# How Does ML Work?

- Most of the ML systems consist of three main components:

**Hypothesis (Model):** The function that approximates the target.

- E.g. Linear Regression, Logistic Regression, SVM, Decision Trees, NN,...

**Optimizer:** The mechanism for improving predictions of our model.

**Loss Function:** The measure of how wrong the predictions are.

# How Does ML Work?

- How are they related to each other? 🤔

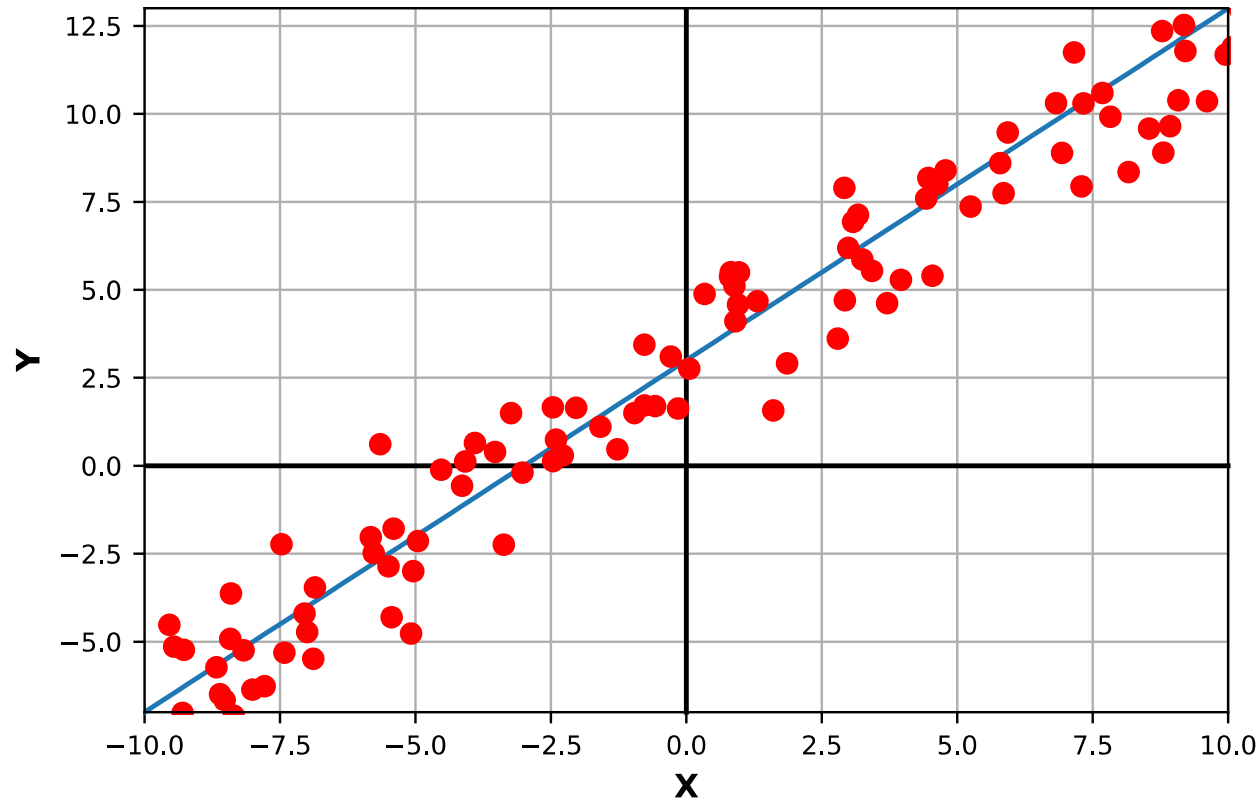
# How Does ML Work?

- We firstly define our task (classification/regression) then choose an appropriate **model**.
- We will use an **optimization method** to minimize the **loss function**.
- Reached a minima?
  - = Model is making the least possible number of mistakes.
  - = **Model trained** 🎉 .

# Linear Regression: 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

- **Hypothesis:**

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

- **Input:** data  $(x_i, y_i)$ ,  $i \in \{1, 2, \dots, N\}$

- (e.g., house size  $x$  and price  $y$ )

- **Goal:** learn values of variable  $(m, b)$

# Notation

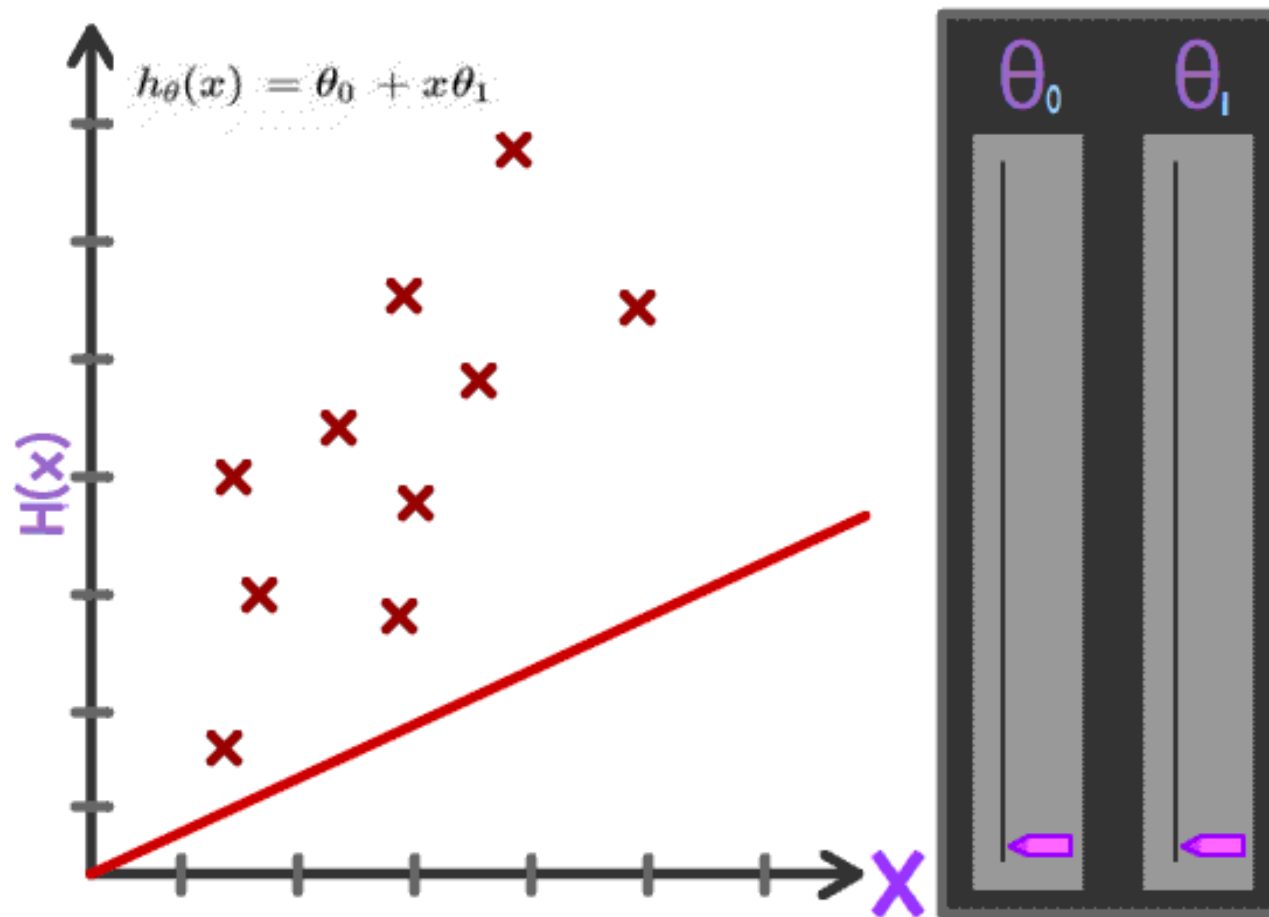
- Some clarification about the notation we will use for this course

$$X_i^{j,[k]}$$

- $i$  is the index of the data,  $j$  is the feature number, and  $k$  is the power.

# Solution Strategy for Solving the Problem

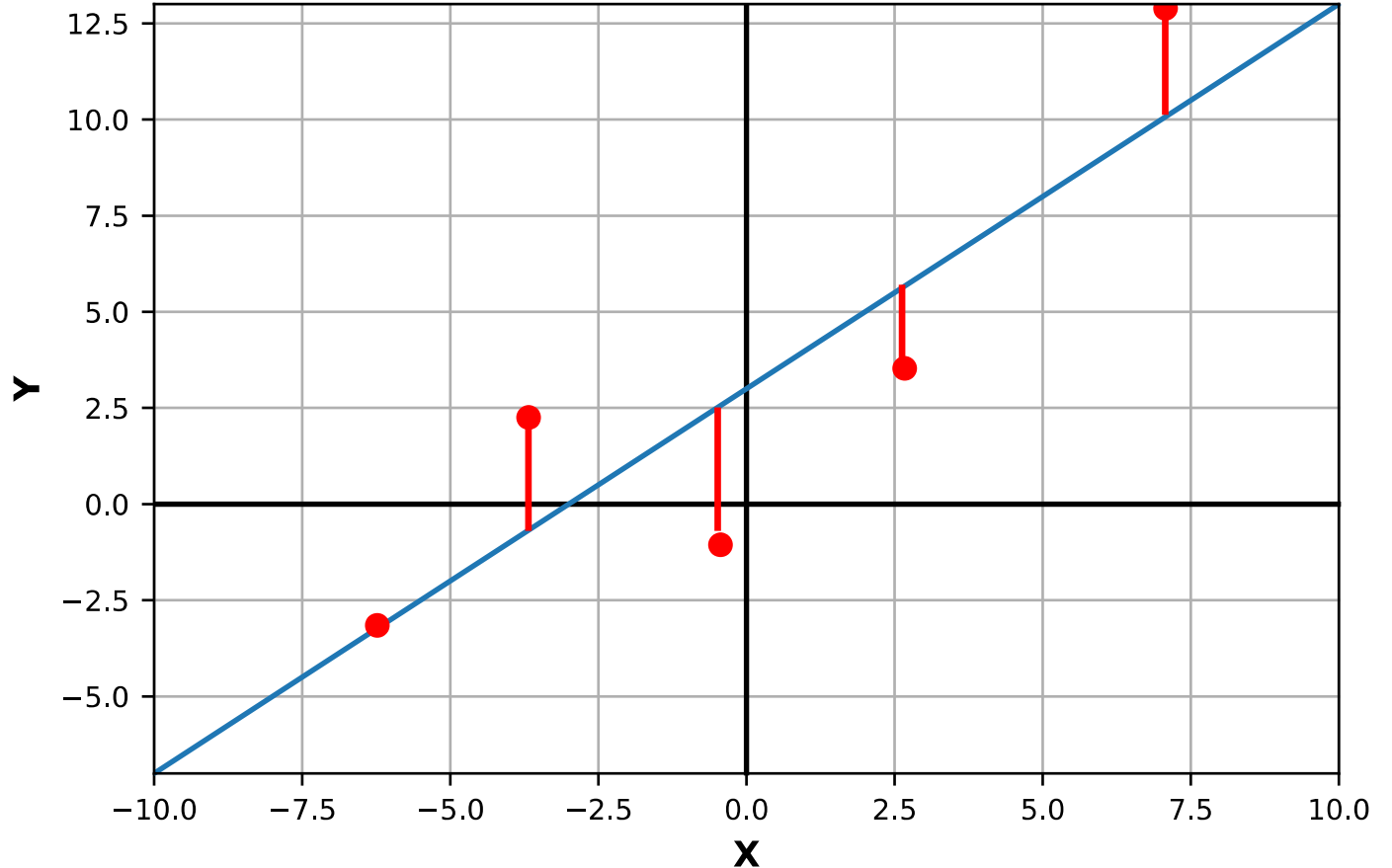
- There are countless possible lines.
- We want a line which is in some sense the “average line” that represents the data.
- Any ideas as to how we can do it?





# Optimization

- To find the "best line," we should minimize the distances between our line's predictions and all the data points.
- How to define that mathematically?



# Loss Function

- For one sample, this can be represented mathematically by:

$$(y_i - \hat{y}_i) \quad (\text{Error})$$

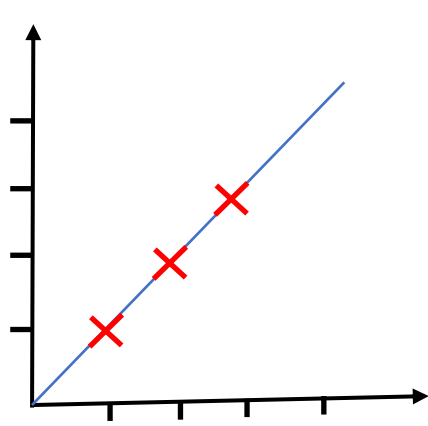
- But this could result in negative value if  $\hat{y} > y$ . Let's square it to remove the negative sign:

$$(y_i - \hat{y}_i)^2 \quad (\text{Squared Error})$$

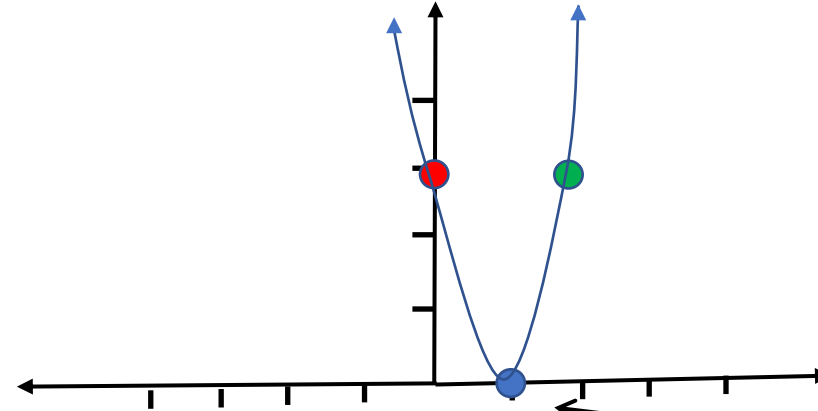
- But we have N samples, not only one. So, let's sum the errors and take the average:

$$\boxed{Loss (MSE) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} \quad (\text{Mean Squared Error})$$

# Intuition of Loss Function



Hypothesis



Loss Function

$$h(x) = mx$$

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

Notice: Lower is better.

$$J(m = 0) = 14$$

$$J(m = 1) = 0$$

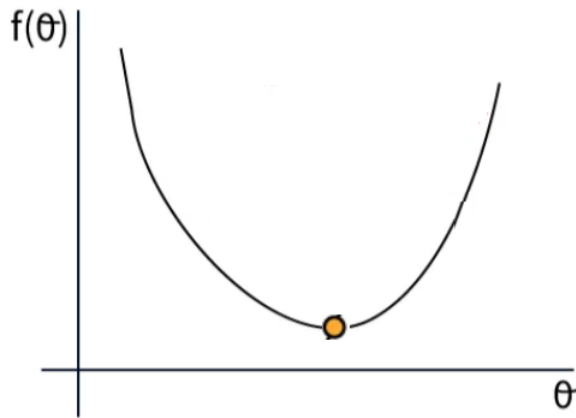
$$J(m = 2) = 14$$

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

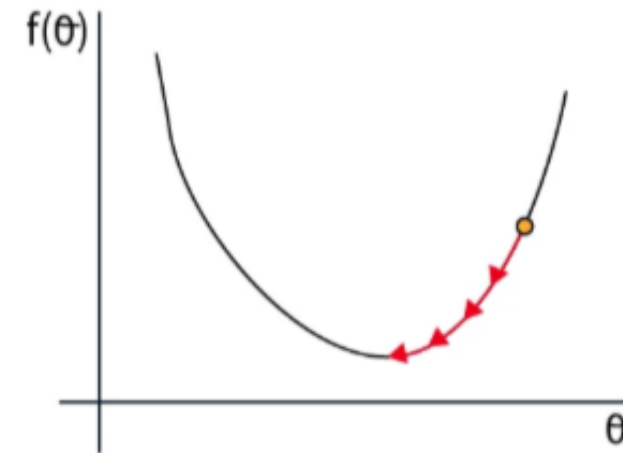
- There are two approaches to find the minima:
  - **Exact (Closed-form):** Directly calculates the solution mathematically by solving for  $f'(x) = 0$ .  
Important Note: can be used only with a very limited number of algorithms.
  - **Approximation (Iterative approach):** Gradually improves the solution step by step.  
Done by optimizers (e.g. Gradient Descent, ADAM,...etc).

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

Closed-form:



Iterative:



- Example:  $y = x^2$  (Solution:  $x = 0$ )
  - Closed-form Final Result:  $x = 0$
  - Iterative Final Result:  $x = 0.00001$  (close enough)

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

- Let's try to solve this using the closed-form here (Assume  $\hat{y} = mx$ ):

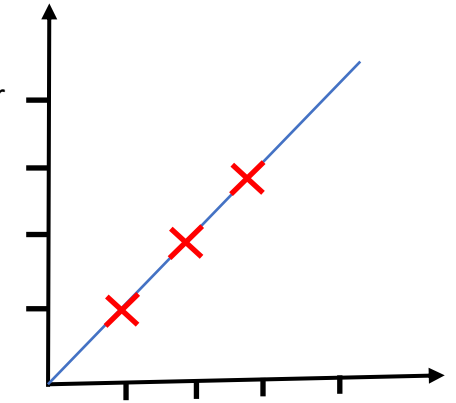
$$J(m) = \sum_{i=1}^3 (y_i - mx_i)^2$$

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

(Notice that  $y = x$  for our 3 points).

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

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



# 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.  $m, b$

# Gradient of the loss 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 loss function

- Simplifying the above expressions, we get:

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

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



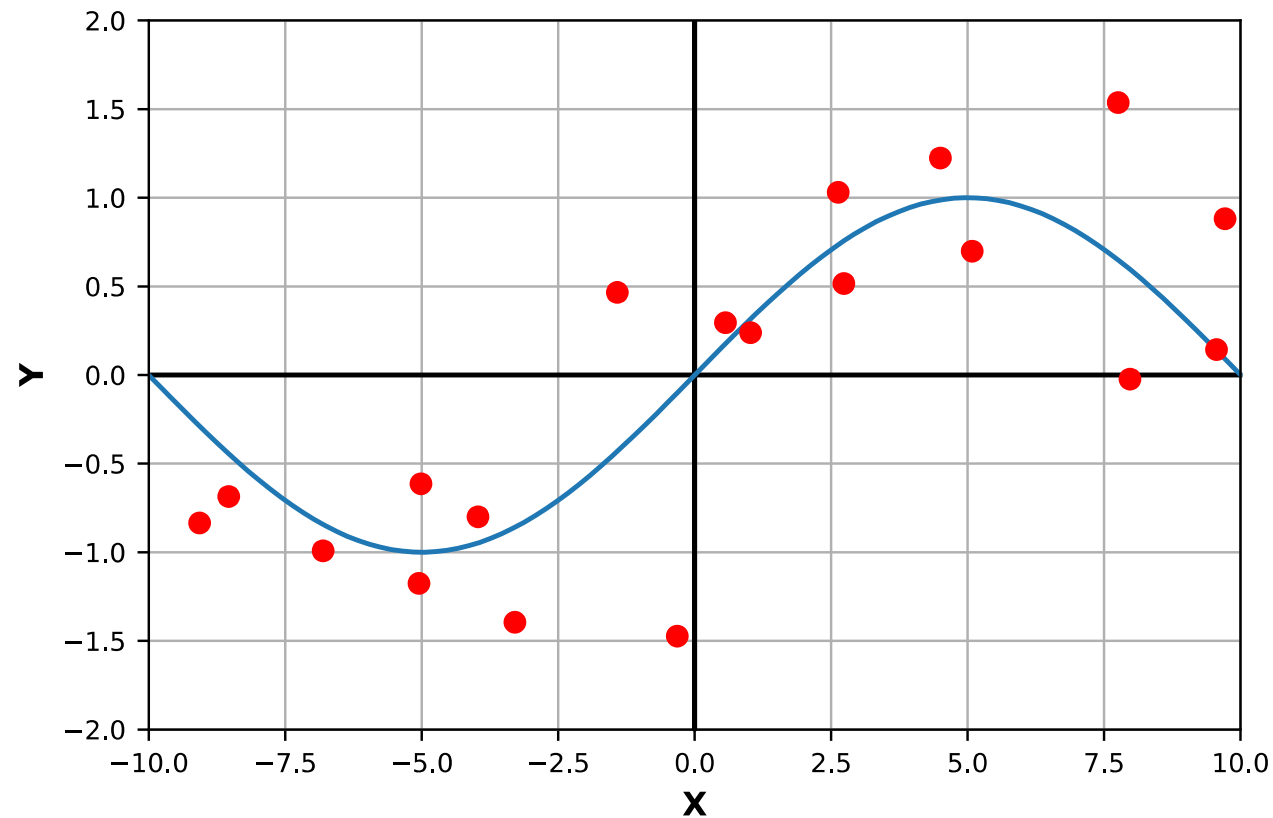
# Gradient of the loss function

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

$$\underbrace{\begin{bmatrix} \frac{\sum_i x_i^2}{N} & \frac{\sum_i x_i}{N} \\ \frac{\sum_i x_i}{N} & 1 \end{bmatrix}} \underbrace{\begin{bmatrix} m \\ b \end{bmatrix}} = \underbrace{\begin{bmatrix} \frac{\sum_i x_i y_i}{N} \\ \frac{\sum_i y_i}{N} \end{bmatrix}}$$

# Fitting Non-linear Data

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



# Transforming the Feature Space (Feature Engineering)



- We can transform features  $x_i$

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

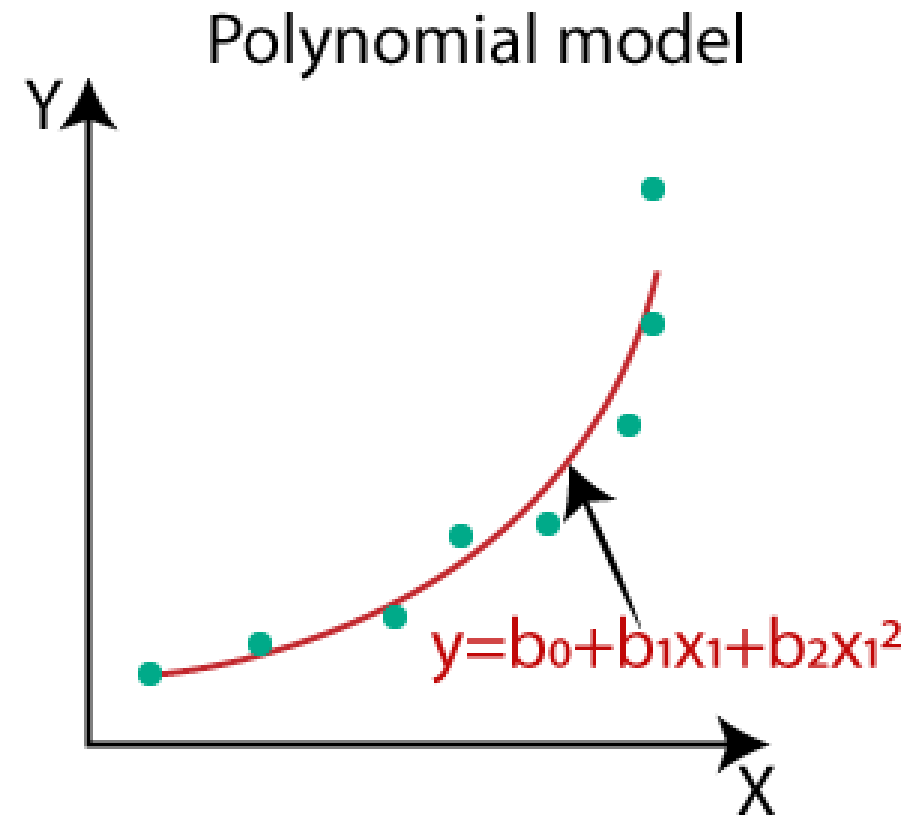
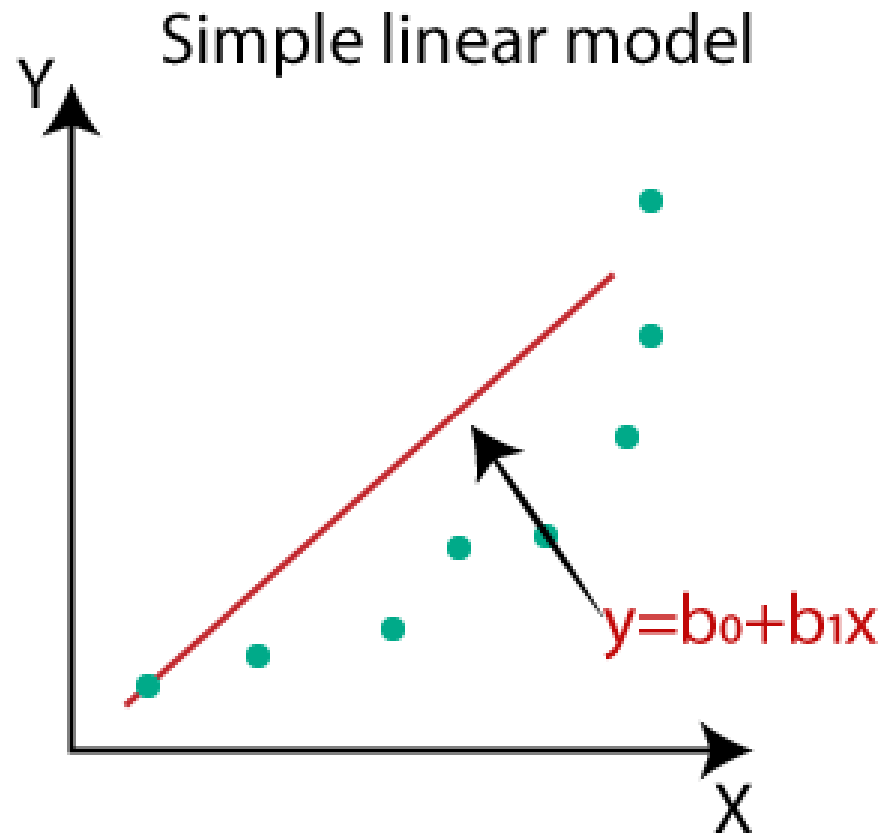
- We will apply some non-linear transformation  $\phi$ :

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

- For example, Polynomial transformation:

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

# Transforming the Feature Space (Feature Engineering)



# Transforming the Feature Space (Feature Engineering)



Example: assume you have:

$$x_i^1: Length$$

$$x_i^2: Width$$

You can add  $x_i^3: Area = x_i^1 * x_i^2$  to the dataset.

Other types:

- Cosine, splines, radial basis functions, etc.
- Encoding (Label encoding, One-hot,...)
- Domain-related features (e.g. financial measures)
- Time-related features (Day, month, year,...)
- Group-level features (e.g., average income per household, total sales per region, median age per team). Often called “Aggregation features”.

# Gradient of the loss function

- Let's get back to the gradients...

# Issues with the Approach

- Assume we have 100 variables instead of 2.
- Calculating gradients like this can quickly become tedious
- **Notice:** 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** (Writing equations as matrices).



# Vectorization

- To truly appreciate the power of vectorization. Let's make the problem a little more complex. The hypothesis function is now

$$\hat{y}_i = w_0 + w_1 x_i^1 + w_2 x_i^2 + \cdots + w_M x_i^M$$

- Where  $w_j$  ( $j = 0, 1, \dots, M$ ) are the unknown weights of the data, and  $x_i^j$  is the  $j$ th feature of the  $i$ th input.
- Next, we denote the discrepancy between  $y_i$  and  $\hat{y}_i$  as  $\epsilon_i$

$$y_i = \hat{y}_i + \epsilon_i$$

# Vectorization

- Now let's collect the above equation for all  $N$  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  $\hat{y}$ , 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:

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

# Vectorization

- 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:

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

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



# Cost function for the Vectorized form

- Notice that we are using the MSE cost function:

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

- 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:

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

# Optimization

- The optimization problem is now:

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

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

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

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



# Linear Least Squares

- We get:

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

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

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

**Closed-form solution for Linear Regression**

# Perspectives on ML Algorithms:

- We can look at ML algorithms from two perspectives:
  - **Loss Minimization** Problem (like what we did).
  - **Probability Maximization** Problem (using Maximum Likelihood Estimation).
- For linear regression, under particular assumptions, these two approaches yield equivalent solutions.

# 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

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

- Now assume that  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ , then  $y_i | x_i \sim \mathcal{N}(\mathbf{x}_i^T \boldsymbol{\theta}, \sigma^2)$
- 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

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

- Using the definition we get

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

# Maximum Likelihood Estimator

- Let's try to maximize:

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

- Note that

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

# Maximum Likelihood Estimator

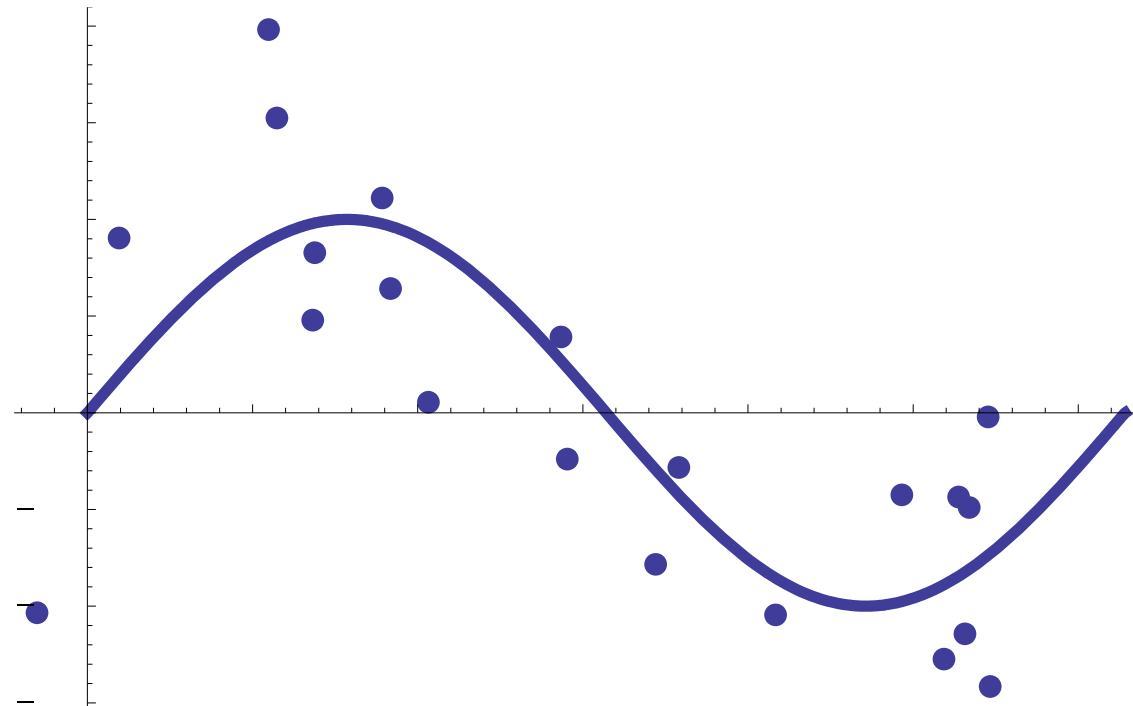
- Furthermore, since the right-hand side of the above equation is monotonic in  $\theta$  (theta), the arg max will not change if we take the logarithm of the expression.

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

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

# Bias and Variance

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

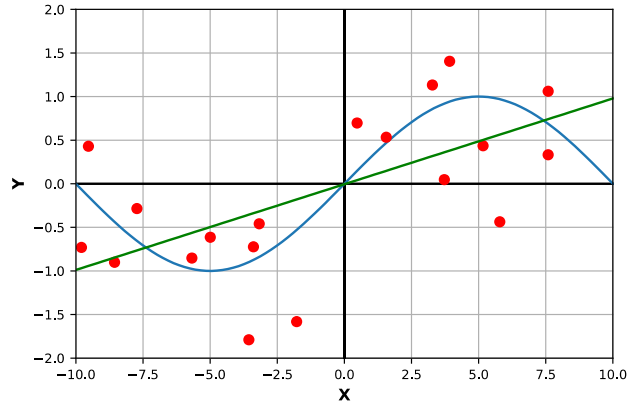


- Can we still use a linear model?

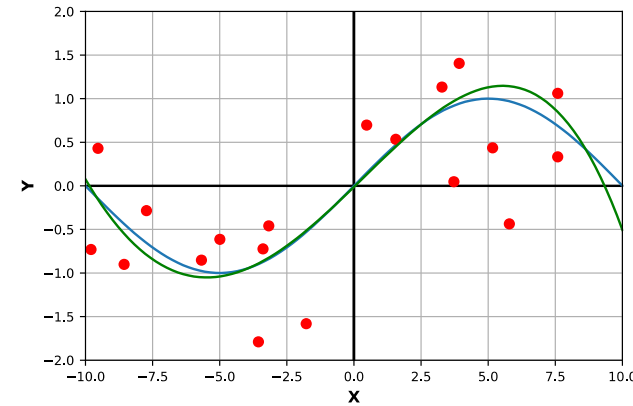


# What is Bias and Variance?

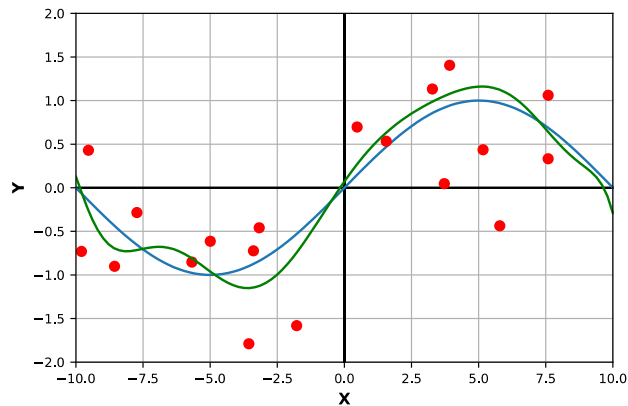
$$\{1, x\}$$



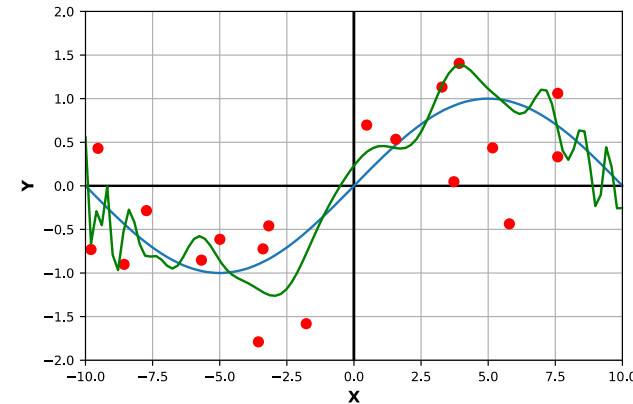
$$\{1, x, x^{[2]}, x^{[3]}, x^{[4]}\}$$



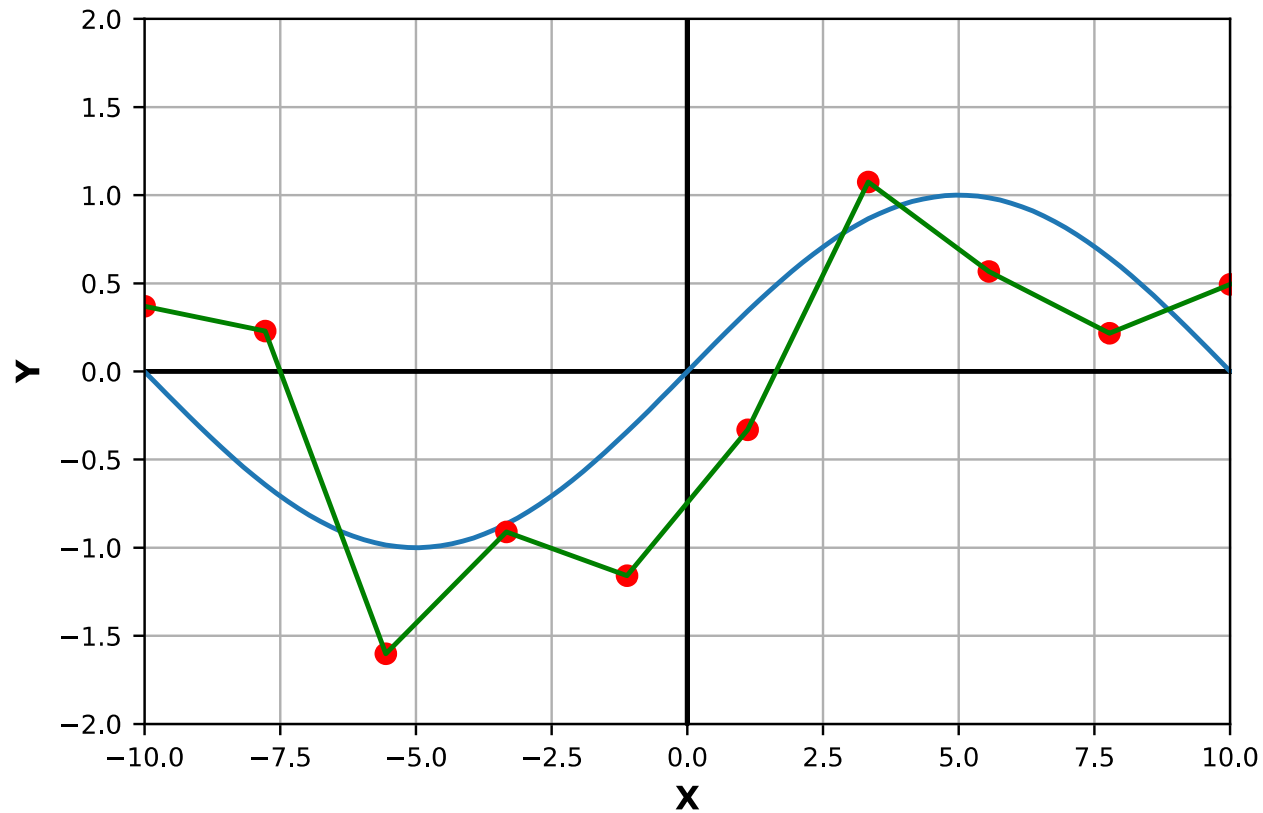
$$\{1, x, x^{[2]}, \dots, x^{[9]}, x^{[10]}\}$$



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



# Real Bad Overfit?



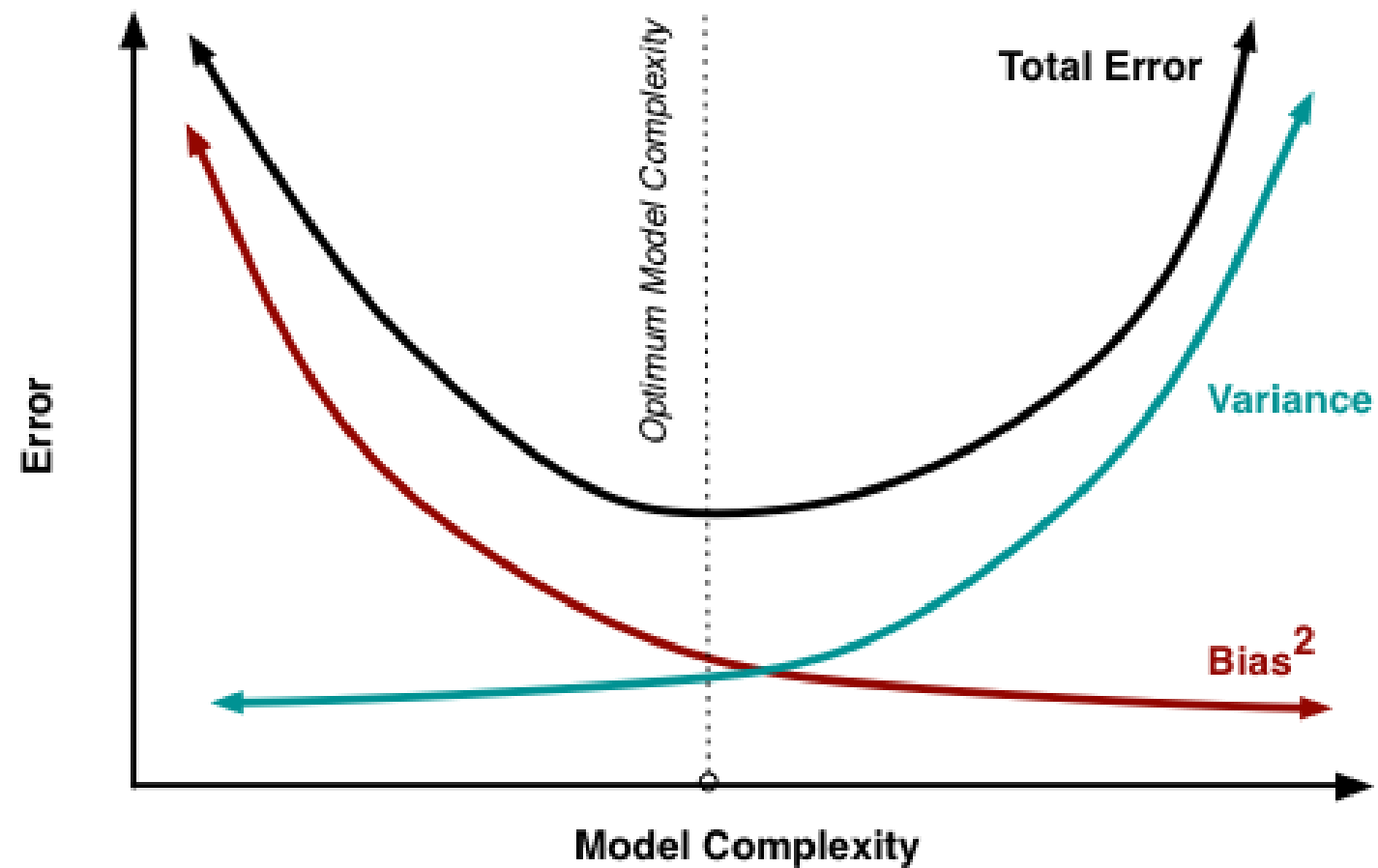
# Bias-Variance Tradeoff



- 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:

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

# Bias Variance Plot



# Evaluating models and Improving them



## R-squared ( $R^2$ or $r^2$ ) — “Coefficient of Determination”

- A **statistical metric** used to measure how well the independent variables explain the variability in the dependent variable.
- It provides a measure of the goodness-of-fit (**performance**) for a regression model:

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

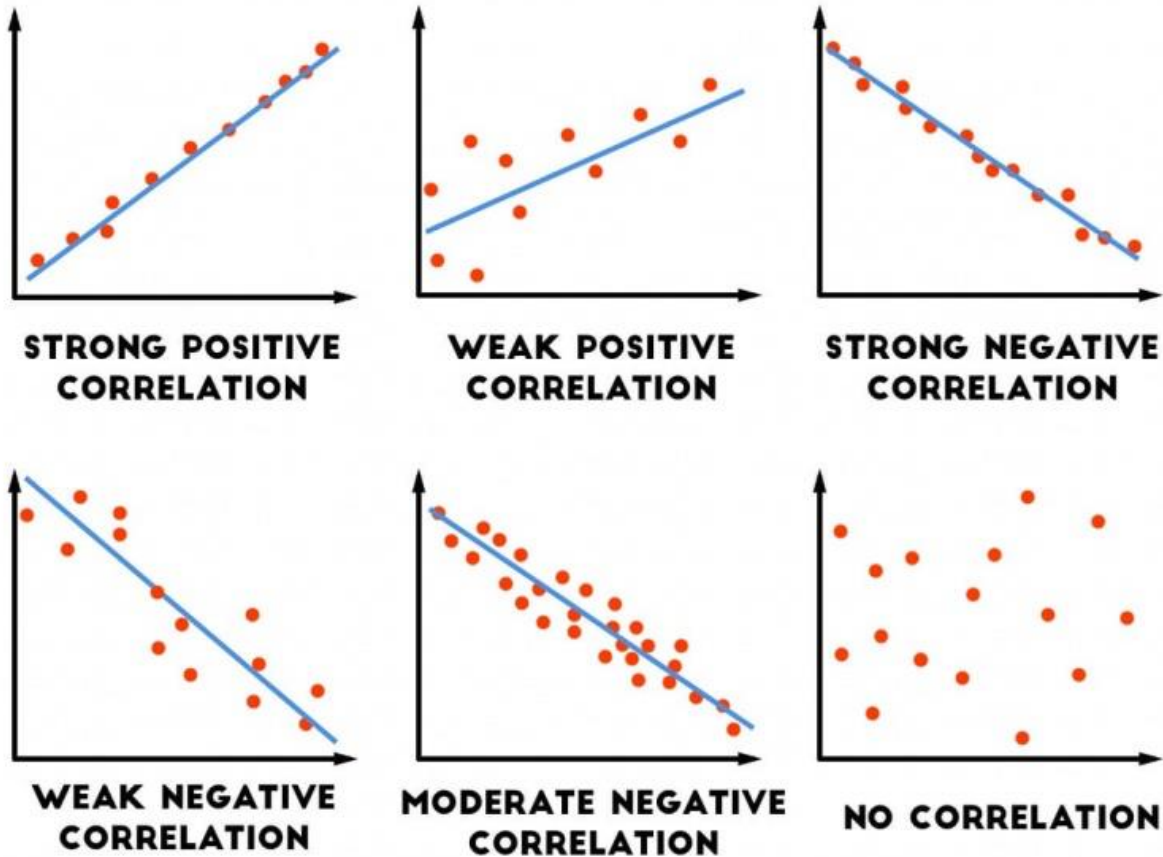
Where:

$SS_{\text{residual}} = \sum (y_i - \hat{y}_i)^2$ : The sum of squared residuals (difference between observed and predicted values).

$SS_{\text{total}} = \sum (y_i - \bar{y})^2$  : (Total Sum of Squares): The total variation in the data around the mean.

[Learn more!](#)

# Visualisation of $R^2$



## Interpretation:

- $R^2=1$ : The model perfectly explains the data.
- $R^2=0$ : The model explains none of the variability (as good as guessing the mean).
- $R^2<0$ : The model performs worse than a simple horizontal line at the mean of the target variable.

# Properties of $R^2$



## Range:

- $0 \leq R^2 \leq 1$ 
  - $R^2=1$ : Perfect model; predictions perfectly match the observations.
  - $R^2=0$ : Model does no better than the mean of the dependent variable.

## Interpretability:

- $R^2$  indicates the proportion of the variance in the dependent variable that is predictable from the independent variables.

## Limitations:

- $R^2$  increases with the addition of independent variables, even if they don't improve model prediction significantly.
- Does not account for overfitting or the complexity of the model.

## Negative Values:

- In rare cases,  $R^2$  can be negative when the model fits the data worse than a horizontal line representing the mean of the dependent variable.

## Adjusted $R^2$ :

- Adjusted  $R^2$  penalizes for the addition of non-significant predictors and is often a better metric for comparing models (here  $n$  is the number of data points and  $p$  is the number of predictors):

$$R_{\text{adj}}^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - p - 1}$$



## R-squared vs. Cost Function:



- **Cost Function:** A mathematical function used to optimize a model during training by minimizing the prediction error (e.g., Mean Squared Error or MSE).
  - Example: Gradient Descent minimizes MSE for linear regression.
- **R-squared:** A metric to evaluate how well the model fits the data **after training**. It is not used during model optimization.

## Limitations of $R^2$ :

1. **Doesn't Indicate Causation:** A high  $R^2$  doesn't mean the predictors cause the dependent variable.
2. **Sensitive to Overfitting:** A complex model might have a high  $R^2$  but poor generalization.
3. **Adjusted  $R^2$ :** For models with many predictors, use adjusted  $R^2$ , which accounts for the number of predictors to avoid overestimating performance.

# Data Split

- To ensure your model doesn't overfit to the training data, you should have another subset called **testing data**.
- You will evaluate your model against this subset, and based on its **metric score (e.g. accuracy)** you will decide if it's overfitting or not.
- But how should I split my data?

# Data Split

- **Hold-out set:**

- A portion of the dataset set aside and not used during training.
- E.g. 80% for training and 20% for testing.

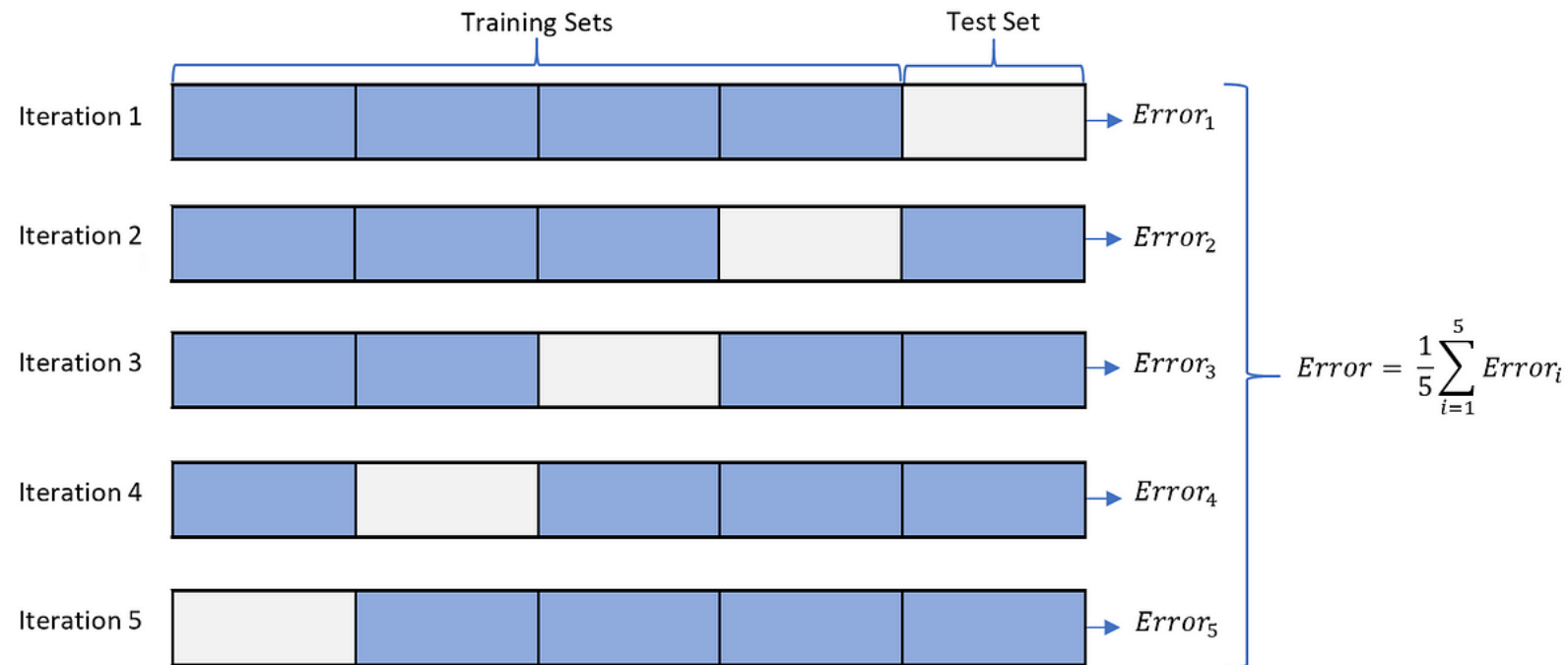
## Issues:

- Imagine you have these labels: [1, 1, 2, 2, 2, 3, 3, 3, 3, 3] and you took last 30% as test: [3,3,3]. You didn't include 1 and 2 in test!  
**Solution:** Always shuffle before split: [3, 2, 3, 1, 3, 2, 3, 1, 3, 2] → test: [1,3,2]
- My dataset is small. Taking 20% as test would not be representative!  
**Solution:** Use KFold.

# Data Split

- **K-Fold Cross Validation (CV):**

- Split data into  $k$  parts (folds), trains on  $k - 1$  folds, test on the remaining fold, and repeats  $k$  times then average the scores.



# Regularization

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  $l_1$  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  $l_2$  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.

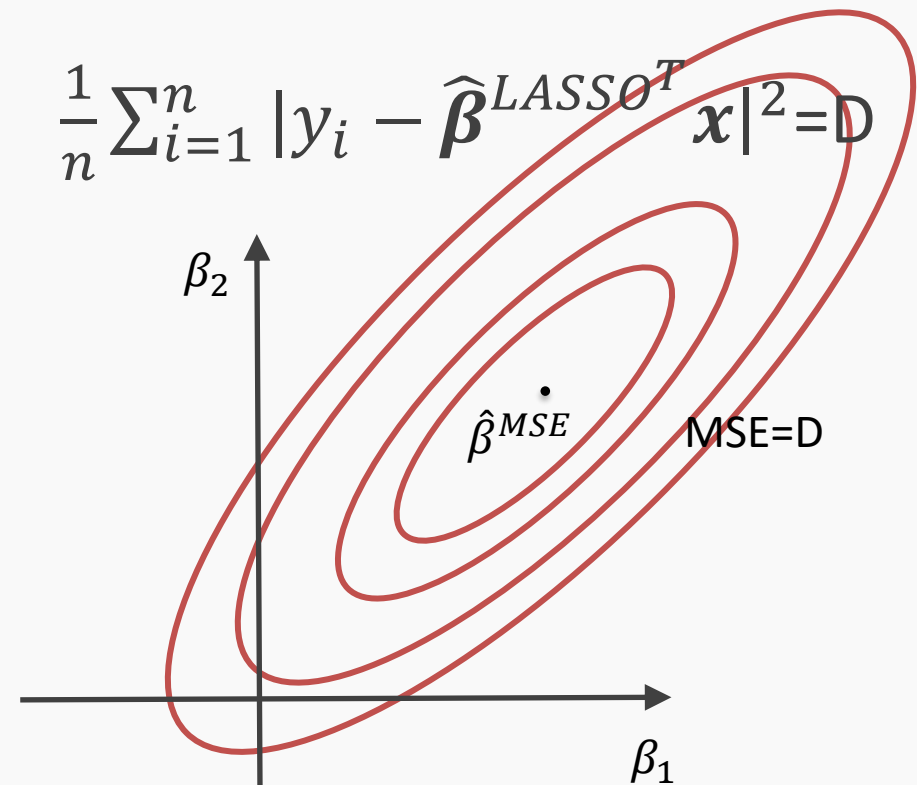
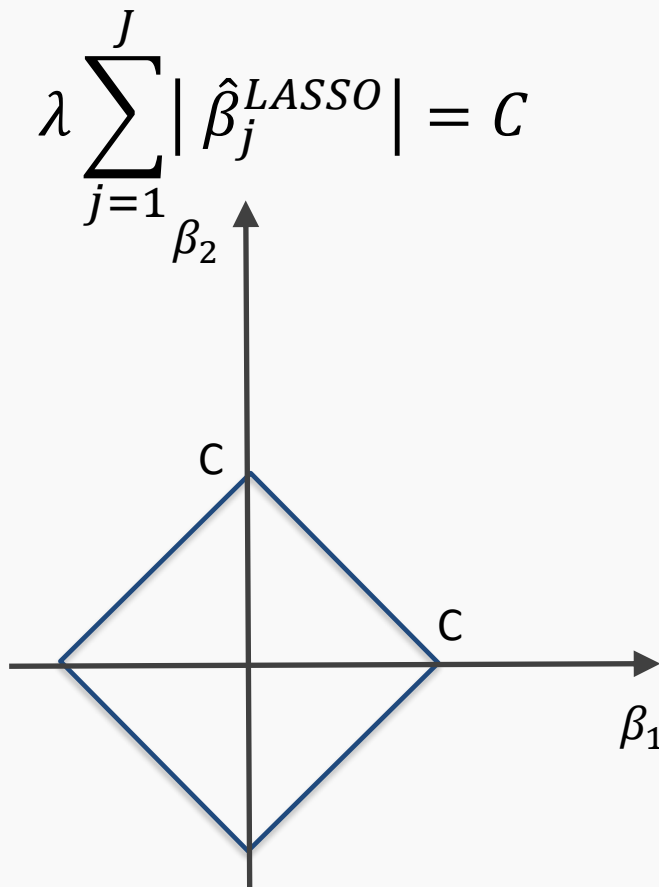
The solution of the Ridge/Lasso regression involves three steps:

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

# The Geometry of Regularization (LASSO)

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

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

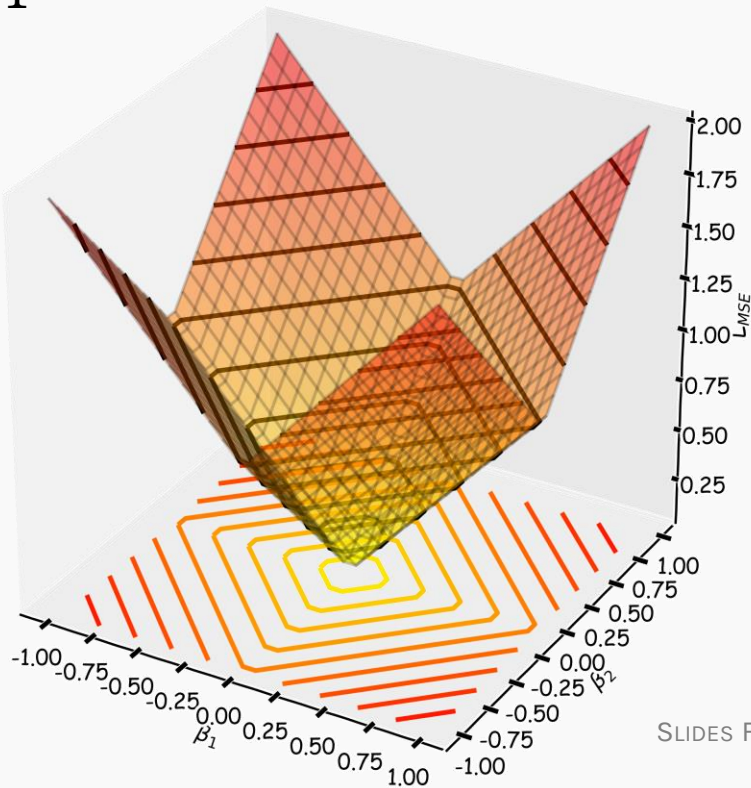


# The Geometry of Regularization (LASSO)

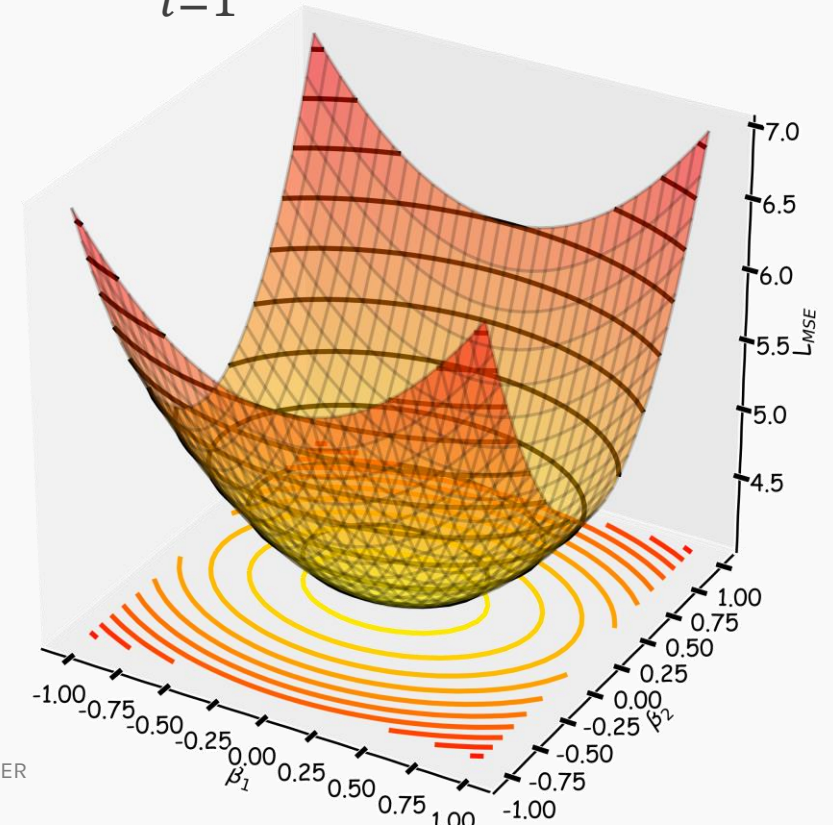
$$L_{LASSO}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n |y_i - \boldsymbol{\beta}^T \mathbf{x}|^2 + \lambda \sum_{j=1}^J |\beta_j|$$

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

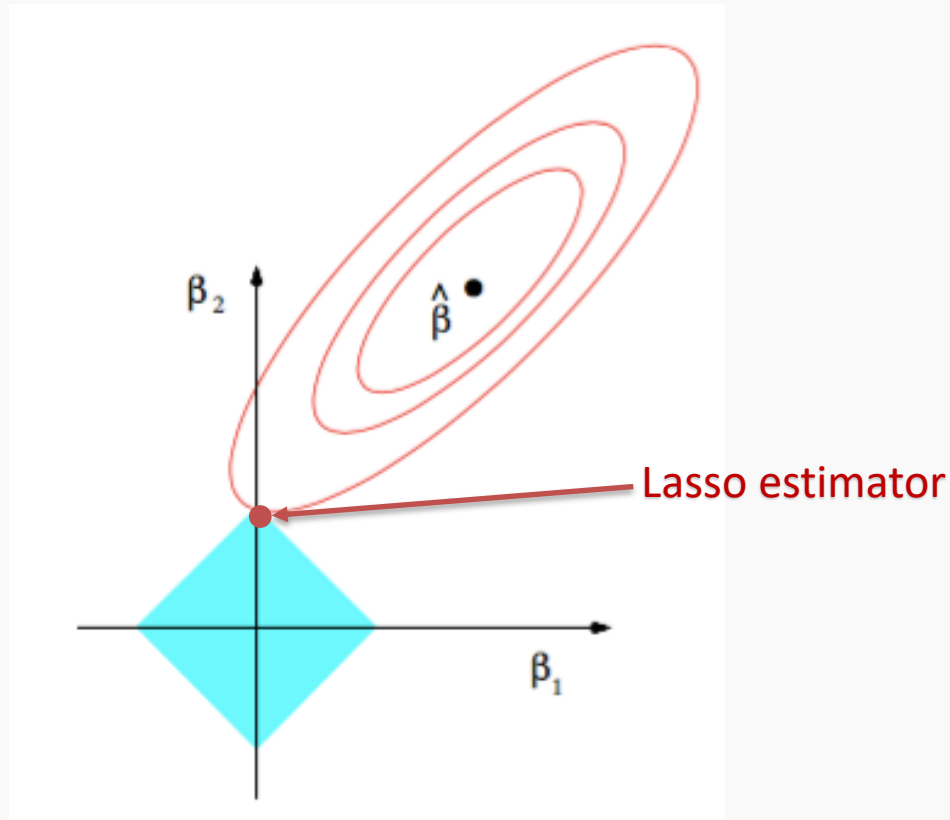
$$L_1 = \lambda \sum_{j=1}^J |\hat{\beta}_j^{LASSO}|$$



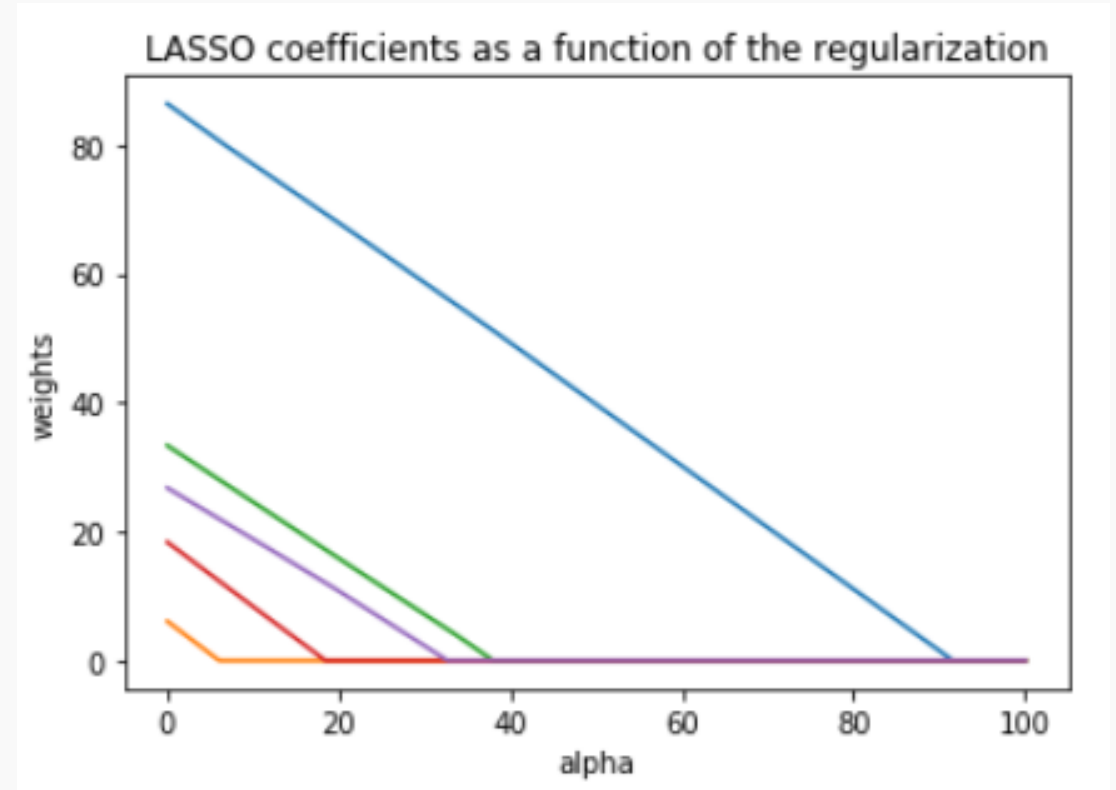
$$L_{MSE}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n |y_i - \boldsymbol{\beta}^T \mathbf{x}|^2$$



# 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.



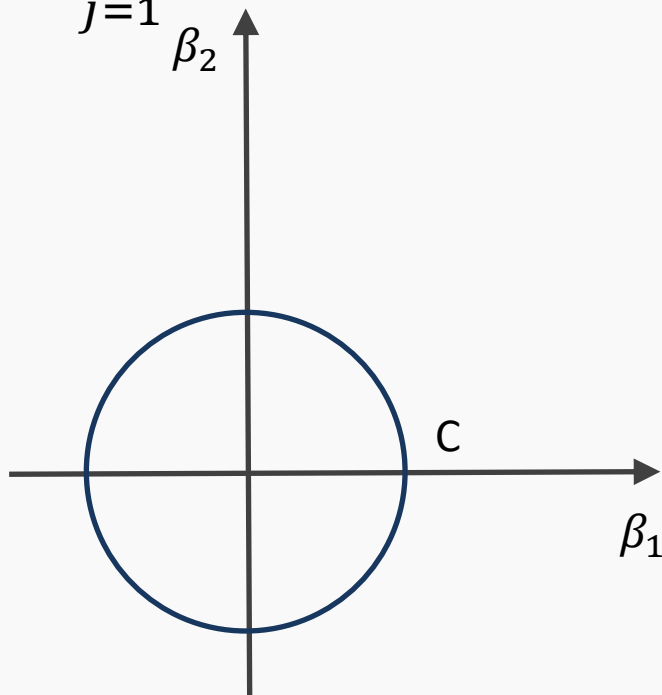
The values of the coefficients decrease as lambda increases, and are nullified fast.

# The Geometry of Regularization (Ridge)

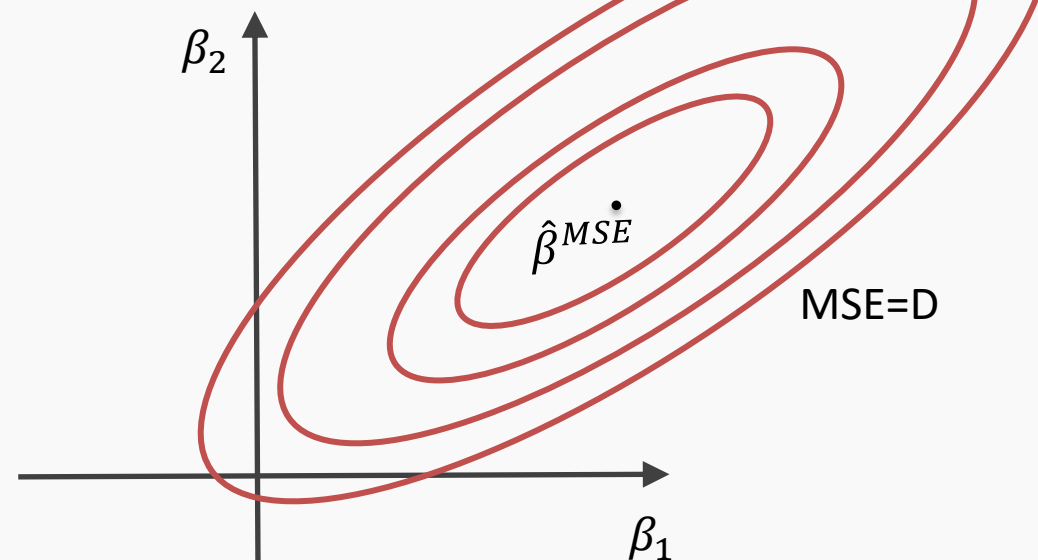
$$L_{\text{Ridge}}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n |y_i - \boldsymbol{\beta}^T \mathbf{x}|^2 + \lambda \sum_{j=1}^J (\beta_j)^2$$

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

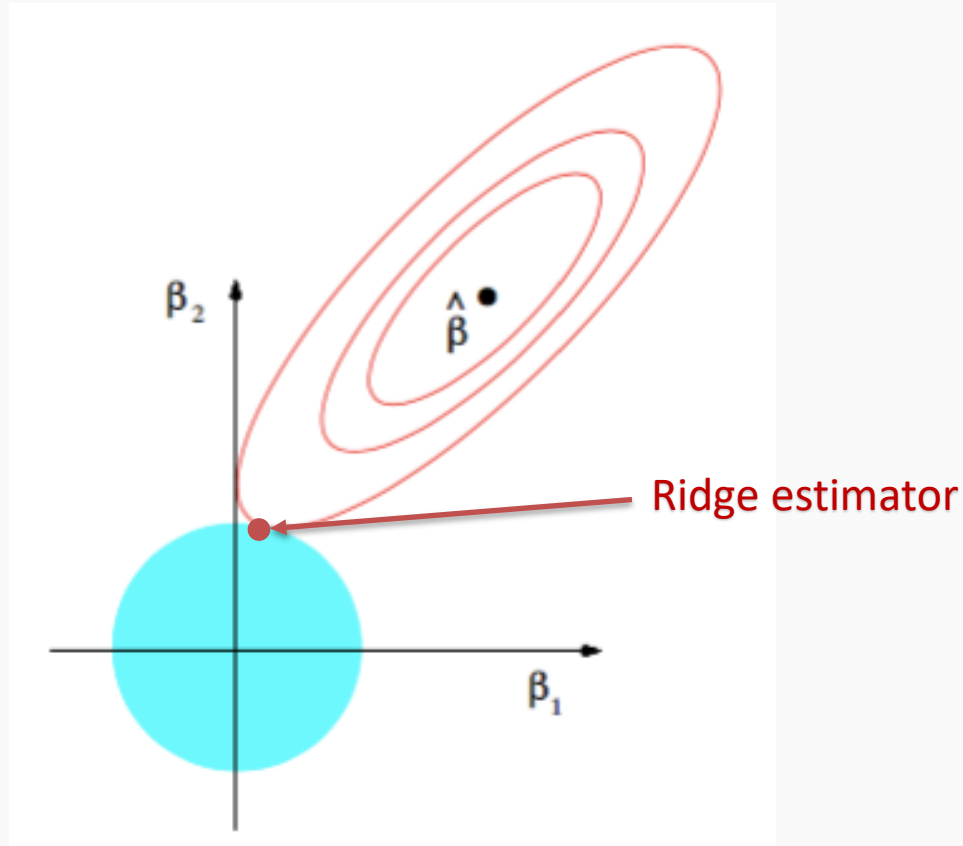
$$\lambda \sum_{j=1}^J |\hat{\beta}_j^{\text{Ridge}}|^2 = C$$



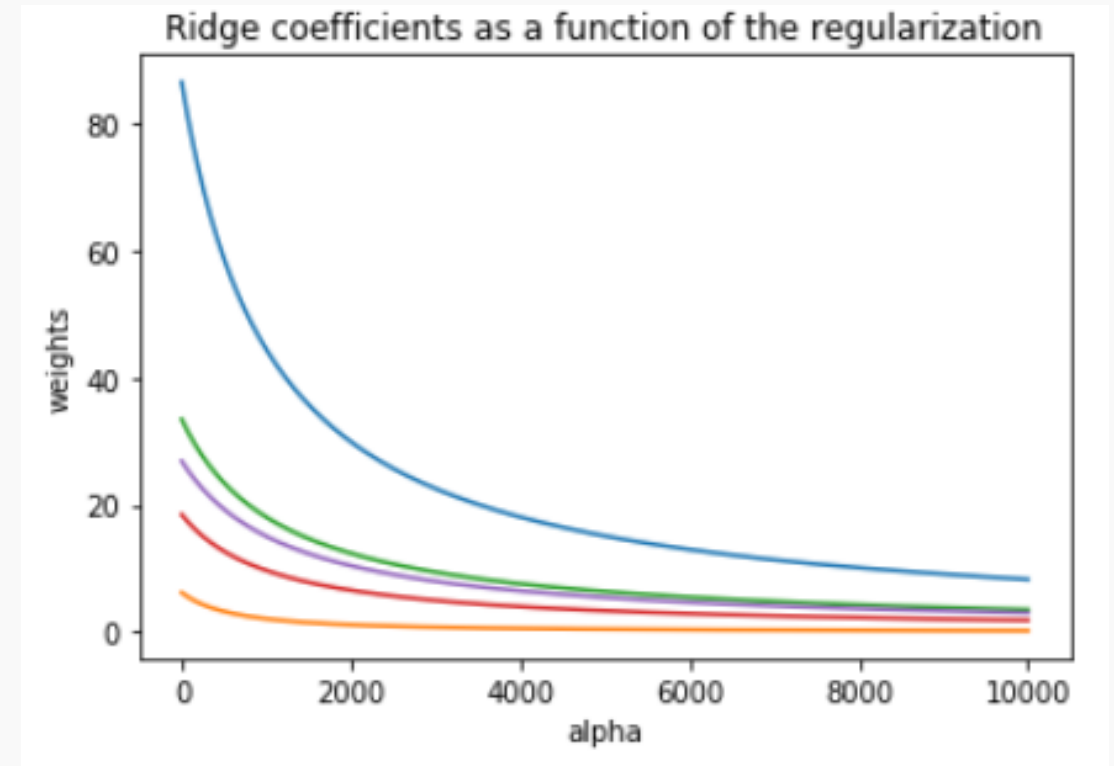
$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{\boldsymbol{\beta}}^{\text{Ridge}^T} \mathbf{x}|^2 = D$$



# Ridge visualized



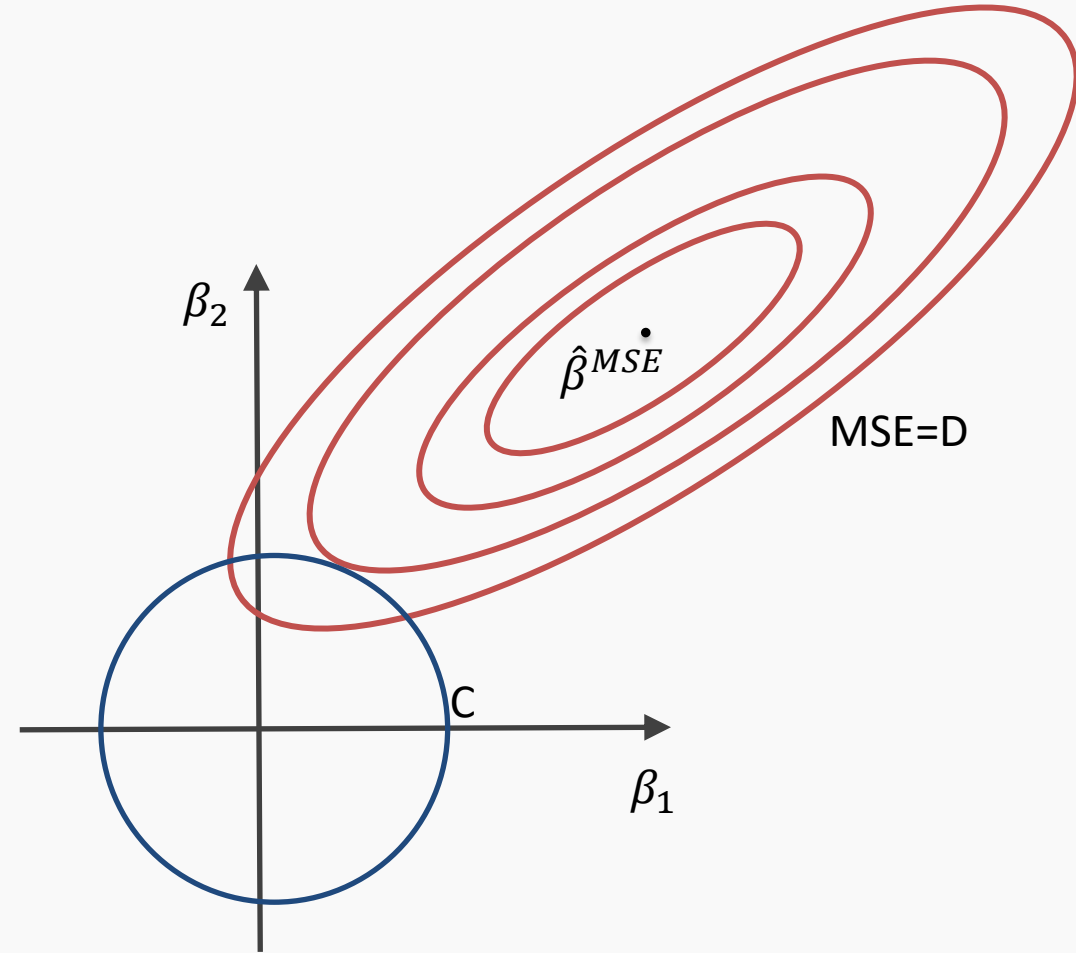
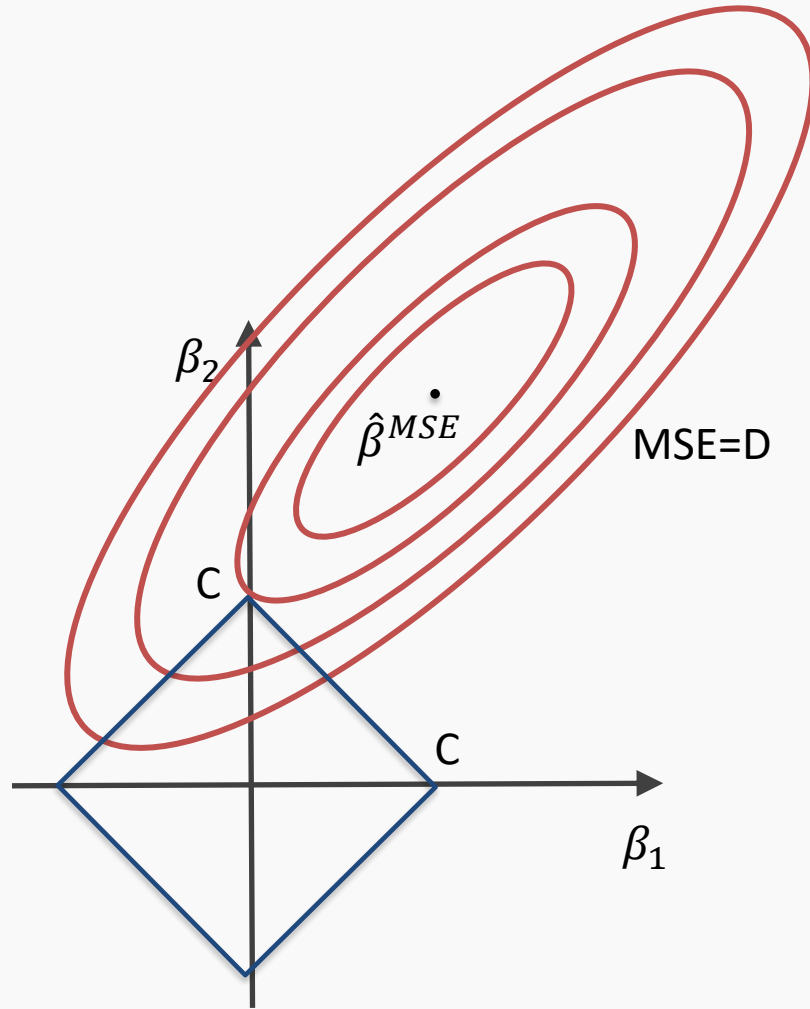
The ridge estimator is where the constraint and the loss intersect.



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



# The Geometry of Regularization



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

1. split data into  $\{\{X, Y\}_{train}, \{X, Y\}_{validation}, \{X, Y\}_{test}\}$
2. for  $\lambda$  in  $\{\lambda_{min}, \dots, \lambda_{max}\}$ :
  1. determine the  $\beta$  that minimizes the  $L_{ridge}$ ,  $\hat{\beta}_{Ridge}(\lambda) = (X^T X + \lambda I)^{-1} X^T Y$ , using the train data.
  2. record  $L_{MSE}(\lambda)$  using validation data.
3. select the  $\lambda$  that minimizes the loss on the validation data,  
 $\lambda_{ridge} = \operatorname{argmin}_{\lambda} L_{MSE}(\lambda)$
4. Refit the model using both train and validation data,  
 $\{\{X, Y\}_{train}, \{X, Y\}_{validation}\}$ , resulting to  $\hat{\beta}_{ridge}(\lambda_{ridge})$
5. report MSE or  $R^2$  on  $\{X, Y\}_{test}$  given the  $\hat{\beta}_{ridge}(\lambda_{ridge})$

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



1. split data into  $\{\{X, Y\}_{train}, \{X, Y\}_{validation}, \{X, Y\}_{test}\}$
2. for  $\lambda$  in  $\{\lambda_{min}, \dots, \lambda_{max}\}$ :
  - A. determine the  $\beta$  that minimizes the  $L_{lasso}$ ,  $\hat{\beta}_{lasso}(\lambda)$ , using the train data. **This is done using a solver.**
  - B. record  $L_{MSE}(\lambda)$  using validation data
3. select the  $\lambda$  that minimizes the loss on the validation data,  
 $\lambda_{lasso} = \operatorname{argmin}_{\lambda} L_{MSE}(\lambda)$
4. Refit the model using both train and validation data,  
 $\{\{X, Y\}_{train}, \{X, Y\}_{validation}\}$ , resulting to  $\hat{\beta}_{lasso}(\lambda_{lasso})$
5. report MSE or  $R^2$  on  $\{X, Y\}_{test}$  given the  $\hat{\beta}_{lasso}(\lambda_{lasso})$

# 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
```

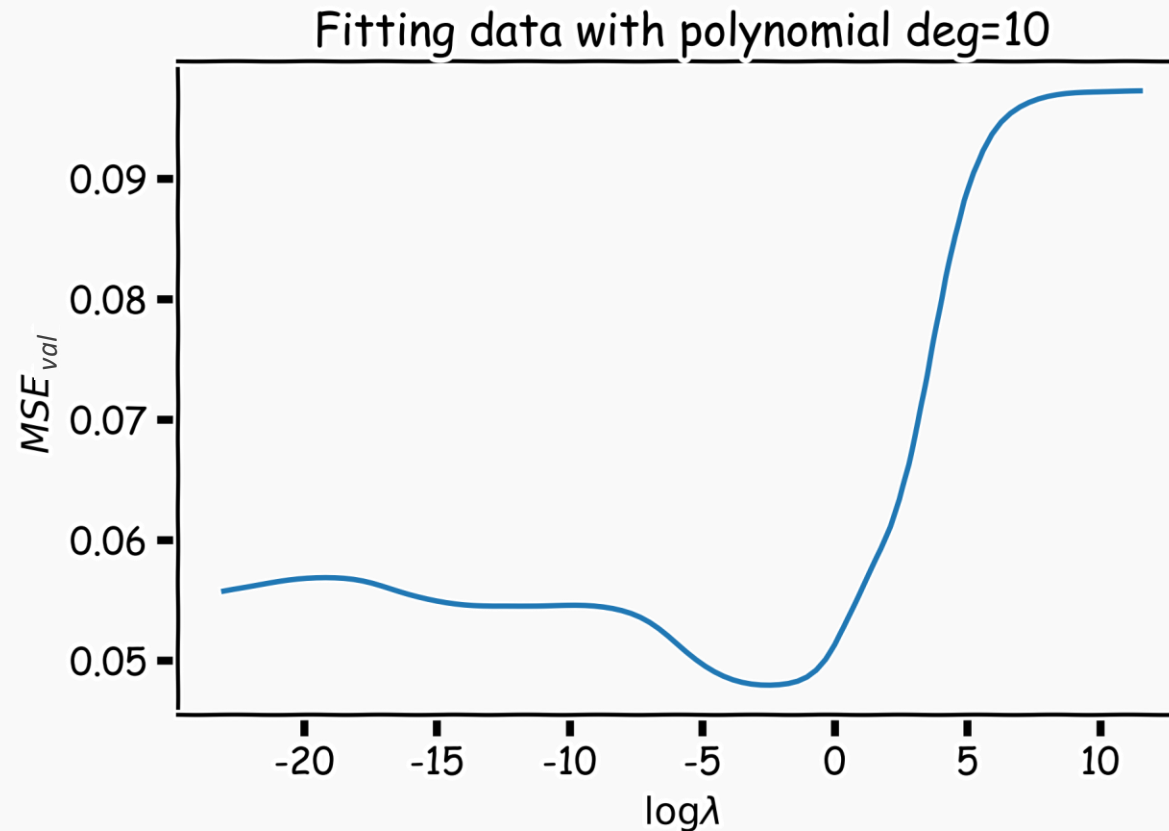
```
In [20]: X_train = train[all_predictors].values
X_val = validation[all_predictors].values
X_test = test[all_predictors].values

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



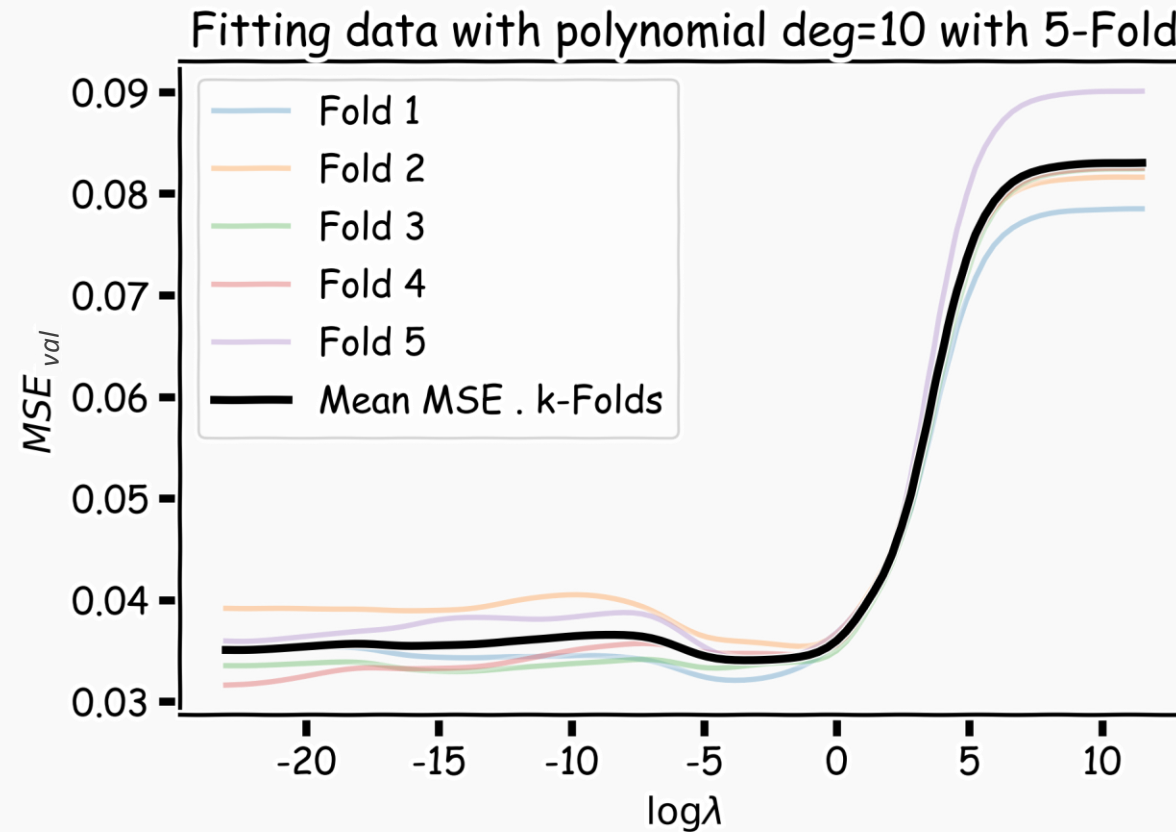
# Ridge regularization with CV: step by step

	$\lambda_1$	$\lambda_2$	...	$\lambda_n$
$k_1$	$L_{11}$	$L_{12}$	..	...
$k_2$	$L_{21}$	...	..	...
...	..	...	..	...
$k_n$	...	...	...	...
$E[]$	$\bar{L}_1$	$\bar{L}_2$	...	$\bar{L}_n$

1. remove  $\{X, Y\}_{test}$  from data
2. split the rest of data into K folds,  $\{\{X, Y\}_{train}^{-k}, \{X, Y\}_{val}^k\}$
3. for  $k$  in  $\{1, \dots, K\}$ 
  1. for  $\lambda$  in  $\{\lambda_0, \dots, \lambda_n\}$ :
    - A. determine the  $\beta$  that minimizes the  $L_{ridge}$ ,  $\hat{\beta}_{ridge}(\lambda, k) = (X^T X + \lambda I)^{-1} X^T Y$ , using the train data of the fold,  $\{X, Y\}_{train}^{-k}$ .
    - B. record  $L_{MSE}(\lambda, k)$  using the validation data of the fold  $\{X, Y\}_{val}^k$

At this point we have a 2-D matrix, rows are for different k, and columns are for different  $\lambda$  values.
4. Average the  $L_{MSE}(\lambda, k)$  for each  $\lambda$ ,  $\bar{L}_{MSE}(\lambda)$  .
5. Find the  $\lambda$  that minimizes the  $\bar{L}_{MSE}(\lambda)$  , resulting to  $\lambda_{ridge}$ .
6. Refit the model using the full training data,  $\{\{X, Y\}_{train}, \{X, Y\}_{val}\}$ , resulting to  $\hat{\beta}_{ridge}(\lambda_{ridge})$
7. report MSE or  $R^2$  on  $\{X, Y\}_{test}$  given the  $\hat{\beta}_{ridge}(\lambda_{ridge})$

# 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?