# Ridge and Lasso - Hyperparameters
## CS109A Introduction to Data Science
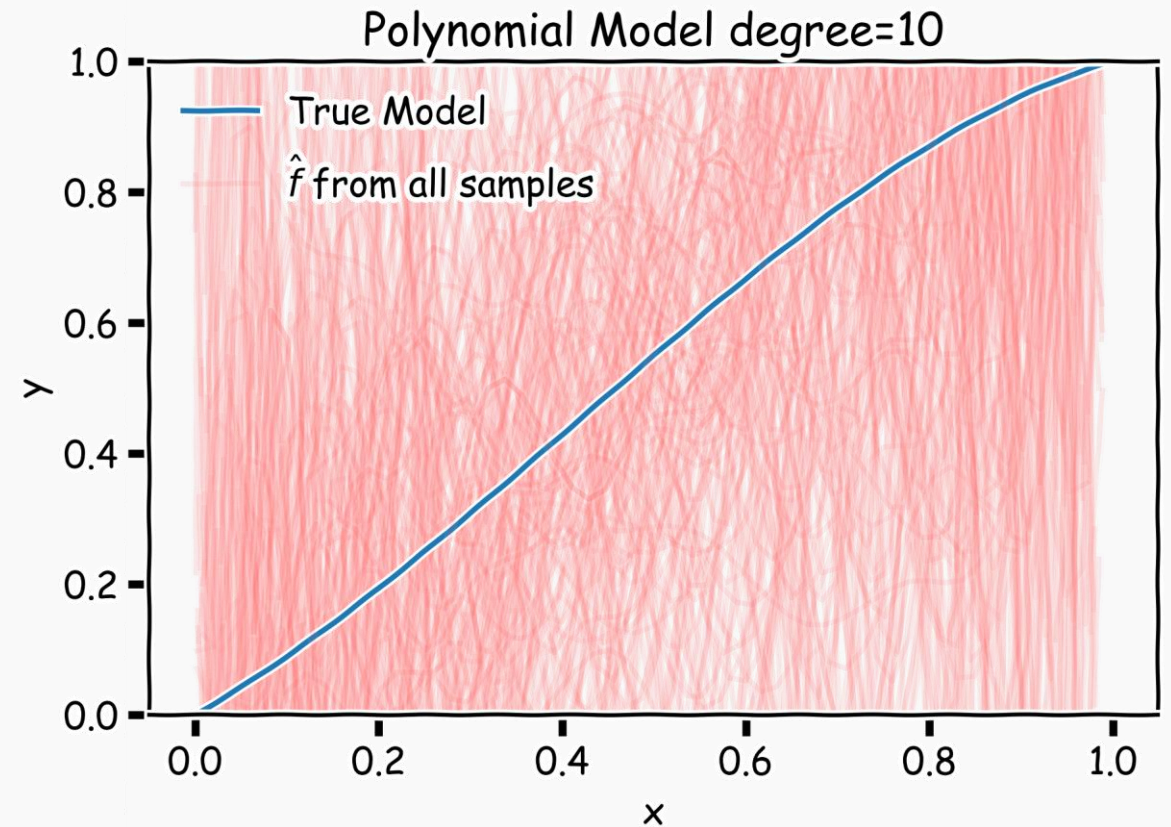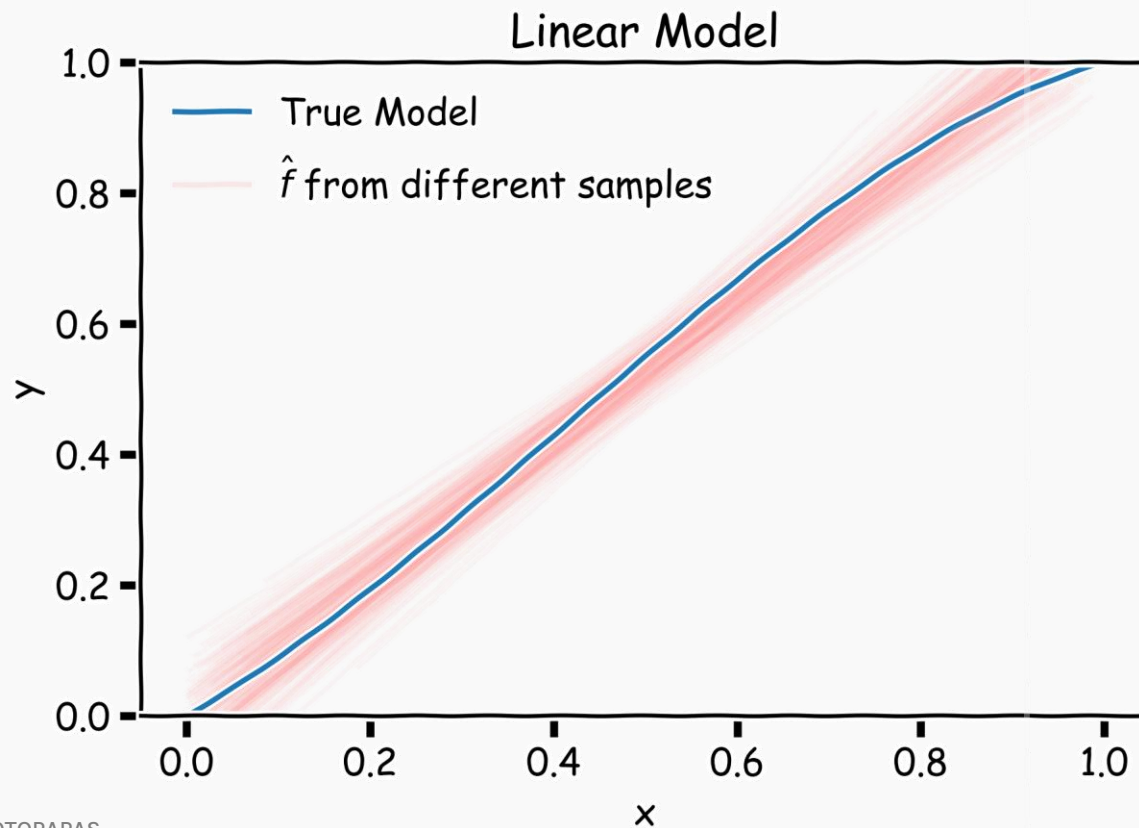Pavlos Protopapas, Kevin Rader, and Chris Gumb

# Outline

- Recap – Model Selection

- Generalization Error, Bias Variance Tradeoff

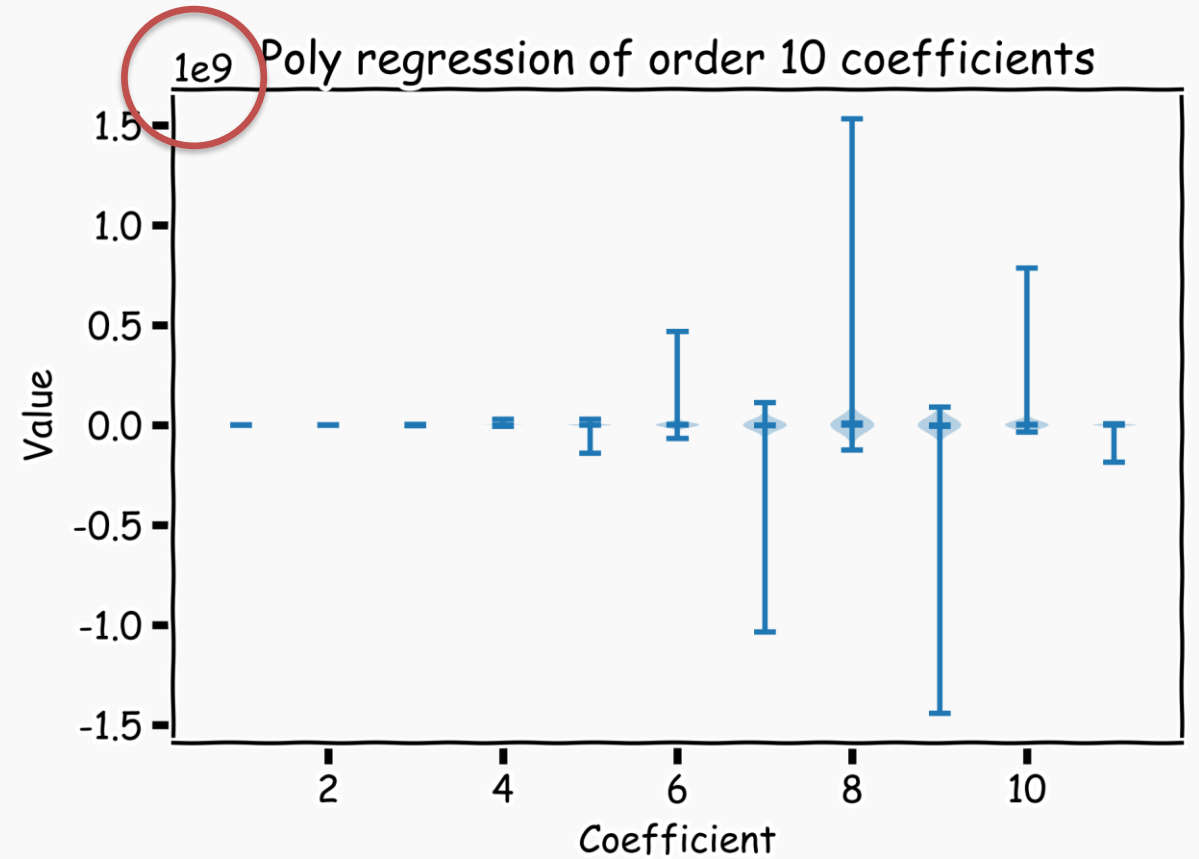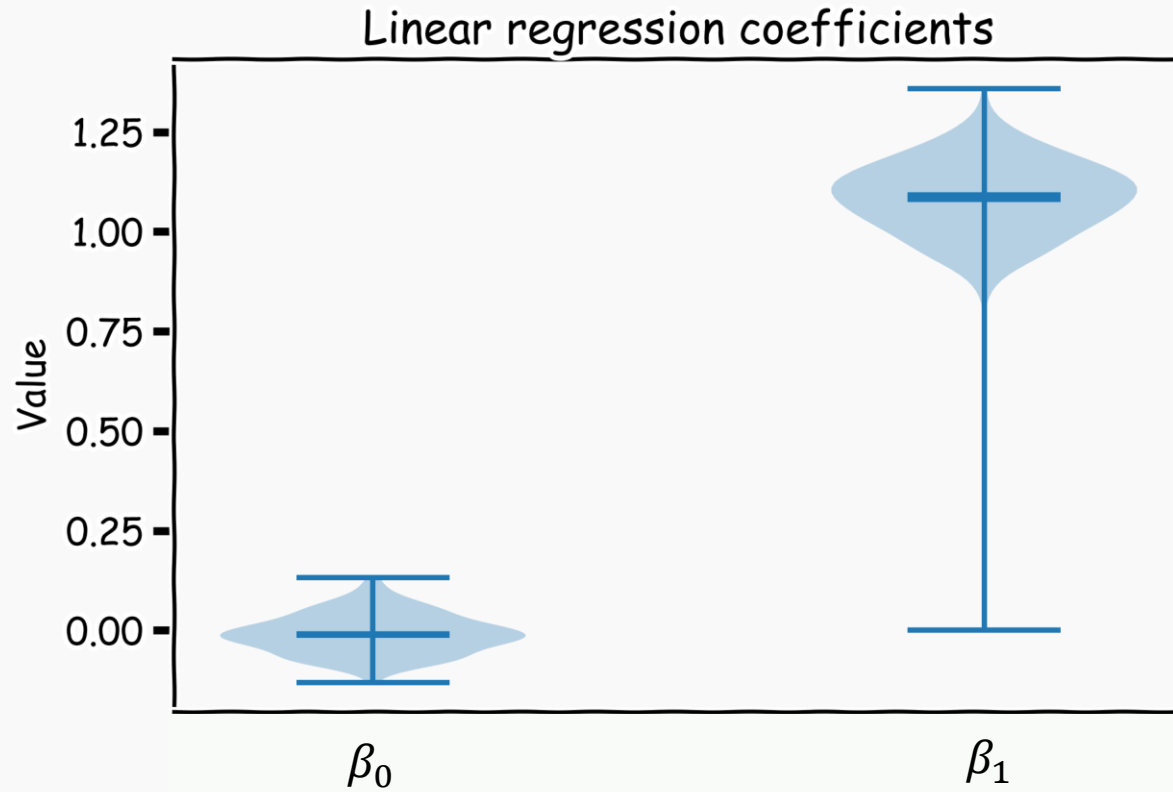- **Regularization Techniques: Lasso, Ridge**

# Bias vs Variance

**Left**: 2000, best fit straight lines, each fitted on a different 20-point training set.

**Right**: Best-fit models using degree-10 polynomial

# Bias vs Variance

# Model Selection

**Model selection** is the application of a principled method to determine the complexity of the model, e.g., choosing a subset of predictors, choosing the degree of the polynomial model etc.

A strong m‌‌‌‌‌‌‌‌oid **overfitting,**

**How do we discourage extreme values in the model parameters?**

- there ar‌‌
  - the f‌‌‌‌‌‌‌‌h dimensionality
  - the polynomial degree is too high
  - too many cross terms are considered

- the coefficients values are too **extreme**

# Game Time

How would you discourage extreme values in the model parameters

**Options:**

    A.  Divide all model parameters by a large number

    B.  Make sure the causal relationship between predictors and response variable is true

    C.  Discard any model with model parameter value larger than 1

    D.  Penalize the model with a penalty that is proportional to the value its parameters

# Regularization

## What we want

### Low model error

Minimize:

$$\frac{1}{n} \sum_{i=1}^{n} \left| y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i \right|^2$$

### Discourage extreme values in model parameters

Minimize:

# Regularization

## What we want

Low model error

Minimize:

$$\frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i\right|^2$$

Discourage extreme values in model parameters

Minimize:

$$L_{reg} = \begin{cases} \displaystyle\sum_{j=1}^{J}\beta_j^2 \\ \displaystyle\sum_{j=1}^{J}|\beta_j| \end{cases}$$

## What we want

Low model error

Discourage extreme values in model parameters

Minimize:

Minimize:

$$\frac{1}{n} \sum_{i=1}^{n} \left| y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i \right|^2$$

$$L_{reg} = \begin{cases} \displaystyle\sum_{j=1}^{J} \beta_j^2 \\ \\ \displaystyle\sum_{j=1}^{J} |\beta_j| \end{cases}$$

How do we combine these two objectives?

## What we want

Low mod... ...e values in

Minimiz...

$$\frac{1}{n} \sum_{i=1}^{n} |y_i - \quad \quad \beta_j^2$$

$$|\beta_j|$$



MSE

Regularization Terms

MAKE GIFS AT GIFSOUP.COM

# Regularization

## What we want

Low model error

<span style="color:#4a90d2">Minimize</span>:

Discourage extreme values in model parameters

<span style="color:#4a90d2">Minimize</span>:

$$\mathcal{L}_{REG} \;=\; \frac{1}{n}\sum_{i=1}^{n}\left| y_i - \boldsymbol{\beta}^{\top} \boldsymbol{x}_i \right|^2 \;+\; L_{reg}$$

# Regularization

## What we want

Low model error

Discourage extreme values in model parameters

Minimize ~~error~~ ... ~~imize~~:

$\lambda$ is the **regularization parameter**. It controls the relative importance between model error and the regularization term

$$\mathcal{L}_{REG} = \frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i\right|^2 + \lambda \; L_{reg}$$

# Regularization

What we want

Low model error

Discourage extreme values in model parameters

mize:

$\lambda = 0$: equivalent to simple linear regression

$\lambda = \infty$: yields a model with $\beta's$ =0

$$\mathcal{L}_{REG} = \frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i\right|^2 + \lambda\, L_{reg}$$

**?**

What we want

Low model error

Discourage extreme values in model parameters

Minimize:                         Minimize:

How do we determine $\lambda$?

$$\mathcal{L}_{REG} = \frac{1}{n} \sum_{i=1}^{n} \left| y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i \right|^2 + \lambda\, L_{reg}$$
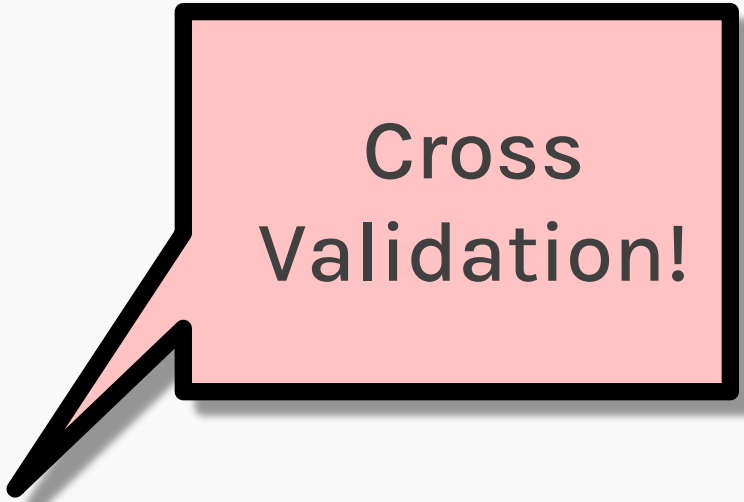
# Regularization

## What we want

**Low model error**

Minimize:

**Discourage extreme values in model parameters**

Minimize:

Cross Validation!

$$\mathcal{L}_{REG} = \frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i\right|^2 + \lambda\, L_{reg}$$

# Regularization: **LASSO** Regression

## What we want

Low model error

Minimize:

Discourage extreme values in model parameters

Minimize:

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

$$\mathcal{L}_{LASSO} = \frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^{\top}\boldsymbol{x}_i\right|^2 + \lambda\sum_{j=1}^{J}|\beta_j|$$

**?**

## What we want

Low model error

Discourage extreme values in
~~del~~ parameters

~~nize~~:

No need to regularize the bias, $\beta_0$
Why?

$$\mathcal{L}_{LASSO} = \frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i\right|^2 + \lambda \sum_{j=1}^{J}|\beta_j|$$

# Regularization: **LASSO** Regression

**Lasso** regression: minimize $\mathcal{L}_{LASSO}$ with respect to $\beta's$

$$\mathcal{L}_{LASSO} = \frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i\right|^2 + \lambda \sum_{j=1}^{J}|\beta_j|$$

**Ridge** regression: minimize $\mathcal{L}_{RIDGE}$ with re...

Note that $\sum_{j=1}^{J} \beta_j^2$ is the $L_2$ norm square of the vector $\boldsymbol{\beta}$

$$\mathcal{L}_{RIDGE} = \frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i\right|^2 + \lambda \sum_{j=1}^{J} \beta_j^2$$

**Ridge** regression: minimize $\mathcal{L}_{RIDGE}$ with respect to $\beta's$

$$\mathcal{L}_{RIDGE} = \frac{1}{n}\sum_{i=1}^{n}\left|y_i - \boldsymbol{\beta}^\top \boldsymbol{x}_i\right|^2 + \lambda \sum_{j=1}^{J}\beta_j^2$$

No need to regularize the bias, $\beta_0$, since it is not connected to the predictors.

For ridge regression there exist an analytical solution for the coefficients:

$$\hat{\beta}_{Ridge}(\lambda) = \left(X^TX + \lambda I\right)^{-1} X^T Y$$



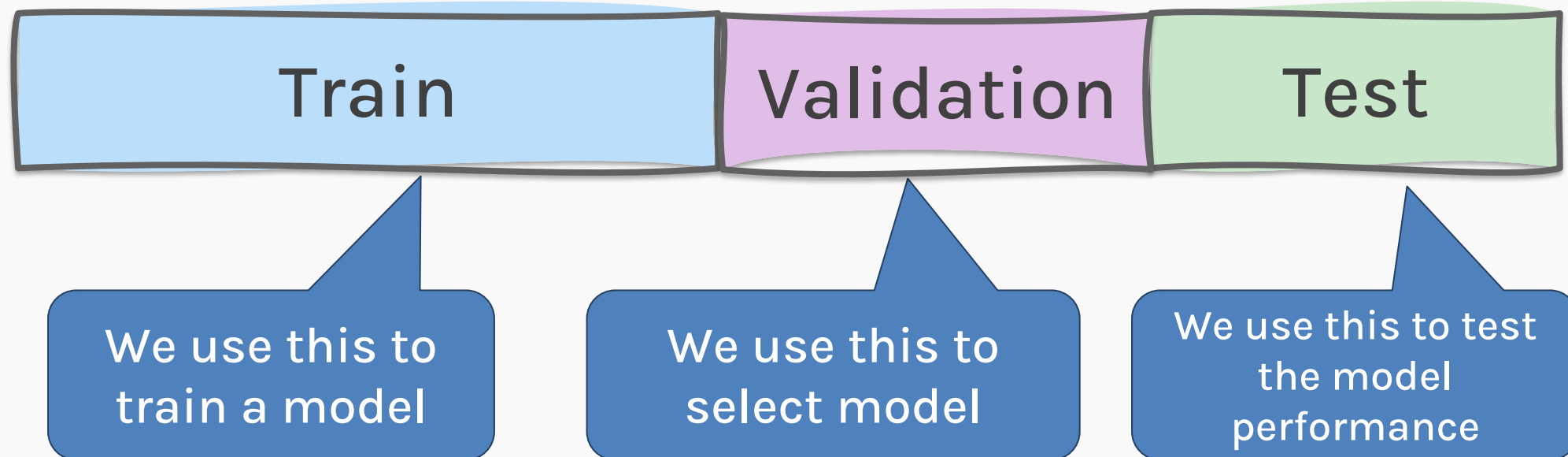So, what are the steps to determine $\lambda$ using validation?

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

For ridge regression there exist an analytical solution for the coefficients:
$$\hat{\beta}_{Ridge}(\lambda) = \left(X^TX + \lambda I\right)^{-1} X^T Y$$

## Step 1 : Split Data

Split data into train, validation and test data.

| Train | Validation | Test |
|---|---|---|

We use this to train a model

We use this to select model

We use this to test the model performance

## Step 2 : Select a range of possible $\lambda$ values

| $\lambda_{min}$ | $\lambda_{min-1}$ | .... | $\lambda_{max-1}$ | $\lambda_{max}$ |
|---|---|---|---|---|

⚠️ The values of $\lambda$ are not fixed—you can choose the range yourself based on the problem.

## Step 2 : Select a range of possible $\lambda$ values

As an example, we will take a range of values from 0.00001 to 10

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

## Step 3 : Determine $\beta$

Train

For each $\lambda$ value, we determine $\beta_{Ridge}(\lambda)$ using train data

$\lambda\ values\ =$

| 0.00001 | 0.0001 | 0.001 | 0.01 | 0.1 |
|---|---|---|---|---|

••••

| 1 | 2 | 3 | 4 | 10 |
|---|---|---|---|---|

$$\beta_{Ridge}(\lambda) = \left(\mathrm{X^T X} + \lambda I\right)^{-1} X^T Y$$

$\beta\ values\ =$



Keep in mind that each $\lambda$ value gives us a different vector of $\beta$ coefficients.
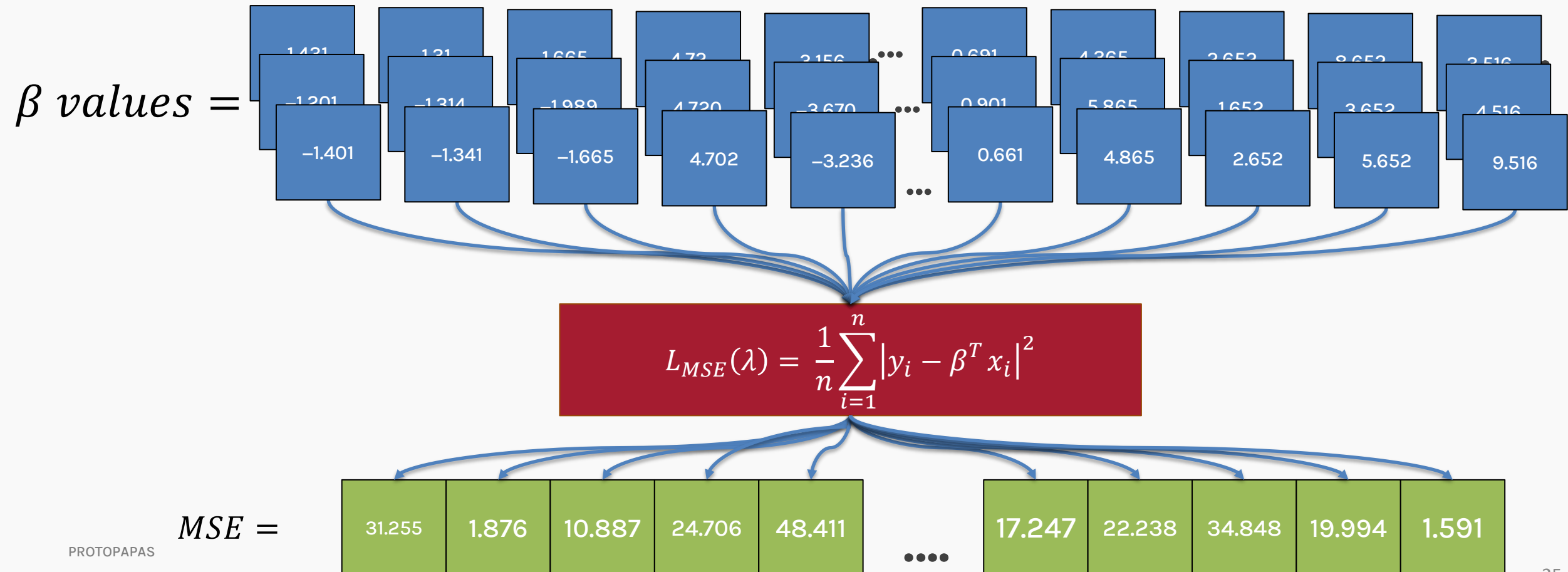
# Ridge regularization with only **validation** : step by step

## Step 4 : Record the $L_{MSE}(\lambda)$

For each vector of $\beta$ coefficients, we record $L_{MSE}(\lambda)$ using the validation data

$$\beta \ values =$$

| 1.431 | 1.31 | 1.665 | 4.72 | 3.156 | ••• | 0.691 | 4.365 | 2.652 | 8.652 | 3.516 |
| -1.201 | -1.314 | -1.989 | 4.720 | -3.670 | ••• | 0.901 | 5.865 | 1.652 | 3.652 | 4.516 |
| -1.401 | -1.341 | -1.665 | 4.702 | -3.236 | | 0.661 | 4.865 | 2.652 | 5.652 | 9.516 |

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

$MSE =$

| 31.255 | 1.876 | 10.887 | 24.706 | 48.411 | •••• | 17.247 | 22.238 | 34.848 | 19.994 | 1.591 |

## Step 5 : Select the $\lambda_{Ridge}$

Validation

Select the $\lambda$ that minimizes the MSE loss on the validation data.

$\lambda\ values\ =$

| 0.00001 | 0.0001 | 0.001 | 0.01 | 0.1 | | 1 | 2 | 3 | 4 | 10 |
|---------|--------|-------|------|-----|----|---|---|---|---|----|

....

$MSE\ =$

| 31.255 | 1.876 | 10.887 | 24.706 | 48.411 | | 17.247 | 22.238 | 34.848 | 19.994 | 1.591 |
|--------|-------|--------|--------|--------|----|--------|--------|--------|--------|-------|

....

## Let's visualize and see it!

# Step 5 : Select the $\lambda_{Ridge}$

Validation

Select the $\lambda$ that minimizes the MSE loss on the validation data.

## Step 6 : Refit the model

Refit the model using both train and validation data using $\lambda_{Ridge}$.

## Step 6 : Refit the model

Refit the model using both train and validation data using $\lambda_{Ridge}$.

This gives us $\widehat{\beta}_{Ridge}(\lambda)$



Training the Model

$$\widehat{\beta}_{Ridge}(\lambda)$$

## Step 7 : Record MSE/R2

Report MSE or R2 on test data given the $\widehat{\beta}_{Ridge}(\lambda)$

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

For Lasso regression, there is **no** analytical solution for the coefficients, so we use a **solver**.



Now, the steps to determine $\lambda$ using validation will still be the same!

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

For Lasso regression, there is **no** analytical solution for the coefficients, so we use a **solver**.
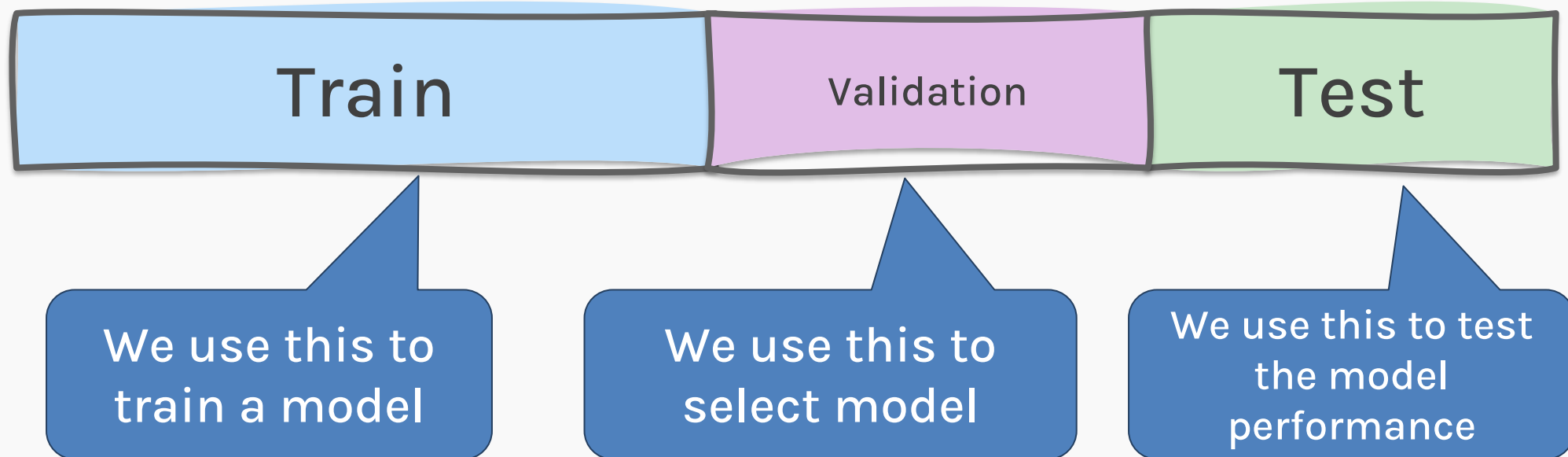
## Step 1 : Split Data

Split data into train, validation and test data.

| Train | Validation | Test |
|-------|------------|------|

We use this to train a model

We use this to select model

We use this to test the model performance

**Step 2 : Select a range of possible $\lambda$ values**

| $\lambda_{min}$ | $\lambda_{min-1}$ | .... | $\lambda_{max-1}$ | $\lambda_{max}$ |
|---|---|---|---|---|

The values of $\lambda$ are not fixed—you can choose the range yourself based on the problem.

## Step 2 : Select a range of possible $\lambda$ values

As an example, we will take a range of values from 0.00001 to 10



$\lambda \, values =$ | 0.00001 | 0.0001 | 0.001 | 0.01 | 0.1 | .... | 1 | 2 | 3 | 4 | 10 |

$\lambda_{min}$          $\lambda_{max}$

## **Step 3 : Determine** $\beta$

Train

For each $\lambda$ value, we determine $\beta_{Lasso}(\lambda)$ using train data

$\lambda\ values\ =$

| 0.00001 | 0.0001 | 0.001 | 0.01 | 0.1 |
|---------|--------|-------|------|-----|

.....

| 1 | 2 | 3 | 4 | 10 |
|---|---|---|---|----|

*Solver*

$\beta\ values\ =$



Keep in mind that each $\lambda$ value gives us a different vector of $\beta$ coefficients.

PROTOPAPAS

35

## **Step 4 : Record the** $L_{MSE}(\lambda)$

Validation

For each vector of $\beta$ coefficients, we record $L_{MSE}(\lambda)$ using the validation data

$$\beta\ values =$$



| 1.431 | 1.31 | 1.665 | 4.72 | 3.156 | $\cdots$ | 0.691 | 4.365 | 2.652 | 8.652 | 3.516 |
| -1.201 | -1.314 | -1.989 | 4.720 | -3.670 | | 0.901 | 5.865 | 1.652 | 3.652 | 4.516 |
| -1.401 | -1.341 | -1.665 | 4.702 | -3.236 | | 0.661 | 4.865 | 2.652 | 5.652 | 9.516 |

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

$$MSE =$$

| 31.255 | 1.876 | 10.887 | 24.706 | 48.411 | $\cdots$ | 17.247 | 22.238 | 34.848 | 19.994 | 1.591 |

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

## Step 5 : Select the $\lambda_{Lasso}$

Validation

Select the $\lambda$ that minimizes the MSE loss on the validation data.

$\lambda\ values =$

| 0.00001 | 0.0001 | 0.001 | 0.01 | 0.1 | .... | 1 | 2 | 3 | 4 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|

$MSE =$

| 31.255 | 1.876 | 10.887 | 24.706 | 48.411 | .... | 17.247 | 22.238 | 34.848 | 19.994 | 1.591 |
|---|---|---|---|---|---|---|---|---|---|---|

## Let's visualize and see it!

## Step 5 : Select the $\lambda_{Lasso}$

Validation

Select the $\lambda$ that minimizes the MSE loss on the validation data.

## Step 6 : Refit the model

Refit the model using both train and validation data using $\lambda_{Lasso}$.

## **Step 6 : Refit the model**

Refit the model using both train and validation data using $\lambda_{Lasso}$.

This gives us $\hat{\beta}_{Lasso}(\lambda)$



Training Data

Training the Model

$$\hat{\beta}_{Lasso}(\lambda)$$

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

## Step 7 : Record MSE/R2

Report MSE or R2 on test data given the $\widehat{\boldsymbol{\beta}}_{Lasso}(\lambda)$

# Lasso regularization with only validation : step by step

## Step 7 : Record MSE/R2

Report MSE or R2 on test data given the $\hat{\beta}_{Lasso}(\lambda)$

Instead of relying on a single validation set, we can use cross-validation to get a more reliable estimate of the best $\lambda$

# Regularization with **CV**: step by step



Train | Test

## Step 1 : Split Training Data

Split training data into K folds.



Train | Train | Train | Train | Train | Test

# Regularization with **CV**: step by step

## Step 2: Select a fold as Validation

Select one fold as validation and the rest as training data

| Train | Train | Train | Train | Valid. |
|-------|-------|-------|-------|--------|

## Step 3: Select a range of possible $\lambda$ values

$\lambda \; values \; =$

| 0.00001 | 0.0001 | 0.001 | 0.01 | 0.1 |
|---------|--------|-------|------|-----|

....

| 1 | 2 | 3 | 4 | 10 |
|---|---|---|---|----|

$\lambda_{min}$

$\lambda_{max}$

## Step 3 : Determine $\beta$

Train

For each $\lambda$ value, we determine $\beta_{Lasso/Ridge}(\lambda)$ using train data

$\lambda\ values\ =$

| 0.00001 | 0.0001 | 0.001 | 0.01 | 0.1 |
|---|---|---|---|---|

••••

| 1 | 2 | 3 | 4 | 10 |
|---|---|---|---|---|

*Ridge or Lasso Regularization*

$\beta\ values\ =$

| 1.431 | 1.31 | 1.665 | 472 | 3.156 | ••• | 0.691 | 4.365 | 3.652 | 8.652 | 3.516 |
|---|---|---|---|---|---|---|---|---|---|---|
| –1.201 | –1.314 | –1.989 | 4.730 | –3.670 | ••• | 0.901 | 5.865 | 1.652 | 3.652 | 4.516 |
| –1.401 | –1.341 | –1.665 | 4.702 | –3.236 | | 0.661 | 4.865 | 2.652 | 5.652 | 9.516 |

•••

Keep in mind that each $\lambda$ value gives us a different vector of $\beta$ coefficients.

## Step 4 : Record the $L_{MSE}(\lambda)$

For each vector of $\beta$ coefficients, we record $L_{MSE}(\lambda)$ using the validation data



$\beta\ values =$

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

$MSE\ values =$

| 31.255 | 1.876 | 10.887 | 24.706 | 48.411 | | 17.247 | 22.238 | 34.848 | 19.994 | 1.591 |

## Step 6: Average MSE, $\overline{L}_{MSE}(\lambda)$

Calculate the average MSE, $\overline{L}_{MSE}(\lambda)$ for each $\lambda$ by averaging $L_{MSE}(\lambda, k)$ over k. folds.

| $MSE_1 =$ | 31.255 | 1.876 | 10.887 | 24.706 | 48.411 | •••• | 17.247 | 22.238 | 34.848 | 19.994 | 1.591 |
|---|---|---|---|---|---|---|---|---|---|---|---|

| $MSE_2 =$ | 33.255 | 1.676 | 6.887 | 4.706 | 38.411 | •••• | 7.247 | 2.238 | 3.848 | 1.994 | 3.591 |
|---|---|---|---|---|---|---|---|---|---|---|---|

| $MSE_k =$ | 36.255 | 3.876 | 4.887 | 5.706 | 8.411 | •••• | 1.247 | 5.238 | 4.848 | 9.994 | 4.591 |
|---|---|---|---|---|---|---|---|---|---|---|---|

| $Average\ MSE =$ | 33.521 | 2.476 | 7.550 | 11.706 | 31.744 | •••• | 8.58 | 9.904 | 14.515 | 10.660 | 3.257 |
|---|---|---|---|---|---|---|---|---|---|---|---|

## Step 7: Select $\lambda$

Find the $\lambda$ that minimizes the $\overline{L}_{MSE}(\lambda)$

$Average\ MSE =$

| 33.521 | 2.476 | 7.550 | 11.706 | 31.744 |

•••• 

| 8.58 | 9.904 | 14.515 | 10.660 | 3.257 |

$\lambda\ values =$

| 0.00001 | 0.0001 | 0.001 | 0.01 | 0.1 |

•••• 

| 1 | 2 | 3 | 4 | 10 |

## Let's visualize and see it!

## Step 7: Select $\lambda$

Find the $\lambda$ that minimizes the $\bar{L}_{MSE}(\lambda)$

## Step 8: Refit the model

Refit the model using both train and validation data using $\lambda$.

## Step 8: Refit the model

Refit the model using both train and validation data using $\lambda$.

This gives us $\widehat{\boldsymbol{\beta}}(\lambda)$



Training Data

Training the Model

$\widehat{\boldsymbol{\beta}}(\lambda)$

## Step 9 : Record MSE/R2

Report MSE or R2 on test data given the $\widehat{\beta}(\lambda)$

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

For ridge regression there exist an analytical solution for the coefficients:
$$\hat{\beta}_{Ridge}(\lambda) = \left(X^T X + \lambda I\right)^{-1} X^T Y$$

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 $L_{ridge}$, $\beta_{Ridge}(\lambda) = \left(X^T X + \lambda I\right)^{-1} X^T Y$, using the train data.

    2. record $L_{MSE}(\lambda)$ using validation data.

3. select the $\lambda$ that minimizes the *MSE* loss on the validation data,
$$\lambda_{ridge} = \text{argmin}_\lambda \, L_{MSE}(\lambda)$$

4. Refit the model using both train and validation data, $\{\{X,Y\}_{train}, \{X,Y\}_{validation}\}$, now using $\lambda_{ridge}$, 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

> For Lasso regression, there is **no** analytical solution for the coefficients, so we use a **solver**.

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 $L_{lasso}$, $\beta_{lasso}(\lambda)$, using the train data. **This is done using a solver.**

   B. record $L_{MSE}(\lambda)$ using the validation data.

3. select the $\lambda$ that minimizes the **_MSE_ loss** on the validation data,

$$\lambda_{lasso} = \text{argmin}_\lambda L_{MSE}(\lambda)$$

4. Refit the model using both train and validation data, $\{\{X,Y\}_{train}, \{X,Y\}_{validation}\}$, now using $\lambda_{Lasso}$, resulting to $\hat{\beta}_{lasso}(\lambda_{lasso})$

5. Report MSE or $R^2$ on $\{X,Y\}_{test}$ given the $\hat{\beta}_{lasso}(\lambda_{lasso})$

# 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 $\{1, ..., K\}$
   for $\lambda$ in $\{\lambda_0, ..., \lambda_n\}$:

   A. determine the $\beta$ that minimizes the $L_{ridge}$, $\beta_{ridge}(\lambda, k) = \left(X^TX + \lambda I\right)^{-1} X^T Y$, **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. Calculate the average MSE, $\bar{L}_{MSE}(\lambda)$  the for each $\lambda$ by averaging $L_{MSE}(\lambda, k)$ over $k$ folds.

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 $\hat{\beta}_{ridge}(\lambda_{ridge})$