

# **STAT 215A Fall 2019**

## **Week 8c**

Tiffany Tang

10/18/19

# Plan for Today:

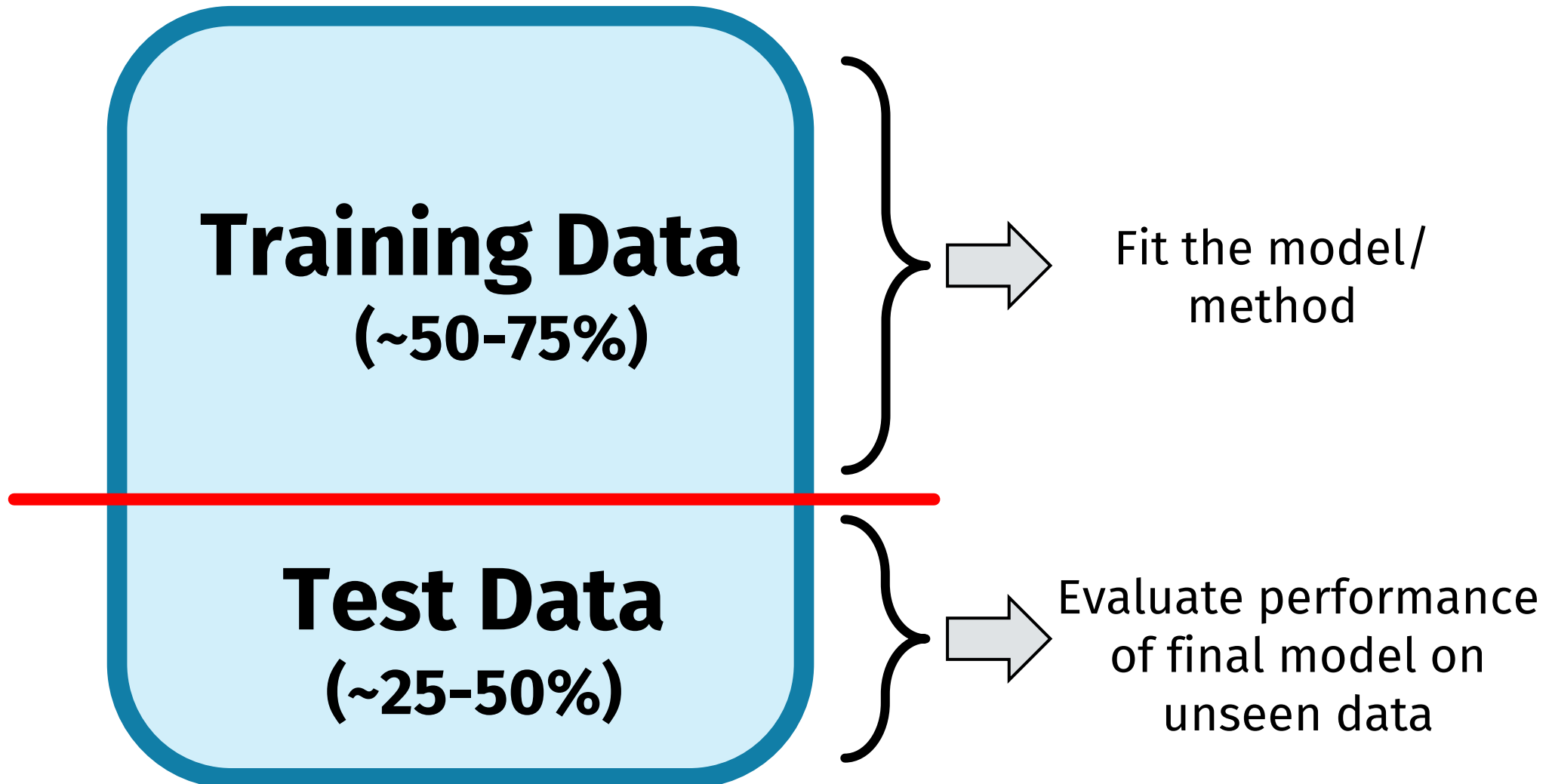
- ▶ Data splitting and the statistical learning pipeline
- ▶ Bias-variance Tradeoff
- ▶ Ordinary Least Squares
- ▶ Regularized Regression (Ridge, Lasso)

# Statistical (aka Machine) Learning

- ▶ **Goal:** to “learn” insights about the current data that are **generalizable** to future data
- ▶ Typically, have one of the following tasks in mind:
  - ▶ **Prediction** – have an outcome of interest
  - ▶ **Data-driven discoveries** – pattern recognition, feature selection, what’s important in the data
- ▶ Two main types of statistical learning:
  - ▶ **Unsupervised learning** – no outcomes/response data
    - ▶ Ex. clustering, dimension reduction, pattern recognition
  - ▶ **Supervised learning** – have outcomes/response data
    - ▶ **Classification** categorical responses
    - ▶ **Regression** – continuous responses

# How to assess generalizability: data splitting

## ► The simplest case:



# How to assess generalizability: data splitting

- ▶ Can use data splitting idea for both unsupervised and supervised learning methods
- ▶ **IMPORTANT:** Do NOT touch the test set until the very end when you are ready to *evaluate* the generalization performance of your final method
  - ▶ **Purpose of the test set:** to obtain an unbiased estimate of the prediction error for your statistical learning pipeline on future data
- ▶ In more realistic pipelines, we will want to compare multiple methods, tune hyper-parameters, etc. We will return to data splitting in this case later.

# How to evaluate?

- ▶ For unsupervised learning methods, this is difficult and typically depends on the method itself
  - ▶ Ex. NMF, kmeans
- ▶ For supervised methods, there are several possible metrics:
  - ▶ **Mean squared error**
  - ▶ Absolute loss
  - ▶ Correlation
  - ▶ AUC
  - ▶ Classification error – 0/1 loss

# Mean Squared Error

## ► Notation:

- $x_i$  = observed feature measurements from sample  $i$
  - $y_i$  = observed response for sample  $i$
  - For simplicity, assume the underlying truth is given by  $y = f(x)$
  - $\hat{y} = \hat{f}(x)$  = prediction from statistical method or algorithm
- Want prediction function to give small MSE

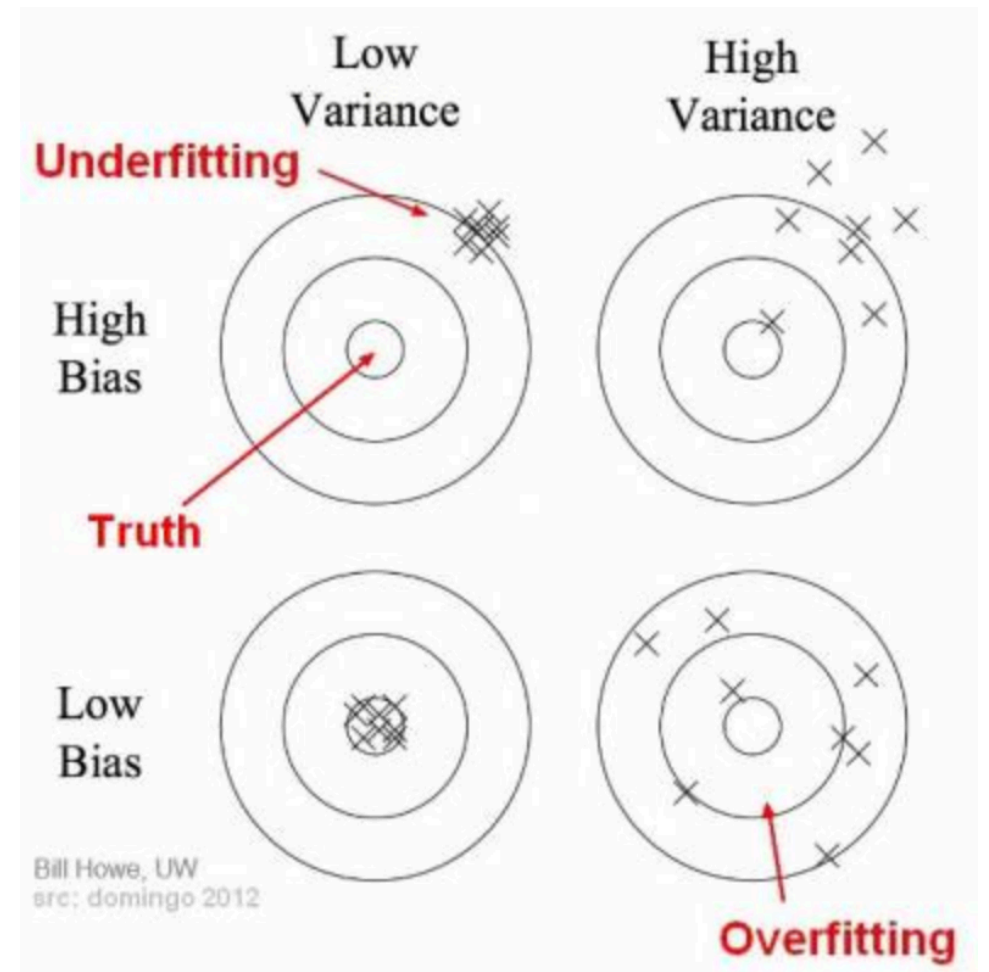
$$\text{MSE} = \mathbb{E} \left[ (Y - \hat{f}(X))^2 \right] \quad (\text{Population})$$

$$\widehat{\text{MSE}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 = \frac{1}{n} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 \quad (\text{Empirical})$$

# Bias-Variance Tradeoff

$$\text{MSE} = \mathbb{E} \left[ (Y - \hat{f}(X))^2 \right] = \text{Var}(\hat{f}(X)) + \text{Bias}^2(\hat{f}(X))$$

- ▶ **Bias:** on average, how wrong is your prediction from the truth
- ▶ **Variance:** if you obtain a new, but similar dataset, how much does this change your predictions

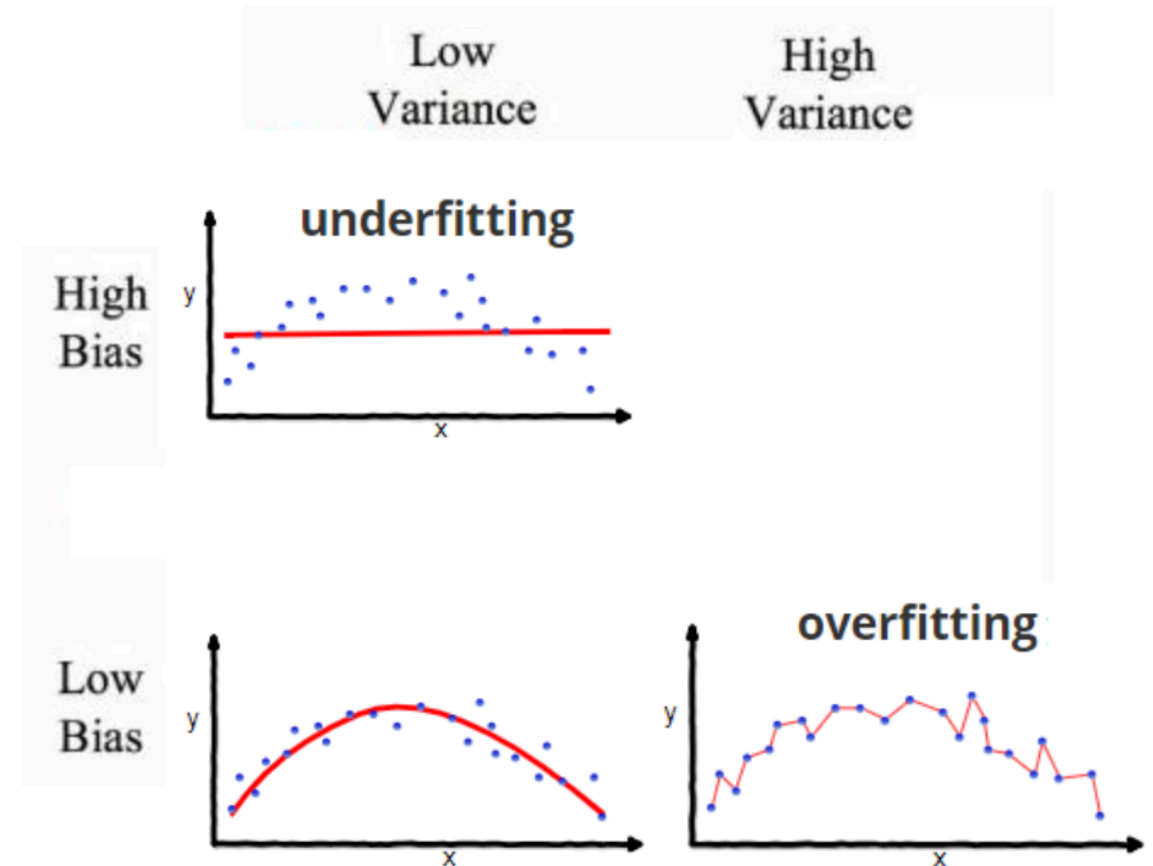




# Bias-Variance Tradeoff

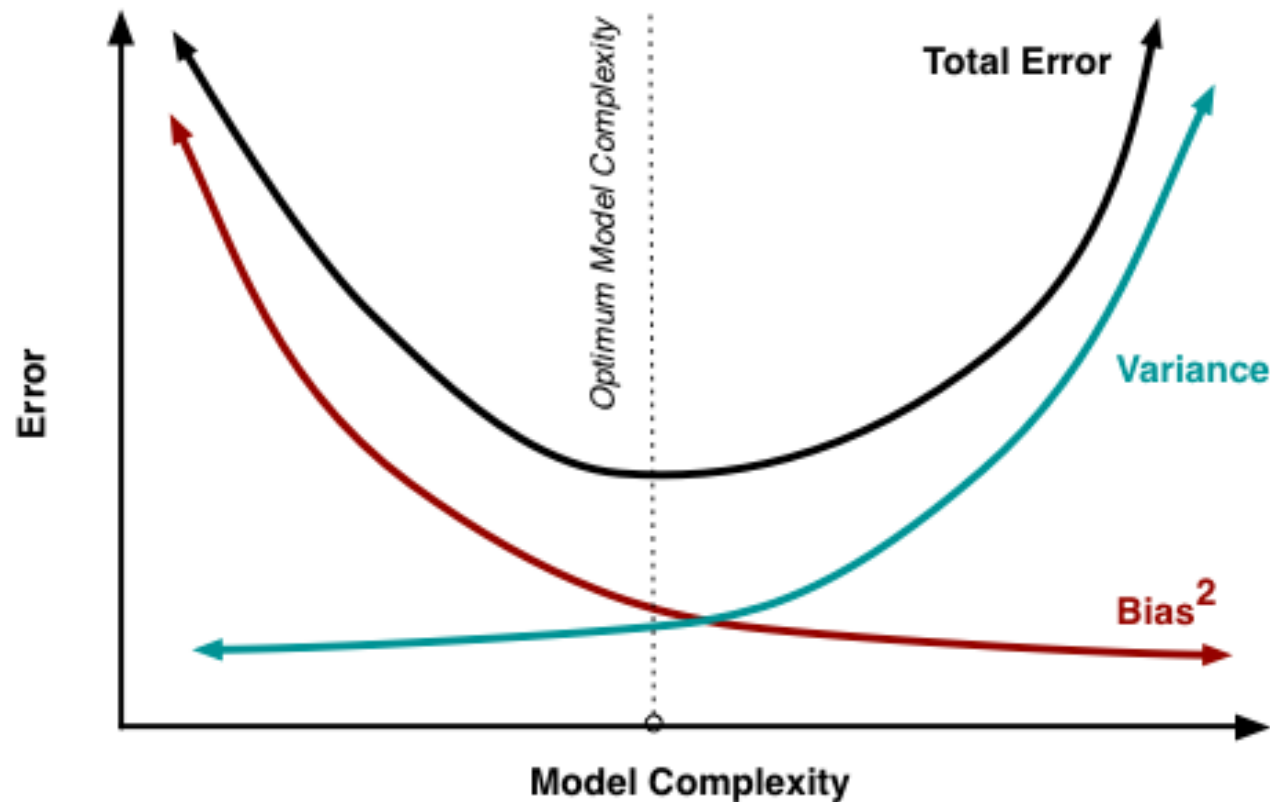
$$\text{MSE} = \mathbb{E} \left[ (Y - \hat{f}(X))^2 \right] = \text{Var}(\hat{f}(X)) + \text{Bias}^2(\hat{f}(X))$$

- ▶ **Bias:** on average, how wrong is your prediction from the truth
- ▶ **Variance:** if you obtain a new, but similar dataset, how much does this change your predictions



# Bias-Variance Tradeoff

$$\text{MSE} = \mathbb{E} \left[ (Y - \hat{f}(X))^2 \right] = \text{Var}(\hat{f}(X)) + \text{Bias}^2(\hat{f}(X))$$



# Supervised Learning: A general framework

1. Split the data into training/test sets
2. On the training data, try to minimize some **loss** function, which quantifies the error between the observed values and the predicted values

$$\min_{\hat{f} \in \mathcal{F}} \mathcal{L}(\mathbf{y}, \hat{f}(\mathbf{X}))$$

► Ex. MSE:

$$\min_{\hat{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

3. Report test set error as a measure of generalizability

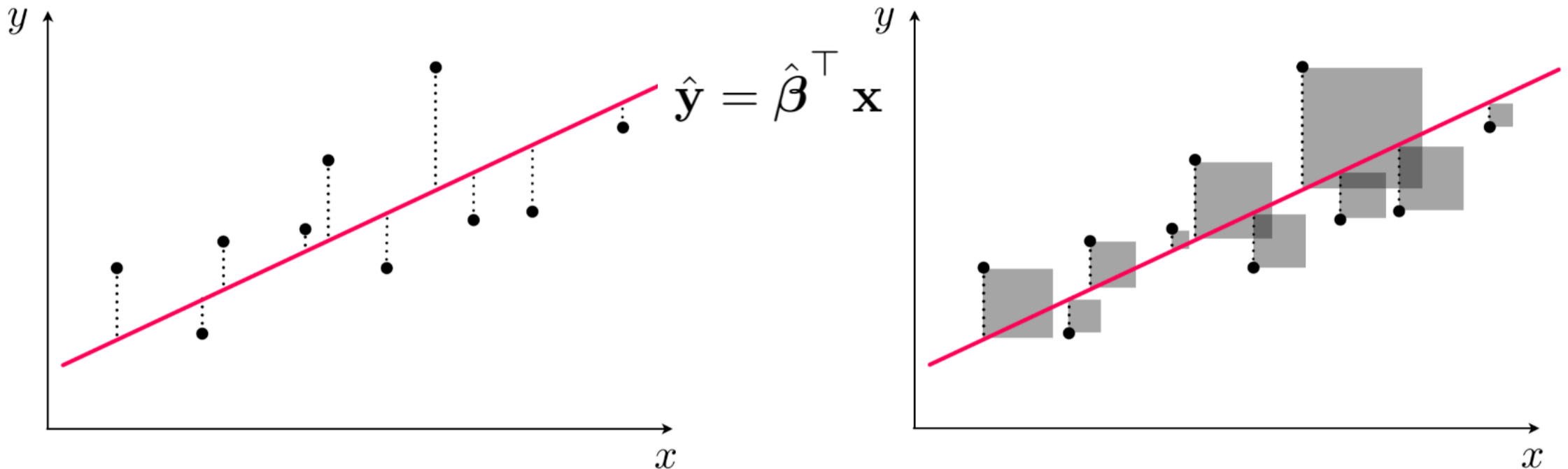
# Ordinary Least Squares (OLS)

- ▶ Find best linear approximation of the data, as measured by MSE

$$\hat{\beta}_{OLS} = \operatorname{argmin}_{\beta} \underbrace{\|\mathbf{y} - \mathbf{X}\beta\|_2^2}_{\text{Loss function}}$$

Loss function

- ▶ Closed-form solution:  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$



# Ordinary Least Squares (OLS)

## Advantages

- ▶ Simple
- ▶ Closed-form solution
- ▶ Interpretable (?)
- ▶ Under some modeling assumptions, OLS has some desirable properties

# Ordinary Least Squares (OLS) + Model

- ▶ Assume that

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where  $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$ .

- ▶ Then  $\widehat{\boldsymbol{\beta}}_{OLS}$  is an unbiased estimator for  $\boldsymbol{\beta}$ .
- ▶ If, in addition,  $\text{Var}(\varepsilon_i) = \sigma^2$  for all  $i$  and  $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$  for all  $i \neq j$ , then OLS is **BLUE** (the “best” linear unbiased estimator)
  - ▶ Out of all linear unbiased estimators,  $\widehat{\boldsymbol{\beta}}_{OLS}$  has the minimum variance (Gauss-Markov Theorem)

$$\text{MSE} = \mathbb{E} \left[ (Y - \hat{f}(X))^2 \right] = \text{Var}(\hat{f}(X)) + \text{Bias}^2(\hat{f}(X))$$

# Ordinary Least Squares

## Advantages

- ▶ Simple
- ▶ Closed-form solution
- ▶ Interpretable (?)
- ▶ Under some modeling assumptions, OLS has some desirable properties

## Disadvantages

- ▶ Too simple
- ▶ When  $p > n$ 
  - ▶ Non-unique solution; high variance
  - ▶ Massive overfitting since training error = 0

# One possible solution

$$\text{MSE} = \mathbb{E} \left[ (Y - \hat{f}(X))^2 \right] = \text{Var}(\hat{f}(X)) + \text{Bias}^2(\hat{f}(X))$$

- ▶ Let's sacrifice some bias for a larger reduction in variance
- ▶ One way to reduce the variance of your predictions is to restrict the class of functions you search over
- ▶ In the linear setting, this motivates regularized linear regression methods such as ridge, Lasso, and elastic net



# Ridge Regression

$$\hat{\beta}^r = \underset{\beta}{\operatorname{argmin}} \underbrace{\|\mathbf{y} - \mathbf{X}\beta\|_2^2}_{\text{Loss function}} \quad \text{subject to} \quad \underbrace{\|\beta\|_2^2 \leq \tau}_{\text{Constraint}}$$

or equivalently,

$$\hat{\beta}^r = \underset{\beta}{\operatorname{argmin}} \underbrace{\|\mathbf{y} - \mathbf{X}\beta\|_2^2}_{\text{Loss function}} + \underbrace{\lambda \|\beta\|_2^2}_{\text{Penalty}}$$

- ▶  $\lambda \geq 0$  is called a tuning or penalty parameter and regulates how much *shrinkage* is introduced
  - ▶  $\lambda \rightarrow 0 \Rightarrow \hat{\beta}^r \rightarrow \text{OLS}$  (no bias, high variance)
  - ▶  $\lambda \rightarrow \infty \Rightarrow \hat{\beta}^r \rightarrow 0$  (high bias, but no variance)
  - ▶  $\lambda$  somewhere in between  $\Rightarrow$  some shrinkage
  - ▶ Choose  $\lambda$  via CV

## But still, why is ridge regression a good idea?

$$\hat{\boldsymbol{\beta}}^r = \operatorname{argmin}_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X} \boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$$

1. The solution exists and is unique even when  $p > n$  (unlike OLS)

$$\hat{\boldsymbol{\beta}}^r = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

## But still, why is ridge regression a good idea?

1. The solution exists and is unique even when  $p > n$  (unlike OLS)

$$\hat{\boldsymbol{\beta}}^r = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

2. Under the model from before,

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \text{where} \quad \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}, \quad \text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$$

there **always** exists a  $\lambda$  such that the Ridge( $\lambda$ ) MSE is less than the OLS MSE:

$$\text{MSE} \left( \mathbf{X} \hat{\boldsymbol{\beta}}^r(\lambda) \right) < \text{MSE} \left( \mathbf{X} \hat{\boldsymbol{\beta}}^{OLS} \right)$$

# But still, why is ridge regression a good idea?

1. The solution exists and is unique even when  $p > n$  (unlike OLS)

$$\hat{\boldsymbol{\beta}}^r = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

2. Under the model from before,

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \text{where} \quad \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}, \quad \text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$$

there **always** exists a  $\lambda$  such that the Ridge( $\lambda$ ) MSE is less than the OLS MSE:

$$\text{MSE} \left( \mathbf{X} \hat{\boldsymbol{\beta}}^r(\lambda) \right) < \text{MSE} \left( \mathbf{X} \hat{\boldsymbol{\beta}}^{OLS} \right)$$

3. Handles correlated features very well – Features that are highly correlated tend to be shrunk together, i.e. they are given equal contribution to the linear model

# Ridge Regression in Practice

- ▶ Don't want to penalize the intercept term
  - ▶ Before applying any regularized regression such as ridge regression, center columns of  $\mathbf{X}$  and center  $\mathbf{y}$
- ▶ Most of the time, should also scale columns of  $\mathbf{X}$  so that we don't penalize coefficients more than others simply because of differing scales
- ▶ These practical guidelines apply to ALL regularization methods

# Feature Selection

- ▶ Want to select the best subset of  $s$  variables for prediction (typically  $s \ll p$ )
- ▶ Looking for data-driven discoveries or interpretable insights
- ▶ Ridge regression does not perform *automatic* feature selection
  - ▶ Can use ridge regression to order the importance of features using the magnitude of their coefficients (which is often done in practice)
  - ▶ But typically there are no exact 0's in the coefficients vector
- ▶ Methods that more directly solve the feature selection problem:
  - ▶ Best subsets ( $L_0$ ) regression
  - ▶ Lasso

# Best Subsets ( $L_0$ ) Regression

$$\hat{\boldsymbol{\beta}}^{best} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \quad ||\mathbf{y} - \mathbf{X}\boldsymbol{\beta}||_2^2 \quad \text{subject to} \quad ||\boldsymbol{\beta}||_0 \leq s$$

- ▶ Want to find the subset of  $s$  features that give the smallest MSE
- ▶ **Major problem:**
  - ▶ Computationally infeasible for moderate-sized problems
  - ▶ Need to do an exhaustive search over  $\binom{p}{s}$  regressions

# Lasso ( $L_1$ ) Regression

$$\hat{\beta}^l = \underset{\beta}{\operatorname{argmin}} \quad ||\mathbf{y} - \mathbf{X}\beta||_2^2 \quad \text{subject to} \quad ||\beta||_1 \leq \tau$$

or equivalently,

$$\hat{\beta}^l = \underset{\beta}{\operatorname{argmin}} \quad ||\mathbf{y} - \mathbf{X}\beta||_2^2 + \lambda ||\beta||_1$$

- ▶ “Best convex relaxation” of the best subsets ( $L_0$ ) problem (i.e., the next best thing that we can actually compute in a reasonable amount of time)
- ▶ This  $L_1$  penalty results in **sparsity** (lots of non-zero entries in  $\hat{\beta}^l$ )
- ▶ More interpretable than ridge but can be worse in prediction if the true underlying process is not exactly sparse



# Lasso ( $L_1$ ) Regression

$$\hat{\beta}^l = \underset{\beta}{\operatorname{argmin}} \quad ||\mathbf{y} - \mathbf{X}\beta||_2^2 \quad \text{subject to} \quad ||\beta||_1 \leq \tau$$

or equivalently,

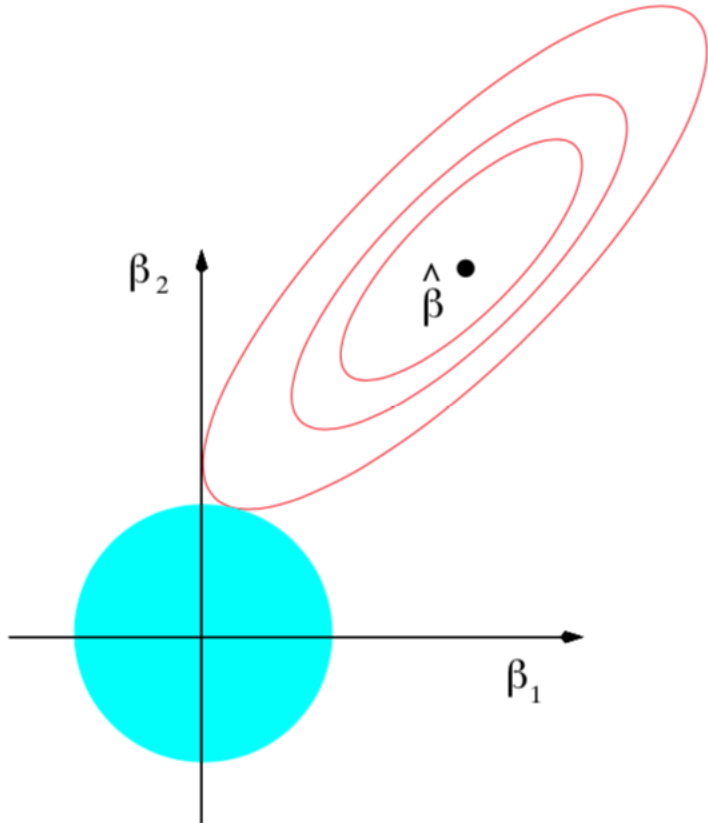
$$\hat{\beta}^l = \underset{\beta}{\operatorname{argmin}} \quad ||\mathbf{y} - \mathbf{X}\beta||_2^2 + \lambda ||\beta||_1$$

- ▶  $\lambda \geq 0$  is called a tuning or penalty parameter and regulates how the *sparsity* level
  - ▶  $\lambda \rightarrow 0 \Rightarrow \hat{\beta}^l \rightarrow \text{OLS}$  (no bias, high variance)
  - ▶  $\lambda \rightarrow \infty \Rightarrow \hat{\beta}^l \rightarrow 0$  (high bias, but no variance)
  - ▶  $\lambda$  somewhere in between  $\Rightarrow$  some sparsity
  - ▶ Choose  $\lambda$  via CV

# Why does the Lasso induce sparsity?

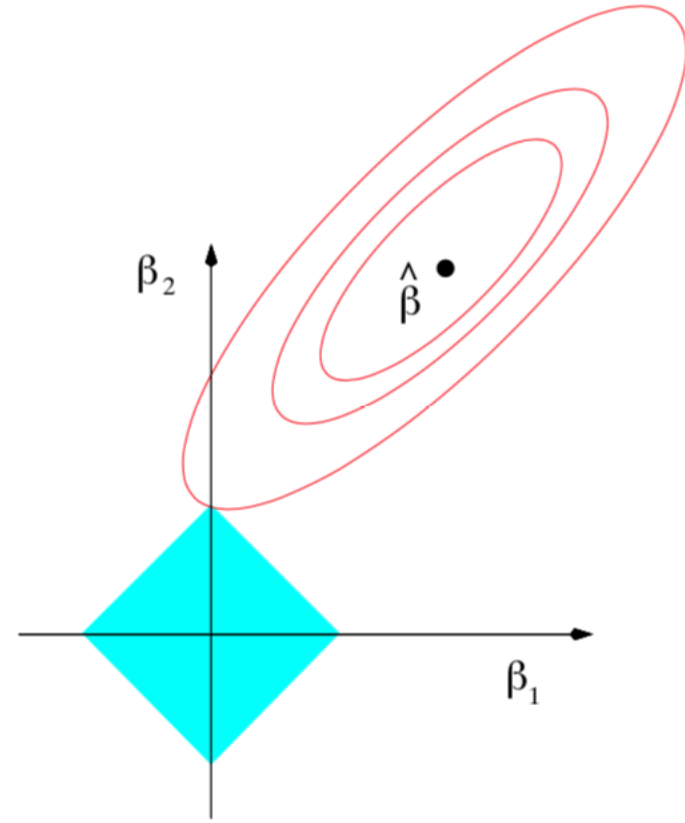
$$\hat{\beta}^r = \underset{\beta}{\operatorname{argmin}} \quad ||\mathbf{y} - \mathbf{X} \beta||_2^2$$

subject to  $||\beta||_2^2 \leq \tau$



$$\hat{\beta}^l = \underset{\beta}{\operatorname{argmin}} \quad ||\mathbf{y} - \mathbf{X} \beta||_2^2$$

subject to  $||\beta||_1 \leq \tau$



# Lasso in Practice

- ▶ Great for feature selection and interpretability
- ▶ **Warning:**
  - ▶ Non-zero Lasso coefficient  $\neq$  statistically significant
  - ▶ Cannot fit OLS after doing Lasso and interpret the p-values in the same way; OLS inference is no longer valid after doing model selection
- ▶ Tends to select one representative from a group of correlated features
  - ▶ Good if goal is a parsimonious model
  - ▶ Bad for data-driven discoveries
- ▶ Like ridge regression, don't forget to center and scale data
- ▶ For large  $\lambda$ , the coefficients of selected variables are heavily biased

# Other Alternatives/Modifications to the Lasso

- ▶ **Weighted Lasso:** use weighted penalty =  $\lambda \sum_{j=1}^p w_j |\beta_j|$ 
  - ▶ If true  $\beta_j$  is large, don't want to penalize it, so want  $w_j$  small
  - ▶ If true  $\beta_j$  is small, want to penalize it more, so want  $w_j$  large
- ▶ **Adaptive Lasso ([Zou \(2006\)](#)):**  $w_j = \left( \frac{1}{|\hat{\beta}_j^{OLS/ridge}|} \right)^\gamma \quad (\gamma \in \{0.5, 1, 2\})$
- ▶ **Non-convex Penalties:**
  - ▶ **SCAD ([Fan \(2001\)](#))**
  - ▶ **MCP ([Zhang \(2009\)](#))** - variables that are non-zero in this model are essentially unbiased
- ▶ **Randomized Lasso ([Meinshausen \(2010\)](#)):** combines stability via bootstrapping and randomized weights to handle correlated variables

# Elastic Net

$$\hat{\boldsymbol{\beta}}^{el} = \operatorname{argmin}_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X} \boldsymbol{\beta}\|_2^2 + \lambda (\alpha \|\boldsymbol{\beta}\|_1 + (1 - \alpha) \|\boldsymbol{\beta}\|_2^2)$$

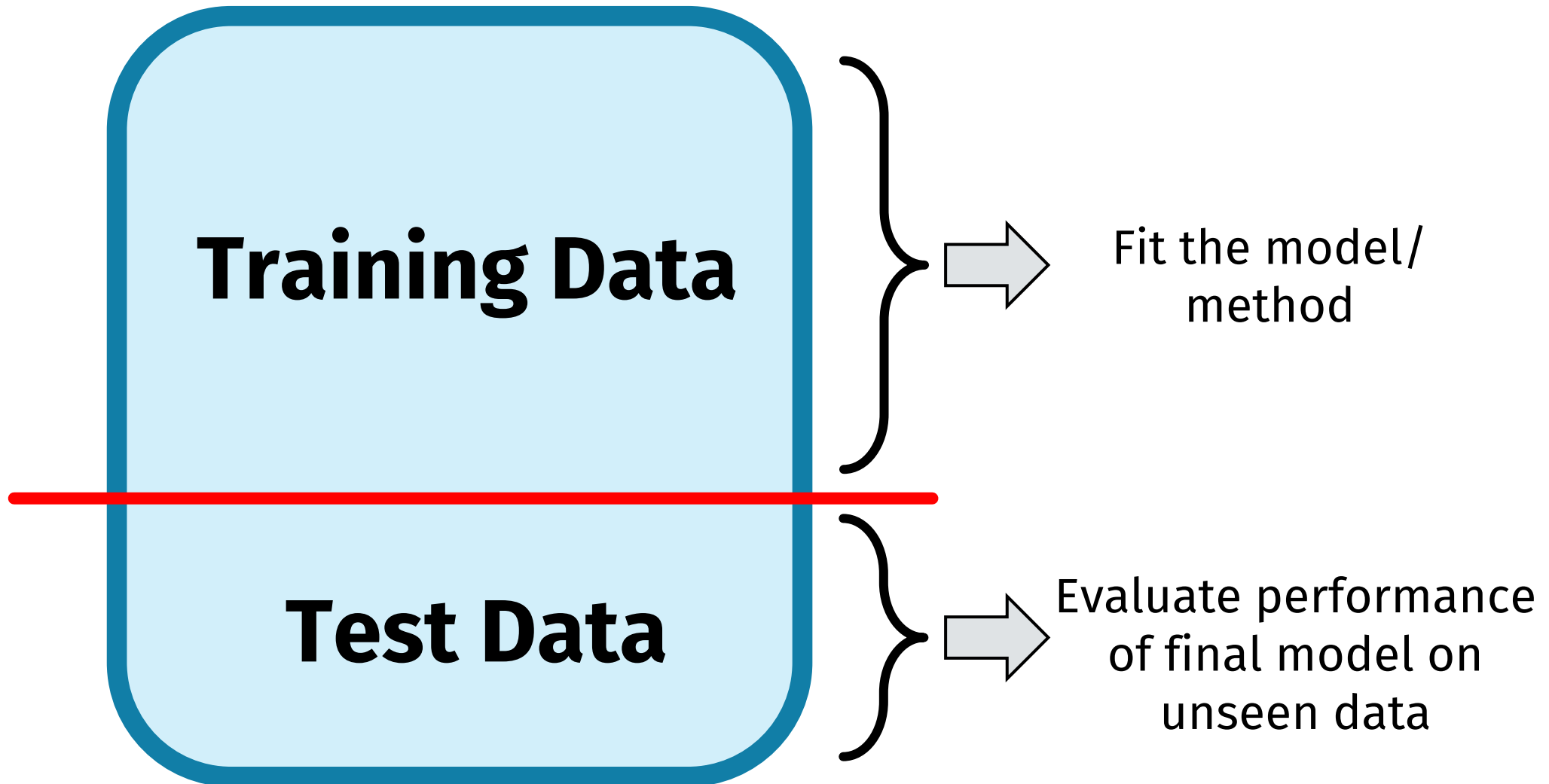
- ▶ Compromise between ridge and Lasso depending on choice of  $\alpha \in [0, 1]$ 
  - ▶ Ridge:  $\alpha = 0$
  - ▶ Lasso:  $\alpha = 1$
- ▶ Can handle correlated features better than the Lasso and enforces more sparsity than ridge
- ▶ Difficult to tune both  $\alpha$  and  $\lambda$  in practice

# Review: OLS + Regularized Regression

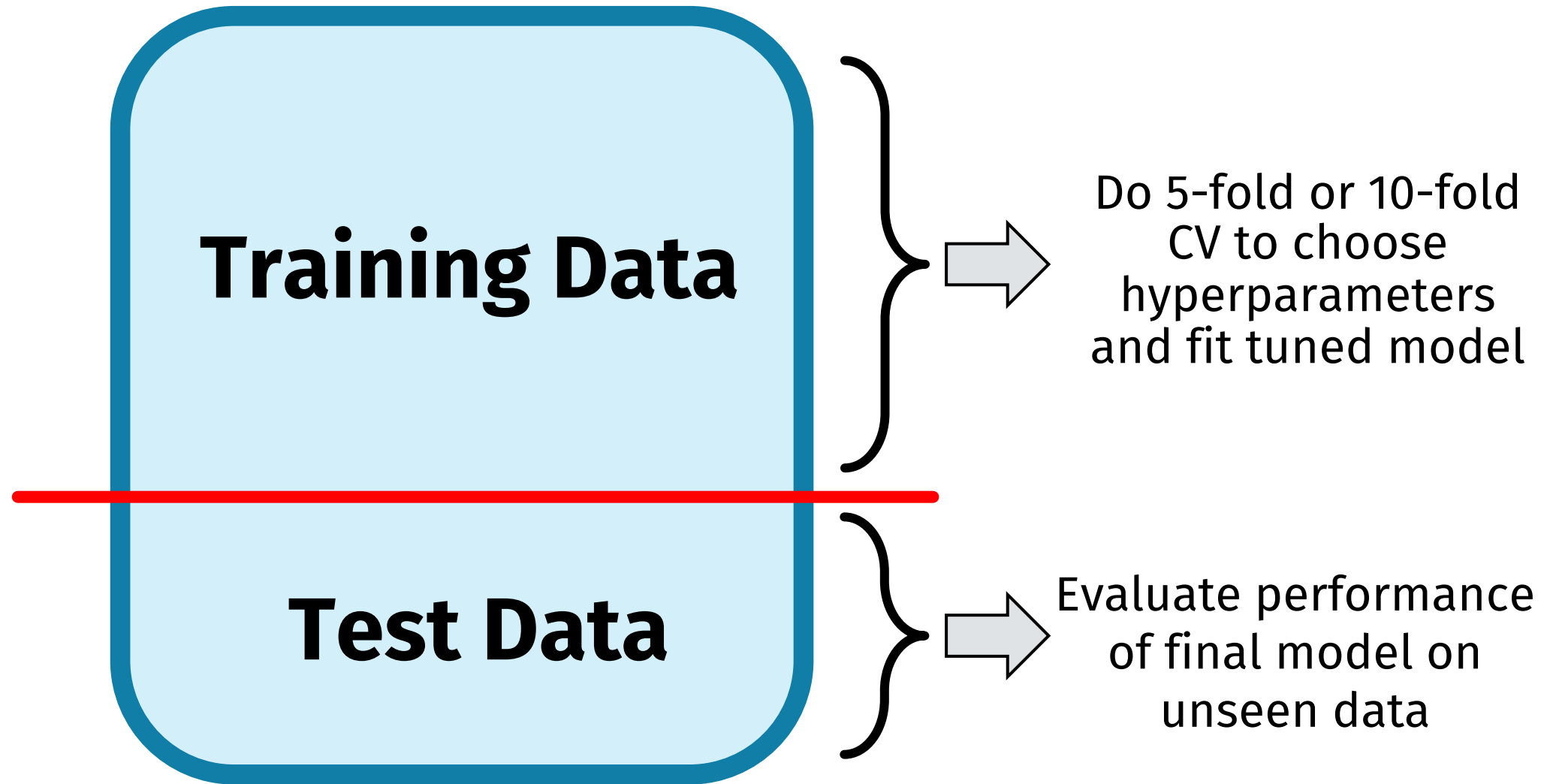
	OLS	Ridge	Lasso	Elastic Net
Advantages	<ul style="list-style-type: none"><li>• Simple</li><li>• Can do inference (with all the caveats)</li></ul>	<ul style="list-style-type: none"><li>• Good with correlated features</li><li>• MSE Existence Theorem (good for prediction)</li></ul>	<ul style="list-style-type: none"><li>• Feature selection and sparsity</li><li>• Interpretability</li></ul>	<ul style="list-style-type: none"><li>• Compromise between ridge and Lasso</li></ul>
Disadvantages	<ul style="list-style-type: none"><li>• Major problems when <math>p &gt; n</math></li></ul>	<ul style="list-style-type: none"><li>• Dense model is not interpretable</li><li>• No automatic feature selection</li></ul>	<ul style="list-style-type: none"><li>• Tends to choose 1 feature from correlated group</li><li>• Can give worse prediction results if truth is not sparse</li></ul>	<ul style="list-style-type: none"><li>• Difficult to tune two parameters in practice</li></ul>

# Back to Data Splitting

## ► The simplest case:

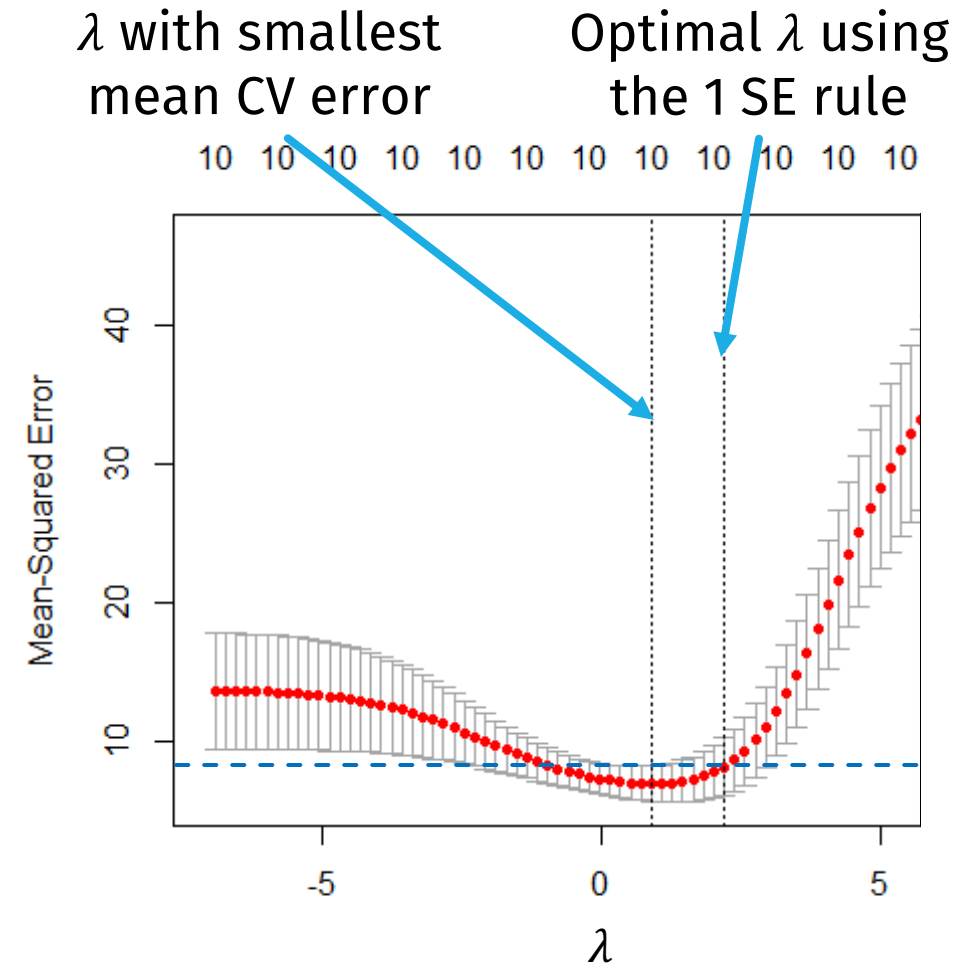


# Data Splitting + Tuning Hyperparameters

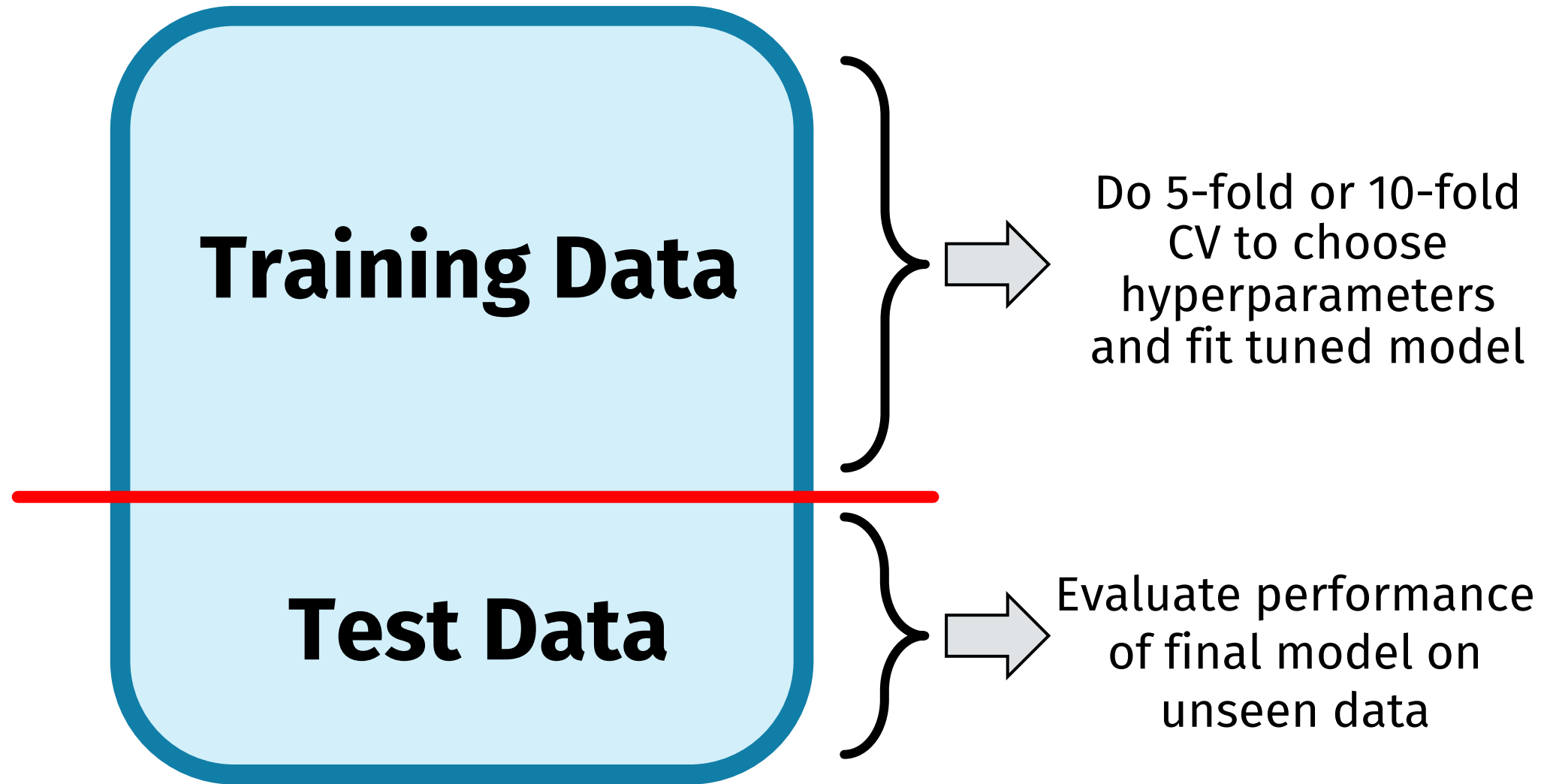




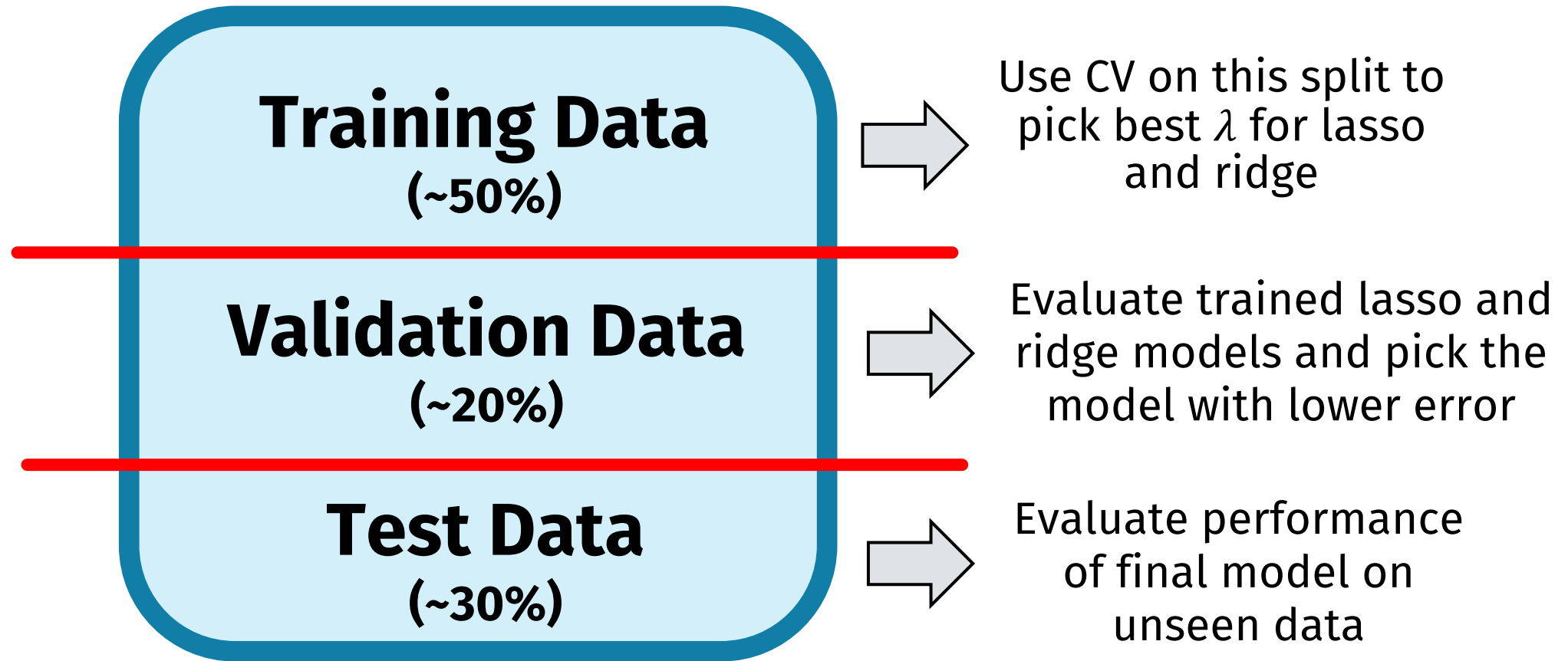
# K-fold Cross Validation for choosing hyperparameters



# Data Splitting + Tuning Hyperparameters



# Data Splitting + Tuning Hyperparameters + Multiple Methods



- ▶ Can repeat this data splitting B times to get a variance estimate of the test error
- ▶ This gives you an unbiased estimate of the prediction error for the **statistical learning process**, NOT a specific model

## Let's implement this pipeline in R using the ozone data

1. Pick two (or more) regression methods and an error metric
  - ▶ Ex: ridge, lasso, elastic net, SCAD, MCP, etc.
2. Split the data into training/validation/test splits
3. On the training set, use CV to tune hyperparameters for each method under consideration
4. Evaluate the tuned methods on the validation set
5. Report the test error of your “best” method
6. Time permitting, repeat steps 2-5 B times to get a variance estimate
7. Time permitting, perform EDA to “diagnose” the errors – are there any observations that are particularly difficult to predict? Any common mistakes?