

# Lecture 17: Regularization, Feature Selection, Cross-Validation Tests

GENOME 560, Spring 2017

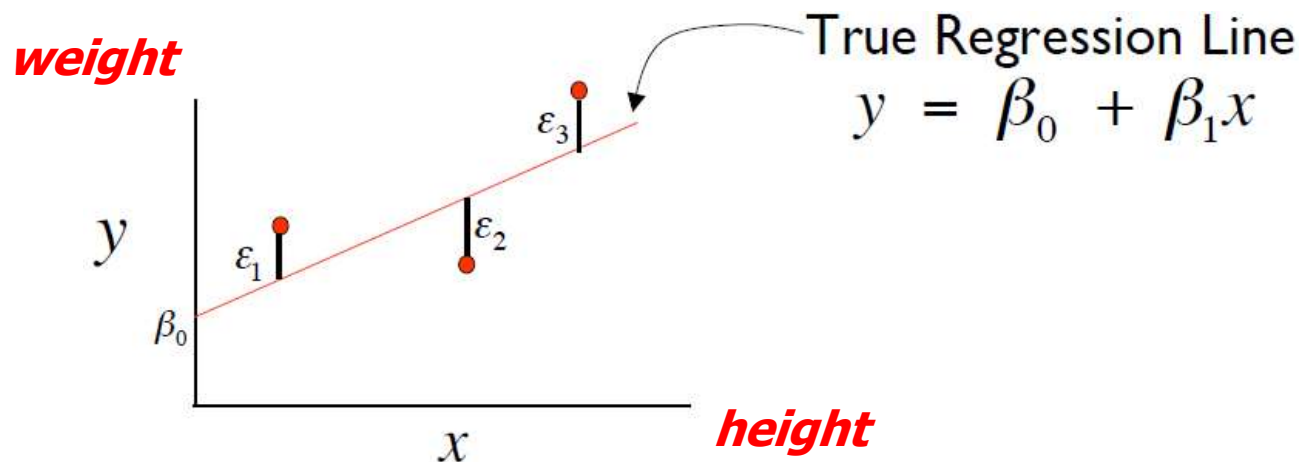
Su-In Lee, CSE & GS (suinlee@uw.edu)

# Review of Last Lecture

- Linear regression
  - Linear regression is a probabilistic model.
- **Definition:** There exists **parameters  $\beta_0, \beta_1$  and  $\sigma^2$** , such that for any fixed value of the predictor variable  $X$ , the outcome variable  $Y$  is related to  $X$  through the model equation:

$$Y = \beta_0 + \beta_1 X + \varepsilon ,$$

where  $\varepsilon$  is a RV assumed to be  $N(0, \sigma^2)$



# Review of Last Lecture

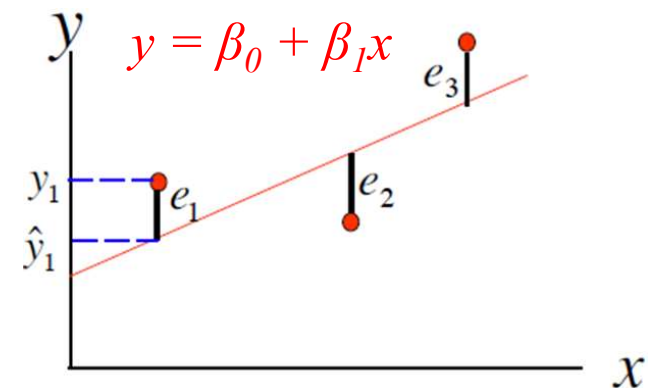
- Parameter estimation

- **Least squares**

- Find  $\beta_0$  and  $\beta_1$  that minimizes SSE.

$$f(\beta_0, \beta_1) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

- Denote the solutions by  $\hat{\beta}_0$  and  $\hat{\beta}_1$ .



- Coefficient of determination ( $R^2$ )

# Outline

- High-dimensional regression problems



- Challenges with high-dimensional data
- Feature selection

- Model selection, Cross validation

- Cross validation test
- L1 vs L2 regularization

- R-session

- Ridge regression using `lm.ridge`

# When there are many predictors

- Say that there are  $p$  predictor variables and the input data have  $n$  samples ( $p > n$ )
- Example
  - Outcome variable  $Y$ : systolic blood pressure
  - Predictor variables  $X_1, \dots, X_p$ : expression levels of  $p$  genes
  - Input data:  $(y, x_1, x_2, \dots, x_p)$  from each of  $n$  patients
    - $(y_1, x_{11}, x_{21}, \dots, x_{p1})$  from the 1<sup>st</sup> patient,
    - $(y_2, x_{12}, x_{22}, \dots, x_{p2})$  from the 2<sup>nd</sup> patient,
    - :
    - $(y_n, x_{1n}, x_{2n}, \dots, x_{pn})$  from the  $n^{\text{th}}$  patient

# Challenges

- Least squares problems have an infinite number of solutions
- Explanation without using linear algebra

- There are *more unknowns than equations*

- Input data:  $(y, x_1, x_2, \dots, x_p)$  from each of  $n$  patients

$$y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{21} + \dots + \beta_p x_{p1} + \varepsilon_1 \quad \longrightarrow \quad 1^{\text{st}} \text{ patient}$$

$$y_2 = \beta_0 + \beta_1 x_{12} + \beta_2 x_{22} + \dots + \beta_p x_{p2} + \varepsilon_2 \quad \longrightarrow \quad 2^{\text{nd}} \text{ patient}$$

⋮

$$y_n = \beta_0 + \beta_1 x_{1n} + \beta_2 x_{2n} + \dots + \beta_p x_{pn} + \varepsilon_n \quad \longrightarrow \quad n^{\text{th}} \text{ patient}$$

- There are  $(p+1)$  unknowns,  $\beta_0, \dots, \beta_p$ , and there are  $n$  equations
  - If  $p > n$ , There is an infinite number of ways to perfectly fit to the linear model ( $\varepsilon_i=0$ )

# Regularization

- We need more equations or *constraints*

- “Regularized” least squares

- The cost function is defined as:

$$f(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})]^2 + \lambda \sum_{j=1}^p \beta_j^2$$

- *L2 regularization term* encourages to choose  $\beta_1, \dots, \beta_p$  that have small magnitude
    - $\lambda$  : tuning parameter
    - With a large enough  $\lambda$ , there is not an infinite number of solutions when  $p > n$
- *L2 regularized linear regression* is also called *Ridge regression*

# Feature selection

- Sometimes, we have *many potential predictors* and want to *select relevant predictors* among them
- Each predictor is called a “feature”
- Motivating example
  - Outcome variable  $Y$  : systolic blood pressure
  - Predictor variables  $X_1, \dots, X_p$  : expression levels of  $p$  genesWe don't need all genes' expression levels to predict systolic blood pressure.



# Feature selection via L1 Regularization

- We want to select a small number of non-zero features
- L1 regularized linear regression

- The cost function is defined as:

$$f(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})]^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- *L1 regularization term* encourages many of  $\beta_1, \dots, \beta_p$  to be set to zero
    - $\lambda$  : tuning parameter
  - *L1 regularized linear regression* is also called *LASSO regression*\*

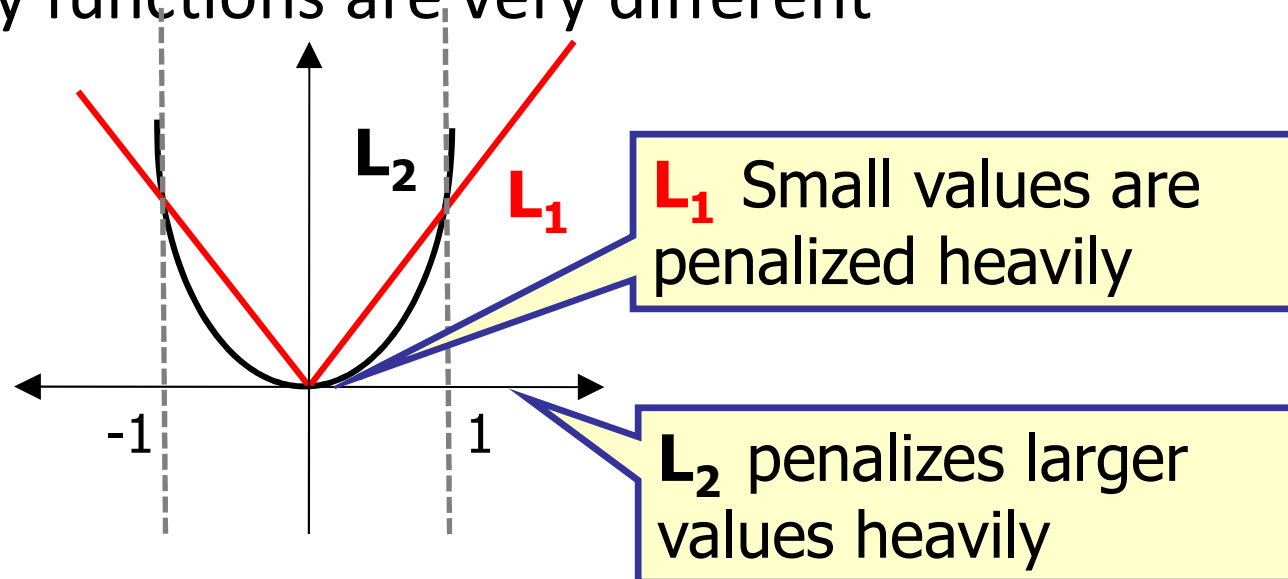
# L1 vs. L2 regularization

- Find  $\beta$ 's that minimize

$$\mathbf{L1:} \quad \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})]^2 + \lambda \sum_{j=1}^p |\beta_j|$$

$$\mathbf{L2:} \quad \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})]^2 + \lambda \sum_{j=1}^p \beta_j^2$$

- Penalty functions are very different



# How to determine $\lambda$ ?

- L1 regularized linear regression (LASSO)

- Find  $\beta$  values that minimizes the cost function:

$$f(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi})]^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- Large  $\lambda$ : Large penalty term in the objective will force  $\beta$  weights to be set to zero
    - Small  $\lambda$ : Small penalty term will force  $\beta$  weights to be set to zero less strongly than when  $\lambda$  is large
- The tuning parameter  $\lambda$  determines how much  $\beta$  weights can be **sparse (having many zeros)**
- How do we select the “optimal” tuning parameter  $\lambda$  ?

# Outline

- High-dimensional regression problems
  - Challenges with high-dimensional data
  - Feature selection
- Model selection, Cross validation
  - Cross validation test
  - L1 vs L2 regularization
- R-session
  - Ridge regression using `lm.ridge`



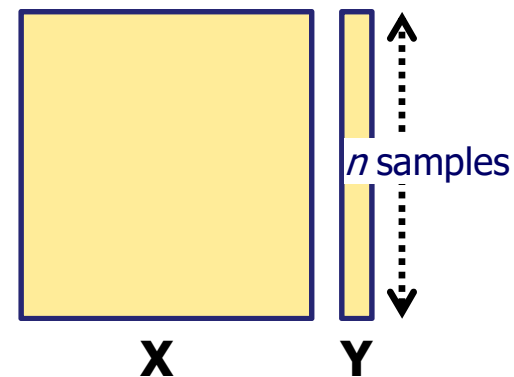
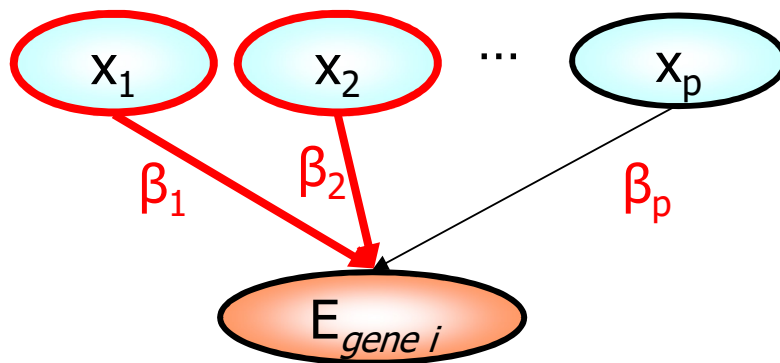
# Cross validation test

- Model selection

- Which regularization method do we want to use?
- How do we select the “optimal” tuning parameter  $\lambda$  ?

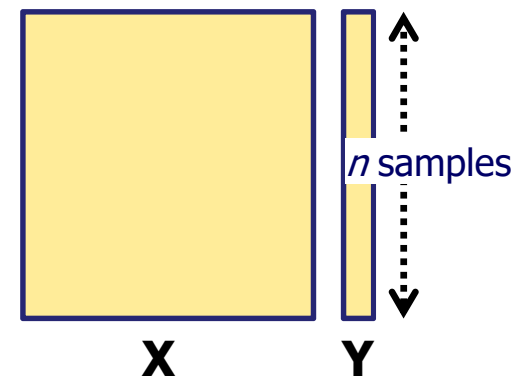
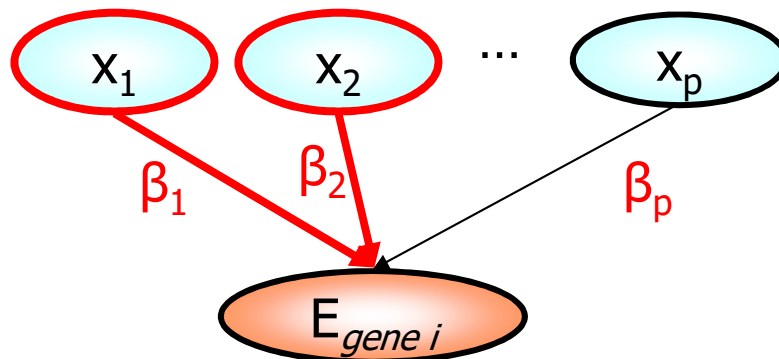
- Criteria

- How well each model fits data?
- How do we estimate the model’s “true” error rate?



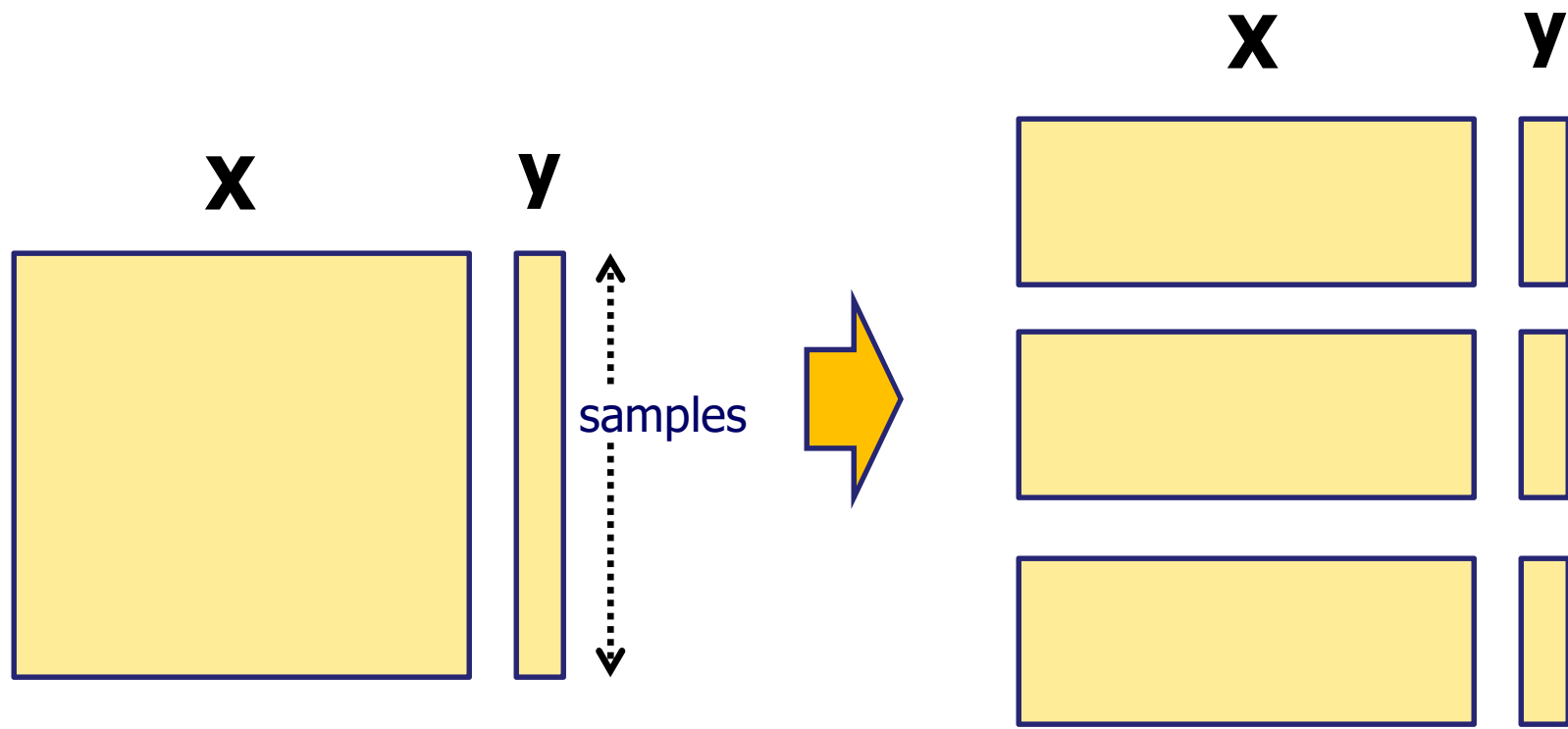
# Cross validation test

- If we had access to *an unlimited number of examples*, these questions would have a straightforward answer
  - Choose the model that provides the lowest error rate on the entire population
  - And, of course, that error rate is the true error rate
- However, in real applications *only a finite set of examples is available*
  - This number is usually smaller than we would hope for!
  - Why? Data collection is a very expensive process



# Cross validation test

- How well the estimated weight values explain the *left out* data (unseen data)?
- Let's divide the samples into  $k$  ( $=3$ ) groups



# K-fold cross validation test

- Let's divide the samples into  $k$  ( $=3$ ) groups

In each of  $k$  folds,

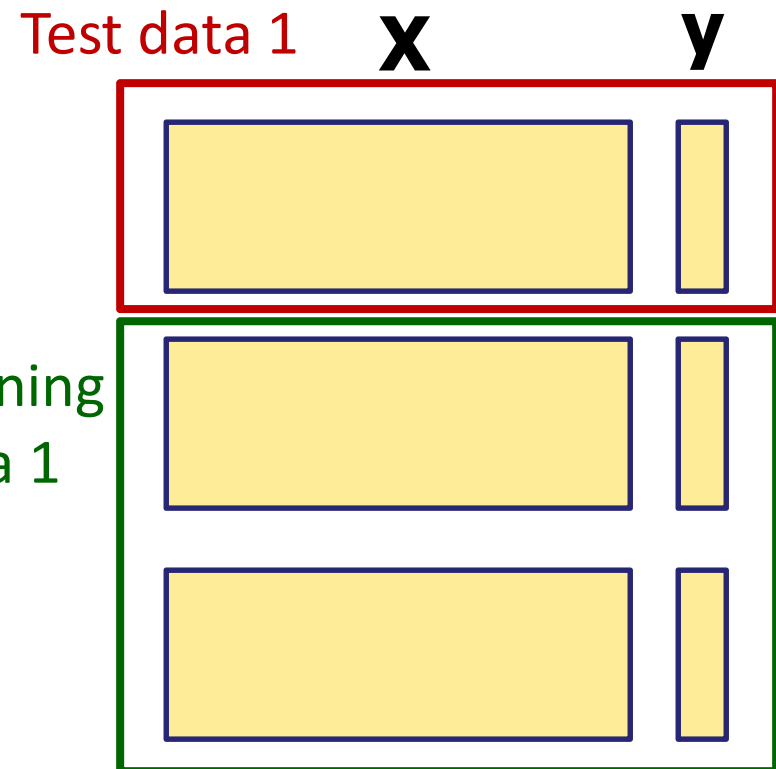
- **Training:** learning the  $\beta$  values using the **Training data 1**
- **Model:** the resulting linear regression model with the learned  $\beta$  values
- **Test:** compute the error in held-out **test data 1**
- We call it “test error”

$$f(\beta_0, \beta_1, \dots, \beta_p)$$

$$= \sum_{i=1}^n \left[ y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi}) \right]^2$$

actual y value

predicted y value  
based on the learned  $\beta$  values





# K-fold cross validation test

- Let's divide the samples into  $k$  ( $=3$ ) groups

In each of  $k$  folds,

- **Training:** learning the  $\beta$  values using the **Training data 1**
- **Model:** the resulting linear regression model with the learned  $\beta$  values
- **Test:** compute the error in held-out **test data 1**
- We call it “test error”

$$f(\beta_0, \beta_1, \dots, \beta_p)$$

$$= \sum_{i=1}^n \left[ y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi}) \right]^2$$

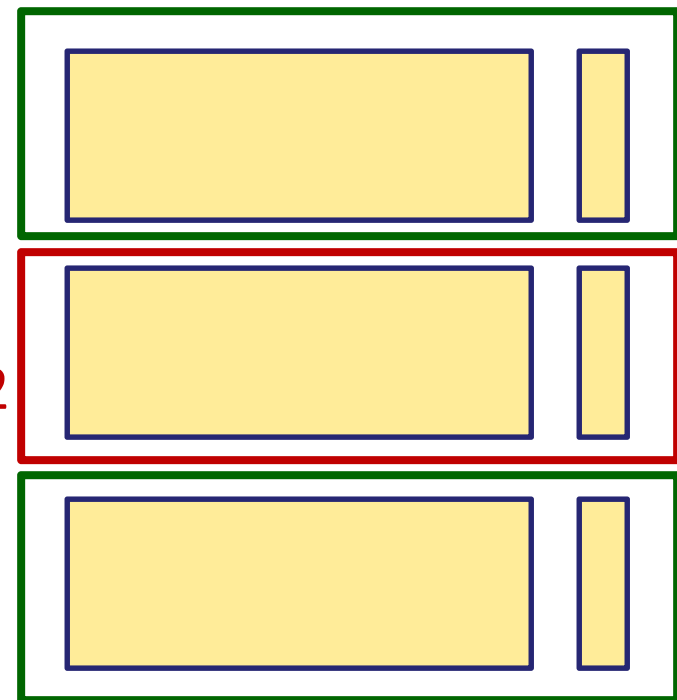
actual y value

predicted y value  
based on the learned  $\beta$  values

Training  
data 2

**X**

**y**



# K-fold cross validation test

- Let's divide the samples into  $k$  ( $=3$ ) groups

In each of  $k$  folds,

- **Training:** learning the  $\beta$  values using the **Training data 1**
- **Model:** the resulting linear regression model with the learned  $\beta$  values
- **Test:** compute the error in held-out **test data 1**
- We call it "test error"

$$f(\beta_0, \beta_1, \dots, \beta_p)$$

$$= \sum_{i=1}^n \left[ y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi}) \right]^2$$

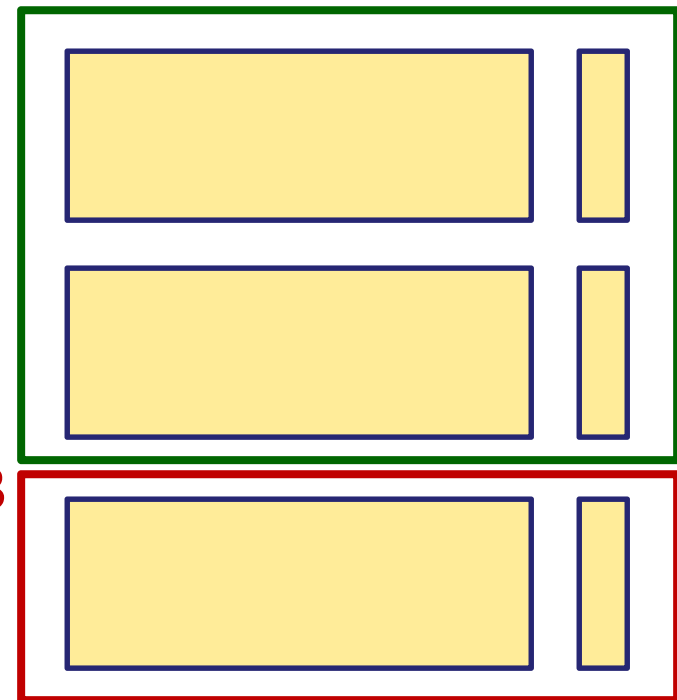
actual y value

predicted y value  
based on the learned  $\beta$  values

Training  
data 3

**X**

**y**



# k-fold cross-validation test

- In k-fold cross-validation the data is first *partitioned into  $k$  equally (or nearly equally) sized segments* or folds.
- Subsequently  *$k$  iterations of training and validation* are performed such that within each iteration a different fold of the data is held-out for validation while the remaining  $k - 1$  folds are used for learning.
- The figure in the previous slide demonstrates an example with  $k = 3$ .
- In data mining and machine learning 10-fold cross-validation ( $k = 10$ ) is the most common.

# Leave-One-Out Cross Validation

- **LOOCV**: a special case of k-fold cross-validation where k equals the number of samples in the data.
- In other words in each iteration nearly all the data except for a single observation (sample) are used for training and the model is tested on that single observation.
- An accuracy estimate obtained using LOOCV is known to be almost unbiased.
- It is widely used when the available data are very rare, especially in biology where only dozens of data samples are available.

# Cross validation test and model selection

- Model selection
  - Which regularization method do we want to use?
  - How do we select the “optimal” tuning parameter  $\lambda$  ?
- Criteria
  - How well each model fits data?
  - Cross validation is one way to estimate the error
- L1 vs. L2
  - L1 (LASSO):  $\beta$  values have many zeros
  - L2 (Ridge):  $\beta$  values tend to be small but not exactly set to zero
- We can decide whether to use L1 or L2 based on the prediction error estimation through cross-validation tests
- We can also determine the tuning parameter  $\lambda$

# Outline

- High-dimensional regression problems
  - Challenges with high-dimensional data
  - Feature selection
- Model selection, Cross validation
  - Cross validation test
  - L1 vs L2 regularization
- R-session
  - Ridge regression using `lm.ridge`




# Input Data

- <http://homes.cs.washington.edu/~suinlee/genome560/data/mice.txt>
- Data on fluctuating proportions of marked cells in marrow from heterozygous Safari cats
- Proportions of cells of one cell type in samples from cats (taken in our department many years ago).  
Column 1 is the ID number of the particular cat. You will want to plot the data from one cat.
  - For example cat 40004 is rows 1:17, 40005a is 18:31, 40005b is 32:47, 40006 is 48:65, 40665 is 66:83 and so on.

# Input Data

- <http://homes.cs.washington.edu/~suinlee/genome560/data/mice.txt>
  - 1<sup>nd</sup> column: mouse ID
  - 2<sup>rd</sup> column: sex
  - 3<sup>th</sup> column: weight
  - 4<sup>th</sup> column: length
  - :

mice



F2	sex	weight_g		length_cm		Trigly	Total_Chol	FFA	Insulin_log
F2_1	2	42.8	9.6	14	1646	132	2.974971994		
F2_2	2	38	10.5	14	1646	132	2.974971994		
F2_3	2	33.5	10.8	109	1216	96	2.974971994		
F2_4	1	51.5	11.3	156	1382	147	3.968482949		
F2_5	1	59.2	11.5	261	1563	140	3.968482949		
F2_6	1	51.5	11	134	1823	161	3.968482949		
F2_7	1	56.5	11.2	192	1589	172	4.31255815		
F2_8	1	48.5	10.7	79	1265	85	4.31255815		
F2_9	1	52.7	11.3	96	1375	92	3.646893624		
F2_10	1	58.4	12	134	1567	153	4.347661709		



# Our Goals

- We want to learn a Ridge regression model for the insulin level (8<sup>th</sup> column) using other variables (sex, length, weight, etc.) as predictors.
- We will learn how to diagnose the pairwise correlations between variables.
- We will choose the optimal lambda based on CV errors.