

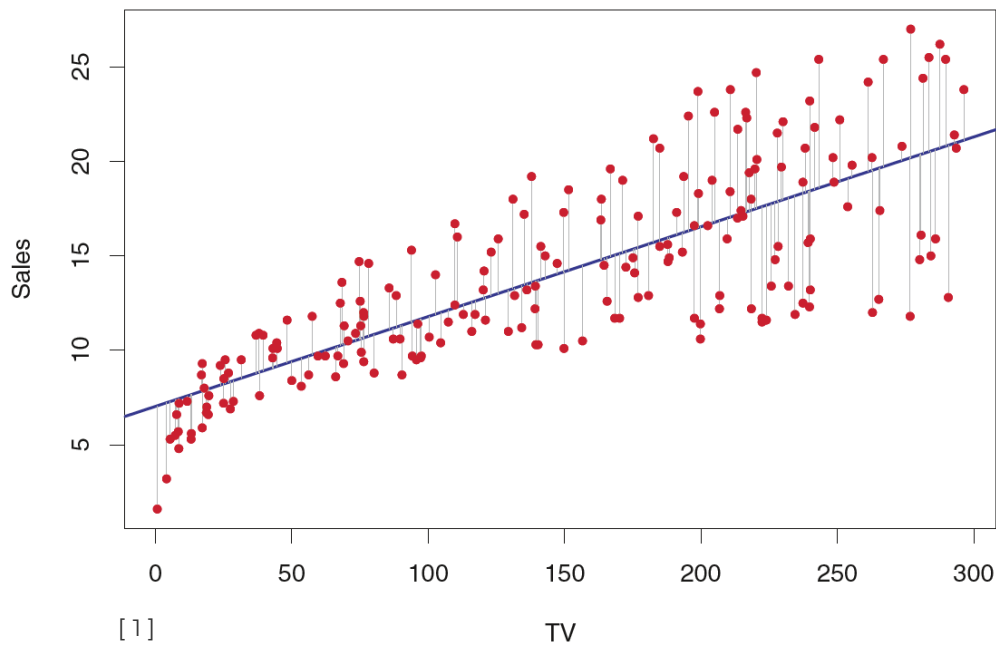
# Supervised Machine Learning

## Module 3



# Kahoot!

# Simple Linear Regression



- Quantitative predictions
- Single input variable
- $Y \approx \beta_0 + \beta_1 X$
- $\beta_0, \beta_1$ :
  - Constant and unknown
  - Coefficients/parameters
- $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$
- $\hat{\beta}_0, \hat{\beta}_1$ :
  - Calculated from training data
  - Reduce closeness

# Estimate Coefficients

## LEAST SQUARES

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

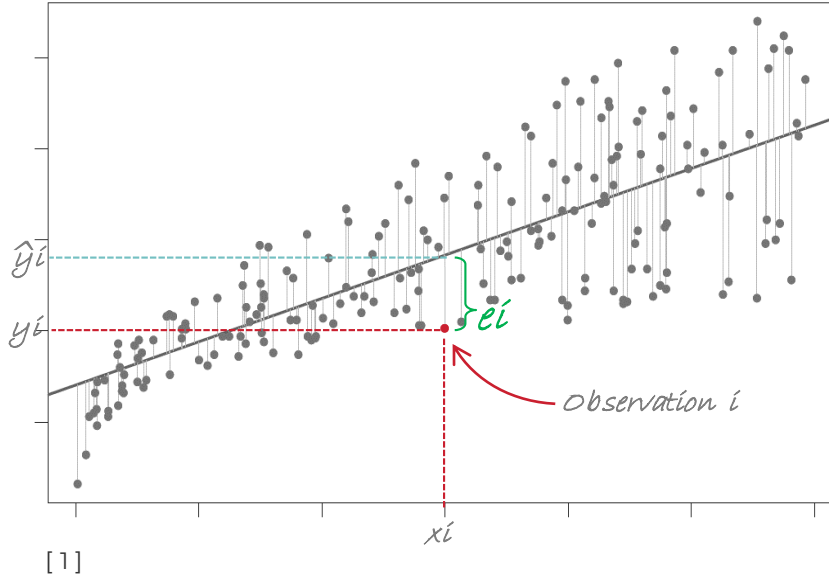
ith residual  $\triangleright e_i = y_i - \hat{y}_i$

### Residual Sum of Squares:

$$\begin{aligned} \text{RSS} &= e_1^2 + e_2^2 + \dots + e_n^2 \\ &= (y_1 - \hat{\beta}_0 + \hat{\beta}_1 x_1)^2 + \dots + (y_n - \hat{\beta}_0 + \hat{\beta}_1 x_n)^2 \end{aligned}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

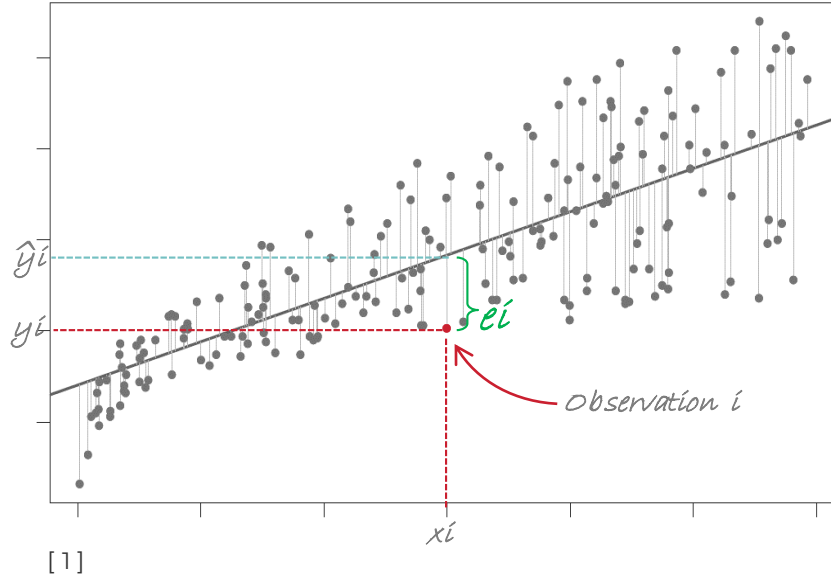
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$



# Estimate Coefficients



## LEAST SQUARES MATRIX APPROACH

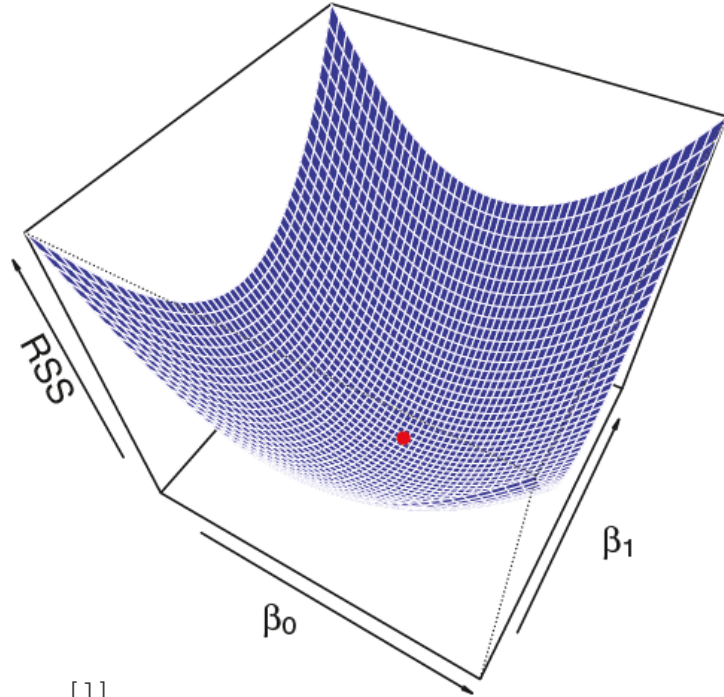


$$X\beta = y$$

$$X = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

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

# Estimate Coefficients

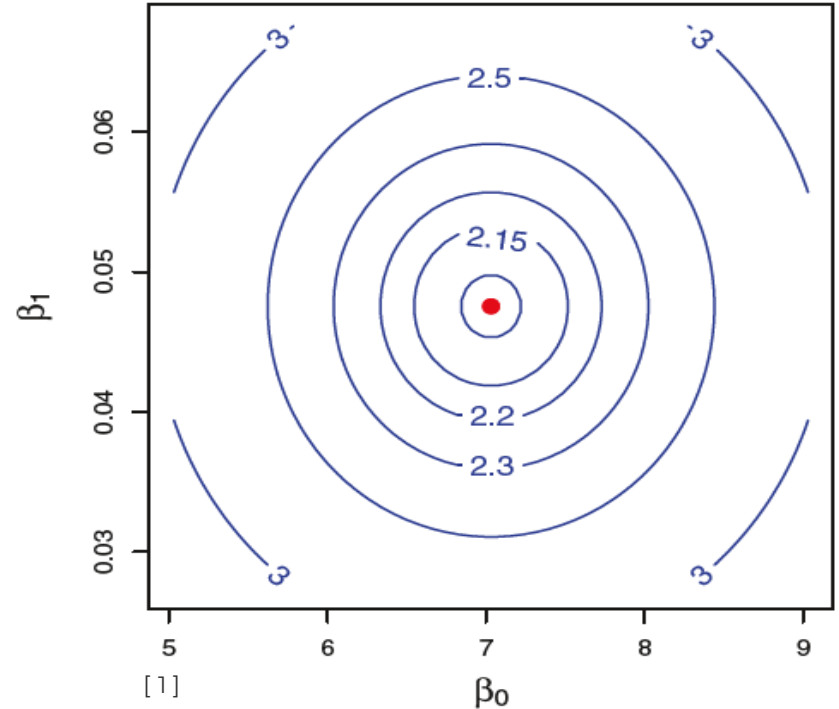


[1]

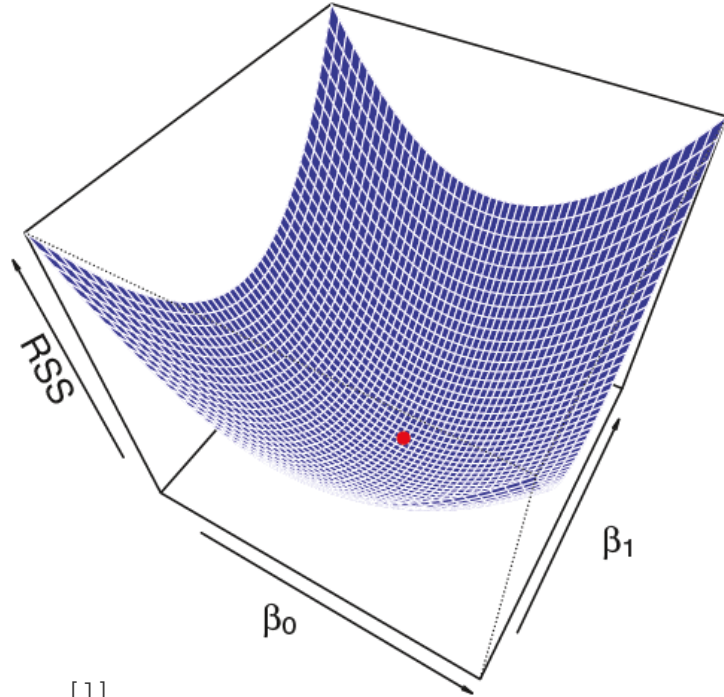
- This diagram shows how different values for each regression coefficient determine RSS value.
- We can see how there is a single solution for the global minimum of the loss function

# Estimate Coefficients

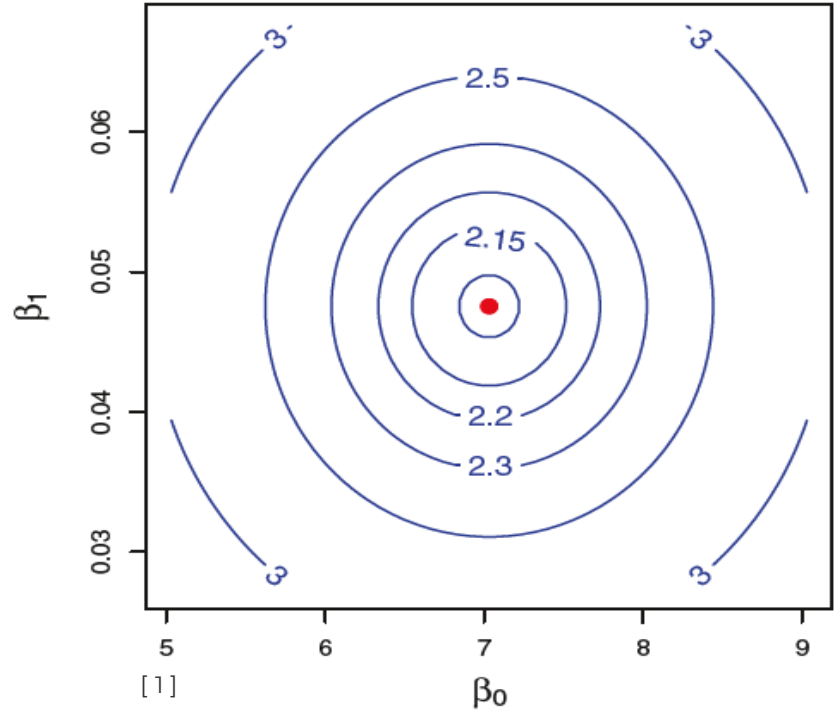
This diagram shows different ellipses where different values of  $\beta$  lead to same RSS value



# Estimate Coefficients



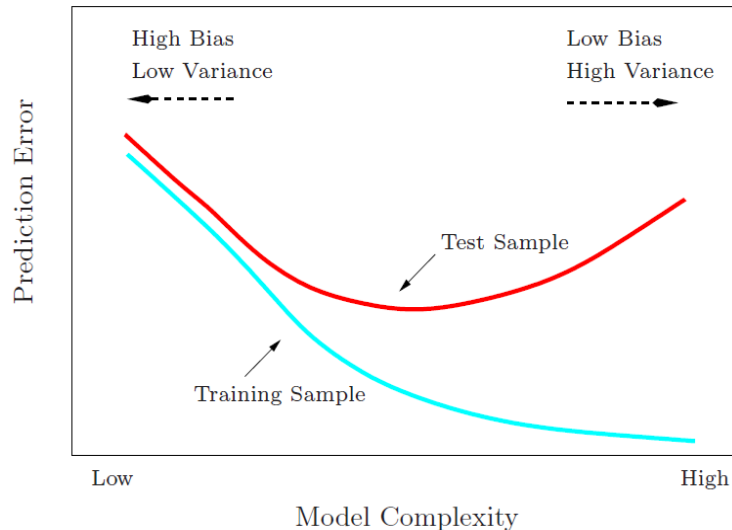
[1]





# Regularization

- Techniques for avoiding overfitting and increasing model interpretability
- Alternative to least squares
- Constrain/regualrize/shrink regression coefficients
- Discourage learning more complex models
- Ridge regression and Lasso Regression



# Ridge Regression

- Similar to OLS but coefficients are slightly minimized by penalty term:

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^p \beta_j^2}_{\text{penalty term}} = \text{RSS} + \underbrace{\lambda \sum_{j=1}^p \beta_j^2}_{\text{penalty term}}$$

- Tuning parameter  $\lambda$  -> increased for bigger shrinkage penalty
- Predictor standardization:

$$X_{\text{changed}} = \frac{X - \mu}{\sigma}$$

What would be  
coefficient value for  
 $\alpha = 0$ ?

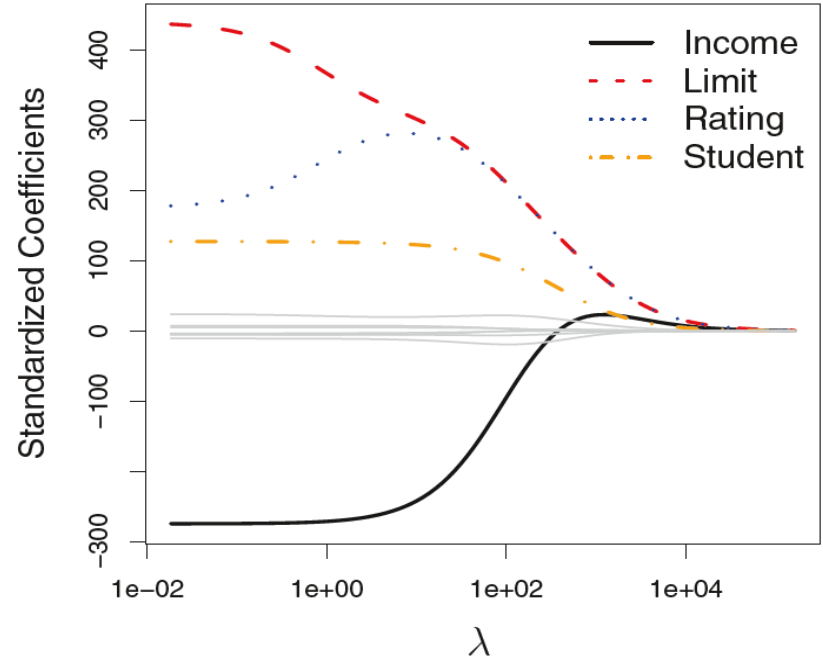


What would be  
coefficient value for  
 $\alpha = \infty$ ?



# Ridge Regression

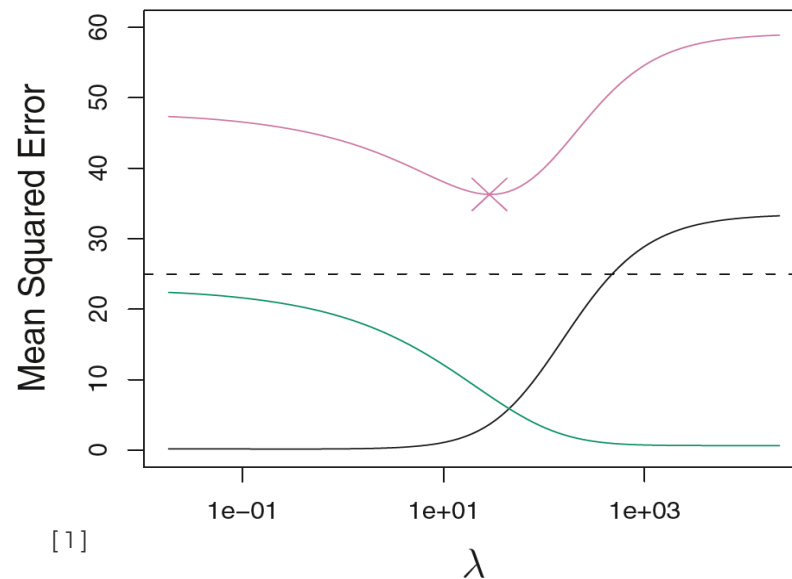
- We see how different coefficient values change when we increase alpha value
- As alpha increases, we expect coefficients to decrease
- We may observe sudden increases in some coefficients due to global minimization



[1]

# Ridge Regression

- We see how global MSE (pink), variance MSE (green) and bias MSE (black) vary when we try different Alpha values
- We expect (see graph)
- Our objective is to find the Alpha that corresponds to the minimum global MSE





# Lasso Regression

- Ridge regression reduces coefficients, but does not set them to zero -> all predictors are kept in the model
- Lasso coefficients minimize the quantity:

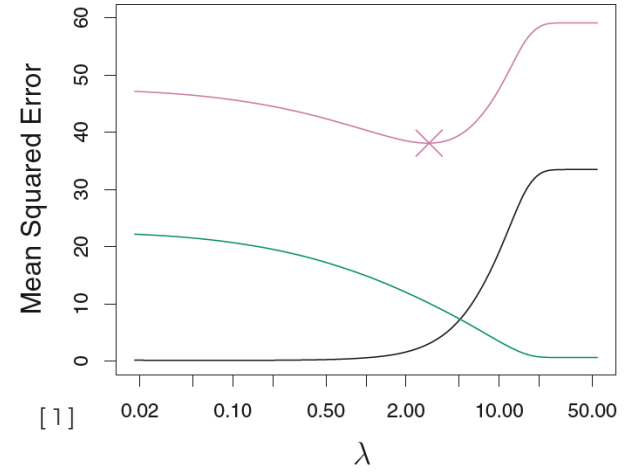
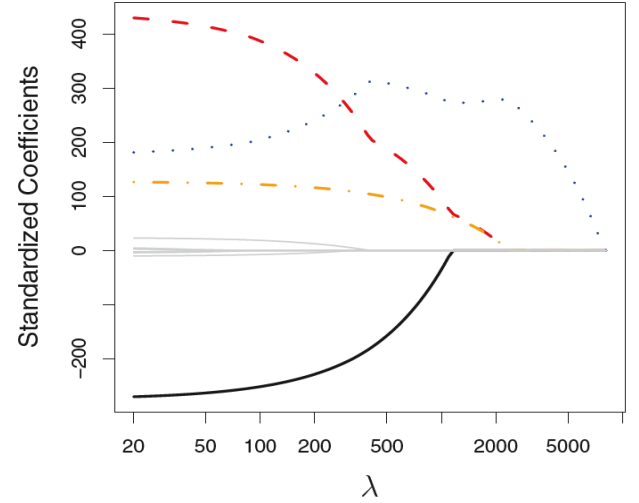
$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^p |\beta_j|}_{\text{penalty}} = \text{RSS} + \underbrace{\lambda \sum_{j=1}^p |\beta_j|}_{\text{penalty}}$$

- Lasso forces some coefficients to be exactly equal to zero -> variable selection
- Lasso models are easier to interpret



# Lasso Regression

- As can be now observed, some coefficients yield zero values as Alpha increases
- On the MSE side, we see again the effects of bias-variance tradeoff for different Alpha values





# Recall: Norms

$$\|x\|_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}}$$

- $L^1$  Norm

Useful when 0 and non-zero have to be distinguished

- $L^2$  Norm

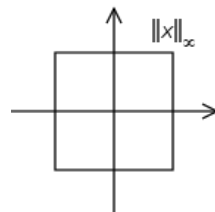
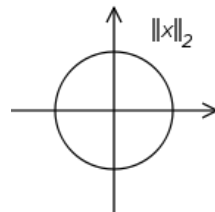
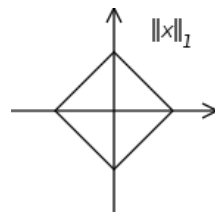
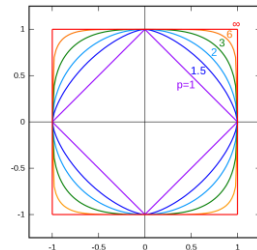
Called Euclidean norm:

- Simply the Euclidean distance between the origin and the point  $x$

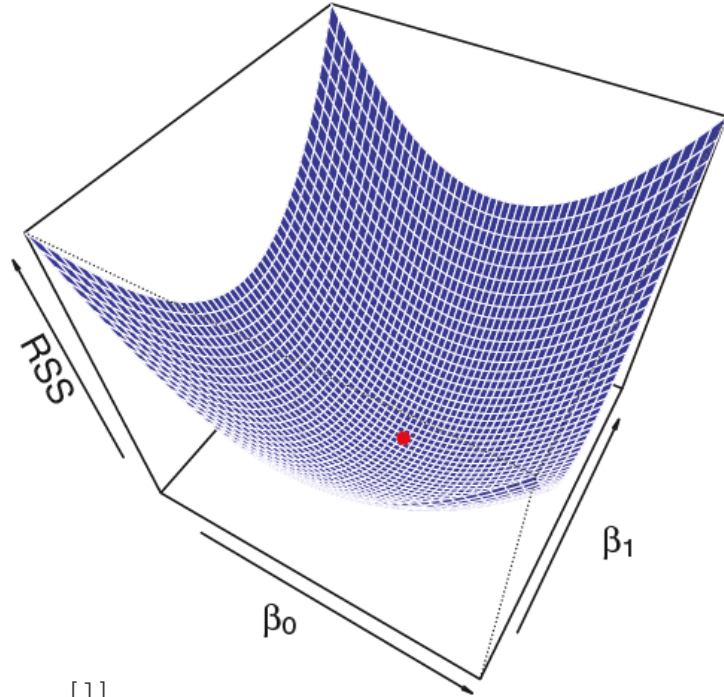
- $L^\infty$  Norm

$$\|x\|_\infty = \max_i |x_i|$$

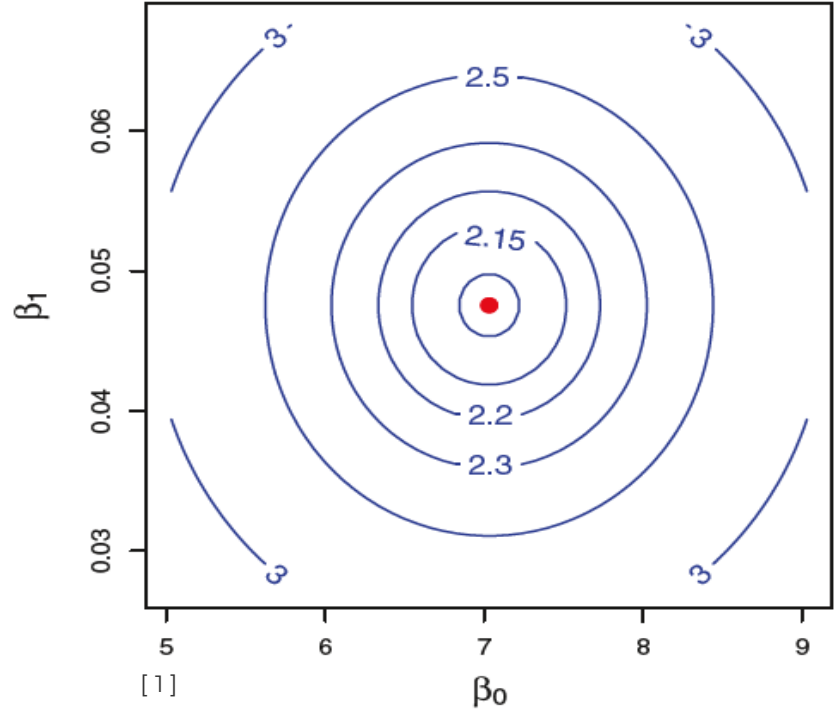
Called max norm



# Recall: Estimate Coefficients



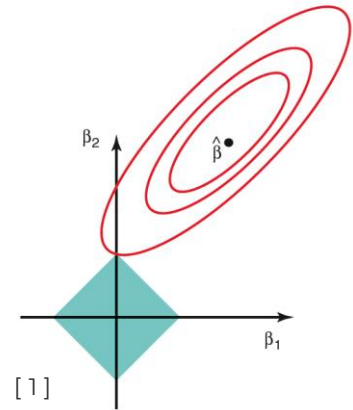
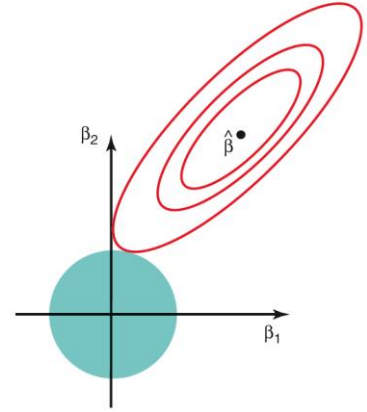
[1]



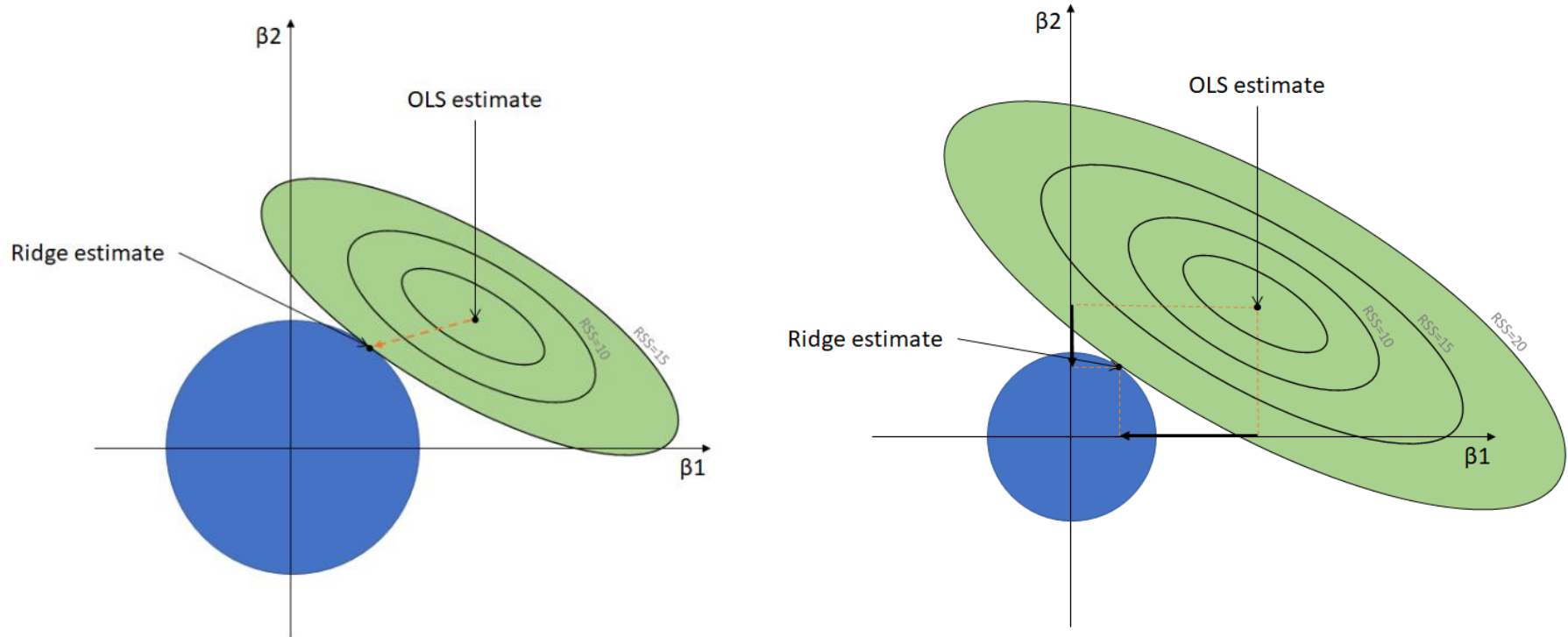
[1]

# Geometric understanding of regularization

- Different regularization terms lead to different regression solutions
- Ridge regression adds a **L2 regularization term** (circle on the right). This shape make it tough to yield values of beta equal to zero
- Lasso adds a **L1 regularization term** (diamond on the right). This sharp shape makes it easier to select beta values equal to zero



# Geometric understanding of regularization



# Ridge vs. Lasso

- Lasso has a major advantage on interpretability due to feature selection
- Regarding prediction accuracy, neither ridge regression nor the lasso will universally dominate the other.
- In general:
  - **Lasso:** small number of predictors have substantial coefficients, remaining predictors have coefficients that are very small or that equal zero.
  - **Ridge:** response is a function of many predictors, all with coefficients of roughly equal size.

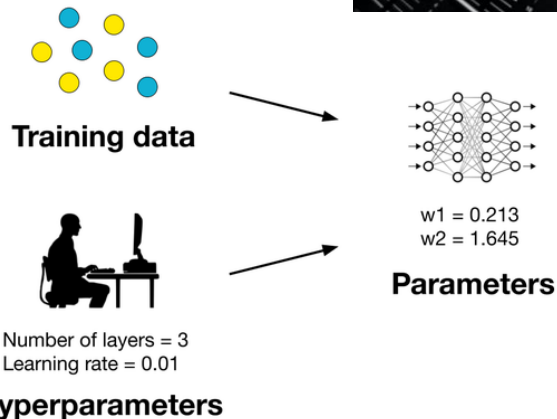
# How would you choose?





# Parameters and hyperparameters

- Hyperparameters: parameter whose value is set before the learning process begins.
- Parameters: its value is derived via training.





$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = \text{RSS} + \lambda \sum_{j=1}^p |\beta_j|$$

# References

- [1] G. James, D. Witten, T. Hastie, R. Tibshirani. An Introduction to Statistical Learning with Applications in R. Springer, 2017.
- [2] T. Hastie, R. Tibshirani, J. Friedman. The Elements of Statistical: Data Mining, Inference and Prediction. Springer, 2009.



IMMUNE  
CODING INSTITUTE