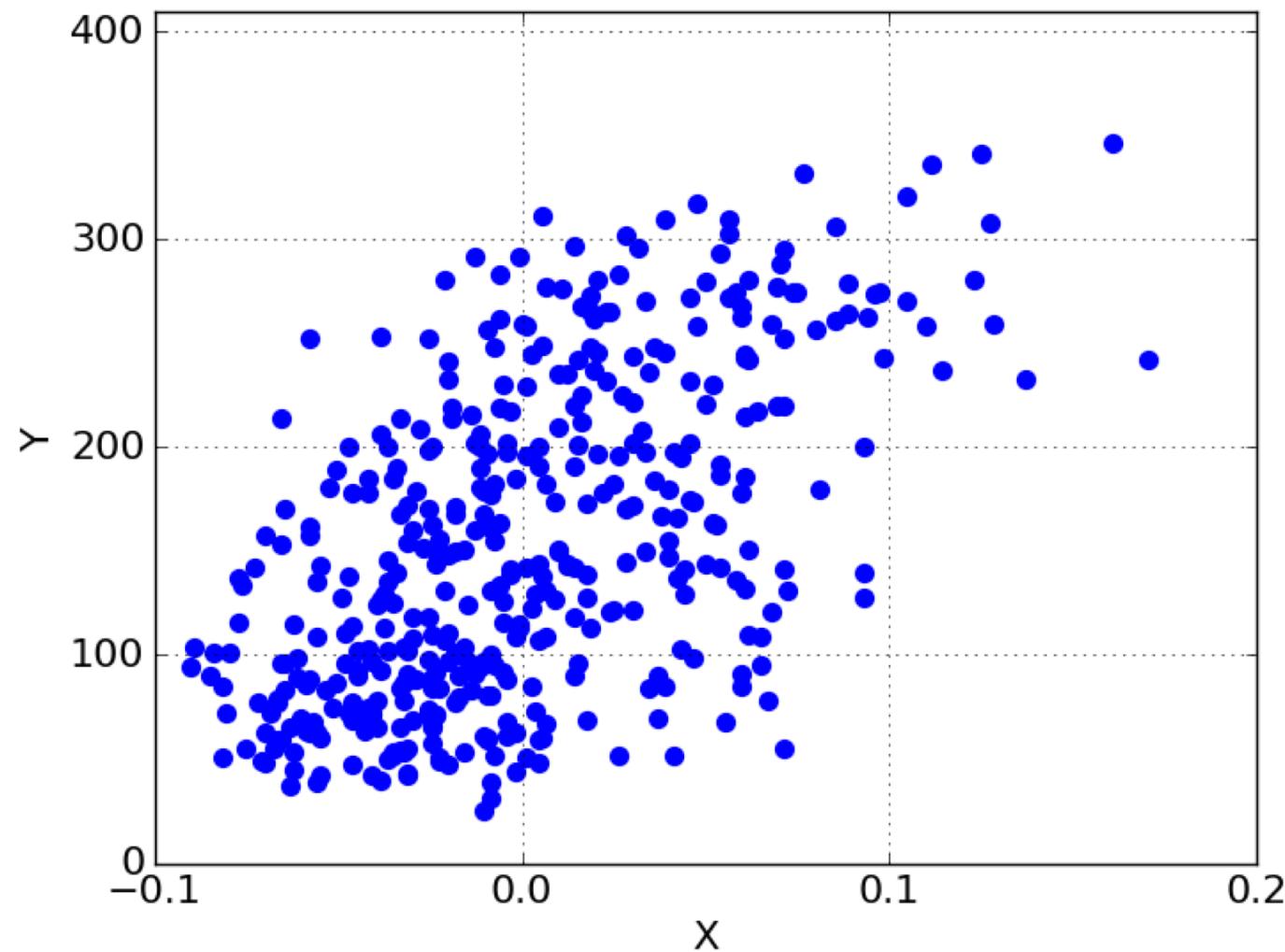




# Regression

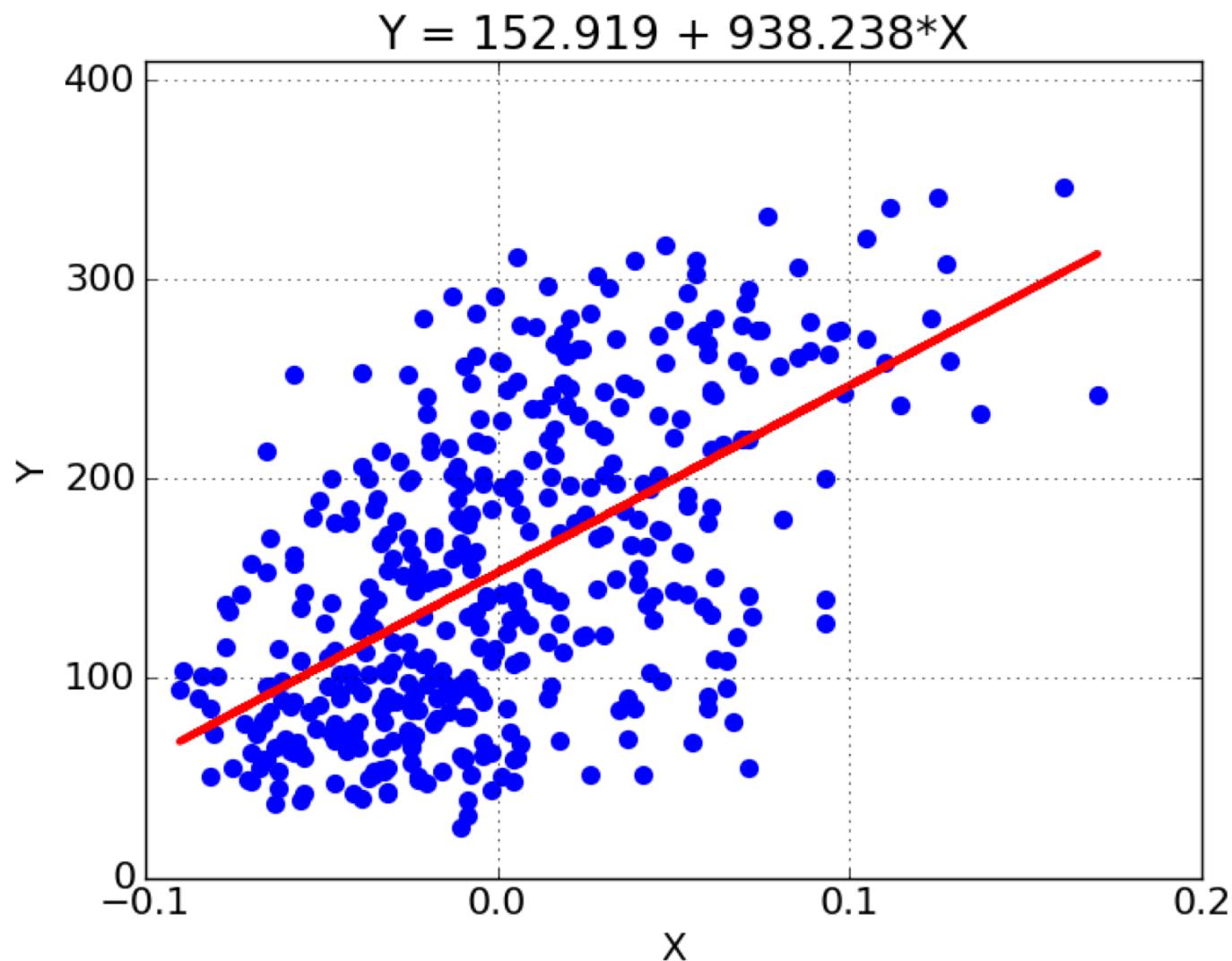
Data Mining and Text Mining (UIC 583 @ Politecnico di Milano)

# Simple Linear Regression



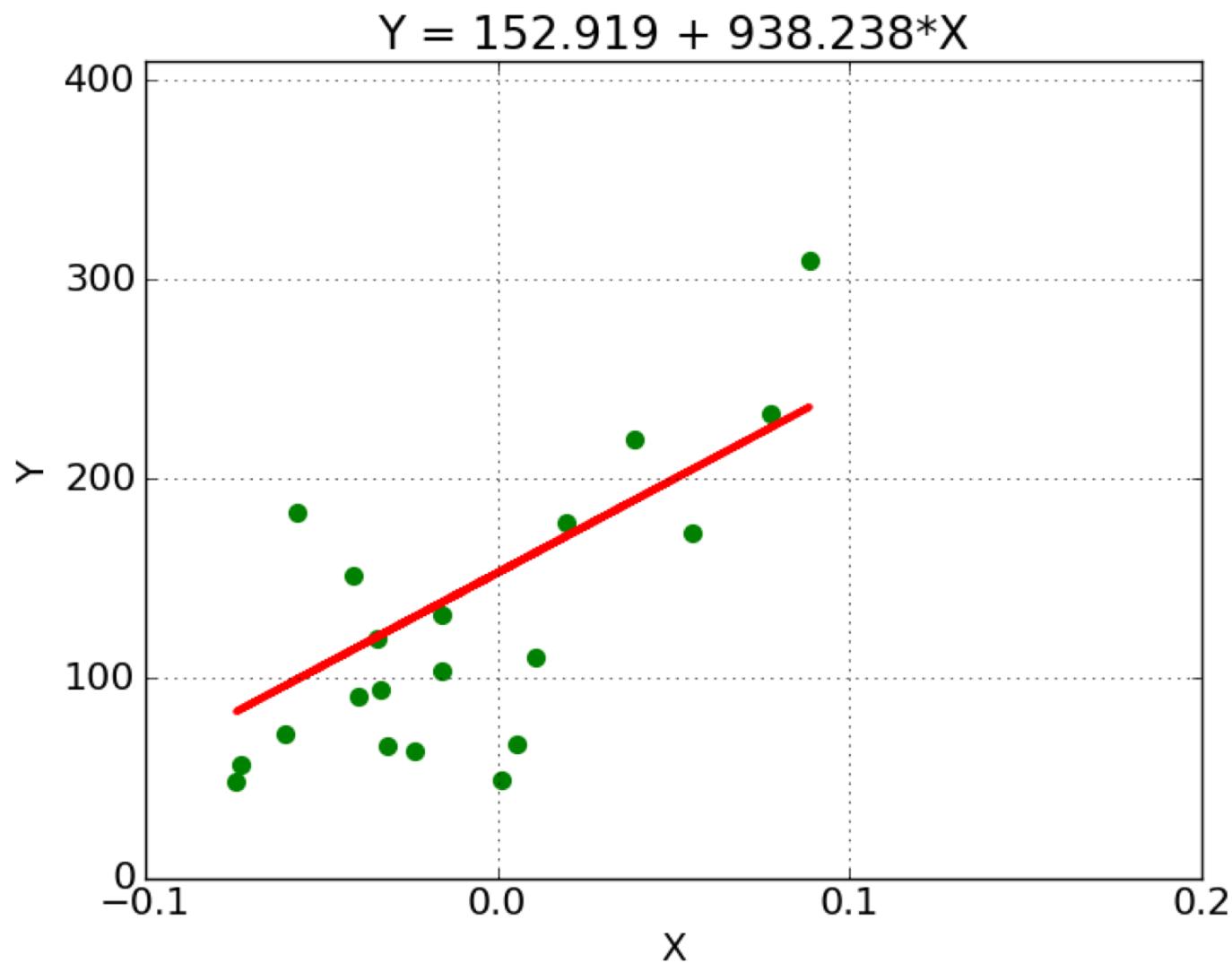
The training data points

Can we predict the value of  $y$  from  $x$ ?



Training data points with a simple linear regression model

But the model we built will not be used not on the same data ...



The previous model applied to the model learned from the training data

regression = model building + model usage

## How Do We Evaluate a Regression Model?

- Given N examples, pairs  $x_i, y_i$ , linear regression computes a model

$$y = w_0 + w_1 x$$

- So that for each point,

$$y_i = w_0 + w_1 x_i + \varepsilon_i$$

- We evaluate the model by computing the Residual Sum of Squares (RSS) computed as,

$$RSS(w_0, w_1) = \sum_{i=1}^N (y_i - (w_0 + w_1 x_i))^2$$

The goal of linear regression is thus to find the weights that minimize RSS

$$w_0, w_1 = \arg \min_{w_0, w_1} RSS(w_0, w_1)$$

$$= \arg \min_{w_0, w_1} \sum_{i=1}^N (y_i - (w_0 + w_1 x_i))^2$$

# How Do We Compute the Best Weights?

- Approach 1
  - Set the gradient of  $RSS(w_0, w_1)$  to zero; but the approach is infeasible in practice
- Approach 2
  - Apply gradient descent

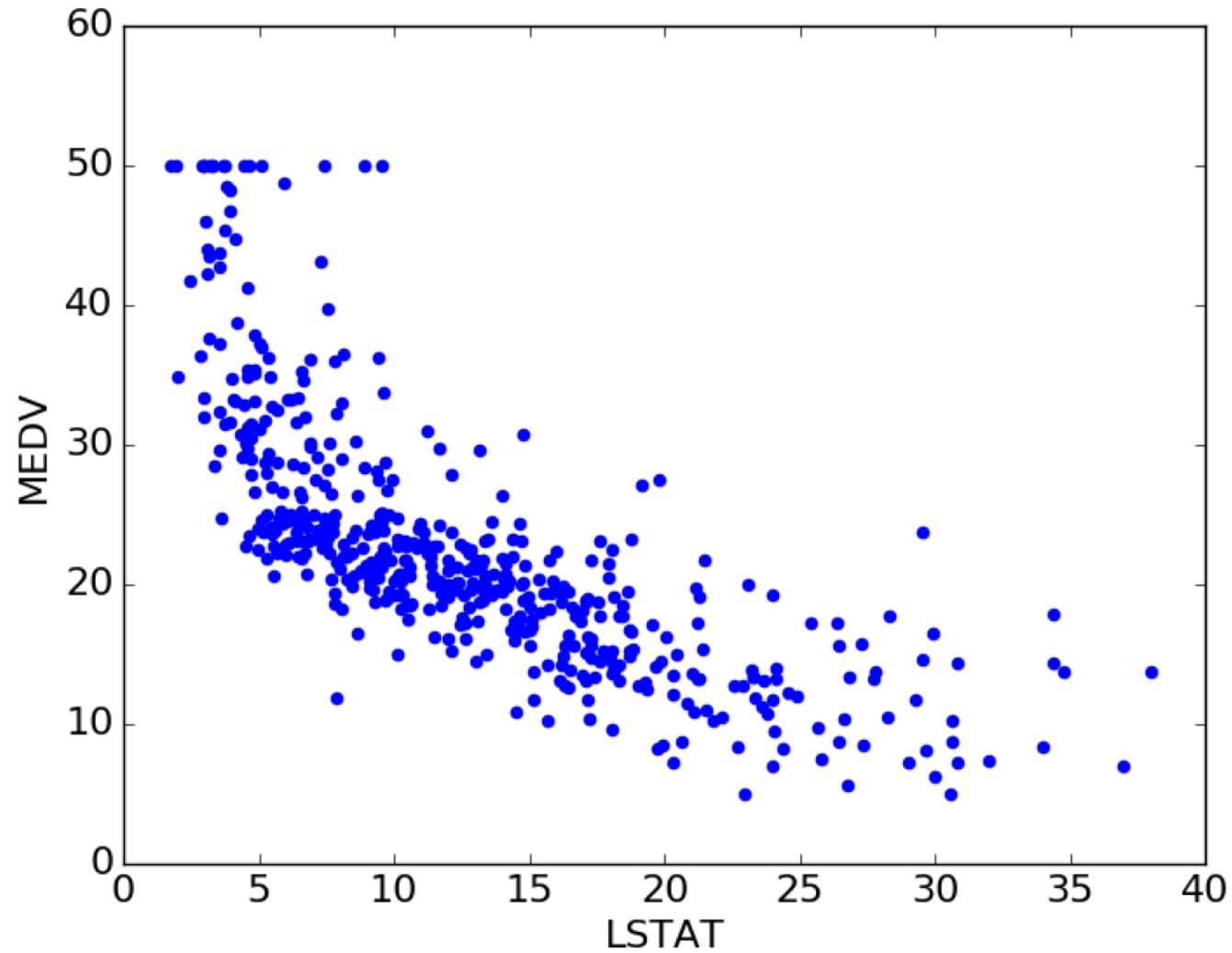
while not converged

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \nabla RSS(\vec{w}^{(t)})$$
- If  $\eta$  is large, we are making large steps but might not converge, if  $\eta$  is small, we might be very slow. Typically,  $\eta$  adapts over time, e.g.,  $\eta(t) = \alpha/t$  or  $\alpha/\sqrt{t}$

- For simple linear regression the gradient of RSS has only two components one for  $w_0$  and one for  $w_1$

$$\nabla RSS(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N (y_i - (w_0 + w_1 x_i)) \\ -2 \sum_{i=1}^N (y_i - (w_0 + w_1 x_i)) x_i \end{bmatrix}$$

# Multiple Linear Regression



Input variable: LSTAT - % lower status of the population

Output variable: MEDV - Median value of owner-occupied homes in \$1000's

Can we predict the property value  
using other variables?

- Given a set of examples associating LSTAT<sub>i</sub> values to MEDV<sub>i</sub> values, simple linear regression find a function f(.) such that

$$MEDV_i = f(LSTAT_i) + \epsilon_i$$

- Where  $\epsilon_i$  is the error to be minimized
- Typically, we assume a model and fit the model into the data
- With linear model we assume that f(.) is computed as

$$f(LSTAT_i) = w_0 + w_1 \times LSTAT_i$$

- A polynomial model would fit the data points with a function,

$$f(LSTAT_i) = w_0 + \sum_{j=1}^D w_j \times LSTAT_i^j$$

- Given, D input variables, assumes that the output  $y$  can be computed as,

$$y = w_0 + \sum_{j=1}^D w_j x_j + \varepsilon$$

- The model cost is computed using the residual sum of square,  $RSS(w)$  as,

$$RSS(\vec{w}) = \sum_{i=1}^N \left( y_i - w_0 - \sum_{j=1}^D w_j x_{i,j} \right)^2$$

- Total sum of squares

$$TSS = \sum_{i=1}^N (y_i - \bar{y})^2$$

- Coefficient of determination

$$R^2 = 1 - \frac{RSS}{TSS}$$

- R<sup>2</sup> measures of how well the regression line approximates the real data points. When R<sup>2</sup> is 1, the regression line perfectly fits the data.

- In general, given a set of input variables  $x$ , a set of  $N$  examples  $x_i$ ,  $y_i$  and a set of  $D$  features  $h_j$  computed from the input variables  $x_i$ , multiple linear regression assumes a model,

$$y_i = \sum_{j=0}^D w_j h_j(\vec{x}_i) + \varepsilon_i$$

- $h_j(\cdot)$  identify variables derived from the original inputs
- $h_j(\cdot)$  could be the squared value of an existing variable, a trigonometric function, the age given the date of birth, etc.

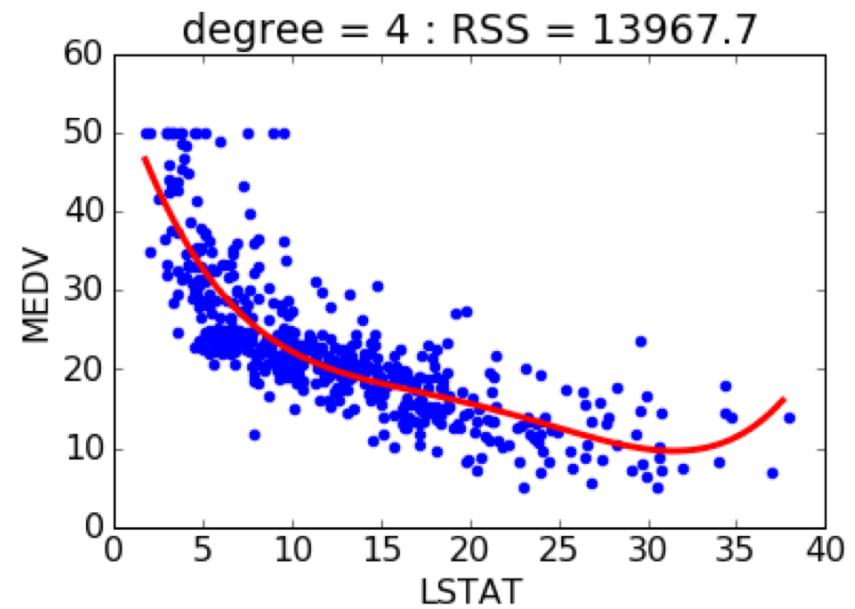
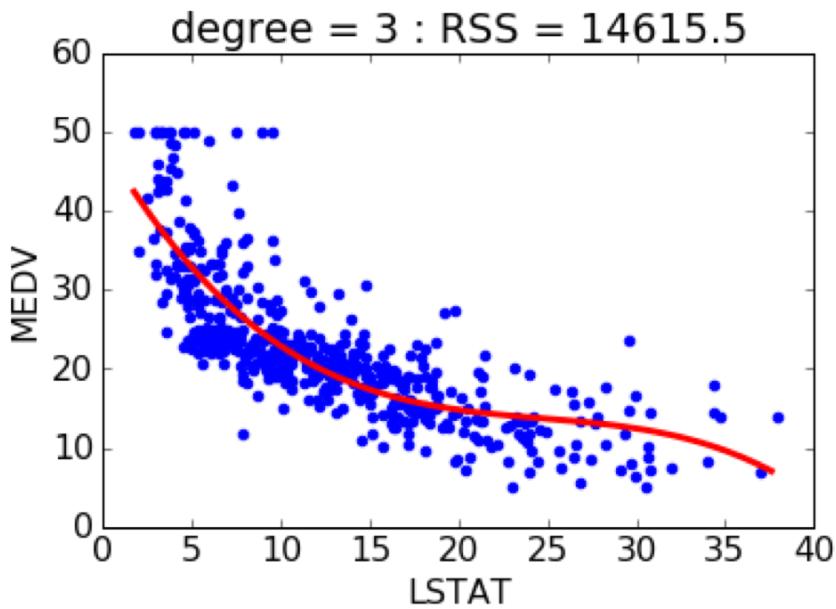
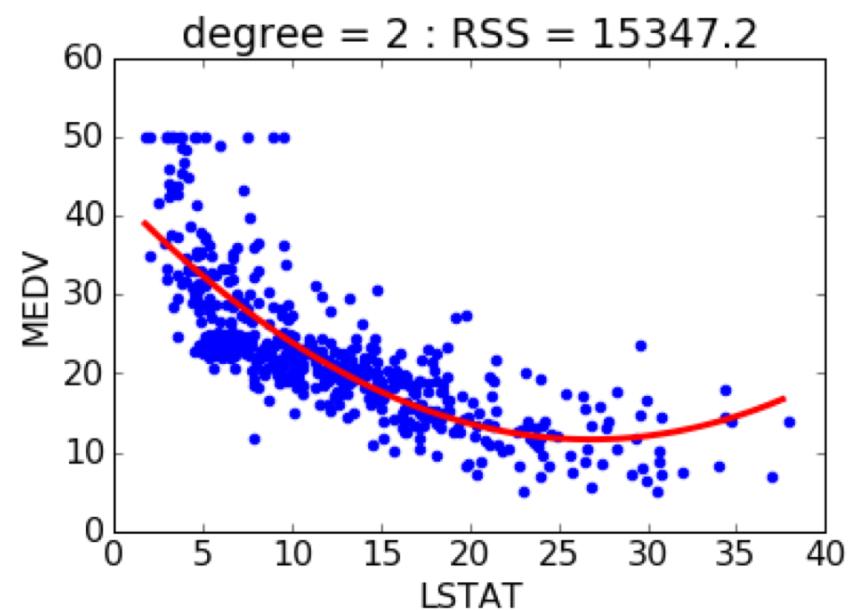
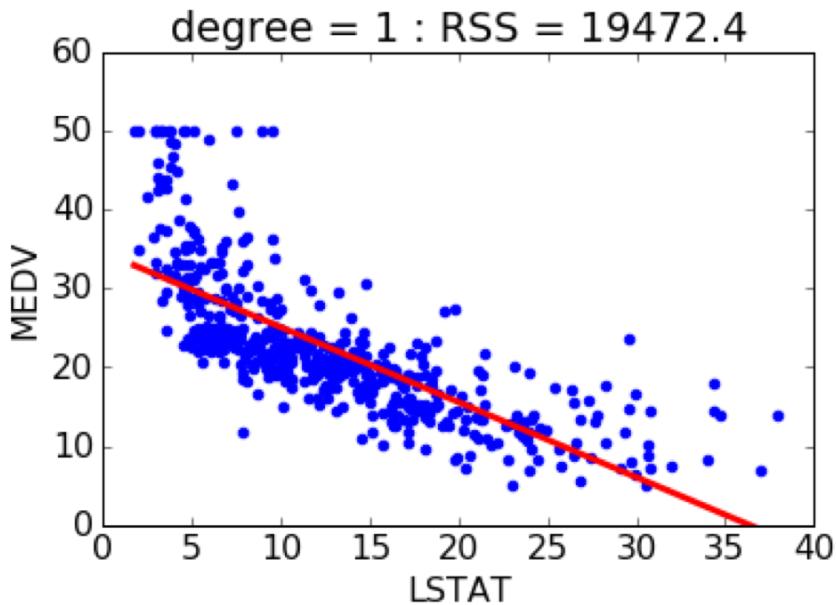
## Multiple Linear Regression: General Formulation

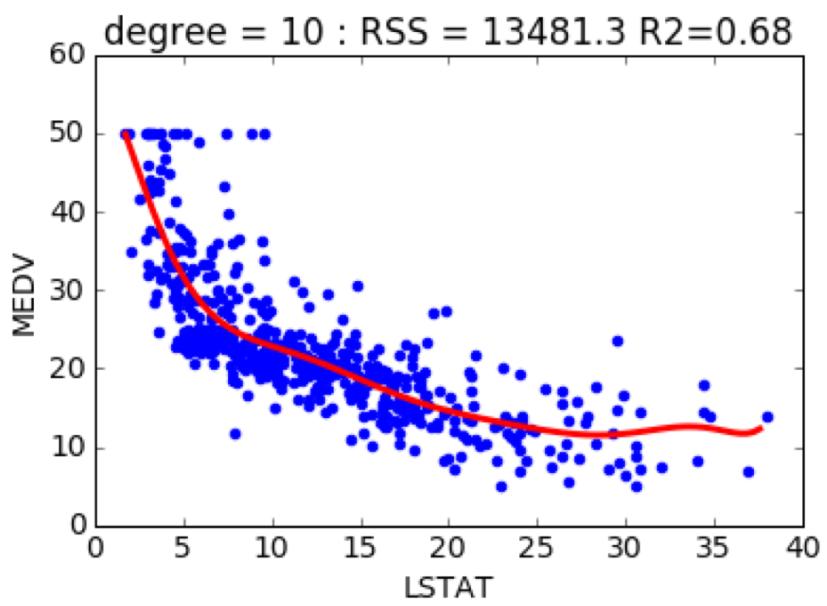
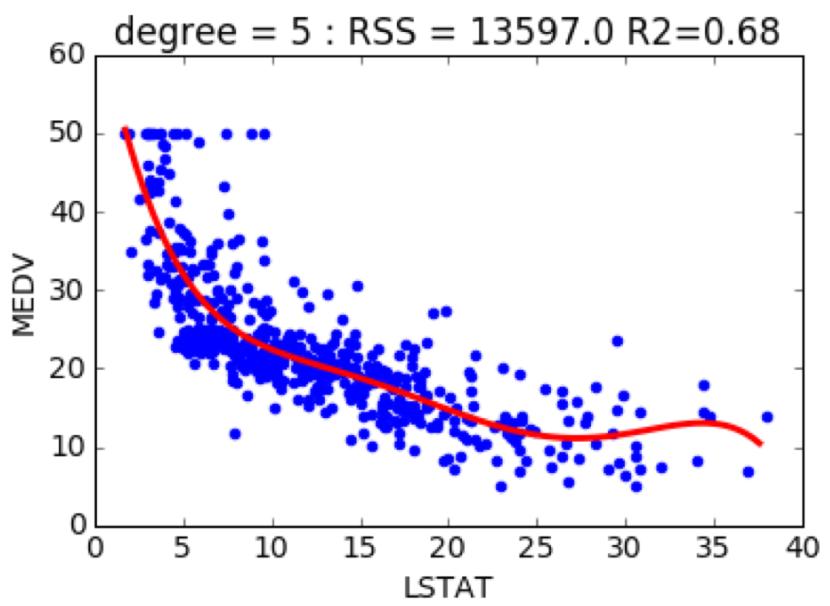
- Multiple linear regression aims at minimizing,

$$RSS(\vec{w}) = \sum_{I=1}^N \left( y_i - \sum_{j=0}^D w_j h_j(\vec{x}_i) \right)^2$$

- For this purpose, it can apply gradient descent to update the weights as,

$$w_j^{(t+1)} = w_j^{(t)} + 2\eta \sum_{i=1}^N h_j(\vec{x}_i)(y_i - \hat{y}_i(\vec{w}^{(t)}))$$





$t = 0$

init  $w^{(0)}$

while not converged

for ( $j = 0; j \leq D; j = j + 1$ )

$$\Delta w_j = -2 \sum_{i=1}^N h_j(\vec{x}_i)(y_i - \hat{y}_i(\vec{w}^{(t)}))$$

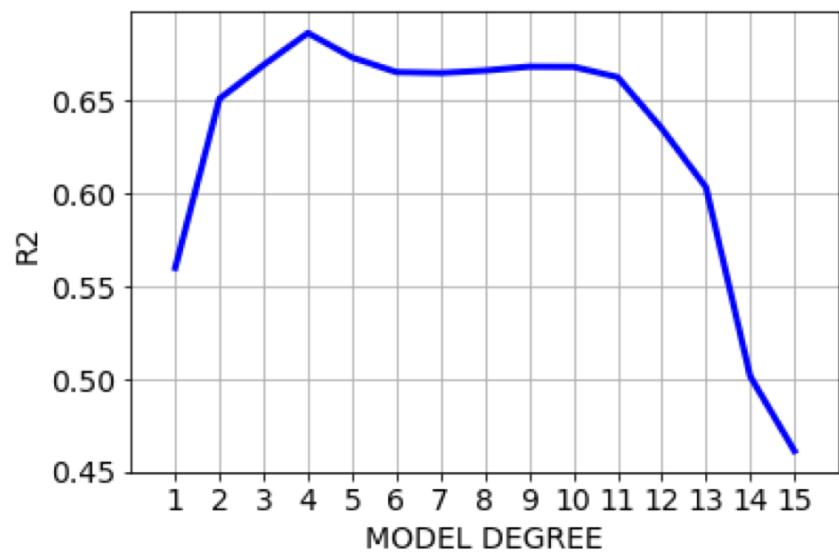
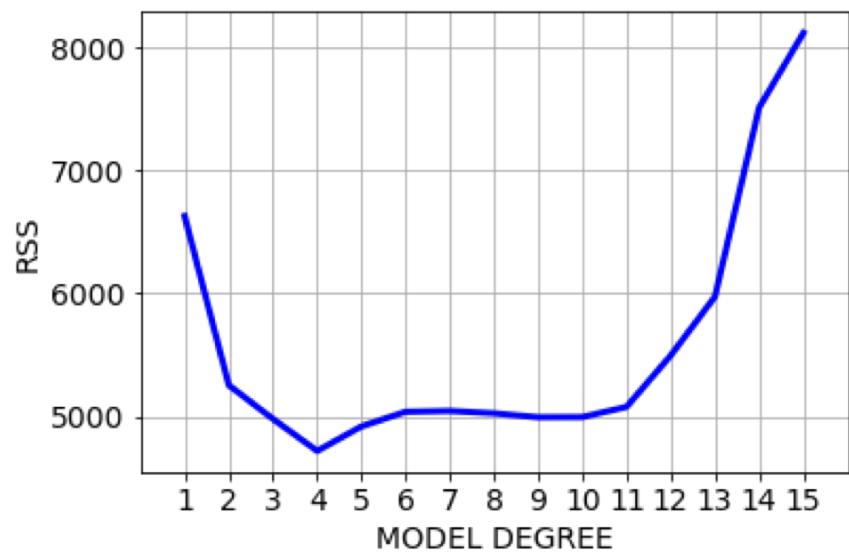
$$w_j^{(t+1)} = w_j^{(t)} - \eta \Delta w_j$$

$t = t + 1$

# Model Evaluation

- Models should be evaluated using data that have not been used to build the model itself
- For example, would be feasible to evaluate students using exactly the same problems solved in class?
- The available data must be split between training and test
- Training data will be used to build the model
- Test data will be used to evaluate the model performance

- Reserves a certain amount for testing and uses the remainder for training
  - Too small training sets might result in poor weight estimation
  - Too small test sets might result in a poor estimation of future performance
- Typically,
  - Reserve  $\frac{1}{2}$  for training and  $\frac{1}{2}$  for testing
  - Reserve  $\frac{2}{3}$  for training and  $\frac{1}{3}$  for testing
- For small or “unbalanced” datasets, samples might not be representative



Given the original dataset, we split the data into 2/3 train and 1/3 test and then apply multiple linear regression using polynomials of increasing degree. The plot show how RSS and  $R^2$  vary.

## Holdout Evaluation using the Housing Data

- Given the original dataset, we split the data into 2/3 train and 1/3 test and then apply multiple linear regression using polynomials of increasing degree.
- RSS initially decreases as polynomials better approximate the data but then higher degree polynomials overfit
- The same is shown by the R<sup>2</sup> statistics



- First step
  - Data is split into  $k$  subsets of equal size
- Second step
  - Each subset in turn is used for testing and the remainder for training
- This is called  $k$ -fold cross-validation and avoids overlapping test sets
- Often the subsets are stratified before cross-validation is performed
- The error estimates are averaged to yield an overall error estimate

## Ten-fold Crossvalidation

30



$p_1$

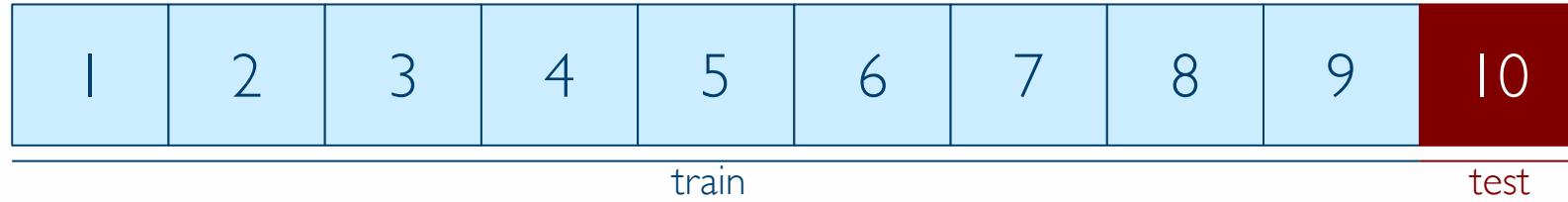


$p_2$

...

...

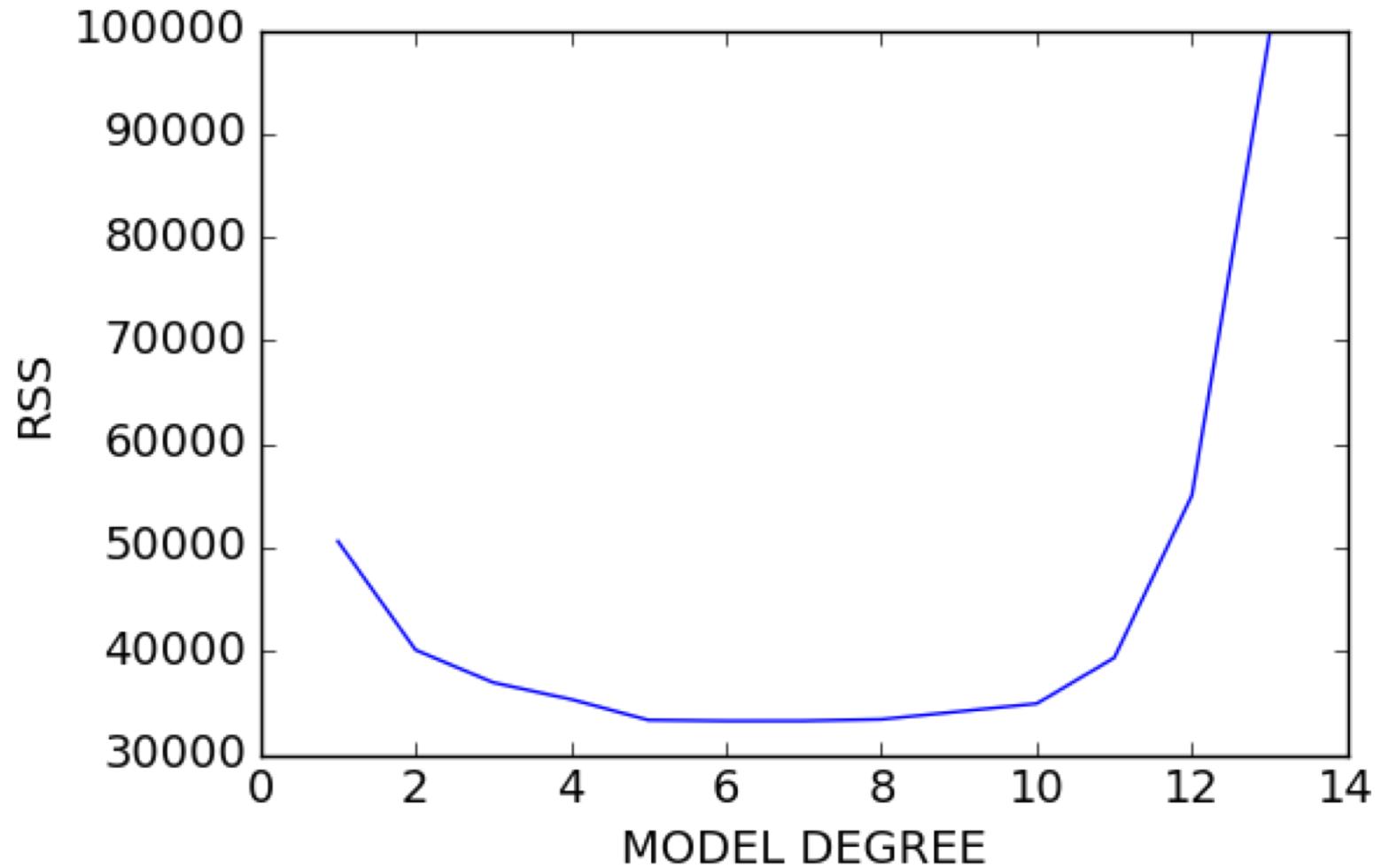
...



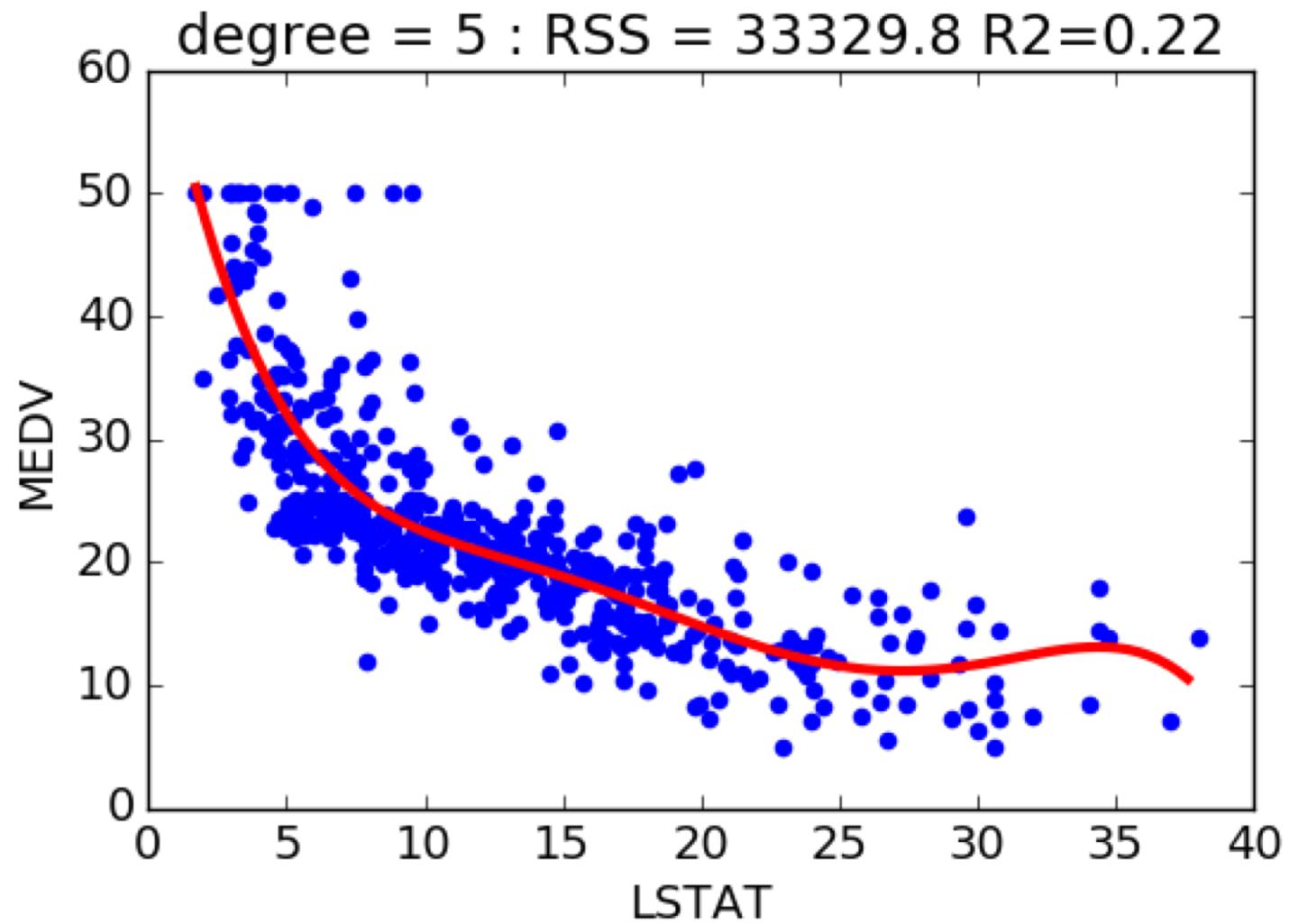
$p_{10}$

The final performance is computed as the average  $p_i$

- Standard method for evaluation stratified ten-fold cross-validation
- Why ten? Extensive experiments have shown that this is the best choice to get an accurate estimate
- Stratification reduces the estimate's variance
- Even better: repeated stratified cross-validation
- E.g. ten-fold cross-validation is repeated ten times and results are averaged (reduces the variance)
- Other approaches appear to be robust, e.g., 5x2 crossvalidation



As we increase the degree of the fitting polynomial, the crossvalidation error starts to increase, because the model starts to overfit! Best performance for the 5<sup>th</sup> degree polynomial.



Fitting using the 5<sup>th</sup> degree polynomial

# Overfitting

# What is Overfitting?

Very good performance  
on the training set

Terrible performance  
on the test set

In regression, overfitting is often associated  
to large weights estimates

Add to the usual cost (RSS) a term  
to penalize large weights to avoid overfitting

Total cost = Measure of Fit + Magnitude of Coefficients

- Minimizes the cost function,

$$\begin{aligned} Cost(\vec{w}) &= RSS(\vec{w}) + \alpha \|\vec{w}\|_2^2 \\ &= \sum_{i=1}^N \left( y_i - \sum_{j=0}^D w_j h_j(x_i) \right)^2 + \alpha \sum_{j=0}^D w_j^2 \end{aligned}$$

- If  $\alpha$  is zero, the cost is exactly the same as before; if  $\alpha$  is infinite, then the only solution corresponds to having all the weights to 0
- In the gradient descent algorithm the update for weight  $j$  becomes,

$$\Delta w_j = -2 \sum_{i=1}^N h_j(\vec{x}_i)(y_i - \hat{y}_i(\vec{w}^{(t)})) + 2\alpha w_j$$

- Minimizes the cost function,

$$Cost(\vec{w}) = RSS(\vec{w}) + \alpha \|\vec{w}\|_1$$

$$= \sum_{i=1}^N \left( y_i - \sum_{j=0}^D w_j h_j(x_i) \right)^2 + \alpha \sum_{j=0}^D |w_j|$$

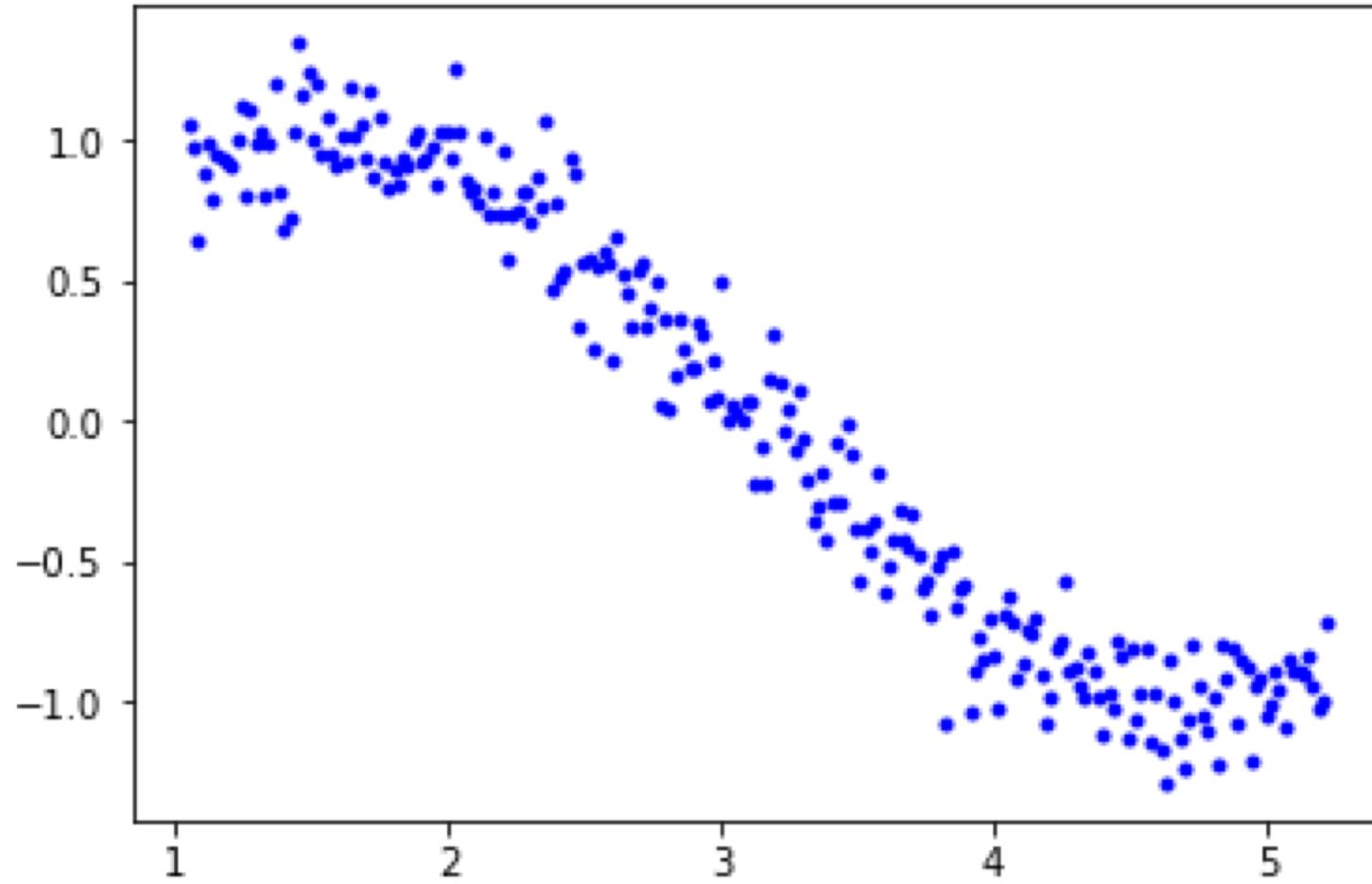
- If  $\alpha$  is zero, the cost is exactly the same as before; if  $\alpha$  is infinite, then the only solution corresponds to having all the weights to 0
- In gradient descent, weight  $j$  is modified as,

$$z_j = \sum_{i=1}^N h_j(\vec{x}_i)^2$$

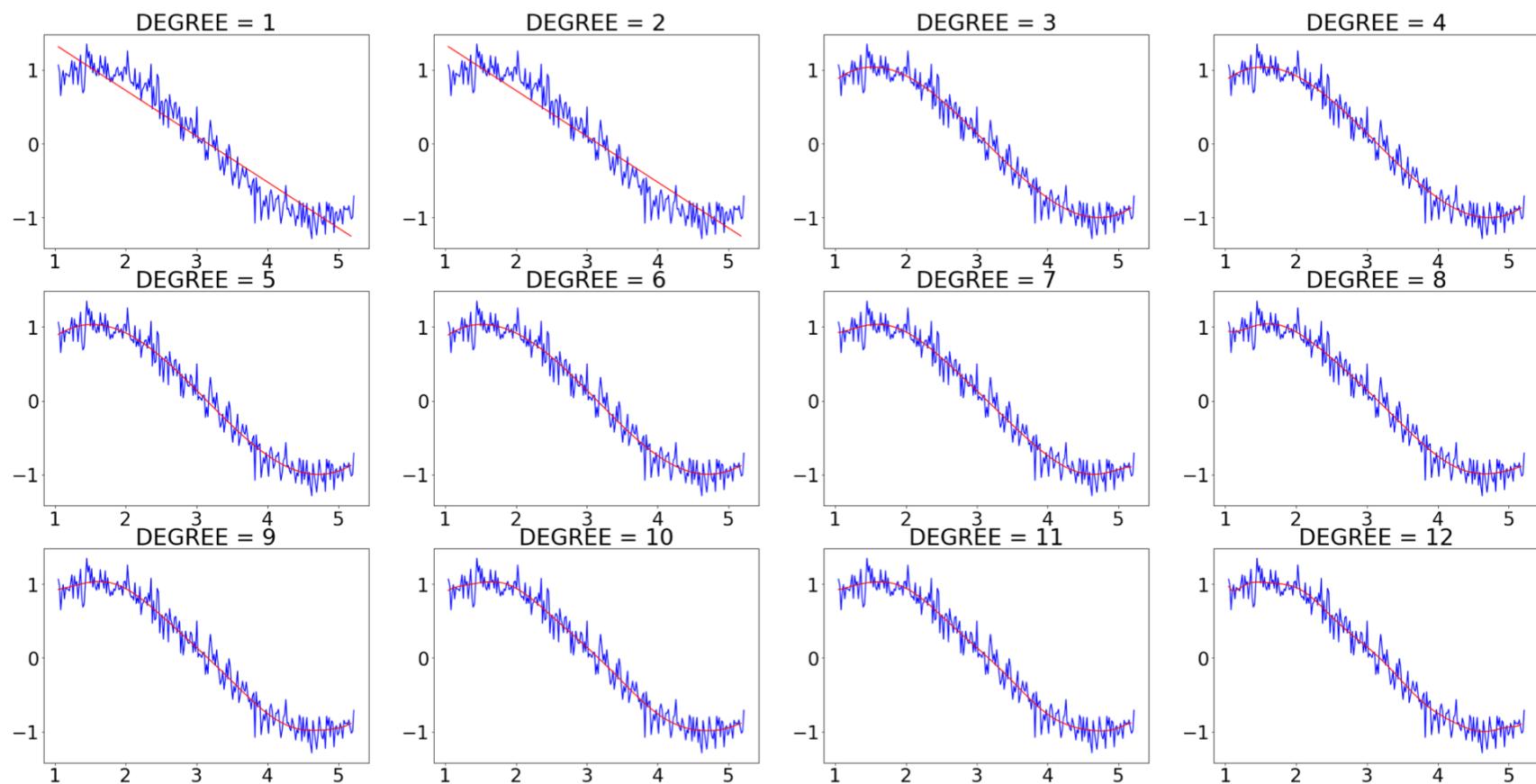
$$\rho_j = \sum_{i=1}^N h_j(\vec{x}_i)(y_i - \hat{y}_i(\vec{w}_{-j}^{(t)}))$$

$$w_j = \begin{cases} (\rho_j + \alpha/2)/z_j & \text{if } \rho_j < -\alpha/2 \\ 0 & \text{if } \rho_j \in [-\alpha/2, \alpha/2] \\ (\rho_j - \alpha/2)/z_j & \text{if } \rho_j > \alpha/2 \end{cases}$$

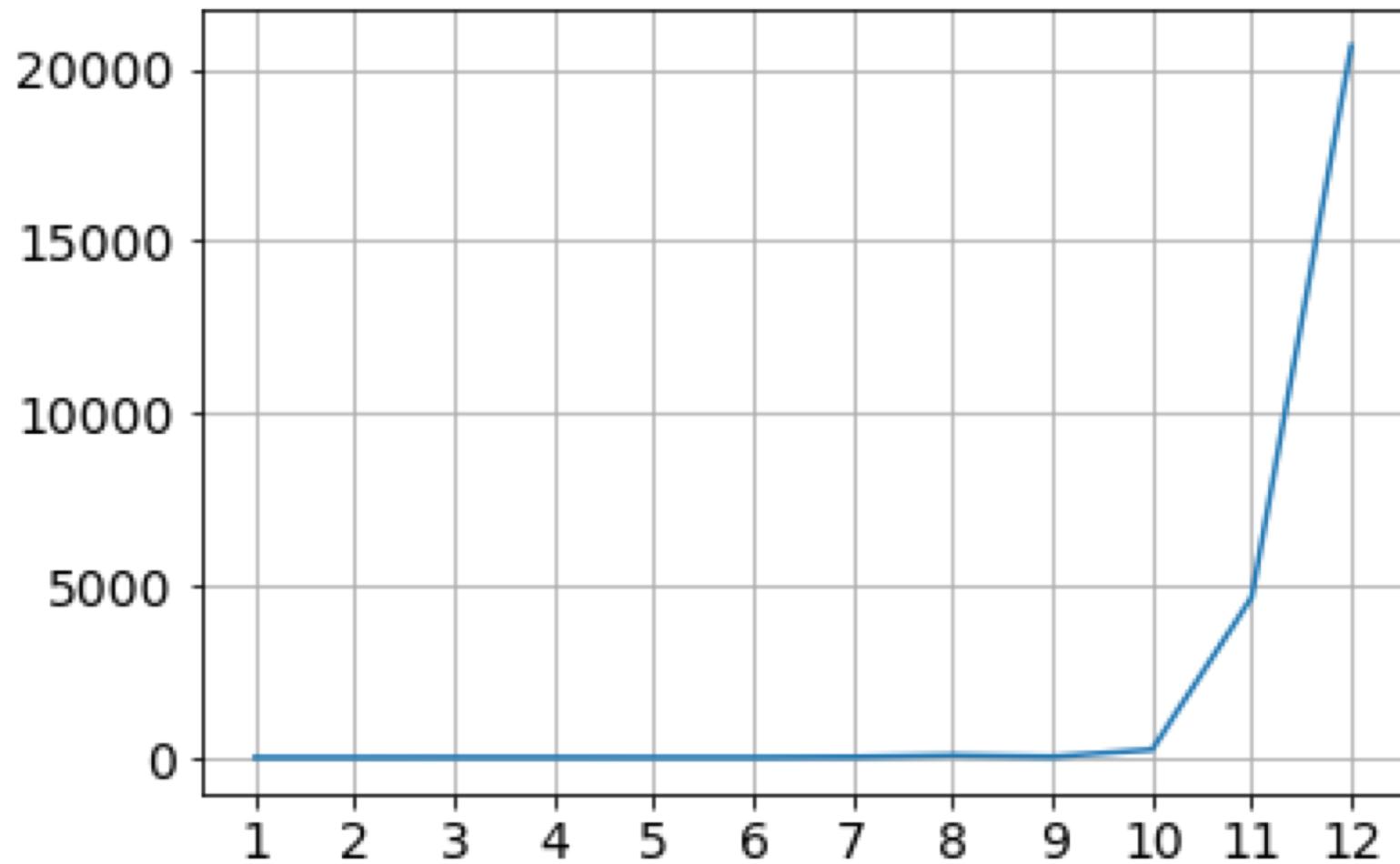
# Example



Simple example data generated using a trigonometric function.

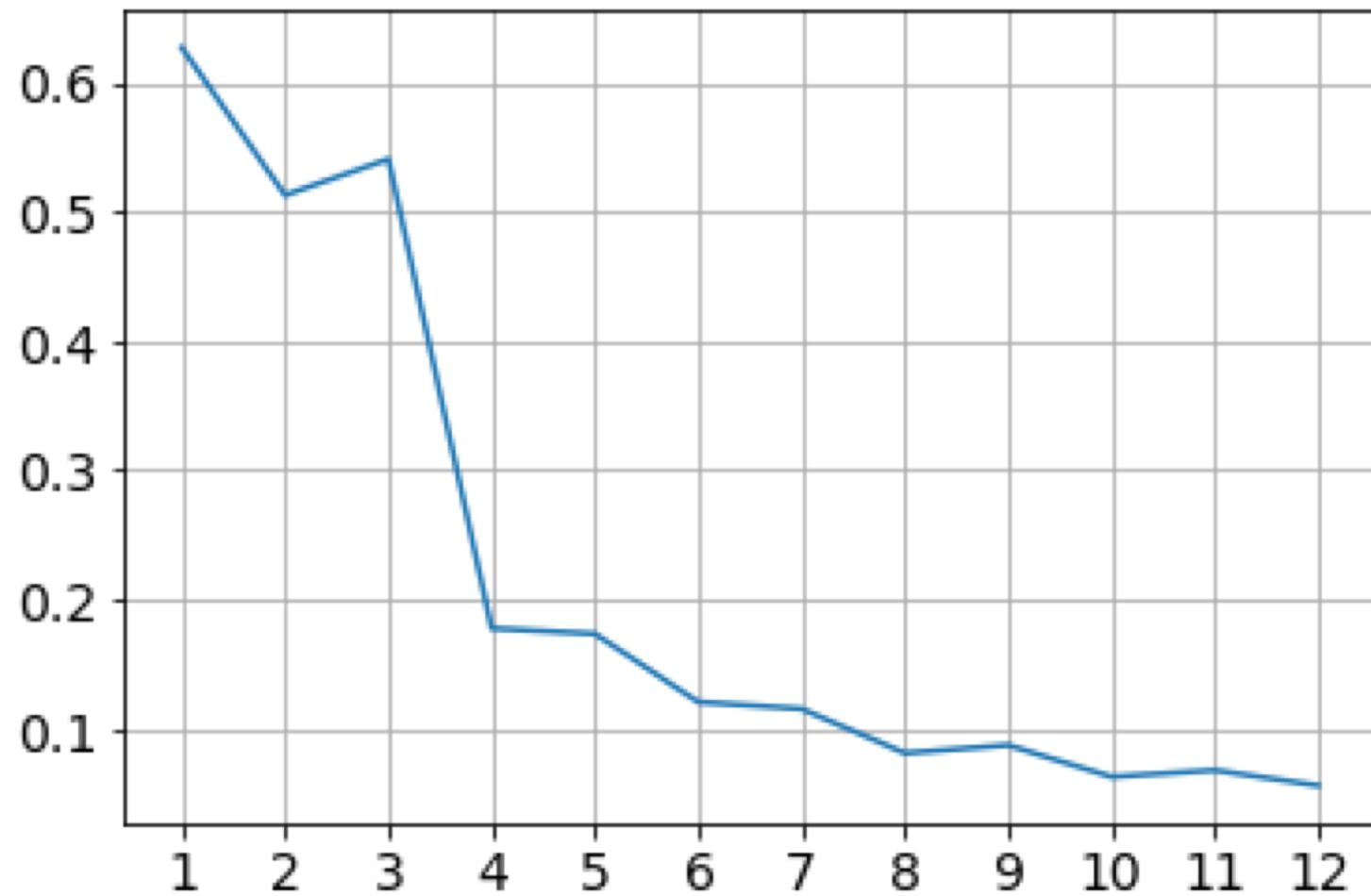


Applying multiple linear regression

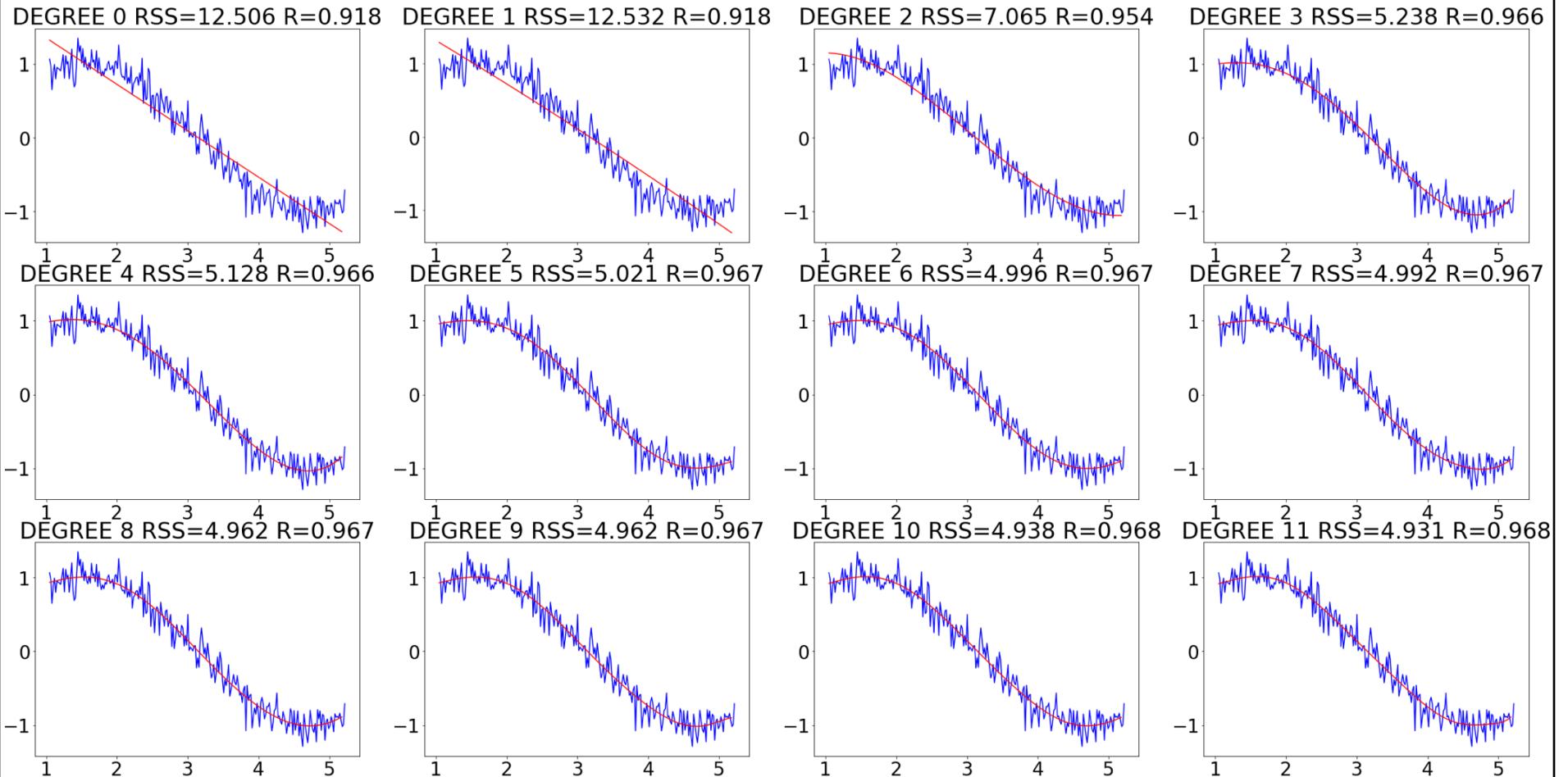


Absolute value of the largest weight computed using simple linear regression with polynomials of increasing degree

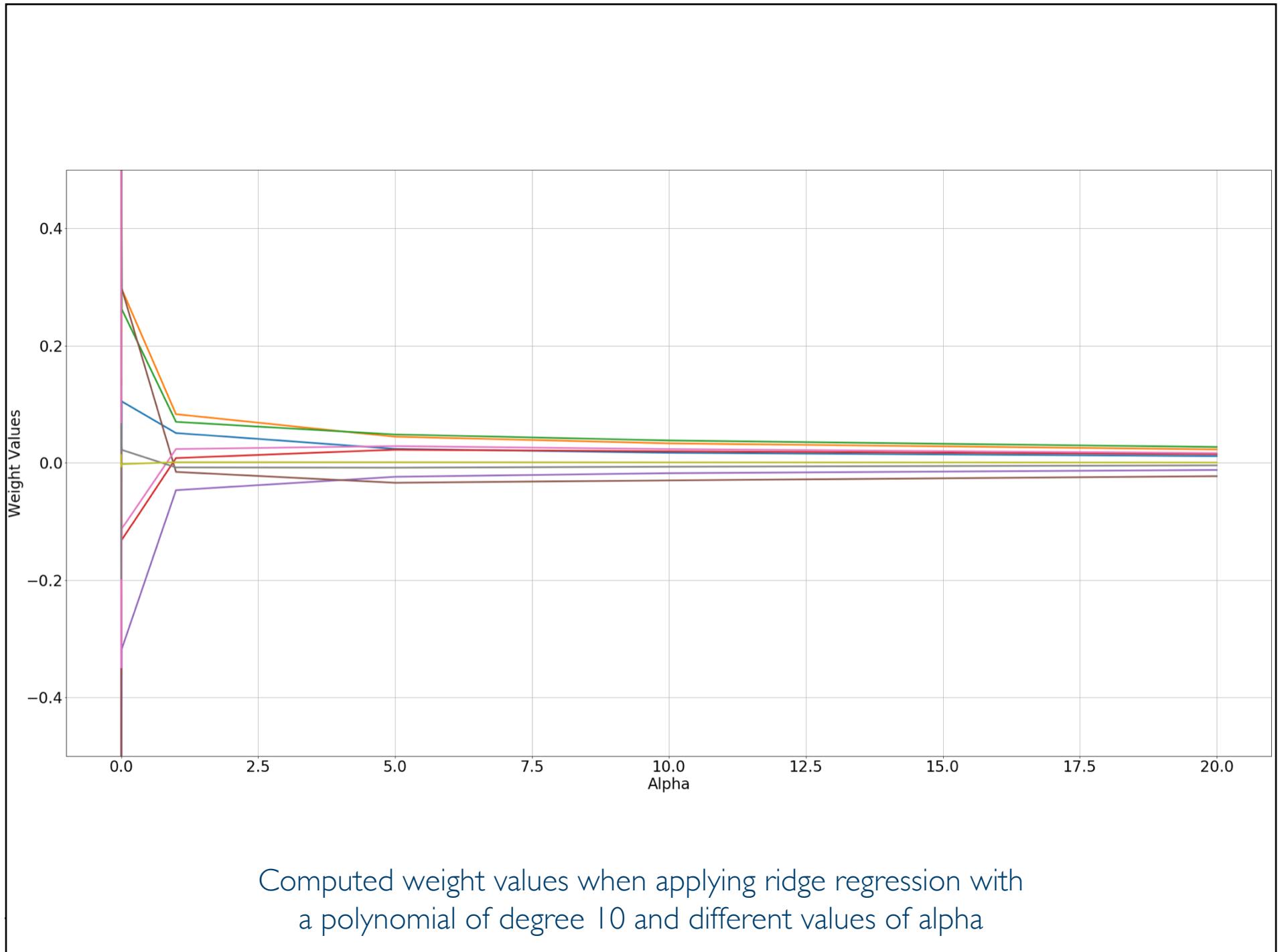
# Ridge Regression



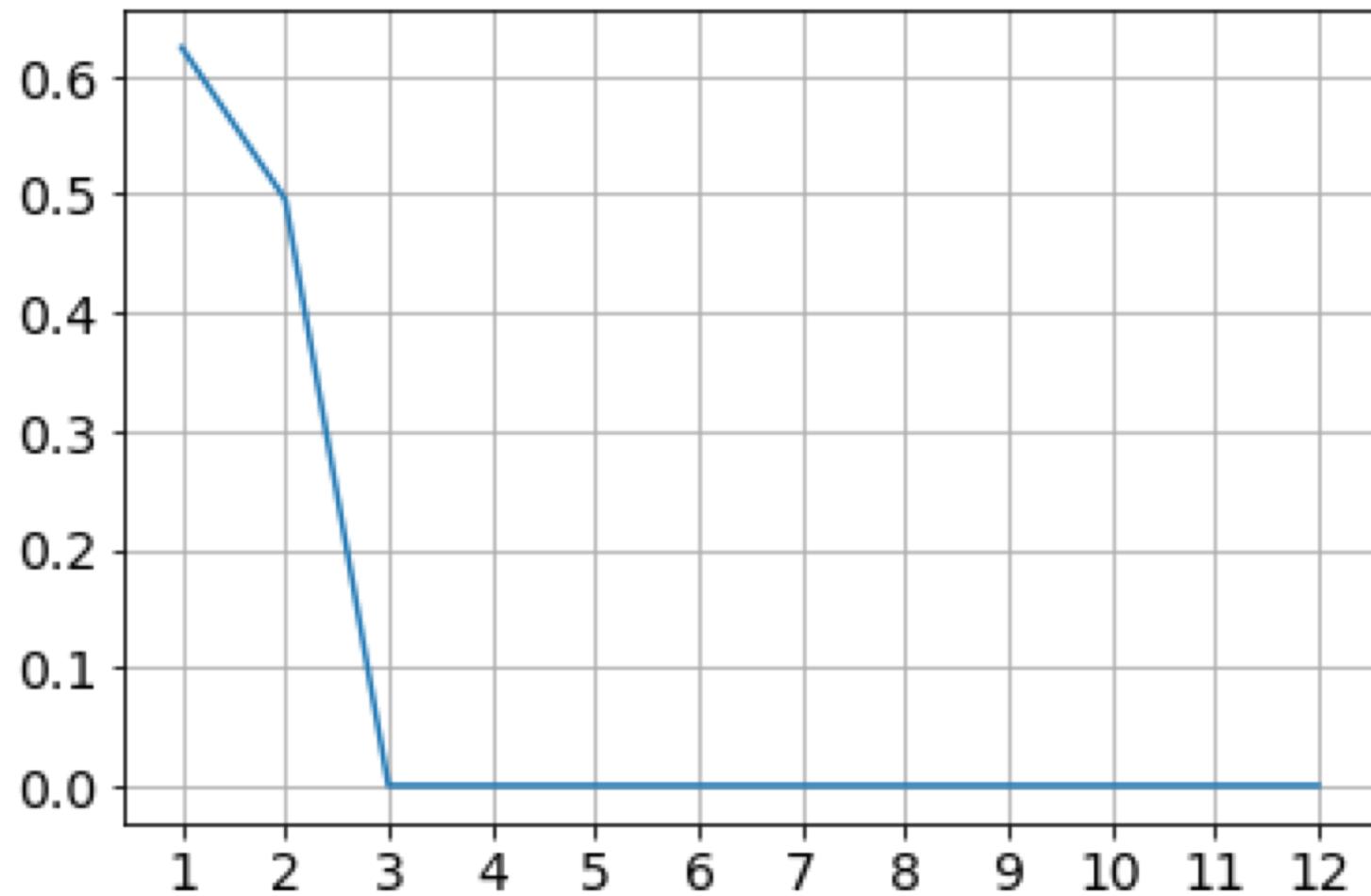
Absolute value of the largest weight computed using  
ridge regression with polynomials of increasing degree



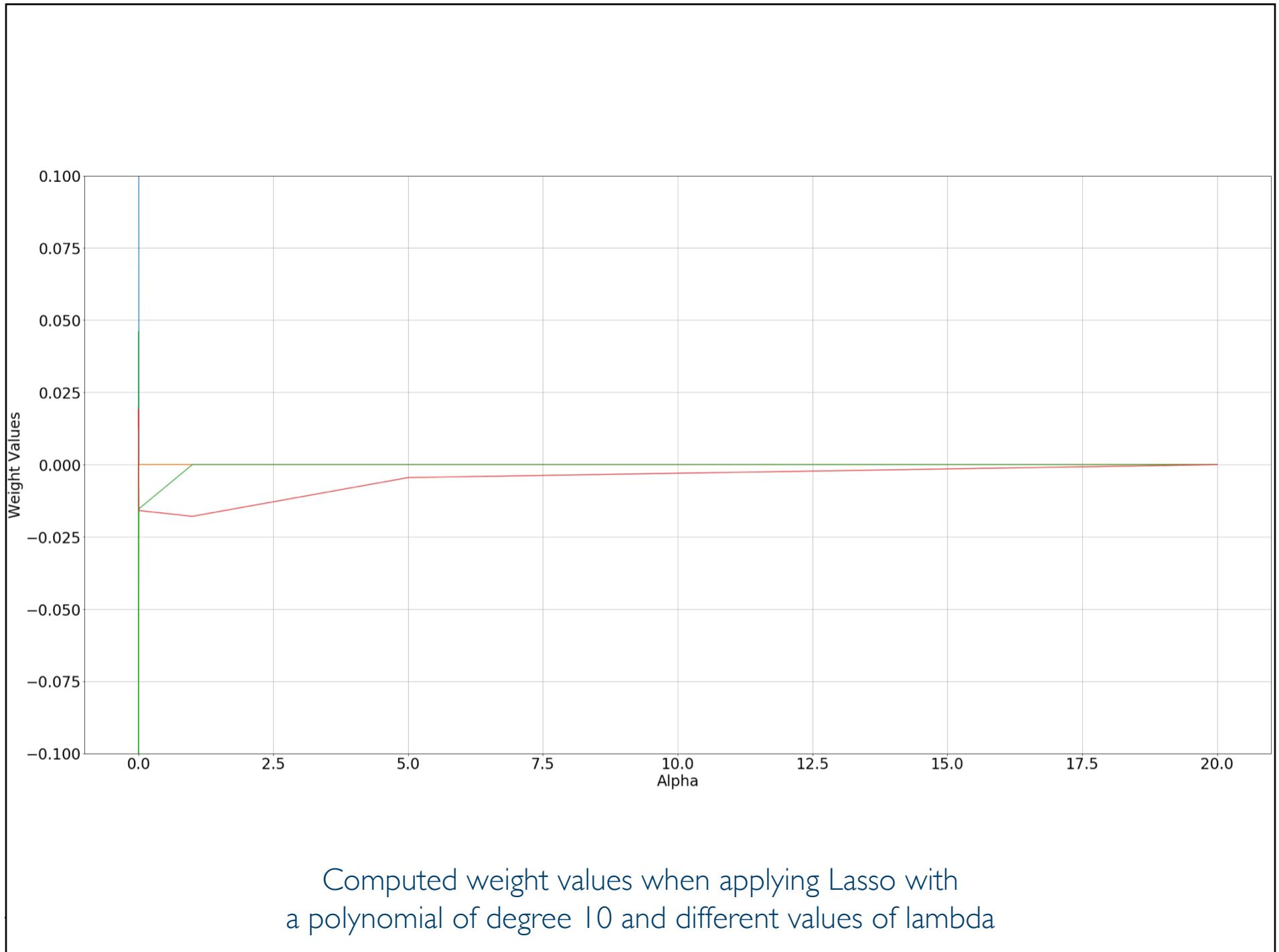
Applying ridge regression with polynomials of increasing degrees



# Lasso



Absolute value of the largest weight computed using  
Lasso with polynomials of increasing degree



Lasso tends to zero out less important features and produces sparser solutions

Basically, by penalizing large weights it also performs feature selection

# Choosing $\alpha$

## Available Data

Training  
model building

Testing  
model evaluation

Training  
model building

Validation  
select  $\alpha$

Testing  
model evaluation

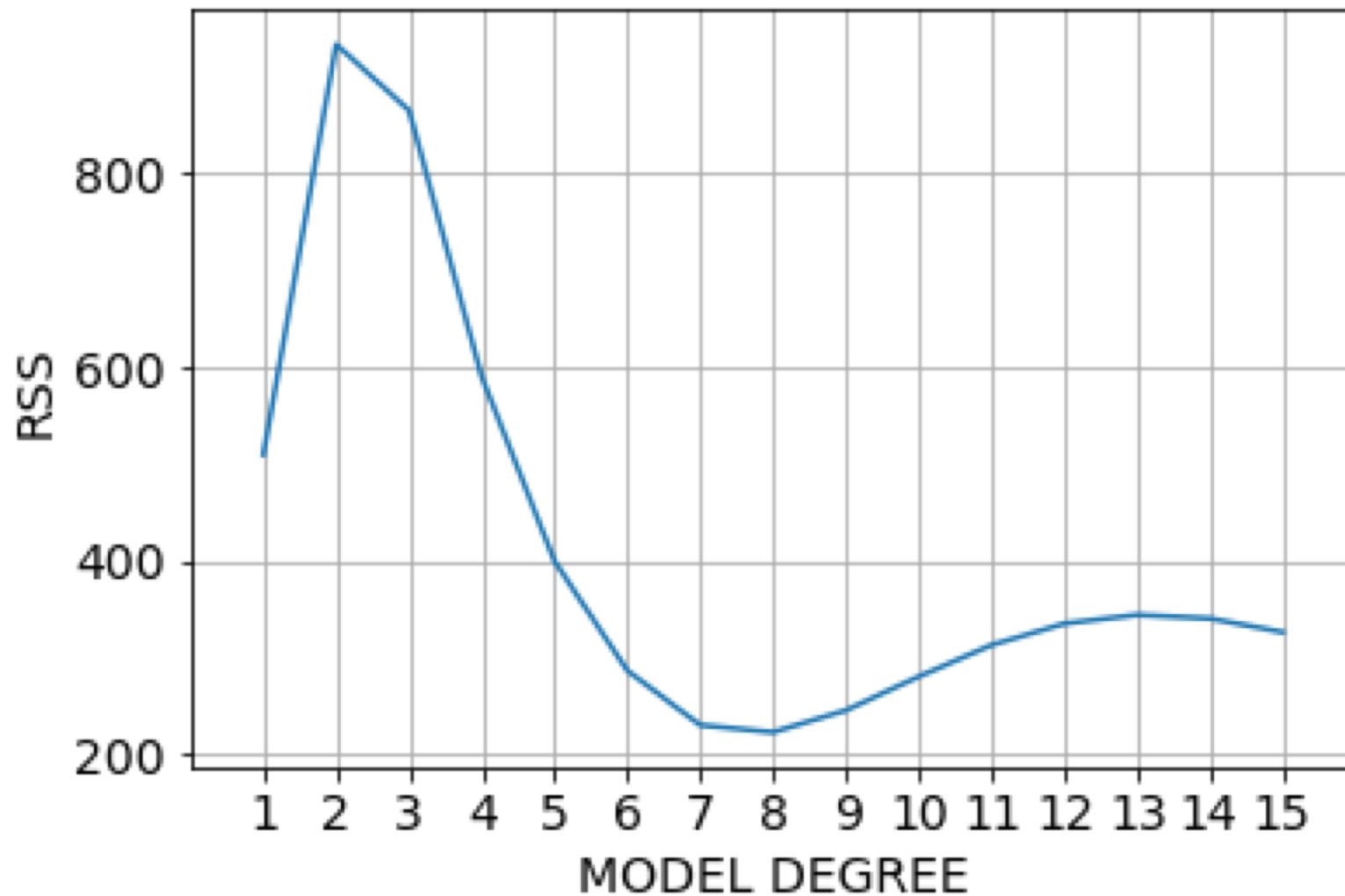
Training &  $\alpha$  Selection  
select the  $\lambda$  with the smallest crossvalidation error then train

Testing  
model evaluation

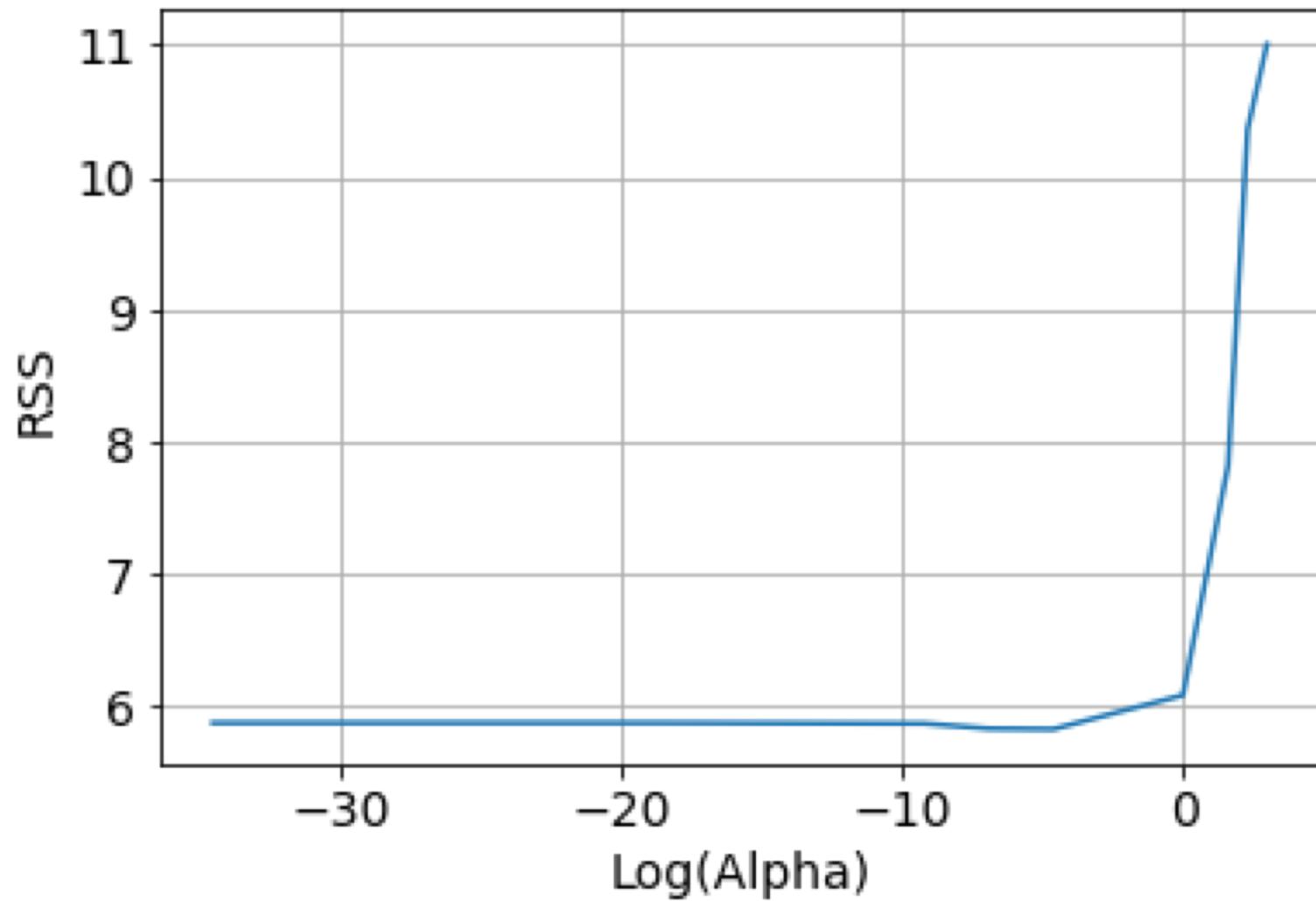
The validation set is also called “dev set”

Because the set is used to develop  
the model before production

- To select the best value of  $\alpha$  we cannot use the test set since it is going to be used for evaluating the final model (which uses  $\alpha$ )
- We need to reserve part of the training data to evaluate possible candidate values of  $\alpha$  and to select the best one
- If we have enough data, we can extract a validation set from the training data which will be used to select  $\alpha$
- If we don't have enough data, we should select  $\alpha$  by applying k-fold crossvalidation over the training data choosing the  $\alpha$  corresponding to the lowest average cost over the k folds



Applying Lasso with a  $\alpha$  of 0.01 with different polynomials



Applying Lasso with differen values of  $\alpha$  – Best  $\alpha$  is 0.01

# Summary

- The goal is to minimize the residual sum of squares (RSS)
- Exact methods
  - Compute the set of weights that minimizes RSS
- Gradient descent (batch and stochastic)
  - Start with a random set of weights and update them based on the direction that minimizes RSS
- Ridge regression/Lasso
  - Compute the cost also using the magnitude of the coefficients
  - The larger the coefficients the more likely we are overfitting

- Check the Python notebooks discussing simple and multiple linear regression, Lasso and Ridge Regression

