



Lecture 6: Regression Introduction

Atsuto Maki
September, 2014

DD2431, CSC/KTH

Part I: we will visit

- Function approximation
- Linear Regression
 - Least squares
 - RANSAC
- KNN Regression

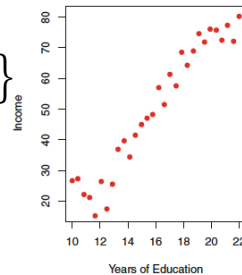
Regression => Real-valued output

Function approximation

- How do we fit this dataset D ?

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

of N pairs of inputs x_i and targets $y_i \in \mathbb{R}$.
 D can be measurements in an experiment.



- Task of regression:
to predict target associated to any arbitrary input

Note: Here we have a single input feature, but inputs to regression tasks are often vectors \mathbf{x} of multiple input features.

Linear Regression

Linear regression tries to estimate the function f and predict the output by

$$\hat{f}(x) = \sum_{i=0}^d w_i x_i = w^T x$$

How to measure the error:

- To see how well $\hat{f}(x)$ approximates $f(x)$, square error is used: $(\hat{f}(x) - f(x))^2$
- Mean Square Error: $E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(x_n) - y_n)^2$ (in-sample)

Minimizing in-sample MSE

E_{in} can be expressed as:

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (w^T x_n - y_n)^2 = \frac{1}{N} \|Xw - Y\|^2$$

where

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix}, \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

We want to compute the parameters w .

Residual sum of squares (RSS)

The sum of squared errors is a convex function of w

$$E_{in}(w) = \|Xw - Y\|^2$$

The gradient with respect to the weights is:

$$\frac{\partial}{\partial w} E_{in}(w) = 2X^T(Xw - Y)$$

The weight vector that sets the gradient to zero minimizes the errors

$$X^T X w = X^T Y$$

$$w = (X^T X)^{-1} X^T Y$$

The linear regression algorithm

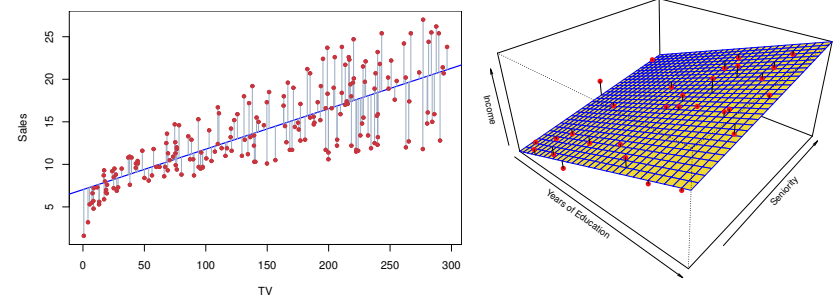
Construct the matrix X and the vector y from the data set $(x_1, y_1), \dots, (x_N, y_N)$ as follows

$$X = \underbrace{\begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad y = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

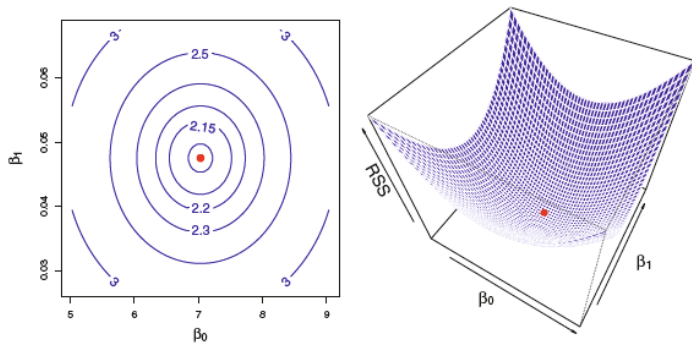
Compute the pseudo-inverse $X^\dagger = (X^T X)^{-1} X^T$.

Return $w = X^\dagger y$.

Examples of least squares fit

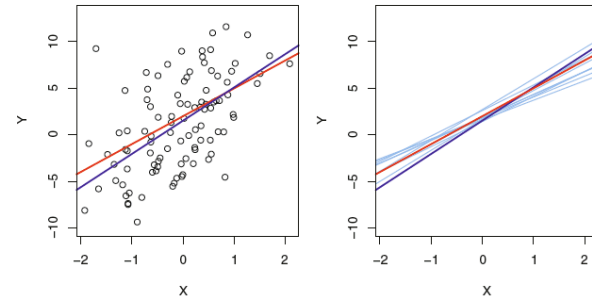


Examples of plots of RSS



Figures from An Introduction to Statistical Learning (G. James et al.)

Least squares line



- Red: the true relationship $f(x) = 2 + 3x$, the **population regression line**
- Blue: the least squares line, estimate based on the observed data
- Light blue (in right): least squares lines, each based on a separate random set of observations

Figures from An Introduction to Statistical Learning (G. James et al.)

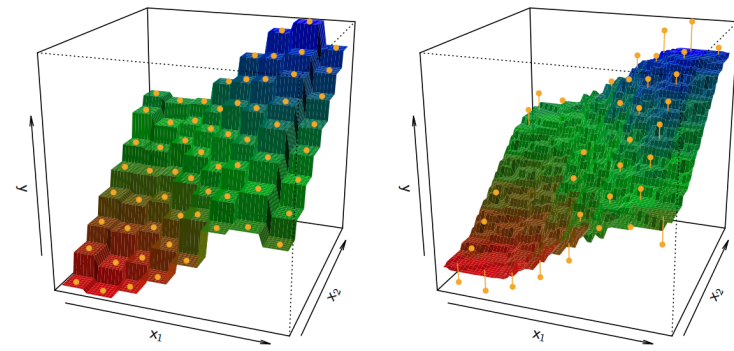
k-NN Regression (non-parametric)

- Similar to the k -NN classifier
- To regress Y for a given value of X , consider k closest points to X in training data and take the average of the responses.

$$f(x) = \frac{1}{k} \sum_{x_i \in N_i} y_i$$

- Larger values of K provide a smoother and less variable fit (lower variance!)

Example plots of $\hat{f}(x)$ with KNN regression

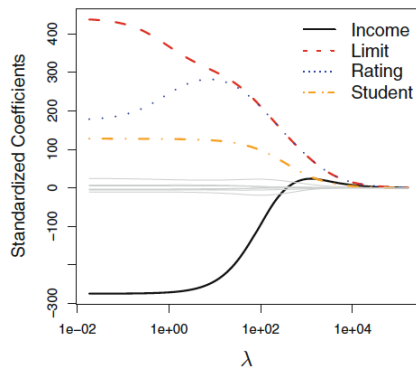


Figures from An Introduction to Statistical Learning (G. James et al.)

Part II: we will visit

- Linear regression + regularization
 - Ridge regression
 - The Lasso (a more recent alternative)

Ridge regression coefficients



As λ increases, the standardized coefficients shrink towards zero (but not exactly forced to zero).

Figures from An Introduction to Statistical Learning (G. James et al.)

Ridge regression

Similar to least squares but minimizes different quantity:

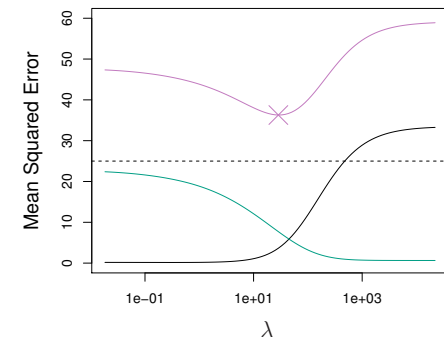
$$RSS + \lambda \sum_{i=1}^d w_i^2$$

The second term is called **shrinkage penalty**

- Shrinkage penalty: small when w_i are close to zero
- The parameter λ : controls the relative impact of the two terms, the selection is critical!

Ridge Regression Bias/Variance

- Purple: MSE
- Black: Bias
- Green: Variance



Increase λ decreases variance while increasing bias

Figures from An Introduction to Statistical Learning (G. James et al.)

The Lasso (Least Absolute Shrinkage and Selection Operator)

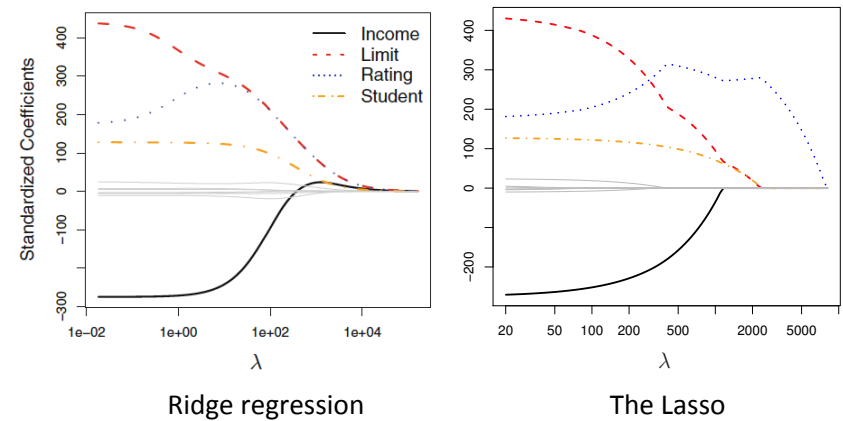
Similar to ridge regression but with slightly different term:

$$RSS + \lambda \sum_{i=1}^d |w_i|$$

The **shrinkage penalty** is now replaced by **l_1 norm**

- Ridge regression: it includes **all features** in the final model, making it harder to interpret – its drawback
- The lasso could be proven mathematically that some coefficients end up being set to **exactly zero**
 - variable selection
 - yielding sparse model

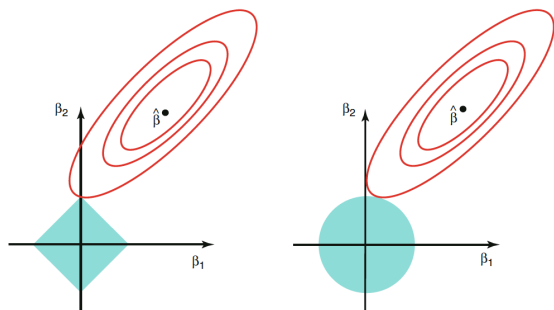
Comparison of estimated coefficients



Figures from An Introduction to Statistical Learning (G. James et al.)

The variable selection property

The regression coefficient estimates: **the first point where an ellipse contacts the constraint region**



The solid blue areas are the constraint regions for
 Left: the Lasso Right: Ridge regression

Figures from An Introduction to Statistical Learning (G. James et al.)