

# Business Analytics II

## Lecture 1

### Regression with Python

Pedro Hernandez Serrano  
p.hernandezserrano@maastrichtuniversity.nl

April 10, 2018

BISS  
Institute

## Today's session

2/46

- ▶ Data Science Applications on Business Cases
  - Data products and type of users
- ▶ Supervised Learning - Regression in Python
  - What is Regression?
  - Use cases: Bias, Spurious correlations, wrong use cases
  - Univariate Statistics
  - Linear Model
  - Multiple regression
  - Overfitting and Complexity
  - Regularization methods: Ridge, Lasso, Elastic-net
- ▶ Hands-on during the session on Jupyter Notebooks
  - Construct a linear regressor from scratch
  - Measuring Overfitting and Performance
  - Regression Models using Turicreate Machine Learning Library

# Data Science Applications on Business Cases

## Data Science Applications

### Business Intelligence vs Business Analytics



Image: <https://www.linkedin.com/pulse/business-intelligence-vs-analytics-what-needs-know-john-vuu/>

## Business Intelligence

- ▶ Recording to DWH
- ▶ Getting new insights from the business data
- ▶ Customized reporting
- ▶ Analyzing the history (past)

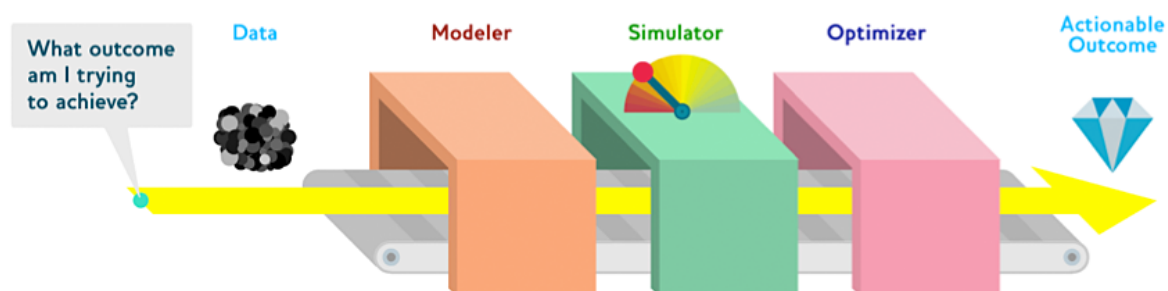
## Business Analytics

- ▶ Exploiting from DWH
- ▶ Understanding the limitations of the business data
- ▶ Customized models
- ▶ Analyzing predictions (future)

Both conduct data-driven decision making.

# Data Science Applications

## Data Product



- ▶ Dynamic reports
- ▶ Static model for decision making
- ▶ Prototype for a new business product
- ▶ Production model for users (consumer, internal. executive, etc.)

Image: <http://blog.kaggle.com/2012/03/28/drivetrain-approach-to-designing-great-data-products/>

## Pros

- ▶ Focuses on data driven decision over politics and gut feelings
- ▶ Automates decisions that might be financially and mentally taxing
- ▶ Improves consistency, accuracy and forces teams to draw out their decisions processes

## Cons

- ▶ If an algorithm is incorrect the team might overly trust it
- ▶ Privacy regulations and data protection might be an issue
- ▶ GDPR and explainability of models

## Example

- ▶ Fire-up Jupyter Notebook via **<https://jupyter.org/try>**
- ▶ Go to the drive folder of the class and make a personal copy
- ▶ Open intro.ipynb
- ▶ Discuss the problem

# Supervised Learning Regression with Python


## Supervised Learning Regression

### What is Regression?

**A technique for determining the statistical relationship between two or more variables where a change in a dependent variable is associated with, and depends on, a change in one or more independent variables.**

- ▶ A Supervised Learning Method of Machine Learning
- ▶ A Statistical Method in Bayesian or Frequentist Inference

#### regression

/rɪˈɡreʃ(ə)n/ 

*noun*

1. a return to a former or less developed state.  
"it is easy to blame unrest on economic regression"
2. **STATISTICS**  
a measure of the relation between the mean value of one variable (e.g. output) and corresponding values of other variables (e.g. time and cost).

[1]<http://www.businessdictionary.com/definition/regression.html>

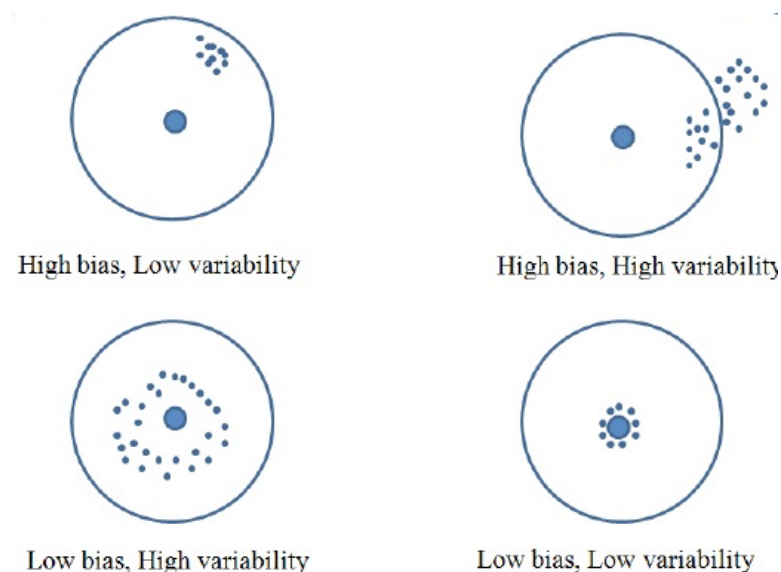
## What is Bias?

**Bias is derived from the ancient Greek word that describes an oblique line (i.e., a deviation from the horizontal). In Data Science, bias is a deviation from expectation in the data.**

- ▶ Several definitions, common usage is decidedly negative.
- ▶ Systematic favoritism of a group
- ▶ Dataset Bias: Sampling bias in a given dataset, give rise to (often unintentional) wrong results when is not representative of the population
- ▶ Social Bias in a Dataset: Humans can introduce bias during data collection.

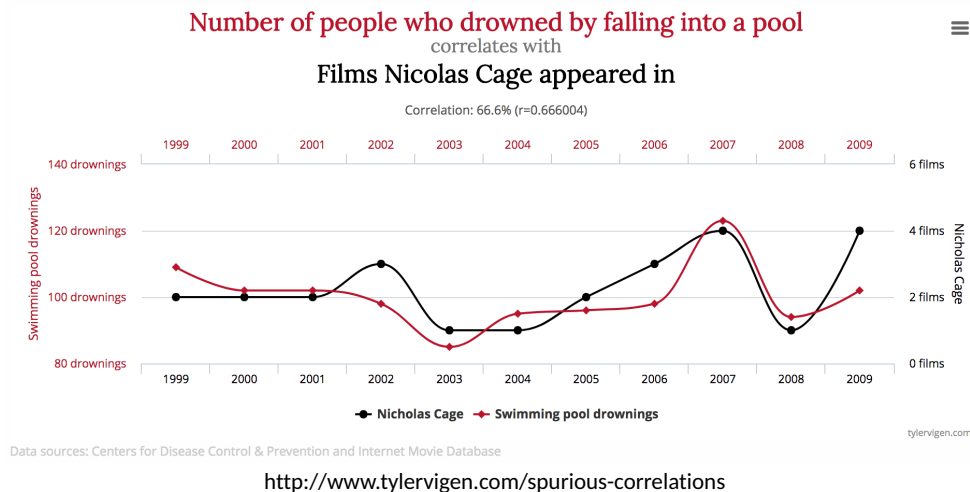
# Supervised Learning Regression

## Bias vs Variability?



## Spurious Correlations

Is a relationship between two variables that appear to have interdependence or association with each other but actually do not.



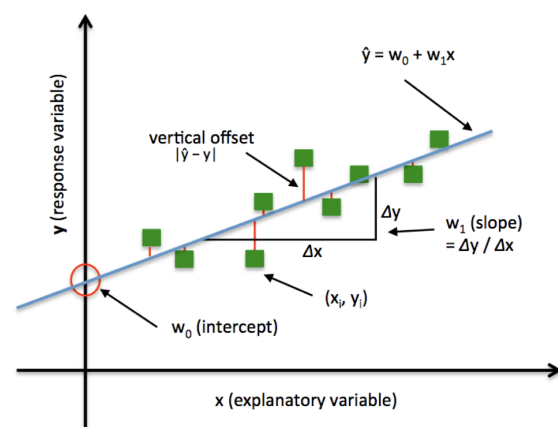
# Supervised Learning Regression

## Linear Model Representation: Statistics vs Machine Learning

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Dependent Variable  $Y_i$  is equal to the Population Y intercept  $\beta_0$  plus the Population Slope Coefficient  $\beta_1$  times the Independent Variable  $X_i$ , plus the Random Error term  $\epsilon_i$ .

The Linear component is  $\beta_0 + \beta_1 X_i$  and the Random Error component is  $\epsilon_i$ .



### Statistics

- ▶  $\beta$  (Coefficient)
- ▶  $Y$  = Dependent Variable
- ▶  $X$  = Independent Variable

### Machine Learning

- ▶  $w$  (Weight)
- ▶  $Y$  = Target/Response Feature
- ▶  $X$  = Input Feature

# Univariate Statistics

## Univariate Statistics

Basics univariate statistics are required to explore a dataset:

- ▶ Discover associations between a variable of interest and potential predictors. It is strongly recommended to start with simple univariate methods before moving to complex multivariate predictors.
- ▶ Assess the prediction performances of machine learning predictors.
- ▶ Most of the univariate statistics are based on the linear model which is one of the main model in machine learning.
- ▶ Check every estimator of the main statistical measures



## Mean

Properties of the expected value operator  $E()$  of a random variable  $X$

$$E(X + c) = E(X) + c \quad (1)$$

$$E(X + Y) = E(X) + E(Y) \quad (2)$$

$$E(aX) = aE(X) \quad (3)$$

The estimator  $\hat{X}$  on a sample of size  $n : x = x_1, \dots, x_n$  is given by

$$\hat{X} = \frac{1}{n} \sum_i x_i$$

$\hat{X}$  is itself a random variable with properties:

$$E(\hat{X}) = \hat{X} \quad \text{Var}(\hat{X}) = \frac{\text{Var}(X)}{n}$$

## Variance

$$\text{Var}(X) = E((X - E(X))^2) = E(X^2) - (E(X))^2$$

The estimator is:

$$\sigma_x^2 = \frac{1}{n-1} \sum_i (x_i - \hat{x})^2$$

## Statistical Bias

$$\text{Bias}(y, \hat{y}) = E((y - \hat{y})^2) - \text{Var}(\hat{y})$$

Where  $\text{MSE} = E((y - \hat{y})^2)$

## Covariance

$$\text{Cov}(X, Y) = E((X - E(X))(Y - E(Y))) = E(XY) - E(X)E(Y)$$

Properties:

$$\text{Cov}(X, X) = \text{Var}(X) \quad (4)$$

$$\text{Cov}(X, Y) = \text{Cov}(Y, X) \quad (5)$$

$$\text{Cov}(cX, Y) = c\text{Cov}(X, Y) \quad (6)$$

$$\text{Cov}(X + c, Y) = \text{Cov}(X, Y) \quad (7)$$

The estimator with  $df = 1$  is

$$\sigma_{xy} = \frac{1}{n-1} \sum_i (x_i - \hat{x})(y_i - \hat{y})$$

# Univariate Statistics

## Correlation

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\text{Std}(X)\text{Std}(Y)}$$

The estimator is:

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

## Standard Error (SE)

As well as the standard deviation (of the sampling distribution) of a statistic:

$$\text{SE}(X) = \frac{\text{Std}(X)}{\sqrt{n}}$$

Commonly considered for the mean with the estimator:  $\text{SE}(\hat{x}) = \sigma_x / \sqrt{n}$

## Testing pairwise associations

Exploring association between pairs of variables.

- ▶ **Categorical variable:** Is a variable that can take on one of a limited, and usually fixed, number of possible values, thus assigning each individual to a particular group or category. The levels are the possible values of the variable (also known as a factor).
- ▶ **Ordinal variable:** Is a categorical variable with a clear ordering of the levels. (none, small, medium and high).
- ▶ **Continuous variable:**  $x \in R$  is one that can take any value in a range of possible values, possibly infinite (salary, experience in years, weight) also known as quantitative.

## Pearson correlation test: Association between two quantitative variables

The test calculates a Pearson correlation coefficient and the  $p$ -value for testing non-correlation.

Let  $x$  and  $y$  two quantitative variables, where samples were observed. The linear correlation coefficient is defined as :

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

Under  $H_0$  the test statistic  $t = \sqrt{n-2} \frac{r}{\sqrt{1-r^2}}$  follow Student distribution with  $n - 2$  degrees of freedom

## Correlation test with Python

```
import numpy as np
import scipy.stats as stats
n = 50
x = np.random.normal(size=n)
y = 2 * x + np.random.normal(size=n)

#Parametric Pearson correlation
cor, pval = stats.pearsonr(x, y)
print("Pearson r, cor: %.4f, pval: %.4f" %(cor, pval))

# Non-Parametric Spearman correlation
cor, pval = stats.spearmanr(x, y)
print("Spearman r, cor: %.4f, pval: %.4f" %(cor, pval))
```

## Independency T-test with Python

```
import numpy as np
import scipy.stats as stats

#Simulate data
nx, ny = 50, 25
x = np.random.normal(loc=1.76, scale=0.1, size=nx)
y = np.random.normal(loc=1.70, scale=0.12, size=ny)

#t-Test Independency
tval, pval = stats.ttest_ind(x, y, equal_var=True)
print("t-Test, tval: %.4f, pval: %.4f" % (tval, pval))
```

# Linear Model

## Linear Model

Given  $n$  random samples  $(y_i, x_i^1, \dots, x_i^p), i = 1, \dots, n$ , the linear regression models the relation between the observations and the independent variables  $x_i^p$  is formulated as

$$y_i = \beta_0 + \beta_1 x_i^1 + \dots + \beta_p x_i^p + \epsilon_i \quad i = 1, \dots, n$$

- **An independent variable.** It is a variable that stands alone and is not changed by the other variables you are trying to measure. In Machine Learning, these variables are also called the predictors or input features.
- **A dependent variable.** It is something that depends on other factors. Usually when you are looking for a relationship between two things you are trying to find out what makes the dependent variable change the way it does. In Machine Learning this variable is called a target variable.

## Case: House Price Problem

How much is my house worth?



Let's take a look at sales in my neighbourhood, How much did they sell for?



## Case: House Price Problem

### Data



$(x_1 = \text{sq.ft.}, y_1 = \$)$



$(x_2 = \text{sq.ft.}, y_2 = \$)$



$(x_3 = \text{sq.ft.}, y_3 = \$)$

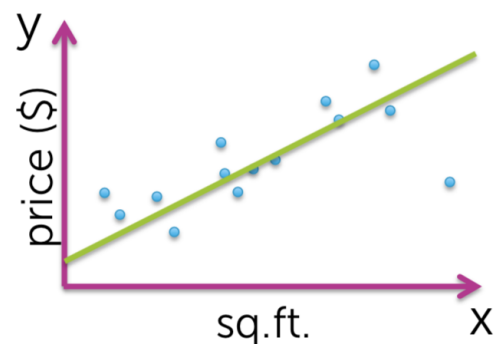


$(x_4 = \text{sq.ft.}, y_4 = \$)$



$(x_5 = \text{sq.ft.}, y_5 = \$)$

$\vdots$



## Case: House Price Problem

- ▶ Go to **simple regression** notebook
- ▶ Load the data of houses **housesdata.csv**
- ▶ Choose **price** as target variable
- ▶ Choose **sq.ft.** as predictor
- ▶ Answer the questions

## Simple Linear regression

In simple linear regression, we attempt to model the relationship between two variables and it has the form:

$$y = \beta_0 + \beta_1 x + \epsilon$$

- ▶  $\beta_1$ : the slope or coefficient or parameter of the model,
- ▶  $\beta_0$ : the **intercept or bias** is the second parameter of the model,
- ▶  $\epsilon_i$ : the  $i$ th error, or residual with  $\epsilon(0, \sigma^2)$ .

The simple regression is equivalent to the Pearson correlation.

## Simple regression

Regression Estimators The goal is to estimate,  $\beta$ , and  $\beta_o$ . We have to minimize the mean squared error (MSE) or the Sum squared error (SSE). The so-called Ordinary Least Squares (OLS)

$$\hat{\beta} = \frac{\frac{1}{n} \sum_i x_i y_i - \bar{y}\bar{x}}{\frac{1}{n} \sum_i (x_i \bar{x})} = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{(\bar{xy}) - (\bar{x})(\bar{y})}{(\bar{x^2}) - (\bar{x})^2}$$

$$\hat{\beta}_o = E(y) - \beta E(x) = \bar{y} - \beta \bar{x}$$

## Fitting

$$\hat{Y} = \hat{\beta}_o + \hat{\beta}x$$

```
def simple_linear_regression(x, y):  
    xy = x * y #create xy variable  
    mean_xy = xy.mean()  
    mean_x = x.mean()  
    mean_y = y.mean()  
    xx = x * x # create xx variable  
    mean_xx = xx.mean()  
    x_sqr = mean_x * mean_x  
    #compute slope  
    cov_xy = mean_xy - mean_x * mean_y  
    var_x = mean_xx - x_sqr  
    slope = cov_xy / var_x  
    #use the slope for the intercept  
    intercept = mean_y - slope * mean_x  
    return intercept, slope
```



## Goodness of fit

The goodness of fit of a statistical model describes how well it fits a set of observations. Measures of goodness of fit typically summarize the discrepancy between observed values and the values expected under the model in question.

- Residual Sum of Squares

$$RSS = \sum_{i=1}^n (y_i - \hat{y})^2$$

- Root Mean Squared Error

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2}$$

# Multiple Regression

Multiple Linear Regression is the most basic supervised learning algorithm. Given: a set of training data  $x_1, \dots, x_n$  with corresponding targets  $y_1, \dots, y_n$ . In linear regression, we assume that the model that generates the data involves only a linear combination of the input variables, i.e

$$y(x_i, \beta) = \beta^0 + \beta^1 x_i^1 + \dots + \beta^P x_i^P,$$

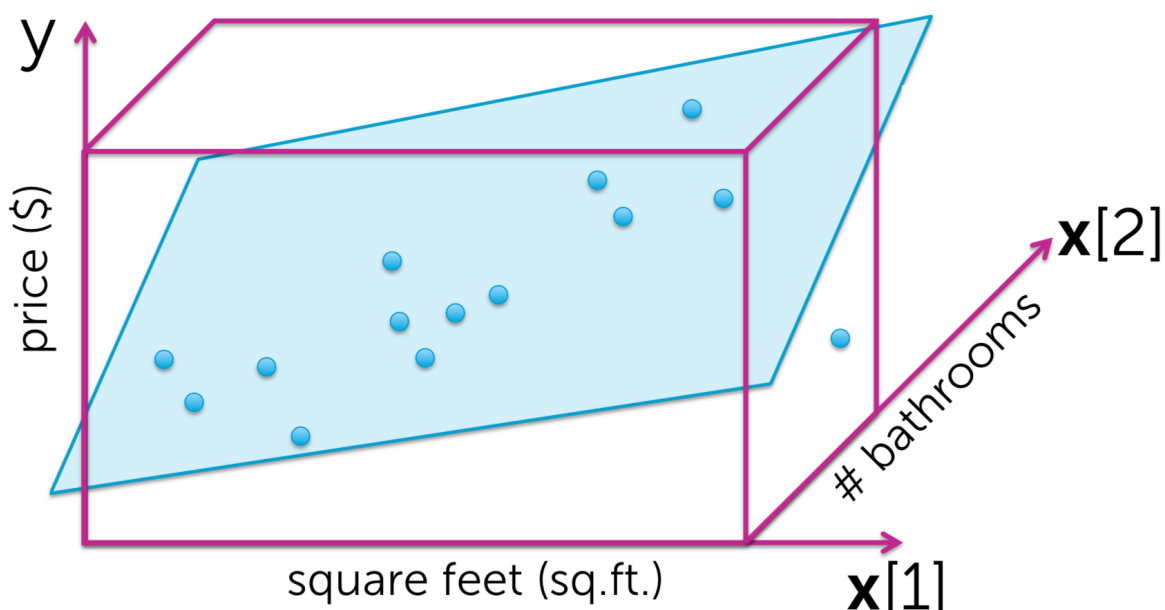
Extending each sample with an intercept  $x_i := [1, x_i] \in \mathbb{R}^{P+1}$  allows us to use a more general notation based on linear algebra and write it as a simple dot product:

$$y(x_i, \beta) = x_i^T \beta,$$

Where  $\beta \in \mathbb{R}^{P+1}$  is a vector of weights that define the  $P + 1$  parameters of the model. From now we have  $P$  regressors + the intercept.

## Multiple regression

- ▶ Linear regression with multiple features
- ▶ More complex functions of a single input



## Case: House Price Problem

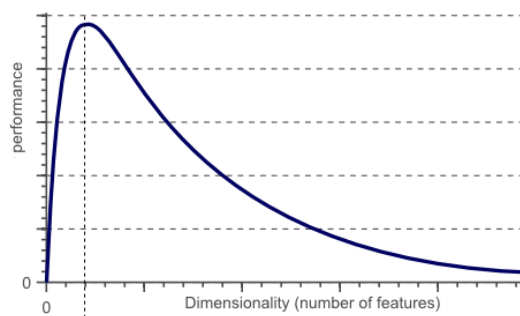
- ▶ Go to **multiple regression** notebook
- ▶ Load the data of houses **housesdata.csv**
- ▶ Check the turicreate API
- ▶ Run the multiple regression example
- ▶ Create the new features
- ▶ Learn multiple models proposed

## Overfitting and Complexity

In statistics and machine learning, overfitting occurs when a statistical model describes random errors or noise instead of the underlying relationships. A model that has been overfitted will generally have poor predictive performance, as it can exaggerate minor fluctuations in the data.

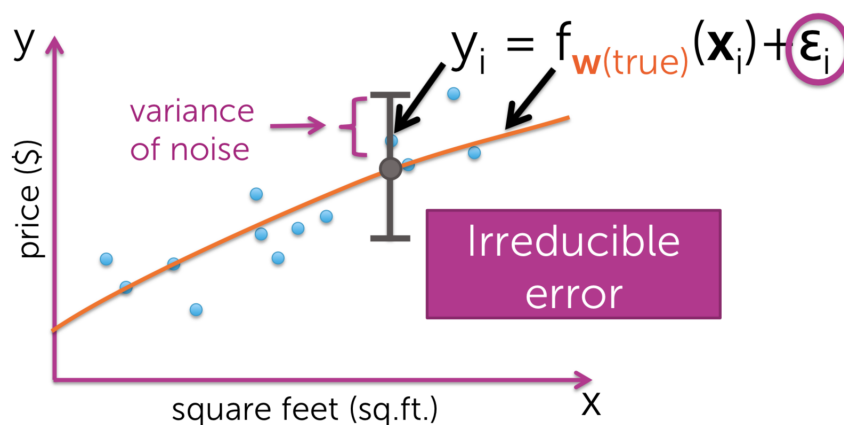
The overfitting phenomenon has three main explanations:

- ▶ Excessively complex models
- ▶ Multicollinearity
- ▶ High dimensionality



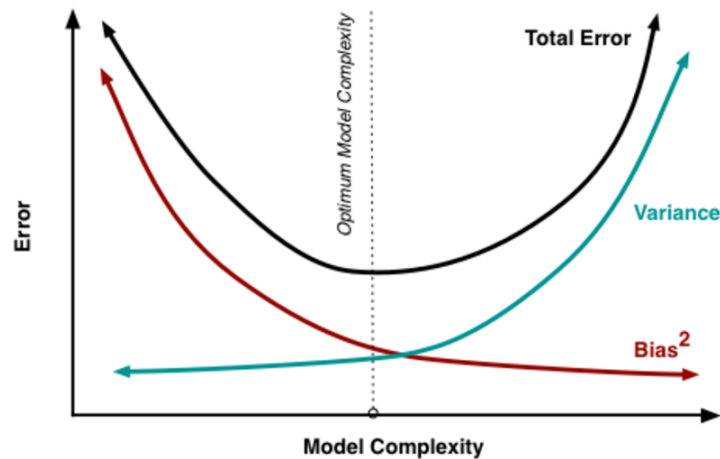
3 sources of error

- ▶ Noise
- ▶ Bias
- ▶ Variance



## MSE in terms of Bias and Variance

$$MSE = Var(\hat{y}) + E((y - \hat{y})^2) = \frac{1}{n-1} \sum_i (y_i - \hat{y})^2 + \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$$



<http://scott.fortmann-roe.com/docs/MeasuringError.html>

Image:

## Regularization Methods

## Ridge regression ( $L_2$ -regularization)

Overfitting generally leads to excessively complex weight vectors, accounting for noise or spurious correlations within predictors. To avoid this phenomenon the learning should constrain the solution in order to fit a global pattern. This constraint will reduce (bias). Adding such a penalty will force the coefficients to be small, i.e. to shrink them toward zeros.

$$\text{Ridge}(\beta) = \|y - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2$$

- ▶ Increasing  $\lambda$  shrinks the coefficients toward 0.
- ▶ Penalizes the objective function by the Euclidian norm of the coefficients such that solutions with large coefficients become unattractive.

# Regularization Methods

## Lasso regression ( $L_1$ -regularization)

Lasso regression penalizes the coefficients by the  $L_1$  norm. This constraint will reduce (bias) the capacity of the learning algorithm. To add such a penalty forces the coefficients to be small, i.e. it shrinks them toward zero. The objective function to minimize becomes:

$$\text{Lasso}(\beta) = \|y - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$$

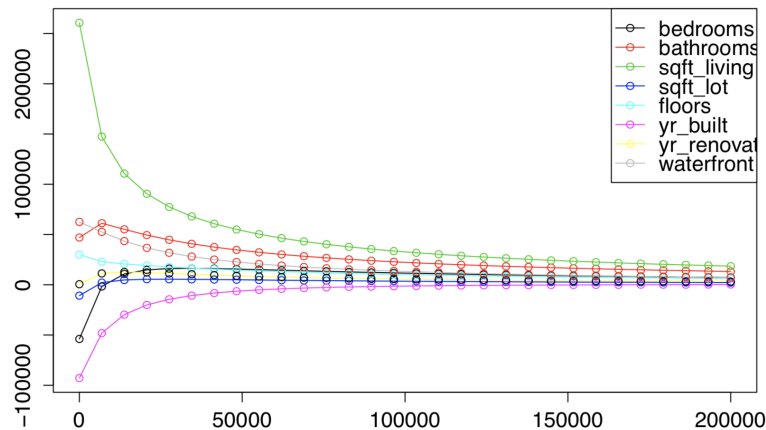
- ▶ Choose the model with the smallest coefficient vector, i.e. smallest  $L_2$  or  $L_1$  norm
- ▶ Choose the model that uses the smallest number of predictors. In other words, choose the model that has many predictors with zero weights.

## Elastic-net regression ( $L_2 - L_1$ -regularization)

The Elastic-net estimator combines the  $L_1$  and  $L_2$  penalties, and results in the problem to:

$$\text{Enet}(\beta) ||y - \mathbf{X}^T \beta||_2^2 + \alpha(\rho ||\beta||_1 + (1 - \rho) ||\beta||_2^2), \quad (8)$$

where  $\alpha$  acts as a global penalty and  $\rho$  as an  $L_1/L_2$  ratio.



# Regularization Methods

## Summary Regularization

- Powerful techniques generally used for creating parsimonious models in presence of a large number of features.
- Large enough to enhance the tendency of a model to overfit (as low as 10 variables might cause overfitting).
- Large enough to cause computational challenges. With modern systems, (millions or billions of features).
- If a feature is too large, means that we are putting a lot of emphasis on that feature.