



ALMA MATER STUDIORUM  
UNIVERSITÀ DI BOLOGNA

DIPARTIMENTO DI  
SCIENZE STATISTICHE "PAOLO FORTUNATI"



# 1 - Introduction to Machine Learning

Giovanni Della Lunga

Halloween Conference in Quantitative Finance

Bologna - October 26-28, 2021

# Outline

- 1 General Ideas
  - What is Machine Learning
  - Feature Scaling
  - Cost Functions
  - Gradient Descent
  - Validation and Testing
  - Bias and Variance
  - Regularization

# Introduction

- These introductory lessons assume a basic level of statistical and mathematical knowledge.
- No previous knowledge of Machine Learning is assumed.
- For this reason I have decided to use two basic texts for the preparation of these lessons that you can consult for further details on the topics we are going to deal with.

*John C. Hull*, **Machine Learning in Business, An Introduction to the World of Data Science**, Amazon (2019)

*Paul Wilmott*, **Machine Learning, An Applied Mathematics Introduction**, Panda Ohana Publishing (2019)

# What is Machine Learning

- Machine learning is a branch of AI
- The idea underlying machine learning is that we give a computer program access to lots of data and let it learn about relationships between variables and make predictions
- Some of the techniques of machine learning date back to the 1950s but improvements in computer speeds and data storage costs have now made machine learning a practical tool
- Machine Learning or data science can be described as the new world of statistics

# Software

- There are several alternatives such as Python, R, MatLab, Spark, and Julia
- Need ability to handle very large data sets and availability of packages that implement the algorithms.
- Python seems to be winning at the moment
- Scikit-Learn has freely available packages for many ML tasks

**A word of caution.** It is tempting to learn a language such as Python and apply various packages to your data without really understanding what the packages are doing or even how the results should be interpreted. This would be a bit like a finance specialist using the Black-Scholes model to value options without understanding where it comes from or its limitations...

# Traditional statistics

- Means, SDs
- Probability distributions
- Significance tests
- Confidence intervals
- Linear regression
- etc

# The new world of statistics

- Huge data sets
- Fantastic improvements in computer processing speeds and data storage costs
- Machine learning tools are now feasible
- Can now develop non-linear prediction models, find patterns in data in ways that were not possible before, and develop multi-stage decision strategies
- New terminology: features, labels, activation functions, target, bias, supervised/unsupervised learning



# Applications of ML

- Credit decisions
- Classifying and understanding customers better
- Portfolio management
- Private equity
- Language translation
- Voice recognition
- Biometrics
- etc

# Types of Machine Learning

- Unsupervised learning (find patterns)
- Supervised learning (predict numerical value or classification)
- Semi-supervised learning (only part of data has values for, or classification of, target)
- Reinforcement learning (multi-stage decision making)

In these introductory lessons we will discuss only the first two types.

# Types of Machine Learning

- **Supervised Learning**

is concerned with using data to make prediction (for example a simple regression model). We can distinguish between supervised learning models that are used to predict a variable than can take a continuum of values and models that are used for classification (for example to learn if a potential borrower can be classified as acceptable or unacceptable credit risk);

- **Unsupervised Learning**

is most concerned with recognizing patterns in data. The main objective usually is not to forecast a particular variable, rather it is to understand the environment represented by the data better. As we will see clustering is one of the most used tool in unsupervised learning.

# Features and Labels

- The data for **Supervised Learning** contains what are referred to as **features** and **labels**;
- *Labels* are the values of the target that is to be predicted;
- *Features* are the variables from which the predictions are to be made (if you are from statistics then think of it as an explanatory variable);

The diagram shows a table with 5 columns and 5 rows. A bracket above the first four columns is labeled 'Features', and a bracket above the last column is labeled 'Label'. A bracket to the left of the rows is labeled 'Rows', and a bracket below the columns is labeled 'Columns'.

	Features				Label
Size	Beds	Baths	Zip	Price	
1100	1	1	64576	1.29	
1900	3	1.5	78321	2.14	
2800	3	3	98712	3.10	
3400	4	3.5	25721	3.75	

# Features and Labels

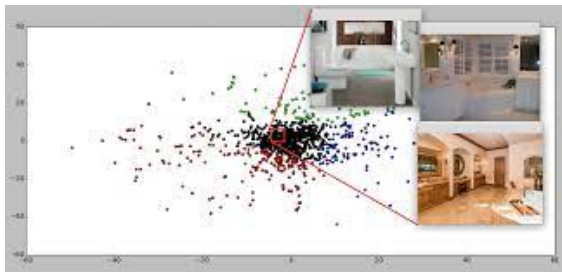
- For example, when predicting the **price of a house**, the *features* could be the *square meters of living space*, the *number of bedrooms*, the *number of bathrooms*, the *size of the garage* and so on.
- The *label* would be the *house price*;

The diagram shows a table with 5 columns and 5 rows. A bracket above the first four columns is labeled 'Features', and a bracket above the last column is labeled 'Label'. A bracket to the left of the rows is labeled 'Rows', and a bracket below the columns is labeled 'Columns'.

Size	Beds	Baths	Zip	Price
1100	1	1	64576	1.29
1900	3	1.5	78321	2.14
2800	3	3	98712	3.10
3400	4	3.5	25721	3.75

# Features and Labels

- The data for **Unsupervised Learning** consists of features but no labels because the model is being used to identify patterns not to forecast something.
- For example we could use an unsupervised model to classify the houses that exist in a certain neighborhood without trying to predict any price.



# Feature Scaling

- Before using many ML algorithms (including those for unsupervised learning), it is important to scale feature values so that they are comparable.
- Z-score scaling involves calculating the mean ( $\mu$ ) and SD ( $\sigma$ ) from the values of each feature from the training set. Scaled feature values for all data sets are then created by subtracting the mean and dividing by the SD.
- The scaled feature values have a mean of zero and SD of one.

$$\text{scaled feature value} = \frac{V - \mu}{\sigma} \quad (1)$$

# Feature Scaling

- Min-max scaling involves calculating the maximum and minimum value of each feature from the training set.
- Scaled feature values for all data sets are then created by subtracting the minimum and dividing by the difference between the maximum and minimum.
- The scaled feature values lie between zero and one:

$$\text{scaled feature value} = \frac{V - V_{min}}{V_{max} - V_{min}} \quad (2)$$



# Cleaning data

- Dealing with inconsistent recording
- Removing unwanted observations
- Removing duplicates
- Investigating outliers
- Dealing with missing items

# Cost Function

- In Machine Learning a **Cost Function** or loss function is used to represent how far away a mathematical model is from the real data;
- One adjust the mathematical model usually by varying parameters within the model so as to minimize the const function;
- Let's take for example a very simple model of the form

$$y = \theta_0 + \theta_1 x$$

where the  $\theta$ s are the parameters that we want to find to give us the best fit to the data;

- Call this function  $h_{\theta}(x)$  to emphasize the dependence on both the variable  $x$  and the two parameters  $\theta_0$  and  $\theta_1$ ;

# Cost Function

- We want to measure how far away the data, the  $y^{(n)}$ s are from the function  $h_{\theta}(x)$ ;
- A common way to do this is via the quadratic cost function

$$J(\theta) = \frac{1}{2N} \sum_{n=1}^N \left[ h_{\theta} \left( x^{(n)} \right) - y^{(n)} \right]^2 \quad (3)$$

- We want the parameters that minimize (3), almost always you are going to have to do this numerically;
- If we have a nice convex function then there is a numerical method that will converge to the solution, it is called **gradient descent**.

# Gradient Descent

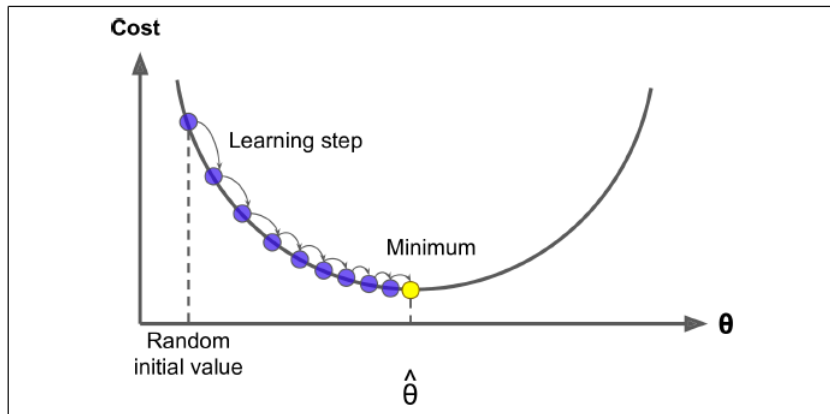
## Linear Regression

- Gradient Descent is a very generic optimization algorithm capable of finding optimal solutions to a wide range of problems.
- Gradient Descent Algorithm measures the local gradient of the error function with regards to the parameter vector  $\theta$ , and it goes in the direction of descending gradient.
- Once the gradient is zero, you have reached a minimum!
- Concretely, you start by filling  $\theta$  with random values (this is called random initialization), and then you improve it gradually, taking one step at a time, each step attempting to decrease the cost function (e.g., the MSE), until the algorithm converges to a minimum

# Gradient Descent

## Linear Regression

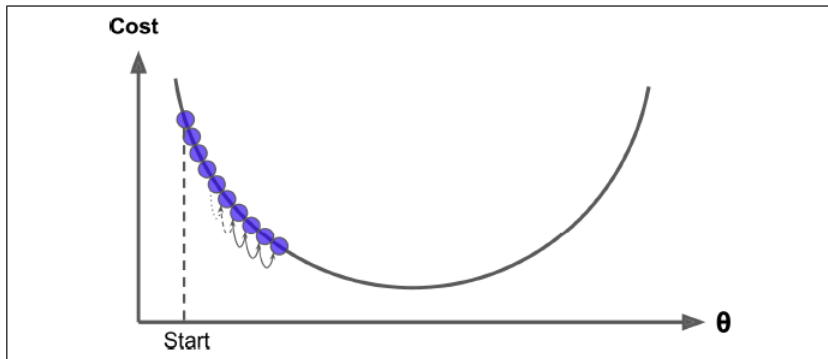
An important parameter in Gradient Descent is the size of the steps, determined by the learning rate hyperparameter.



# Gradient Descent

## Linear Regression

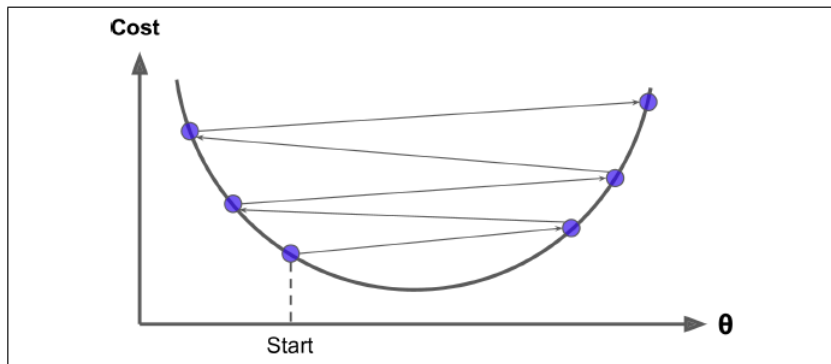
If the learning rate is too small, then the algorithm will have to go through many iterations to converge, which will take a long time...



# Gradient Descent

## Linear Regression

... on the other hand, if the learning rate is too high, you might jump across the valley. This might make the algorithm diverge failing to find a good solution.



# Gradient Descent

Having defined an appropriate learning rate, the scheme works as follow

- Start with an initial guess for each parameter  $\theta_k$ ;
- Move  $\theta_k$  in the direction of the slope

$$\text{New } \theta_k = \text{Old } \theta_k - \beta \frac{\partial J}{\partial \theta_k}$$

where  $\beta$  is our learning rate;

In the above description of gradient descent we have used all of the data points simultaneously. This is called **batch gradient descent**



# Gradient Descent

- Rather than use all of the data in the parameter updating we can use a technique called **stochastic gradient descent**;
- This technique is the same as the batch gradient descent except that you only update using **one** of the data points each time and this data point is chosen **randomly**;
- Especially in high-dimensional optimization problems this reduces the computational burden, achieving faster iterations in trade for a lower convergence rate;
- When the learning rates decrease with an appropriate rate, and subject to relatively mild assumptions, stochastic gradient descent converges almost surely to a global minimum when the objective function is convex or pseudoconvex, and otherwise converges almost surely to a local minimum;

# Validation and Testing

- When data is used for forecasting there is a danger that the machine learning model will work very well for data, but will not generalize well to other data;
- An obvious point is that it is important that the data used in a machine learning model be representative of the situations to which the model is to be applied;
- It is also important to test a model **out-of-sample**, by this we mean that the model should be tested on data that is different from the sample data used to determine the parameters of the model;
- Data scientist refer to the sample data as the **training set** and the data used to determine the accuracy of the model as the **test set**;
- Often a **validation set** is used as well as we explain later;

# Validation and Testing

- We will illustrate the use of a training set and the test data set with a very simple Example (Hull, Chapter 1);
- We suppose that we are interested in forecasting the salaries of people from their age;
- This simple regression model is an example of supervised learning...

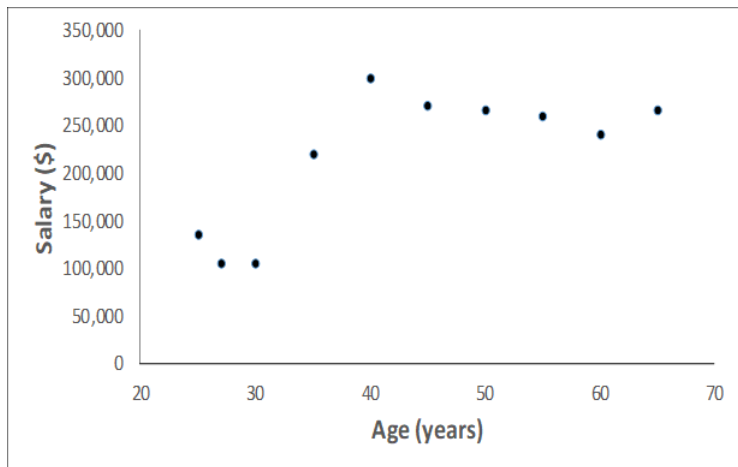
# Validation and Testing: Training Set

**Table 1.** Salary as a function of Age for a certain profession in a certain area)

Age (years)	Salary (\$)
25	135,000
55	260,000
27	105,000
35	220,000
60	240,000
65	265,000
45	270,000
40	300,000
50	265,000
30	105,000

# Validation and Testing: Training Set

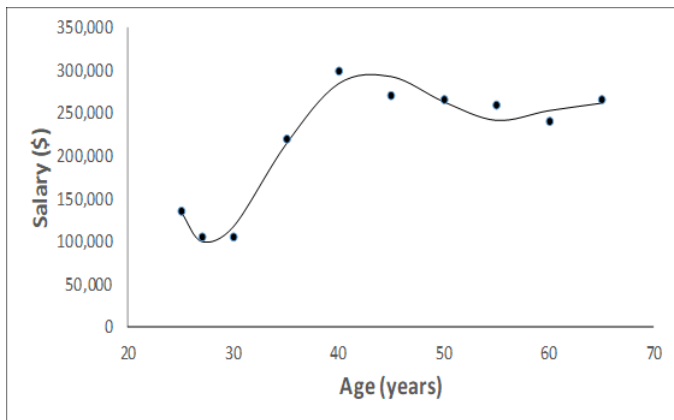
**Figure 1.** Scatter plot of Salary as a function of Age (see Table 1)



## Validation and Testing: Training Set

**Figure 2.** It's tempting to choose a model that fits the data really well for example with a polynomial of degree five ( $Y = \text{Salary}$ ,  $X = \text{Age}$ ):

$$Y = a + b_1X + b_2X^2 + b_3X^3 + b_4X^4 + b_5X^5$$



# Validation and Testing: Discussion of Training Result

- The model provides a good fit to the data;
- The standard deviation of the difference between the salary given by the model and the actual salary for the ten individuals in the training data set (which is referred to as the **root mean square error (rmse)**) is \$12902;
- However common sense would suggest that we may over-fitted the data;
- We need to check the model out-of-sample;
- To use the language of *data science* we need to determine whether the model generalizes well to a new data set that is different from the training set.

# Validation and Testing: Validation Set

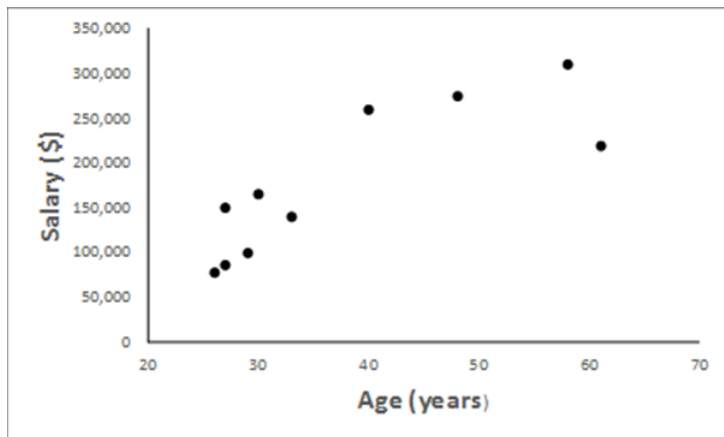
**Table 2.** An Out-of-Sample Validation Set

Age (years)	Salary (\$)
30	166,000
26	78,000
58	310,000
29	100,000
40	260,000
27	150,000
33	140,000
61	220,000
27	86,000
48	276,000



# Validation and Testing: Validation Set

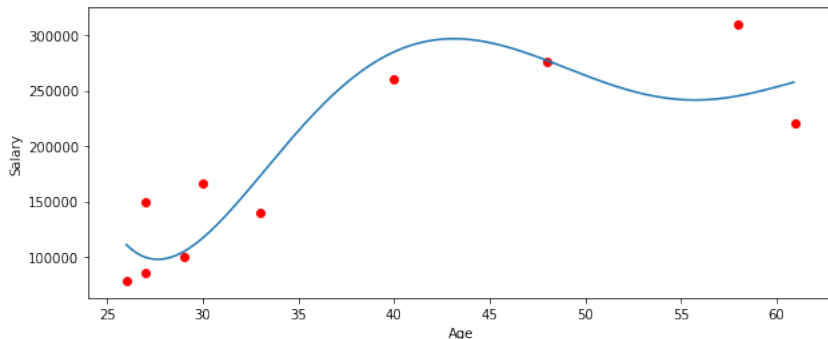
**Figure 3.** Scatter Plot for Validation Set



# Validation and Testing: Discussion of Validation Result

## The Fifth Order Polynomial Model Does Not Generalize Well

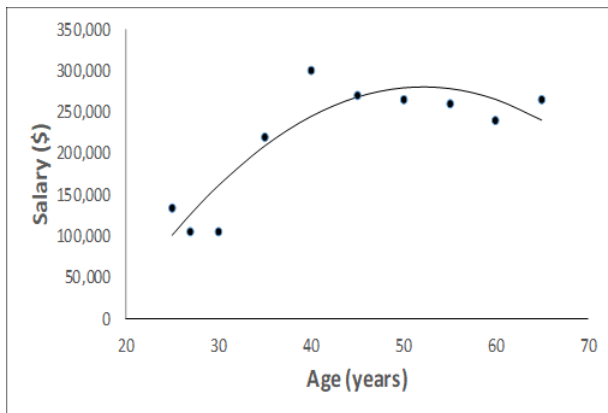
- The root mean squared error (rmse) for the training data set is \$12,902
- The rmse for the test data set is \$38,794
- We conclude that the model overfits the data



# Validation and Testing

**Figure 4.** A Simpler Quadratic Model

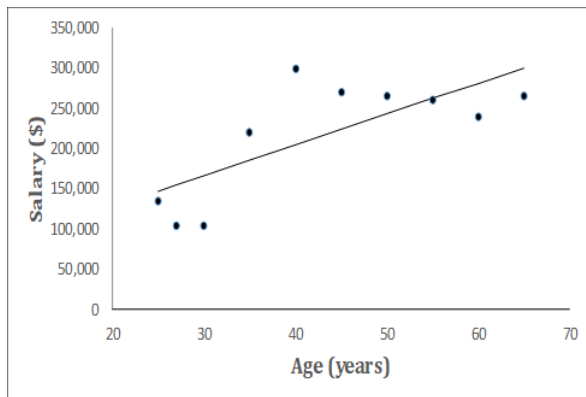
$$Y = a + b_1X + b_2X^2$$



# Validation and Testing

**Figure 5 . Linear Model**

$$Y = a + b_1X$$



# Validation and Testing

**Table 3.** Summary of Results: The linear model under-fits while the 5th degree polynomial over-fits:

	Polynomial of degree 5	Quadratic model	Linear model
Training set	12, 902	32,932	49,731
Validation set	38,794	33,554	49,990

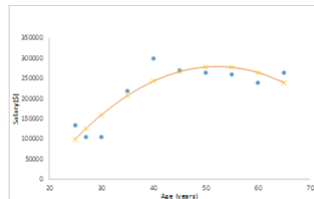
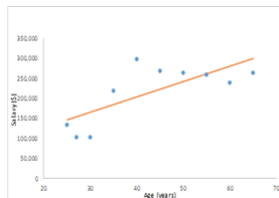
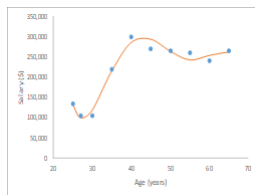
# Validation and Testing

**Table 3.** Summary of Results: The linear model under-fits while the 5th degree polynomial over-fits:

	Polynomial of degree 5	Quadratic model	Linear model
Training set	12,902	32,932	49,731
Validation set	38,794	33,554	49,990

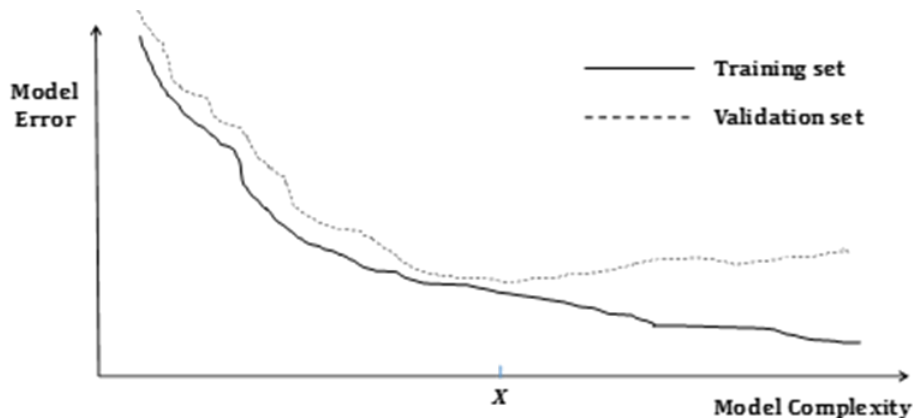
# Validation and Testing

**Figure 6.** Overfitting/Underfitting Example: predicting salaries for people in a certain profession in a certain area (only 10 observations)



Overfitting ————— Underfitting ————— Best model?

# Typical Pattern of Errors for Training Set and Validation Set





# Validation and Testing

## ML Good Practice

- Divide data into three sets
- Training set
- Validation set
- Test set
- Develop different models using the **training set** and compare them using the **validation set**;
- Rule of thumb: increase model complexity until model no longer generalizes well to the validation set;
- The **test set** is used to provide a final out-of-sample indication of how well the chosen model works;

# Bias and Variance

- Suppose there is a relationship between an independent variable  $x$  and a dependent variable  $y$ :

$$y = f(x) + \epsilon \quad (4)$$

where  $\epsilon$  is an error term with mean zero and variance  $\sigma^2$ .

- The error term captures either genuine randomness in the data or noise due to measurement error.
- Suppose we find, with a Machine Learning technique, a deterministic model for this relationship:

$$y = \hat{f}(x) \quad (5)$$

# Bias and Variance

- Now it comes a new data point  $x'$  not in the training set and we want to predict the corresponding  $y'$ ;
- The error we will observe in our model at point  $x'$  is going to be

$$\hat{f}(x') - f(x') - \epsilon \quad (6)$$

- There are two different sources of error in this equation.
- The first one is included in the factor  $\epsilon$ ;
- The second one, more interesting, is due to what is in our training set.
- A robust model should give us the same prediction whatever data we used for training our model.

# Bias and Variance

- Let's look at the average error:

$$E \left[ \hat{f}(x') \right] - f(x') \quad (7)$$

where the expectation is taken over random samples of training data (having the same distribution as the training data).

- This is the definition of the **bias**

$$\text{Bias} \left[ \hat{f}(x') \right] = E \left[ \hat{f}(x') \right] - f(x') \quad (8)$$

# Bias and Variance

- We can also look at the mean square error

$$E \left[ \left( \hat{f}(x') - f(x') - \epsilon \right)^2 \right] = \left[ \text{Bias} \left( \hat{f}(x') \right) \right]^2 + \text{Var} \left[ \hat{f}(x') \right] + \sigma^2$$

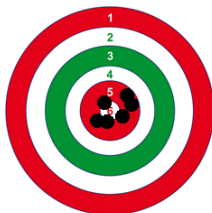
where we remember that  $\hat{f}(x')$  and  $\epsilon$  are independent.

- This show us that there are two important quantities, the **bias** and the **variance** that will affect our results and that we can control to some extent.

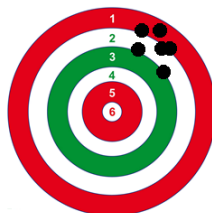
# Bias and Variance

- **What is Bias?** It's the difference between the average prediction of our model and the correct value which we are trying to predict. Model with high bias pays very little attention to the training data and oversimplifies the model. It always leads to high error on training and test data.
- **What is Variance?** It's the variability of model prediction for a given data point or a value which tells us spread of our data. Model with high variance pays a lot of attention to training data and does not generalize on the data which it hasn't seen before. As a result, such models perform very well on training data but has high error rates on test data.

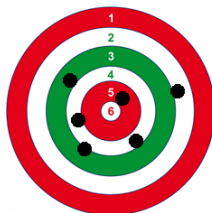
# Bias and Variance



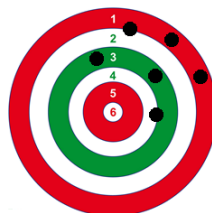
*Low Bias and Low Variance*



*High Bias and Low Variance*



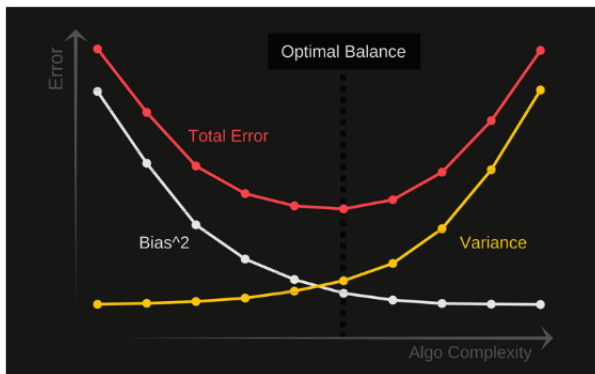
*Low Bias and High Variance*



*High Bias and High Variance*

# Bias and Variance

Unfortunately, we often find that there is a trade-off between bias and variance. As one is reduced, the other is increased. This is the matter of over- and under-fitting.





# Regularization

- Linear regression can over-fit, particularly when there are a large number of correlated features.
- Results for validation set may not then be as good as for training set
- Regularization is a way of avoiding overfitting and reducing the number of features. Alternatives:
  - Ridge
  - Lasso
  - Elastic net
- We must first scale feature values

## Ridge regression (analytic solution)

- Ridge regression is a regularization technique where we change the function that is to be minimize;
- Reduce magnitude of regression coefficients by choosing a parameter  $\lambda$  and minimizing

$$\frac{1}{2N} \sum_{n=1}^N \left[ h_{\theta} \left( x^{(n)} \right) - y^{(n)} \right]^2 + \lambda \sum_{n=1}^N b_i^2 \quad (9)$$

- This change has the effect of encouraging the model to keep the weights  $b_j$  as small as possible;
- The Ridge regression should only be used for determining model parameters using the training set. Once the model parameters have been determined the penalty term should be removed for prediction;
- What happens as  $\lambda$  increases?

# Lasso Regression (must use gradient descent)

- Lasso is short for *Least Absolute Shrinkage and Selection Operator*;
- Similar to ridge regression except we minimize

$$\frac{1}{2N} \sum_{n=1}^N \left[ h_{\theta} \left( x^{(n)} \right) - y^{(n)} \right]^2 + \lambda \sum_{n=1}^N |b_n| \quad (10)$$

- This function cannot be minimized analytically and so a variation on the gradient descent algorithm must be used;
- Lasso regression also has the effect of simplifying the model. It does this by setting the weights of unimportant features to zero. When there are a large number of features, Lasso can identify a relatively small subset of the features that form a good predictive model.

# Elastic Net Regression (must use gradient descent)

- Middle ground between Ridge and Lasso
- Minimize

$$\frac{1}{2N} \sum_{n=1}^N \left[ h_{\theta} \left( x^{(n)} \right) - y^{(n)} \right]^2 + \lambda_1 \sum_{n=1}^N b_n^2 + \lambda_2 \sum_{n=1}^N |b_n| \quad (11)$$

- In Lasso some weights are reduced to zero but others may be quite large. In Ridge, weights are small in magnitude but they are not reduced to zero. The idea underlying Elastic Net is that we may be able to get the best of both by making some weights zero while reducing the magnitude of the others.

# Salary Vs Age Example: The Effect of Regularization

Age (years)	Salary (\$)
25	135,000
55	260,000
27	105,000
35	220,000
60	240,000
65	265,000
45	270,000
40	300,000
50	265,000
30	105,000

We apply regularization to the model:

$$Y = a + b_1X + b_2X^2 + b_3X^3 + b_4X^4 + b_5X^5 \quad (12)$$

where  $Y$  is salary and  $X$  is age.

# Salary Vs Age Example: The Effect of Regularization

Data with Z-score scaling

Observ.	X	$X^2$	$X^3$	$X^4$	$X^5$
1	-1.290	-1.128	-0.988	-0.874	-0.782
2	0.836	0.778	0.693	0.592	0.486
3	-1.148	-1.046	-0.943	-0.850	-0.770
4	-0.581	-0.652	-0.684	-0.688	-0.672
5	1.191	1.235	1.247	1.230	1.191
6	1.545	1.731	1.901	2.048	2.174
7	0.128	-0.016	-0.146	-0.253	-0.333
8	-0.227	-0.354	-0.449	-0.511	-0.544
9	0.482	0.361	0.232	0.107	-0.004
10	-0.936	-0.910	-0.861	-0.803	-0.745

# Salary Vs Age Example: The Effect of Regularization

Ridge Results,  $\lambda = 0.02$  is similar to quadratic model

$\lambda$	$a$	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$
0	216.5	-32,623	135,403	-215,493	155,315	-42,559
0.02	216.5	97.8	36.6	-8.5	35.0	-44.6
0.10	216.5	56.5	28.1	3.7	-15.1	-28.4

# Salary Vs Age Example: The Effect of Regularization

Lasso Results,  $\lambda = 1$  is similar to the quadratic model

$\lambda$	$a$	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$
0	216.5	-32,623	135,403	-215,493	155,315	-42,559
0.02	216.5	-646.4	2,046.6	0.0	-3,351.0	2,007.9
0.1	216.5	355.4	0.0	-494.8	0.0	196.5
1	216.5	147.4	0.0	0.0	-99.3	0.0



# Let's code ...

