

Lecture Note - 01: Introduction to Data Science: Toy Problem, Linear Algebra, Randomness

Dihui Lai

March 29, 2020

Contents

1	Data Science - Toy Problem	2
1.1	Toy Problem	2
1.2	Price Estimation	2
2	Structured Data and Linear Model	3
2.1	Tabular Data and Matrix	3
2.2	Linear Model	3
2.3	Other Models	3
3	Geometric Interpretation, Visualization and Randomness	4
3.1	Pants, Socks and Cost from 100 Friends	4
3.2	Incomplete Data and Visualization	4
3.3	High Dimension Space Projection and Randomness	5

1 Data Science - Toy Problem

1.1 Toy Problem

Suppose you learned from 3 of your friends who went on shopping recently, who bought pants and socks. The number and costs are shown as below:

	Pants	Socks	Cost
John	1	2	23
Lisa	1	2	26
David	1	1	24

According to John and Lisa, the prices of a pant and a sock can be calculated as $P = 20$ and $S = 3$, respectively. However, David should have paid 23 dollar given the inferred prices. Why did David pay 24 dollar instead of 23? It could be due to price variation.

1.2 Price Estimation

To get a good estimation of the prices of socks and pants, we can use the following error function

$$\epsilon = (P + S - 23)^2 + (P + 2S - 26)^2 + (P + S - 24)^2$$

Ideally, we would like to have our estimated socks (S) and pants (P) price as close to the real cost, i.e. minimize ϵ . Use basic calculus knowledge, we know the optimal P and S should satisfy the following equations.

$$\begin{cases} \frac{\partial \epsilon}{\partial P} = 0 \\ \frac{\partial \epsilon}{\partial S} = 0 \end{cases} \quad (1)$$

$$\begin{cases} \frac{\partial \epsilon}{\partial P} = 2(P + S - 23) + 2(P + 2S - 26) + 2(P + S - 24) = 0 \\ \frac{\partial \epsilon}{\partial S} = 2(P + S - 23) + 4(P + 2S - 26) + 2(P + S - 24) = 0 \end{cases} \implies \begin{cases} 9P + 12S - 219 = 0 \\ 8P + 12S - 198 = 0 \end{cases} \quad (2)$$

$$\begin{cases} P = 21 \\ S = 2.5 \end{cases} \quad (3)$$

2 Structured Data and Linear Model

2.1 Tabular Data and Matrix

In general, if we want to consider a model of m types of goods and collect data from n people. The toy model can be generalized to a problem that needs to estimate m variables on n data points

$$\begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{bmatrix} \sim \begin{bmatrix} x_1^1 & x_2^1 & x_1^2 & \dots & x_m^1 \\ x_1^2 & x_2^2 & x_1^2 & \dots & x_m^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^n & x_2^n & x_1^n & \dots & x_m^n \end{bmatrix}$$

Using vector notation, we have

$$Y = \vec{y} = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{bmatrix}$$

,

$$X = [\vec{x}_1, \vec{x}_2, \dots, \vec{x}_m] = \begin{bmatrix} x_1^1 & x_2^1 & x_1^2 & \dots & x_m^1 \\ x_1^2 & x_2^2 & x_1^2 & \dots & x_m^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^n & x_2^n & x_1^n & \dots & x_m^n \end{bmatrix}$$

The vectors \vec{x}_i are called covariates, or predictors. \vec{y} is normally called target variable

2.2 Linear Model

If we assume \vec{y} is linearly dependent on \vec{x} s, we have a linear model

$$\hat{\vec{y}} = \beta_1 \vec{x}_1 + \beta_2 \vec{x}_2 + \dots + \beta_m \vec{x}_m$$

An optimal model should estimate $\hat{\vec{y}}$ as close as \vec{y} . e.g.

$$\epsilon = (\hat{\vec{y}} - \vec{y}) \cdot (\hat{\vec{y}} - \vec{y})$$
$$\frac{\partial \epsilon}{\partial \beta_i} = 0, i=1, 2, 3, \dots, m$$

How can we solve the problem?

2.3 Other Models

In general, \vec{y} could be any function of \vec{x} i.e. $\vec{y} = f(\vec{x})$.

- Kepler's Law: $T^2 \sim r^3$. Note Kepler's law becomes linear if we do a log transformation on both side i.e. $2\log T = 3\log r$
- House price: $P \sim f(\text{size}, \text{location})$

3 Geometric Interpretation, Visualization and Randomness

3.1 Pants, Socks and Cost from 100 Friends

Suppose now you collect data from 100 friends. Everyone of them has bought a few pants and some socks from the same store. In a perfect world, everyone remember the number of pants/socks they bought and the corresponding cost. You end up having a perfect data like this:

	Socks	Pants	Cost
Person1	3	6	129
Person2	8	6	144
Person3	8	5	124
Person4	8	3	84
Person5	3	6	129
...
Person100

By looking at the number closely, you figured out that the price of a pant is 20 dollars and the price is a sock is 3 dollars. And it is consistent across the whole data set.

3.2 Incomplete Data and Visualization

In reality it is very hard to get all information you need to build a pricing model for pants and socks, some people will not tell you the cost and some people can not tell you the number of pants/socks that they bought. Let us assume that all your friends can only tell you the number of pants that they bought and the total cost of their purchase. So you end up having a data set like this:

	Pants	Cost
Person1	6	129
Person2	6	144
Person3	5	124
Person4	3	84
Person5	6	129
...
Person100

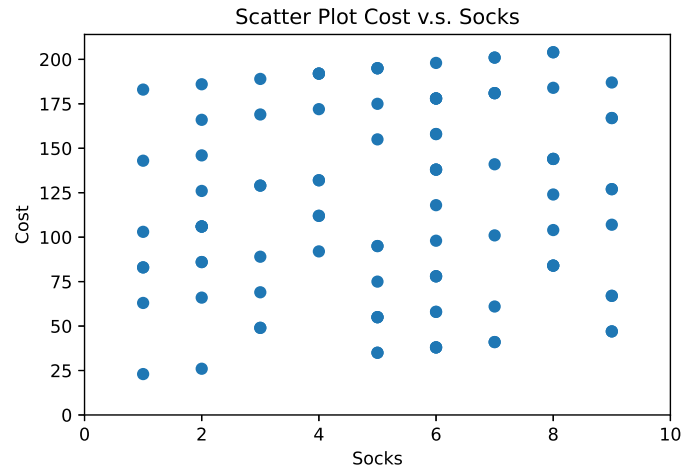


Figure 2: Scatter plot of total cost against the number of socks bought

While dealing with real world data problem, missing information is almost guaranteed. For example, while modeling the house price, it is unlikely that we will know the price that the buyers are willing to pay; while modeling a public company's stock price, it is almost impossible to know all information related to the company. **Therefore, we have to make good assumptions about the information that we do not have a.k.a noise.**