

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH  
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN  
**KHOA CÔNG NGHỆ THÔNG TIN**



BÁO CÁO ĐỒ ÁN CUỐI KỲ MÔN HỌC  
**TOÁN ỨNG DỤNG VÀ THỐNG KÊ**

**ĐỒ ÁN: DATA FITTING**

**Giáo viên hướng dẫn**

Thầy Nguyễn Hữu Toàn

Cô Võ Nam Thực Đoàn

**Lớp**

21\_2

**Sinh viên thực hiện**

21120150 – Nguyễn Song Toàn

Tp. Hồ Chí Minh, tháng 06 năm 2023

**Lời cảm ơn**

Em xin cảm ơn thầy Nguyễn Hữu Toàn đã cung cấp cho chúng em những kiến thức bổ ích, có ý nghĩa thực tiễn để em có thể thực hiện đồ án này. Em cũng xin cảm ơn cô Võ Nam Thục Đoan đã cung cấp những bài thực hành để em vận dụng được những lý thuyết đã học và có những kỹ thuật lập trình phục vụ cho việc thực hiện các đồ án.

## **Mục lục**

<b>Mục lục</b> .....	2
<b>Danh mục hình</b> .....	3
<b>A. Thông tin đồ án</b> .....	4
1. Đề bài.....	4
2. Đánh giá mức độ hoàn thành .....	4
<b>B. Cơ sở lý thuyết</b> .....	4
1. Bài toán bình phương tối thiểu (Least square) .....	4
2. Mô hình hồi quy tuyến tính (Linear regression model) .....	4
<b>C. Mô hình sử dụng ở câu a và câu b</b> .....	5
1. Câu a - mô hình đánh giá giá nhà từ các yếu tố tác động.....	5
2. Câu b - mô hình đánh giá lương nhân viên từ các yếu tố tác động .....	5
<b>D. Mô tả ý nghĩa của các hàm</b> .....	5
<b>E. Kết quả đồ án</b> .....	6
<b>Tài liệu tham khảo</b> .....	8

### **Danh mục hình**

HÌNH 1 - MÔ HÌNH ĐÁNH GIÁ GIÁ NHÀ .....	6
HÌNH 2 - MÔ HÌNH ĐÁNH GIÁ LƯƠNG NHÂN VIÊN .....	7

## A. Thông tin đề án

### 1. Đề bài

- Sử dụng bài toán data fitting trong xây dựng mô hình đánh giá lương nhân viên từ các yếu tố tác động theo dữ liệu được cung cấp.
- Sử dụng bài toán data fitting trong việc xây dựng mô hình đánh giá giá nhà từ các yếu tố tác động từ dữ liệu được cung cấp.

### 2. Đánh giá mức độ hoàn thành

Nội dung	Mức độ hoàn thành
Xây dựng mô hình đánh giá lương nhân viên	100%
Xây dựng mô hình đánh giá giá nhà	100%

## B. Cơ sở lý thuyết

### 1. Bài toán bình phương tối thiểu (Least square)

Xét ma trận  $A$  có kích thước  $m \times n$  ( $m > n$ ) và vector  $b$  có kích thước  $m$ . Bình phương tối thiểu  $\hat{x}$  của  $A$  và  $b$  trên giá trị  $x$  được tính như sau:

$$\hat{x} = \|Ax - b\|^2$$
$$\hat{x} = (A^T A)^{-1} A^T b$$

### 2. Mô hình hồi quy tuyến tính (Linear regression model)

- Nghiệm  $\hat{x}$  của phương trình hồi quy tuyến tính  $Ax = b$  được tính như sau:

$$\hat{x} = (A^T A)^{-1} A^T b$$

- Cách xác định các tham số của mô hình hồi quy tuyến tính:

Cho  $y$  là biến phụ thuộc,  $x_1, x_2, x_3, \dots, x_n$  là biến độc lập.

Tìm 1 hàm số biểu diễn  $y = f(x_1, x_2, \dots, x_n)$

Gọi  $y = b_1x_1 + b_2x_2 + \dots + b_nx_n$

Ta có:  $RSS(\theta) = \|\hat{y}_i - y_i\|^2$

Trong đó:

- $\hat{y}$ : giá trị được tính bằng cách thay biến độc lập vào hàm số
- $y$ : giá trị thực biến phụ thuộc

$$RSS(\theta) = \|A\theta - Y\|^2$$

$$A = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{21} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1n} & x_{2n} & \dots & x_{nn} \end{pmatrix}, \theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{pmatrix}, Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

Tìm  $\theta$  sao cho  $\|A\theta - Y\|^2$  min

Khi đó:  $\theta = (A^T A)^{-1} (A^T Y)$

Từ đó, xác định các tham số của mô hình hồi quy tuyến tính.

- Cách tính chuẩn vector phần dư:

Sau khi tìm được các tham số thì tính:

$$\hat{Y} = \begin{pmatrix} Y(x_1) \\ Y(x_2) \\ \vdots \\ Y(x_n) \end{pmatrix} \Rightarrow r = \hat{Y} - Y = \begin{pmatrix} Y(x_1) - y_1 \\ Y(x_2) - y_2 \\ \vdots \\ Y(x_n) - y_n \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{pmatrix}$$

$$\Rightarrow \|r\| = \sqrt{r_1^2 + r_2^2 + \dots + r_n^2}$$

### C. Mô hình sử dụng ở câu a và câu b

#### 1. Câu a - mô hình đánh giá giá nhà từ các yếu tố tác động

$$\ln(Y) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

- $x_1$  đại diện cho đặc trưng sqft / 100
- $x_2$  đại diện cho đặc trưng Bedrooms
- $x_3$  đại diện cho đặc trưng Baths
- $x_4$  đại diện cho đặc trưng Age

#### 2. Câu b - mô hình đánh giá lương nhân viên từ các yếu tố tác động

$$\ln(Y) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

- $x_1$  đại diện cho đặc trưng educ
- $x_2$  đại diện cho đặc trưng exper
- $x_3$  đại diện cho đặc trưng hrswk

### D. Mô tả ý nghĩa của các hàm

- Hàm **create\_identity\_matrix(n)**: tạo ra ma trận đơn vị  $I_n$ .
- Hàm **multiply\_matrix\_vector(matrix, vector)**: trả ra kết quả của nhân ma trận với vector.
- Hàm **multiply\_two\_matrixes(A, B)**: trả ra kết quả của nhân 2 ma trận.
- Hàm **transpose\_matrix(A)**: trả ra ma trận chuyển vị của ma trận A.
- Hàm **swap\_two\_rows(matrix, r1, r2)**: hoán vị 2 dòng của ma trận với tham số là ma trận matrix, chỉ số của 2 dòng r1 và r2.
- Hàm **inverse(A)**: trả ra ma trận nghịch đảo của ma trận A (nếu có).
- Hàm **subtract\_two\_vectors(A, B)**: trả ra kết quả của trừ 2 vector A và B.
- Hàm **create\_y\_A(data)**: tạo ra vector y và ma trận A dựa trên bộ dữ liệu data được đọc từ file dùng cho câu a.

- Hàm này sử dụng vòng lặp để tạo ra vector  $y$  mang giá trị  $\ln(price)$  và ma trận  $A$  gồm cột đầu tiên chứa toàn số 1 và các cột tiếp theo lần lượt là cột  $sqft/100$ ,  $Bedrooms$ ,  $Baths$ ,  $Age$ , các cột này là giá trị của các cột trong dataframe mà đọc vào từ file `br2.csv`.
- Hàm **vectorTheta(A, y)**: xây dựng vector các tham số của mô hình dựa trên cơ sở lý thuyết đã học.
  - Hàm này sử dụng các hàm hỗ trợ tính toán ma trận, xây dựng ma trận nghịch đảo như: **multiply\_two\_matrixes(A, B)**, **multiply\_matrix\_vector(matrix, vector)**, **inverse(A)** để thực hiện tìm ra vector tham số theo công thức như sau:
 
$$\theta = (A^T A)^{-1} (A^T Y)$$
- Hàm **getModel\_a(Theta, x)**: trả về chuỗi là hàm số của mô hình khi có danh sách biến  $x$  và vector tham số  $\Theta$ .
  - Hàm này sử dụng vòng lặp và xử lý chuỗi để có được chuỗi theo định dạng như sau: “ $Y^* = \ln Y = \theta_0 + \theta_1 * x_1 + \theta_2 * x_2 + \theta_3 * x_3$ ” với  $\theta_i$  là các giá trị tham số cụ thể.
- Hàm **exportModel\_a(A, y, n)**: xây dựng danh sách biến dựa trên số lượng đặc trưng  $n$ , xây dựng vector  $\Theta$  từ ma trận  $A$  và vector  $y$  và trả ra hàm số của mô hình khi có được danh sách biến và danh sách tham số  $\Theta$ .
- Hàm **vectorYhat(A, Theta)**: tạo ra vector  $\hat{Y}$  từ ma trận  $A$  và vector  $\Theta$  để phục vụ cho việc tính sai số của mô hình (tính chuẩn).
- Hàm **get\_error(y, A)**: trả ra chuẩn vector phần dư của mô hình.
- Hàm **create\_y\_A\_b(data)**: tạo ma trận  $A$  và vector  $y$  cho câu b.
  - Hàm này sử dụng vòng lặp để tạo ra vector  $y$  mang giá trị  $\ln(wage)$  và ma trận  $A$  gồm cột đầu tiên chứa toàn số 1 và các cột tiếp theo lần lượt là cột  $educ$ ,  $exper$ ,  $hrswk$ , các cột này là giá trị của các cột trong dataframe mà đọc vào từ file `cps4_small.csv`.
- Trong đồ án này, em sử dụng thư viện *pandas* để đọc dữ liệu từ file `.csv` và dùng thư viện *math* để tính log.

## E. Kết quả đồ án

### - Mô hình đánh giá giá nhà từ các yếu tố tác động:

Mô hình đánh giá giá nhà:  
 $Y^* = \ln Y = 10.91896 + 0.03308 * x_1 - 0.05895 * x_2 + 0.21457 * x_3 - 0.00660 * x_4$   
 Chuẩn vector phần dư là: 9.12

Hình 1 - Mô hình đánh giá giá nhà

**- Mô hình đánh giá lương nhân viên từ các yếu tố tác động:**

Mô hình đánh giá lương nhân viên:

$$Y^* = \ln Y = 1.10054 + 0.09031 * x_1 + 0.00578 * x_2 + 0.00894 * x_3$$

Chuẩn vector phần dư là: 16.21

Hình 2 - Mô hình đánh giá lương nhân viên



### **Tài liệu tham khảo**

- [1] Slide Toán ứng dụng và thống kê, trường đại học Khoa học tự nhiên, ĐHQG-HCM
- [2] Tài liệu hướng dẫn lab Linear Regression, trường đại học Khoa học tự nhiên, ĐHQG-HCM