

**The Pennsylvania State University**  
**The Graduate School**

**DICTIONARY LEARNING FOR SIGNAL CLASSIFICATION**  
**UNDER SPARSITY CONSTRAINTS**

A Ph.D. Dissertation Proposal in  
Electrical Engineering  
by  
Tiep Huu Vu

Submitted in Partial Fulfillment  
of the Requirements  
for the Degree of

Doctor of Philosophy

November 2015

## **Thesis Committee**

Vishal Monga  
Assistant Professor of Electrical Engineering  
Thesis Advisor, Chair of Committee

William E. Higgins  
Distinguished Professor of Electrical Engineering

Kenneth Jenkins  
Professor of Electrical Engineering

Robert T. Collins  
Associate Professor of Computer Science and Engineering

# Contents

<b>1</b>		
	<b>Linear Regression</b>	<b>1</b>
1.1	1. Giới thiệu . . . . .	1
1.2	2. Phân tích toán học . . . . .	2
	1.2.1 Dạng của Linear Regression . . . . .	2
	1.2.2 Sai số dự đoán . . . . .	2
	1.2.3 Hàm mất mát . . . . .	3

# Chapter 1

## Linear Regression

### 1.1 1. Giới thiệu

Quay lại ví dụ đơn giản được nêu trong bài trước: một căn rộng  $x_1\text{m}^2$ , có  $x_2$  phòng ngủ và cách trung tâm  $x_3\text{km}$  có giá là bao nhiêu. Giả sử chúng ta đã có số liệu thống kê từ 1000 căn nhà trong thành phố đó, liệu rằng khi có một căn nhà mới với các thông số về diện tích, số phòng ngủ và khoảng cách tới trung tâm, chúng ta có thể dự đoán được giá của căn nhà đó không? Nếu có thì hàm dự đoán  $y = f(\mathbf{x})$  sẽ có dạng như thế nào. Ở đây  $\mathbf{x} = [x_1, x_2, x_3]$  là một vector hàng chứa thông tin *input*,  $y$  là một số vô hướng (scalar) biểu diễn *output* (tức giá của căn nhà trong ví dụ này).

**Lưu ý về ký hiệu toán học:** trong các bài viết của tôi, các số vô hướng được biểu diễn bởi các chữ cái viết ở dạng không in đậm, có thể viết hoa, ví dụ  $x_1, N, y, k$ . Các vector được biểu diễn bằng các chữ cái thường in đậm, ví dụ  $\mathbf{y}, \mathbf{x}_1$ . Các ma trận được biểu diễn bởi các chữ viết hoa in đậm, ví dụ  $\mathbf{X}, \mathbf{Y}, \mathbf{W}$ .

Một cách đơn giản nhất, chúng ta có thể thấy rằng: i) diện tích nhà càng lớn thì giá nhà càng cao; ii) số lượng phòng ngủ càng lớn thì giá nhà càng cao; iii) càng xa trung tâm thì giá nhà càng giảm. Một hàm số đơn giản nhất có thể mô tả mối quan hệ giữa giá nhà và 3 đại lượng đầu vào là:

$$y \approx f(\mathbf{x}) = \hat{y} \quad (1.1)$$

$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + w_0 \quad (1) \quad (1.2)$$

trong đó,  $w_1, w_2, w_3, w_0$  là các hằng số,  $w_0$  còn được gọi là bias. Mối quan hệ  $y \approx f(\mathbf{x})$  bên trên là một mối quan hệ tuyến tính (linear). Bài toán chúng ta đang làm là một bài toán thuộc loại regression. Bài toán đi tìm các hệ số tối ưu  $\{w_1, w_2, w_3, w_0\}$  chính vì vậy được gọi là bài toán Linear Regression.

**Chú ý 1:**  $y$  là giá trị thực của *outcome* (dựa trên số liệu thống kê chúng ta có trong tập *training data*), trong khi  $\hat{y}$  là giá trị mà mô hình Linear Regression dự đoán được. Nhìn chung,  $y$  và  $\hat{y}$  là hai giá trị khác nhau do có sai số mô hình, tuy nhiên, chúng ta mong muốn rằng sự khác nhau này rất nhỏ.

**Chú ý 2:** *Linear* hay *tuyến tính* hiểu một cách đơn giản là *thẳng, phẳng*. Trong không gian hai chiều, một hàm số được gọi là *tuyến tính* nếu đồ thị của nó có dạng một *đường thẳng*. Trong không gian ba chiều, một hàm số được gọi là *tuyến tính* nếu đồ thị của nó có dạng một *mặt phẳng*. Trong không gian nhiều hơn 3 chiều, khái niệm *mặt phẳng* không còn phù hợp nữa, thay vào đó, một khái niệm khác ra đời được gọi là *siêu mặt phẳng* (*hyperplane*). Các hàm số tuyến tính là các hàm đơn giản nhất, vì chúng thuận tiện trong việc hình dung và tính toán. Chúng ta sẽ được thấy trong các bài viết sau, *tuyến tính* rất quan trọng và hữu ích trong các bài toán Machine Learning. Kinh nghiệm cá nhân tôi cho thấy, trước khi hiểu được các thuật toán *phi tuyến* (non-linear, không phẳng), chúng ta cần nắm vững các kỹ thuật cho các mô hình *tuyến tính*.

## 1.2 2. Phân tích toán học

### 1.2.1 Dạng của Linear Regression

Trong phương trình (1) phía trên, nếu chúng ta đặt  $\mathbf{w} = [w_0, w_1, w_2, w_3]^T$  là vector (cột) hệ số cần phải tối ưu và  $\bar{\mathbf{x}} = [1, x_1, x_2, x_3]$  (đọc là *x bar* trong tiếng Anh) là vector (hàng) dữ liệu đầu vào *mở rộng*. Số 1 ở đầu được thêm vào để phép tính đơn giản hơn và thuận tiện cho việc tính toán. Khi đó, phương trình (1) có thể được viết lại dưới dạng:

$$y \approx \bar{\mathbf{x}}\mathbf{w} = \hat{y} \quad (1.3)$$

Chú ý rằng  $\bar{\mathbf{x}}$  là một vector hàng. (Xem thêm về ký hiệu vector hàng và cột tại đây)

### 1.2.2 Sai số dự đoán

Chúng ta mong muốn rằng sự sai khác  $e$  giữa giá trị thực  $y$  và giá trị dự đoán  $\hat{y}$  (đọc là *y hat* trong tiếng Anh) là nhỏ nhất. Nói cách khác, chúng ta muốn giá trị sau đây càng

nhỏ càng tốt:

$$\frac{1}{2}e^2 = \frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(y - \bar{\mathbf{x}}\mathbf{w})^2 \quad (1.4)$$

trong đó hệ số  $\frac{1}{2}$  (*lại*) là để thuận tiện cho việc tính toán (khi tính đạo hàm thì số  $\frac{1}{2}$  sẽ bị triệt tiêu). Chúng ta cần  $e^2$  vì  $e = y - \hat{y}$  có thể là một số âm, việc nói  $e$  nhỏ nhất sẽ không đúng vì khi  $e = -\infty$  là rất nhỏ nhưng sự sai lệch là rất lớn. phần sau.

### 1.2.3 Hàm mất mát

Điều tương tự xảy ra với tất cả các cặp (*input, outcome*)  $(\mathbf{x}_i, y_i), i = 1, 2, \dots, N$ , với  $N$  là số lượng dữ liệu quan sát được. Điều chúng ta muốn, tổng sai số là nhỏ nhất, tương đương với việc tìm  $\mathbf{w}$  để hàm số sau đạt giá trị nhỏ nhất:

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \bar{\mathbf{x}}_i \mathbf{w})^2 \quad (1.5)$$

Hàm số  $\mathcal{L}(\mathbf{w})$  được gọi là **hàm mất mát** (loss function) của bài toán Linear Regression. Chúng ta luôn mong muốn rằng sự mất mát (sai số) là nhỏ nhất, điều đó đồng nghĩa với việc tìm vector hệ số  $\mathbf{w}$  sao cho giá trị của hàm mất mát này càng nhỏ càng tốt. Giá trị của  $\mathbf{w}$  làm cho hàm mất mát đạt giá trị nhỏ nhất được gọi là *điểm tối ưu* (optimal point), ký hiệu:

sdfadsf