

annotated  
version

Machine Learning Course - CS-433

# Least Squares

Oct 4, 2016

©Mohammad Emtiyaz Khan 2015

minor changes by Martin Jaggi 2016



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

## Motivation

In rare cases, one can compute the optimum of the cost function analytically. Linear regression using MSE is one such case. Here the solution can be obtained explicitly, from using a linear equation system, which is called the normal equations. This method is one of the most popular methods for data fitting, and called least squares.

To derive the normal equations, we use the optimality conditions for convex functions. See the previous lecture notes on optimization.

$$\underline{\nabla \mathcal{L}(\mathbf{w}^*) \stackrel{!}{=} \mathbf{0}}$$

This is a system of  $D$  equations.

*Exercise:* Derive the normal equation for a 1-parameter linear model, for  $\mathcal{L} = \text{MSE}$ .

find  $w$

$$\text{MSE}(w) = \frac{1}{2N} \sum_{n=1}^N (y_n - x_n^T w)^2$$

1-parameter model

$$\mathcal{L}(w_0) = \frac{1}{N} \sum \frac{1}{2} (y_n - w_0)^2$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_0} &= \frac{1}{N} \sum (y_n - w_0)(-1) \\ &= w_0 - \frac{1}{N} \sum y_n \end{aligned}$$

$$\stackrel{!}{=} 0$$

$$\Rightarrow w_0 = \frac{1}{N} \sum y_n = \bar{y}$$

$$\mathcal{L} = \text{MSE}$$

## Normal Equations

Recall the expression of the gradient for multiple linear regression, under mean squared error:

$$\underline{\nabla \mathcal{L}(\mathbf{w}) = -\frac{1}{N} \mathbf{X}^\top \mathbf{e} = -\frac{1}{N} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w})}$$

$\mathbf{e}$

$$\in \mathbb{R}^D$$

Set it to zero to get the **normal equations** for linear regression.

$D$  equations

$D = \# \text{ parameters}$

$\# \text{ columns of } \mathbf{X}$

$$\mathbf{X}^\top \mathbf{e} = \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) \stackrel{!}{=} \mathbf{0}$$

implying that the error is orthogonal to all rows of  $\mathbf{X}^\top$  (= columns of  $\mathbf{X}$ ).

$\mathcal{L}$  convex :

$$\nabla \mathcal{L} = 0 \stackrel{!}{\Leftrightarrow} \text{optimality}$$

# Geometric Interpretation

Denote the  $d$ 'th column of  $\mathbf{X}$  by  $\mathbf{x}_{:d}$ .

$$y \approx x_w$$

↑  
error

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}, \mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1D} \\ x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix}$$

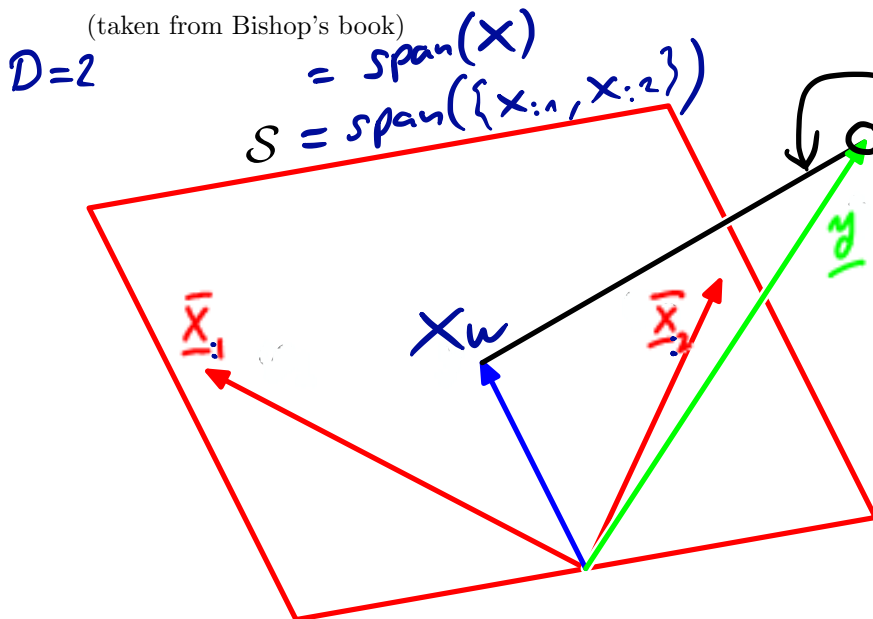
←  $x_1$   
←  $x_n$   
←  $x_N$

The normal equations tell us to choose a vector in the span of  $\mathbf{X}$ , such that the error vector  $\mathbf{e}$  will be orthogonal to the span.

$$\text{span}(\{u, v\}) \subseteq \mathbb{R}^N$$

$$= \{w_1 u + w_2 v \mid w_1 \in \mathbb{R}, w_2 \in \mathbb{R}\}$$

The following figure illustrates this:



Normal Equations

$$\mathbf{X}^T(\mathbf{y} - \mathbf{x}_w) = 0$$

↑  
 $\mathbf{e}$

Q1: what if  $y \in \text{span}(\mathbf{X})$ ?

Q2: what if  $x_{:1} = x_{:2}$ ?

Q3: what if  $y \perp \text{span}(\mathbf{X})$ ?

$$(\Rightarrow \mathbf{x}_w = 0)$$

$$\mathbf{x}_w = w_1 \mathbf{x}_{:1} + w_2 \mathbf{x}_{:2}$$

# Least Squares

When  $\mathbf{X}^\top \mathbf{X}$  is **invertible**, we have a closed-form expression for the minimum.

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

We can use this model to predict a new value for an unseen datapoint (test point)  $\mathbf{x}_m$ :

$$\hat{y}_m := \mathbf{x}_m^\top \mathbf{w}^* = \mathbf{x}_m^\top \boxed{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}} \quad \mathbf{w}^*$$

## Invertibility and Uniqueness

The **Gram matrix**  $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{D \times D}$  is invertible iff  $\mathbf{X}$  has full column **rank**, or in other words  $\text{rank}(\mathbf{X}) = D$ .

*Proof:* Fundamental theorem of linear algebra.

find  $\mathbf{w}$ :

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

$$\mathbf{X}^\top \mathbf{y} = \underbrace{\mathbf{X}^\top \mathbf{X}}_{\mathbb{R}^{D \times D}} \mathbf{w}$$

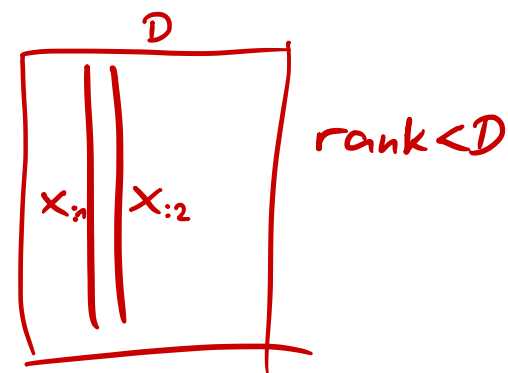
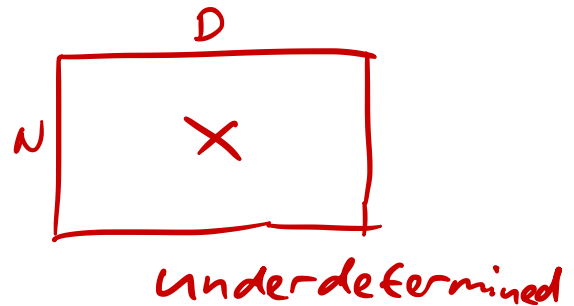
solve a linear system

$$\boxed{\mathbf{X}^\top} \cdot \overset{D}{\boxed{\mathbf{X}}} = \boxed{\phantom{\mathbf{X}^\top \mathbf{X}}} \quad D \times D$$

# Rank Deficiency and Ill-Conditioning

Unfortunately, in practice,  $\mathbf{X}$  is often rank deficient! (meaning  $\text{rank}(\mathbf{X}) < D$ )

- if  $D > N$ , we always have  $\text{rank}(\mathbf{X}) < D$  (since row rank = col. rank)
- if  $D \leq N$ , but some of the columns  $\mathbf{x}_{:d}$  are (nearly) collinear. In the later case, the matrix is ill-conditioned, leading to numerical issues when solving the linear system.



Condition number

$$\frac{\lambda_{\max}(\mathbf{X}^T \mathbf{X})}{\lambda_{\min}(\mathbf{X}^T \mathbf{X})}$$

## Summary of Linear Regression

We have studied three types of methods:

- ① Grid Search
- ② Iterative Optimization Algorithms  
(Stochastic) Gradient Descent

- ③ Least squares  
closed-form solution, for linear MSE

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

↓

$$ND^2 + D^3 + \frac{N \cdot D}{D^2}$$

comp. gram      comp. -1

Computational cost?

$$O(ND^2 + D^3)$$

# Additional Notes

## Closed-form solution for MAE

Can you derive closed-form solution for 1-parameter model when using MAE cost function?

See this short article: <http://www.johnmyleswhite.com/notebook/2013/03/22/modes-medians-and-means-an-unifying-perspective/>.

## Implementation

*solve  $Aw=b$*

There are many ways to solve a linear system, but using the QR decomposition is one of the most robust ways. Matlab's backslash operator and also NumPy's linalg package implement this in just one line:

```
1 w = np.linalg.solve(A, b)
```

For a robust implementation, see Sec. 7.5.2 of Kevin Murphy's book.

## ToDo

1. Revise linear algebra to understand why  $\mathbf{X}$  needs to have full rank. Read the Wikipedia page on the rank of a matrix.
2. Understand the robust implementation of the linear system solver, and play with it during the lab. Read Kevin Murphy's Section 7.5.2 for details.
3. Understand ill-conditioning. Reading about the "condition number" in Wikipedia will help. Also, understanding SVD is essential. Here is another link provided by Dana Kianfar (EPFL) <http://www.cs.uleth.ca/~holzmann/notes/illconditioned.pdf>.
4. Work out the computational complexity of least-squares, when using an existing linear systems solver (use the Wikipedia page on computational complexity).