

Least Squares

Mohammad Emtiyaz Khan
EPFL

Sep 24, 2015



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

©Mohammad Emtiyaz Khan 2015

Motivation

In rare cases, we can compute the minimum of the cost function analytically. Linear regression using MSE is one such case. The solution is obtained using [normal equations](#). This is called [least squares](#).

To derive the equation, we use the optimality conditions. See the lecture notes for Gradient Descent.

$$\frac{\partial \mathcal{L}(\beta^*)}{\partial \beta} = 0$$

Using this, derive the normal equation for 1-parameter model.

Normal equations

Recall the expression of the gradient for multiple linear regression:

$$\frac{\partial \mathcal{L}}{\partial \beta} = -\frac{1}{N} \tilde{\mathbf{X}}^T \mathbf{e} = -\frac{1}{N} \tilde{\mathbf{X}}^T (\mathbf{y} - \tilde{\mathbf{X}}\beta)$$

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

$$\tilde{\mathbf{X}}^T \mathbf{e} = \tilde{\mathbf{X}}^T (\mathbf{y} - \tilde{\mathbf{X}}\beta) = 0$$

implying that the error is orthogonal to rows of $\tilde{\mathbf{X}}^T$ and columns of $\tilde{\mathbf{X}}$.

1-param model

$$\mathcal{L}(\beta_0) = \sum_{n=1}^N (y_n - \beta_0)^2$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \beta_0} &= -2 \sum_n (y_n - \beta_0) \\ &= -2 \sum_n y_n + 2N\beta_0 \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 \mathcal{L}}{\partial \beta_0^2} &= 2N > 0 \quad \beta_0 = \frac{1}{N} \sum y_n \triangleq \bar{y} \end{aligned}$$

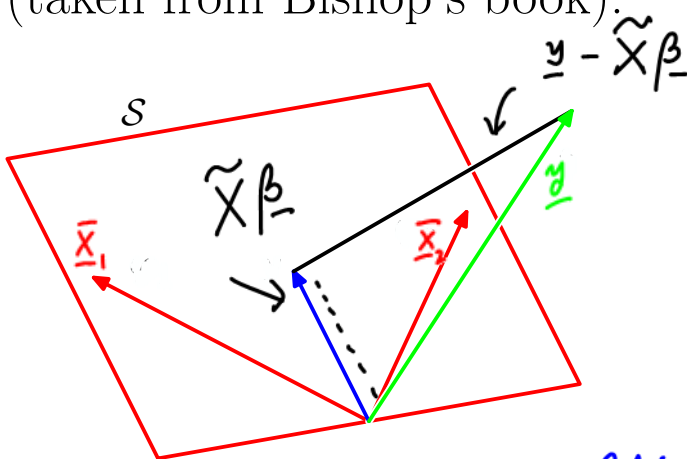
Geometric Interpretation

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

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

$N \times D$
 $\tilde{\mathbf{x}}_1^T$
 $\tilde{\mathbf{x}}_2^T$
 \vdots
 $\tilde{\mathbf{x}}_N^T$

The normal equations suggest to choose a vector in the span of $\tilde{\mathbf{X}}$. The following figure illustrates this (taken from Bishop's book).



An example of "span"

$$\exists \beta_1, \beta_2 \in \mathbb{R}$$

$$\underline{a} = \beta_1 \underline{x}_1 + \beta_2 \underline{x}_2$$

$$\underline{a} \in \text{Span}(X)$$

Q1: What happens when $\underline{y} \in \text{Span}(\tilde{X})$?

Q2: What happens when $\underline{x}_1 = \underline{x}_2$?

Least-squares

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

$$\beta^* = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

We can predict values for a new \mathbf{x}_* .

$$\hat{y}_* = \tilde{\mathbf{x}}_*^T \beta^* = \tilde{\mathbf{x}}_*^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

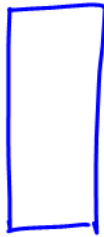
Invertibility and uniqueness

The Gram matrix $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ is invertible iff $\tilde{\mathbf{X}}$ has full column rank.

Proof: Assume $N > D$. The fundamental theorem of linear algebra states that the dimensionality of null space is zero for full column rank. This implies that the Gram matrix is positive definite, which implies invertibility.

$N > D$

$\tilde{\mathbf{X}}$



$N \times D$

$\text{Rank}(\tilde{\mathbf{X}}) = D$

\Rightarrow Nullspace Dimens. = 0

$\tilde{\mathbf{X}} \underline{a} \neq 0$, when $\underline{a} \neq 0$

$\Rightarrow \underline{a}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \underline{a} > 0 \quad \forall \underline{a} \neq 0$

Rank deficiency and ill-conditioning

Unfortunately, $\tilde{\mathbf{X}}$ could often be rank deficient in practice, e.g. when $D > N$, or when the columns $\bar{\mathbf{x}}_d$ are (nearly) collinear. In the later case, the matrix is ill-conditioned, leading to numerical issues.

Condition Number = $\frac{\lambda_{\max}}{\lambda_{\min}}$

SVD

$\tilde{\mathbf{X}} = \begin{bmatrix} \text{U} & \text{S} & \text{V}^T \end{bmatrix}$

$N \times D \quad N \times D \quad D \times D \quad D \times D$

Summary of linear regression

We have studied three methods:

1. Grid search
2. (Stochastic) gradient descent
3. Least squares

Your answers

Correct

$O(NDM^D)$

✓

$O(ND^2)$

$O(NDI)$
where I is # iterations

$O(N^2D^2)$

$O(ND^2 + D^3)$
 $\approx O(ND^2)$
when $N > D$

Additional Notes

Closed-form solution for MAE

Can you derive close-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

There are many ways to implement matrix inversion, but using QR decomposition is one of the most robust ways. Matlab's backslash operator implements this (and much more) in just one line.

```
1 beta = inv(X'*X) * (X'*y)
2 beta = pinv(X'*X) * (X'*y)
3 beta = (X'*X) \ (X'*y)
```

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

To do

1. Revise linear algebra to understand why $\tilde{\mathbf{X}}$ needs to have full rank. Read the Wikipedia page on rank of a matrix.
2. For details on the geometrical interpretation, see Bishop 3.1.2. However, better to read this after the lecture on “basis-function expansion”. Also, note that notation in the book is different. This might make the reading difficult.
3. Understand matrix inversion robust implementation and play with it during the lab. Read Kevin Murphy's section 7.5.2 for details.

4. 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>.

5. Work out the computational complexity of least-squares (use the Wikipedia page on computational complexity).

Link doesn't work. Try these new links!

http://engrwww.usask.ca/classes/EE/840/notes/ILL_Conditioned%20Systems.pdf

<http://www.math.wsu.edu/math/faculty/lih/c220.pdf>