

Annotated
Version

Machine Learning Course - CS-433

Least Squares

Oct 4, 2022

Martin Jaggi

Last updated on: October 3, 2022

credits to Mohammad Emtiyaz Khan & Rüdiger Urbanke

EPFL

Motivation

In rare cases, one can compute the optimum of the cost function analytically. Linear regression using a mean-squared error cost function is one such case. Here the solution can be obtained explicitly, by solving a linear system of equations. These equations are sometimes called the normal equations. This method is one of the most popular methods for data fitting. It is called least squares.

① To derive the normal equations, we first show that the problem is convex. We then use the optimality conditions for convex functions (see the previous lecture notes on optimization). I.e., at the optimum parameter, call it \mathbf{w}^* , it must be true that the gradient of the cost function is $\mathbf{0}$. I.e.,

②
$$\nabla \mathcal{L}(\mathbf{w}^*) = \mathbf{0}.$$

This is a system of D equations.

Find w

$$\min_w \mathcal{L}(w)$$

$$\parallel \frac{1}{N} \sum_{n=1}^N (y_n - w_0)^2 \parallel f_w(x)$$

1-param model

① \mathcal{L} is convex in w ?
yes

②
$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w} &= \frac{1}{N} \sum_{n=1}^N (y_n - w_0)(-1) \\ &= w_0 - \frac{1}{N} \sum_{n=1}^N y_n \\ &\stackrel{!}{=} 0 \\ &\Rightarrow w_0 = \bar{y} \end{aligned}$$

global solution

→ find solution w^*

Normal Equations

Recall that the cost function for linear regression with mean-squared error is given by

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \mathbf{w})^2 = \frac{1}{2N} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}),$$

$f_w(x)$
 \parallel
 $\uparrow \mathbb{R}^D \quad \uparrow \mathbb{R}^D$

$= \frac{1}{2N} \|\mathbf{e}\|^2$
 \mathbf{e}

where

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

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

① We claim that this cost function is convex in the \mathbf{w} . There are several ways of proving this:

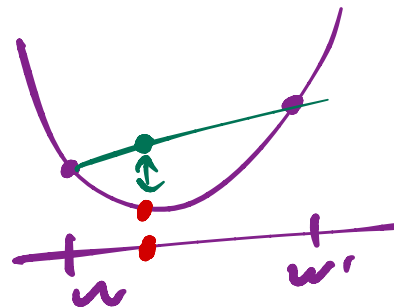
1. Simplest way: observe that \mathcal{L} is naturally represented as the sum (with positive coefficients) of the simple terms $(y_n - \mathbf{x}_n^\top \mathbf{w})^2$. Further, each of these simple terms is the composition of a linear function with a convex function (the square function). Therefore, each of these simple terms is convex and hence the sum is convex.

$\mathcal{L}_n(w)$

$\Rightarrow \text{convex}$

Convex ?

line-segment



2. Directly verify the definition, that for any $\lambda \in [0, 1]$ and \mathbf{w}, \mathbf{w}' ,

$$\mathcal{L}(\lambda \mathbf{w} + (1 - \lambda) \mathbf{w}') - (\lambda \mathcal{L}(\mathbf{w}) + (1 - \lambda) \mathcal{L}(\mathbf{w}')) \leq 0.$$

Computation: LHS =

$$-\frac{1}{2N} \lambda(1 - \lambda) \|\mathbf{X}(\mathbf{w} - \mathbf{w}')\|_2^2, \leq 0 \quad \checkmark$$

which indeed is non-positive.

3. We can compute the second derivative (the Hessian) and show that it is positive semidefinite (all its eigenvalues are non-negative). For the present case a computation shows that the Hessian has the form

$$\frac{1}{N} \mathbf{X}^T \mathbf{X}.$$

This matrix is indeed positive semidefinite since its non-zero eigenvalues are the squares of the non-zero singular values of the matrix \mathbf{X} .

$$\nabla^2 \mathcal{L} = \left(\frac{\partial^2 \mathcal{L}}{\partial w_i \partial w_j} \right)_{ij}$$

MSE

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

$$\nabla^2 \mathcal{L} = \frac{1}{N} \underbrace{\mathbf{X}^T \mathbf{X}}_{D \times D \text{ matrix}}$$

$\nabla^2 \mathcal{L}$ p.s.d. $\forall \mathbf{w}$

$\Rightarrow \mathcal{L}$ convex \checkmark

P.S.D.:

$$\mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} \geq 0 \quad \forall \mathbf{v} \in \mathbb{R}^D \\ = \|\mathbf{X}\mathbf{v}\|^2$$

②

Now where we know that the function is convex, let us find its minimum. If we take the gradient of this expression with respect to the weight vector \mathbf{w} we get

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

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

solve for w

If we set this expression to $\mathbf{0}$ we get the normal equations for linear regression,

$$\mathbf{X}^\top (\underbrace{\mathbf{y} - \mathbf{X}\mathbf{w}}_{\substack{\text{error} \\ e}}) = \mathbf{0}.$$

row of \mathbf{X}^\top
= feature vector
of dataset

① + ② \Rightarrow 1st order opt. condition
global optimum
 $\min_w \mathcal{L}(w)$

Geometric Interpretation

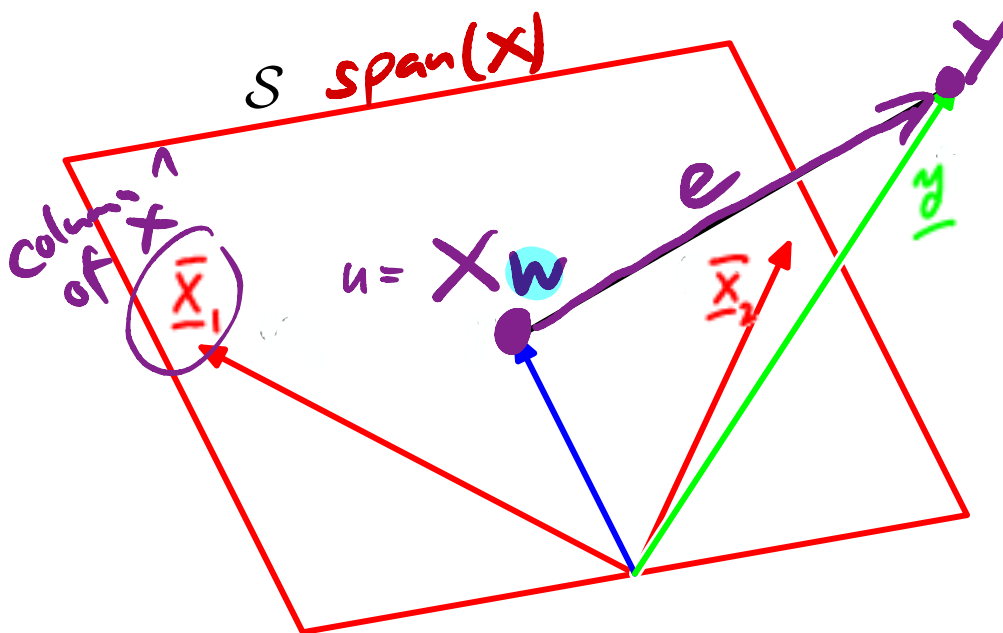
The error e is orthogonal to all columns of \mathbf{X} .

The span of \mathbf{X} is the space spanned by the columns of \mathbf{X} . Every element of the span can be written as $\mathbf{u} = \mathbf{X}\mathbf{w}$ for some choice of \mathbf{w} . Which element of $\text{span}(\mathbf{X})$ shall we take? The normal equations tell us that the optimum choice for \mathbf{u} , call it \mathbf{u}^* , is that element so that $\mathbf{y} - \mathbf{u}^*$ is orthogonal to $\text{span}(\mathbf{X})$. In other words, we should pick \mathbf{u}^* to be equal to the projection of \mathbf{y} onto $\text{span}(\mathbf{X})$.

$$\text{span}(\mathbf{X}) = \{ \mathbf{X}\mathbf{w} \mid \mathbf{w} \in \mathbb{R}^D \}$$

The following figure illustrates this:

(taken from Bishop's book)



$$\min_{\mathbf{w}} \|\mathbf{e}\|^2$$

$$\hat{L}(\mathbf{w})$$

Least Squares

The matrix $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{D \times D}$ is called the **Gram matrix**. If it is invertible, we can multiply the normal equation by the inverse of the Gram matrix from the left to get 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 :

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

test training set

Invertibility and Uniqueness

Note that the Gram matrix $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{D \times D}$ is invertible if and only if \mathbf{X} has full column rank, or in other words $\text{rank}(\mathbf{X}) = D$.

Proof: To see this assume first that $\text{rank}(\mathbf{X}) < D$. Then there exists a non-zero vector \mathbf{u} so that $\mathbf{X}\mathbf{u} = \mathbf{0}$. It follows that $\mathbf{X}^\top \mathbf{X}\mathbf{u} = \mathbf{0}$, and so $\text{rank}(\mathbf{X}^\top \mathbf{X}) < D$. Therefore, $\mathbf{X}^\top \mathbf{X}$ is not invertible.

Conversely, assume that $\mathbf{X}^\top \mathbf{X}$ is not invertible. Hence, there exists a non-zero vector \mathbf{v} so that $\mathbf{X}^\top \mathbf{X}\mathbf{v} = \mathbf{0}$. It follows that

$$\mathbf{0} = \mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} = (\mathbf{X}\mathbf{v})^\top (\mathbf{X}\mathbf{v}) = \|\mathbf{X}\mathbf{v}\|^2.$$

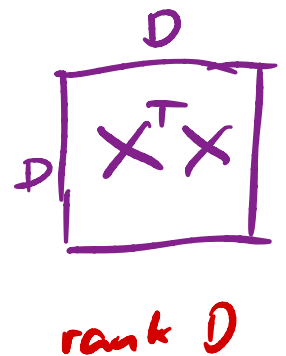
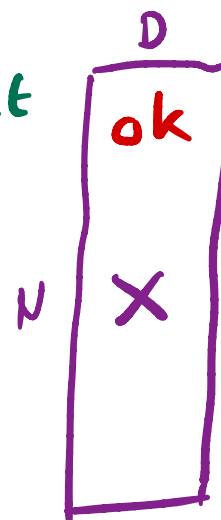
This implies that $\mathbf{X}\mathbf{v} = \mathbf{0}$, i.e., $\text{rank}(\mathbf{X}) < D$.

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

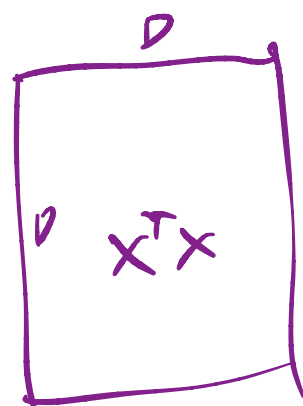
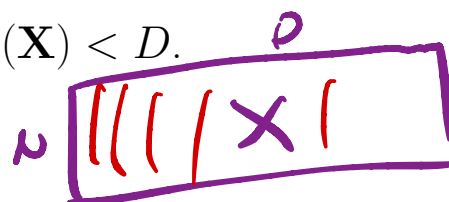
$$\mathbf{X}^\top \mathbf{y} = \underbrace{\mathbf{X}^\top \mathbf{X}}_{D \times D} \mathbf{w}$$

solve linear system

case $D \ll N$



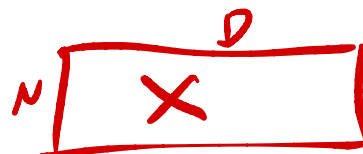
case $D \gg N$



Rank Deficiency and Ill-Conditioning

Unfortunately, in practice, \mathbf{X} is often rank deficient.

- If $D > N$, we always have $\text{rank}(\mathbf{X}) < D$
(since row rank = col. rank)

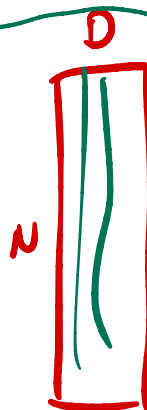


over-parameterized

$\Rightarrow \mathbf{X}^T \mathbf{X}$ not invertible

$\Rightarrow \mathbf{w}^*$ not unique

- If $D \leq N$, but some of the columns $\mathbf{x}_{:,d}$ are (nearly) collinear, then the matrix is ill-conditioned, leading to numerical issues when solving the linear system.



$\Rightarrow \text{rank} < D$

Can we solve least squares if \mathbf{X} is rank deficient? Yes, using a linear system solver.

$$\underset{D \times D}{\mathbf{X}^T \mathbf{X}} \underset{D}{\mathbf{w}} = \underset{D}{\mathbf{X}^T \mathbf{y}}$$

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

condition number

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

zero if rank-deficient

computational cost = $O(N \cdot D^2 + D^3)$

Annotations: "compute $\mathbf{X}^T \mathbf{X}$ " points to $N \cdot D^2$; "inverse" points to 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

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:

```
w = np.linalg.solve(X, y)
```

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