

CX1104: Linear Algebra for Computing

$$\underbrace{\begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mn} \end{bmatrix}}_{A \quad m \times n} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}}_{x \quad n \times 1} = \underbrace{\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}}_{b \quad m \times 1}$$

Chap. No : **7.1.3A**

Lecture : **Least Squares**

Topic : **Least Squares**

Supplementary – Proving Solutions of

Concept : Normal Equation are Least Square
Solutions

Instructor: **A/P Chng Eng Siong**

TAs: **Zhang Su, Vishal Choudhari**

Proving Solutions of Normal Equation are Least Square Solutions

The set of least-squares solutions of $A\mathbf{x} = \mathbf{b}$ coincides with the nonempty set of solutions of the normal equations $A^T A\mathbf{x} = A^T \mathbf{b}$.

Proofs by:

1. Gilbert Strang, MIT (2 Proofs)
2. Jeffrey Chasnov, HKUST
3. Alex Townsend, MIT
4. Alexey Grigorev
5. Quora

Gilbert Strang (Proof 1)

In Strang's Lecture 15 video, Strang: 18.06

Projection onto Subspace

https://youtu.be/Y_Ac6KiQ1t0?t=951

1. 3-D, projecting a point b onto a plane span by A 's columns (2-D) (18:50)

- Vector b is not in the plane (span by A 's columns). (Fig a)
- What is that plane? (19:40)
 - The plane is span by columns of A : a_1 and a_2 and hence, form the basis for the subspace.
 - The columns of A are independent, and need NOT be perpendicular.
- (22:20) $e = b - p$ (is orthogonal to the plane of A) (crucial fact)

2. What is small p (the projected vector of b on subspace spanned by A 's columns)?

- 23:45: p is some linear combination of columns of A . (Fig b)
- $p = A\hat{x}$, we are looking for \hat{x} such that error vector e is orthogonal to plane.

3. Now e is orthogonal to the column space of A . (26:00)

- We have 2 equations to solve for 2 unknowns (fig c)
- Reducing it to $A^T(b - A\hat{x}) = 0$ (fig d) -> The normal equation

Fig (a)

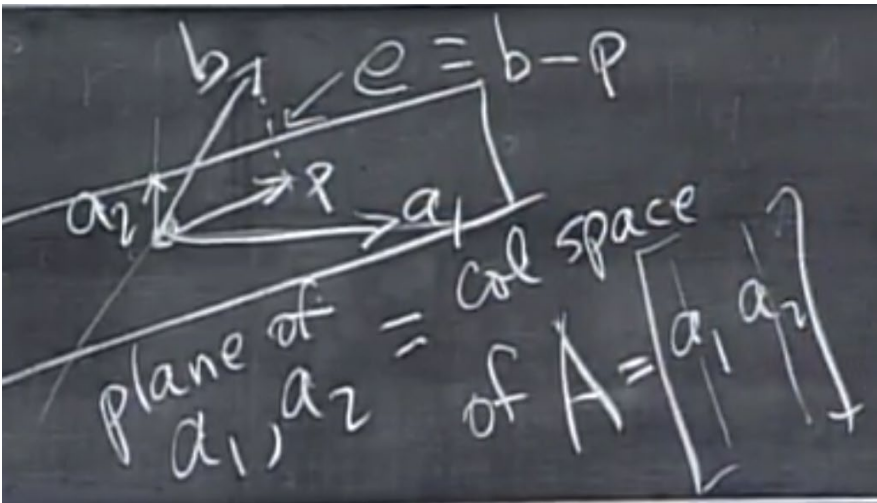


Fig (b)

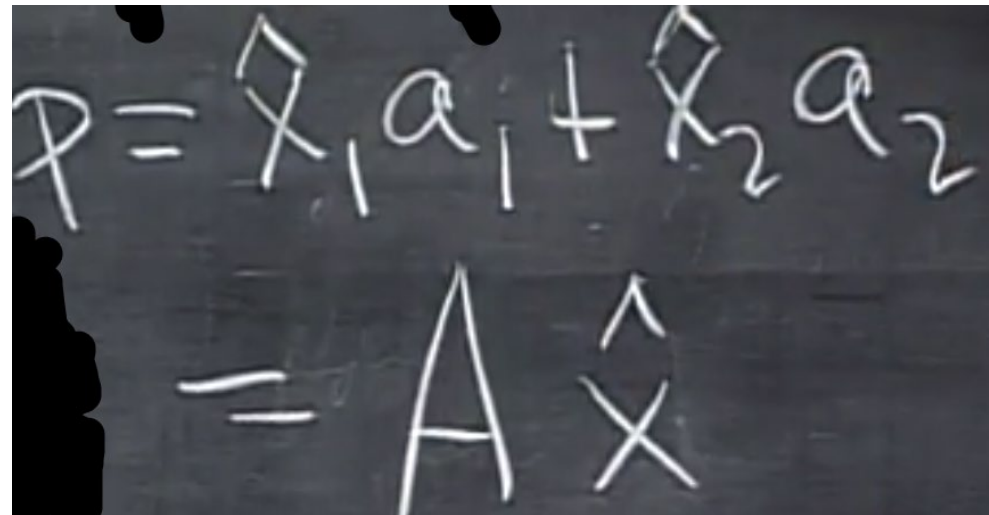


Fig (c)

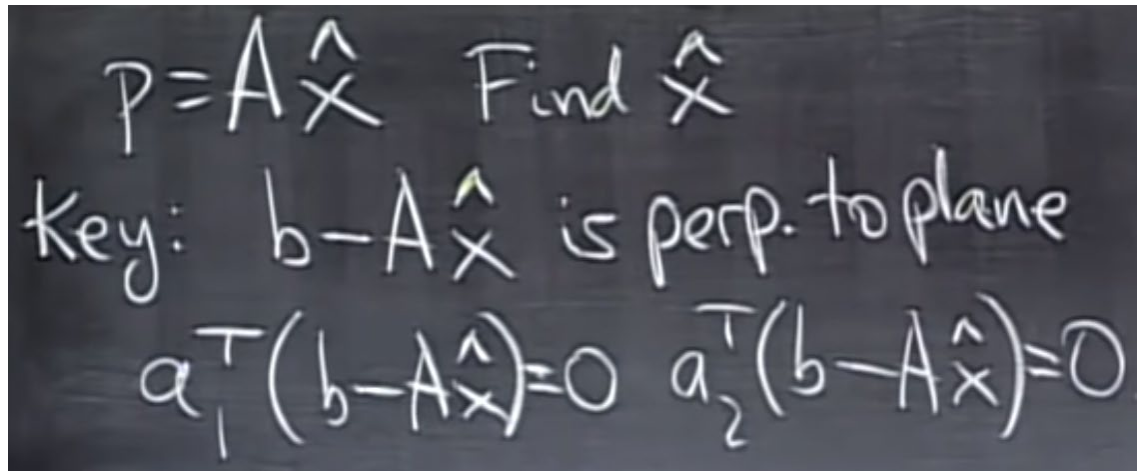
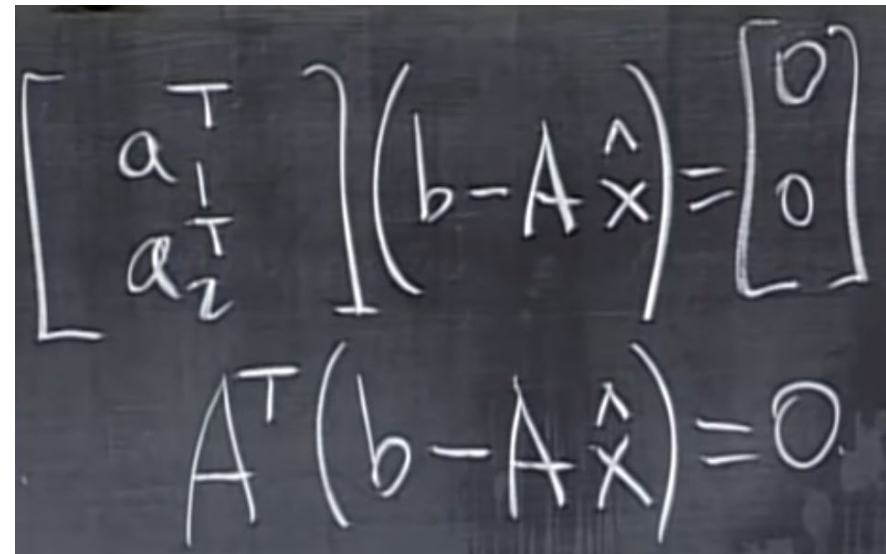


Fig (d)



$$e \in N(A^T)$$
$$e \perp C(A)$$

Gilbert Strang (Proof 2 - Using Dot Product)

Projections onto subspaces

Projections

If we have a vector \mathbf{b} and a line determined by a vector \mathbf{a} , how do we find the point on the line that is closest to \mathbf{b} ?

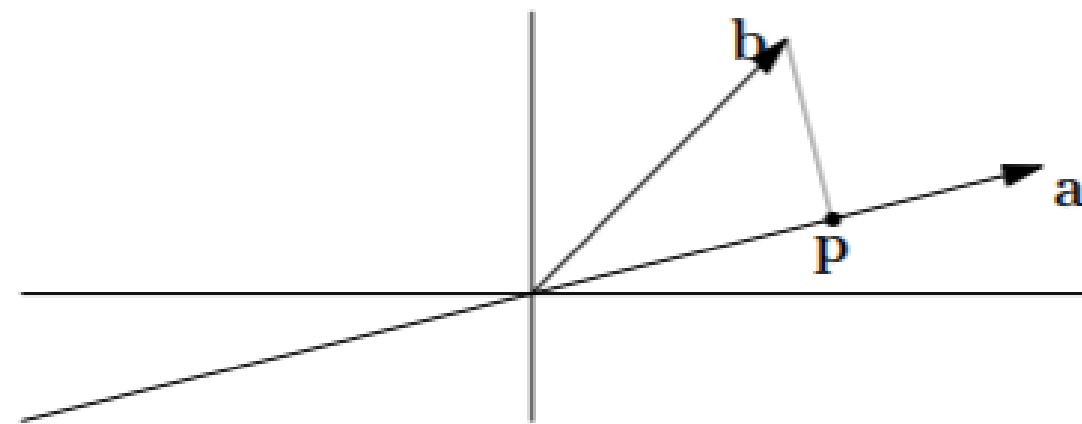


Figure 1: The point closest to \mathbf{b} on the line determined by \mathbf{a} .

We can see from Figure 1 that this closest point \mathbf{p} is at the intersection formed by a line through \mathbf{b} that is orthogonal to \mathbf{a} . If we think of \mathbf{p} as an approximation of \mathbf{b} , then the length of $\mathbf{e} = \mathbf{b} - \mathbf{p}$ is the error in that approximation.

We could try to find \mathbf{p} using trigonometry or calculus, but it's easier to use linear algebra. Since \mathbf{p} lies on the line through \mathbf{a} , we know $\mathbf{p} = x\mathbf{a}$ for some number x . We also know that \mathbf{a} is perpendicular to $\mathbf{e} = \mathbf{b} - x\mathbf{a}$:

$$\begin{aligned}\mathbf{a}^T(\mathbf{b} - x\mathbf{a}) &= 0 \\ x\mathbf{a}^T\mathbf{a} &= \mathbf{a}^T\mathbf{b} \\ x &= \frac{\mathbf{a}^T\mathbf{b}}{\mathbf{a}^T\mathbf{a}},\end{aligned}$$

and $\mathbf{p} = \mathbf{a}x = \mathbf{a}\frac{\mathbf{a}^T\mathbf{b}}{\mathbf{a}^T\mathbf{a}}$. Doubling \mathbf{b} doubles \mathbf{p} . Doubling \mathbf{a} does not affect \mathbf{p} .

Why project?

As we know, the equation $A\mathbf{x} = \mathbf{b}$ may have no solution. The vector $A\mathbf{x}$ is always in the column space of A , and \mathbf{b} is unlikely to be in the column space. So, we project \mathbf{b} onto a vector \mathbf{p} in the column space of A and solve $A\hat{\mathbf{x}} = \mathbf{p}$.

Projection in higher dimensions

In \mathbb{R}^3 , how do we project a vector \mathbf{b} onto the closest point \mathbf{p} in a plane?

If \mathbf{a}_1 and \mathbf{a}_2 form a basis for the plane, then that plane is the column space of the matrix $A = [\mathbf{a}_1 \ \mathbf{a}_2]$.

We know that $\mathbf{p} = \hat{x}_1\mathbf{a}_1 + \hat{x}_2\mathbf{a}_2 = A\hat{\mathbf{x}}$. We want to find $\hat{\mathbf{x}}$. There are many ways to show that $\mathbf{e} = \mathbf{b} - \mathbf{p} = \mathbf{b} - A\hat{\mathbf{x}}$ is orthogonal to the plane we're projecting onto, after which we can use the fact that \mathbf{e} is perpendicular to \mathbf{a}_1 and \mathbf{a}_2 :

$$\mathbf{a}_1^T(\mathbf{b} - A\hat{\mathbf{x}}) = 0 \quad \text{and} \quad \mathbf{a}_2^T(\mathbf{b} - A\hat{\mathbf{x}}) = 0.$$

In matrix form, $A^T(\mathbf{b} - A\hat{\mathbf{x}}) = \mathbf{0}$. When we were projecting onto a line, A only had one column and so this equation looked like: $\mathbf{a}^T(\mathbf{b} - x\mathbf{a}) = 0$.

Note that $\mathbf{e} = \mathbf{b} - A\hat{\mathbf{x}}$ is in the nullspace of A^T and so is in the left nullspace of A . We know that everything in the left nullspace of A is perpendicular to the column space of A , so this is another confirmation that our calculations are correct.

We can rewrite the equation $A^T(\mathbf{b} - A\hat{\mathbf{x}}) = \mathbf{0}$ as:

$$A^T A \hat{\mathbf{x}} = A^T \mathbf{b}.$$

When projecting onto a line, $A^T A$ was just a number; now it is a square matrix. So instead of dividing by $\mathbf{a}^T \mathbf{a}$ we now have to multiply by $(A^T A)^{-1}$.

In n dimensions,

$$\begin{aligned}\hat{\mathbf{x}} &= (A^T A)^{-1} A^T \mathbf{b} \\ \mathbf{p} = A\hat{\mathbf{x}} &= A(A^T A)^{-1} A^T \mathbf{b} \\ P &= A(A^T A)^{-1} A^T.\end{aligned}$$

Jeffrey Chasnov

If you understand Strang's proof,
Chasnov repeats it in the first 4 minutes.



Ref: <https://www.youtube.com/watch?v=WABC6wmuLOk>

Alex Townsend

Townsend's slide:

<https://math.mit.edu/classes/18.085/summer2016/handouts/LeastSquares.pdf>

Another way to show the same thing:

- <https://theclevermachine.wordpress.com/2012/09/01/derivation-of-ols-normal-equations/>
- <https://sites.math.washington.edu/~burke/crs/308/LeastSquares.pdf>

Online tool to perform differentiation of matrix:

<http://www.matrixcalculus.org/>

Matrix calculus reference (Advance):

- <https://atmos.washington.edu/~dennis/MatrixCalculus.pdf>
- <https://www.comp.nus.edu.sg/~cs5240/lecture/matrix-differentiation.pdf>
- <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

NORMAL EQUATIONS:

$$A^T A x = A^T b$$

Why the normal equations? To find out you will need to be slightly crazy and totally comfortable with calculus.

In general, we want to minimize¹

$$f(x) = \|b - Ax\|_2^2 = (b - Ax)^T(b - Ax) = b^T b - x^T A^T b - b^T A x + x^T A^T A x.$$

If x is a global minimum of f , then its gradient $\nabla f(x)$ is the zero vector. Let's take the gradient of f remembering that

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}.$$

We have the following three gradients:

$$\nabla(x^T A^T b) = A^T b, \quad \nabla(b^T A x) = A^T b, \quad \nabla(x^T A^T A x) = 2A^T A x.$$

To calculate these gradients, write out $x^T A^T b$, $b^T A x$, and $x^T A^T A x$, in terms of sums and differentiate with respect to x_1, \dots, x_n (this gets very messy).

Thus, we have

$$\nabla f(x) = 2A^T A x - 2A^T b,$$

just like we saw in the example. We can solve $\nabla f(x) = 0$ or, equivalently $A^T A x = A^T b$ to find the least squares solution. Magic.

Is this the global minimum? Could it be a maximum, a local minimum, or a saddle point? To find out we take the “second derivative” (known as the Hessian in this context):

$$Hf = 2A^T A.$$

Next week we will see that $A^T A$ is a positive semi-definite matrix and that this implies that the solution to $A^T A x = A^T b$ is a global minimum of $f(x)$. Roughly speaking, $f(x)$ is a function that looks like a bowl.

Alexey Grigorev



Alexey Grigorev

A personal page

ML Wiki

[Home](#) [⚙ Page Info ▾](#)

Normal Equation

Contents [hide]

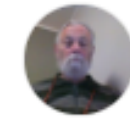
- 1 Normal Equation
- 2 Multivariate Linear Regression Problem
- 3 Least Squares
 - 3.1 Minimization
- 4 Linear Algebra Point of View
 - 4.1 Projection onto $C(A)$
 - 4.2 Invertability of $A^T A$
 - 4.3 \mathbb{R}^2 Case
- 5 Example
 - 5.1 \mathbb{R}^2 Case
- 6 Normal Equation vs Gradient Descent
- 7 Additional
 - 7.1 Orthogonalization
 - 7.2 Singular Value Decomposition
 - 7.3 Regularization
 - 7.4 Implementation
- 8 See Also
- 9 Sources

Projection onto Subspaces

Contents [hide]

- 1 Projections
- 2 Two-Dimensional Case: Motivation and Intuition
 - 2.1 Trigonometry
 - 2.2 Linear Algebra
 - 2.3 Properties
- 3 Projection onto Subspaces
 - 3.1 Motivation
 - 3.2 Projection onto Plane
- 4 Projection Matrix
 - 4.1 Projecting on a Line
 - 4.1.1 Properties
 - 4.2 General Case
 - 4.2.1 Properties
 - 4.3 P as an action of A
- 5 Theorem: $A^T A$ is Invertible
- 6 Projection onto Orthogonal Basis
- 7 Applications
- 8 Sources

Sometimes Quora has the Answer!



Terry Moore, PhD in statistics

Answered February 1, 2018



How complex is the proof / derivation of the normal equation in linear regression?

It's quite simple once you have learned how to differentiate vector equations. If you don't use vectors and matrices it can become a bit messy.

You can also do it by completing the square. I think this is more straightforward.

You have the model $y = X\beta + \epsilon$ where X is a matrix and the rest vectors. We want to minimise $\epsilon^T \epsilon$.

Now

$$\epsilon^T \epsilon = (y - X\beta)^T (y - X\beta) = y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta.$$

We want the terms containing β to be in the form $(\beta - k)^T A(\beta - k)$. Multiply this out: $\beta^T A\beta - k^T A\beta - \beta^T Ak + k^T Ak$.

This works if $A = X^T X$ and $Ak = X^T y$, i.e. $X^T Xk = X^T y$, i.e. $k = (X^T X)^{-1} X^T y$.

Thus

$$\epsilon^T \epsilon = (\beta - k)^T X^T X(\beta - k) + y^T (I - X(X^T X)^{-1} X^T) y.$$

The first term is zero (a minimum because it is the sum of squares of the components of $X(\beta - k)$) if $\beta = k = (X^T X)^{-1} X^T y$. The term $y^T (I - X(X^T X)^{-1} X^T) y$ is the sum of squared residuals.

The normal equations come from this by multiplying $k = (X^T X)^{-1} X^T y$ on the left by $X^T X$. Note that I assumed that the latter matrix is non-singular. If not it is possible to use a generalised inverse (G is a generalised inverse of A if $AGA = A$. Such an inverse always exists but is not unique except if it is an ordinary inverse.)

Yes, pretty trivial.

950 views · View Upvoters



Upvote · 7



Share

