

ECE595 / STAT598: Machine Learning I

Lecture 1.2: Geometry

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline

Mathematical Background

- **Lecture 1: Linear regression: A basic data analytic tool**
- Lecture 2: Regularization: Constraining the solution
- Lecture 3: Kernel Method: Enabling nonlinearity

Lecture 1: Linear Regression

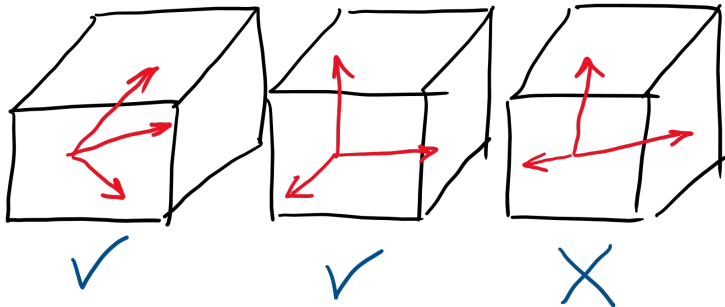
- Linear Regression
 - Notation
 - Loss Function
 - Solving the Regression Problem
- **Geometry**
 - **Projection**
 - **Minimum-Norm Solution**
 - **Pseudo-Inverse**

Linear Span

Given a set of vectors $\{\mathbf{a}_1, \dots, \mathbf{a}_d\}$, the **span** is the set of all possible linear combinations of these vectors.

$$\text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_d\} = \left\{ \mathbf{z} \mid \mathbf{z} = \sum_{j=1}^d \alpha_j \mathbf{a}_j \right\} \quad (1)$$

Which of the following sets of vectors can span \mathbb{R}^3 ?

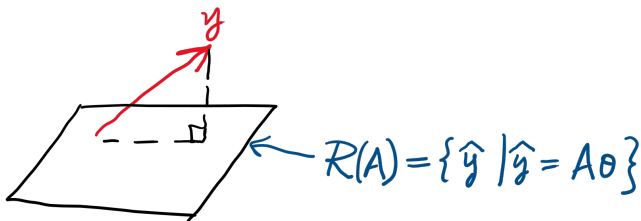


Geometry of Linear Regression

Given θ , the product $\mathbf{A}\theta$ can be viewed as

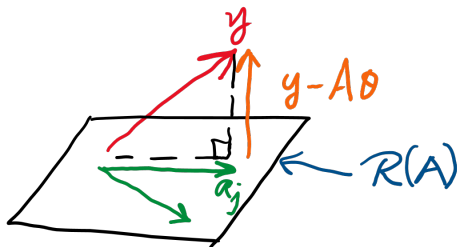
$$\mathbf{A}\theta = \begin{bmatrix} | & | & \dots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_d \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} = \sum_{j=1}^d \theta_j \mathbf{a}_j.$$

So the set of all possible $\mathbf{A}\theta$'s is equivalent to $\text{span}\{\mathbf{a}_1, \dots, \mathbf{a}_d\}$. Define the **range** of \mathbf{A} as $\mathcal{R}(\mathbf{A}) = \{\hat{\mathbf{y}} \mid \hat{\mathbf{y}} = \mathbf{A}\theta\}$. Note that $\mathbf{y} \notin \mathcal{R}(\mathbf{A})$.



Orthogonality Principle

- Consider the error $\mathbf{e} = \mathbf{y} - \mathbf{A}\theta$.
- For the error to minimize, it must be **orthogonal** to $\mathcal{R}(\mathbf{A})$, which is the span of the columns.
- This **orthogonality principle** means that $\mathbf{a}_j^T \mathbf{e} = 0$ for all $j = 1, \dots, d$, which implies $\mathbf{A}^T \mathbf{e} = 0$.



Normal Equation

- The orthogonality principle, which states that $\mathbf{A}^T \mathbf{e} = \mathbf{0}$, implies that $\mathbf{A}^T (\mathbf{y} - \mathbf{A}\boldsymbol{\theta}) = \mathbf{0}$ by substituting $\mathbf{e} = \mathbf{y} - \mathbf{A}\boldsymbol{\theta}$.
- This is called the **normal equation**:

$$\mathbf{A}^T \mathbf{A} \boldsymbol{\theta} = \mathbf{A}^T \mathbf{y}. \quad (2)$$

- The predicted value is

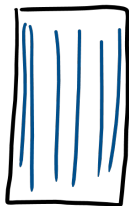
$$\hat{\mathbf{y}} = \mathbf{A} \hat{\boldsymbol{\theta}} = \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

- The matrix $\mathbf{P} \stackrel{\text{def}}{=} \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ is a projection onto the span of $\{\mathbf{a}_1, \dots, \mathbf{a}_d\}$, i.e., the range of \mathbf{A} .
- \mathbf{P} is called the **projection matrix**. It holds that $\mathbf{P} \mathbf{P} = \mathbf{P}$.
- The error $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$ is

$$\begin{aligned} \mathbf{e} &= \mathbf{y} - \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} \\ &= (\mathbf{I} - \mathbf{P}) \mathbf{y}. \end{aligned}$$

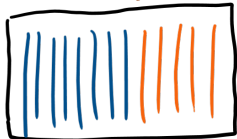
Over-determined and Under-determined Systems

- Assume \mathbf{A} has full column rank.
- Over-determined \mathbf{A} : Tall and skinny. $\hat{\theta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$.
- Under-determined \mathbf{A} : Fat and short. $\hat{\theta} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} \mathbf{y}$.
- If \mathbf{A} is under-determined, then there exists a non-trivial **null space** $\mathcal{N}(\mathbf{A}) = \{\theta \mid \mathbf{A}\theta = 0\}$.
- This implies that if $\hat{\theta}$ is a solution, then $\hat{\theta} + \theta_0$ is also a solution as long as $\theta_0 \in \mathcal{N}(\mathbf{A})$. (Why?)



over-determined

$$\mathcal{N}(\mathbf{A}) = \{\theta \mid \mathbf{A}\theta = 0\}$$



under-determined

Minimum-Norm Solution

- Assume \mathbf{A} is fat and has full row rank.
- Since \mathbf{A} is fat, there exists infinitely many $\hat{\boldsymbol{\theta}}$ such that $\mathbf{A}\hat{\boldsymbol{\theta}} = \mathbf{y}$.
- So we need to pick one in order to be unique.
- It turns out that $\hat{\boldsymbol{\theta}} = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{y}$ is the solution to

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \|\boldsymbol{\theta}\|^2 \quad \text{subject to } \mathbf{A}\boldsymbol{\theta} = \mathbf{y}. \quad (3)$$

(You can solve this problem using Lagrange multiplier. See Appendix.)

- This is called the **minimum-norm** solution.

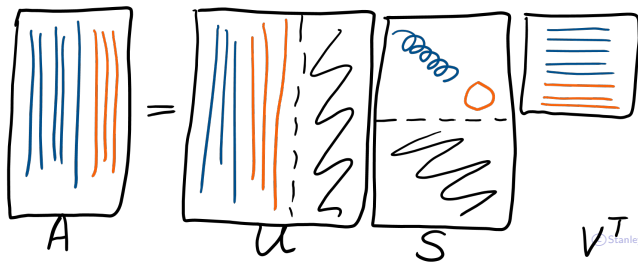


What if Rank-Deficient?

- If \mathbf{A} is rank-deficient, then $\mathbf{A}^T \mathbf{A}$ is not invertible
- Approach 1: **Regularization**. See Lecture 2.
- Approach 2: **Pseudo-inverse**. Decompose $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$.
- $\mathbf{U} \in \mathbb{R}^{N \times N}$, with $\mathbf{U}^T \mathbf{U} = \mathbf{I}$. $\mathbf{V} \in \mathbb{R}^{d \times d}$, with $\mathbf{V}^T \mathbf{V} = \mathbf{I}$.
- The diagonal block of $\mathbf{S} \in \mathbb{R}^{N \times d}$ is $\text{diag}\{s_1, \dots, s_r, 0, \dots, 0\}$.
- The solution is called the **pseudo-inverse**:

$$\hat{\theta} = \mathbf{V}\mathbf{S}^+ \mathbf{U}^T \mathbf{y}, \quad (4)$$

where $\mathbf{S}^+ = \text{diag}\{1/s_1, \dots, 1/s_r, 0, \dots, 0\}$.



Reading List

Linear Algebra

- Gilbert Strang, Linear Algebra and Its Applications, 5th Edition.
- Carl Meyer, Matrix Analysis and Applied Linear Algebra, SIAM, 2000.
- Univ. Waterloo Matrix Cookbook. <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

Linear Regression

- Stanford CS 229 (Note on Linear Algebra)
<http://cs229.stanford.edu/section/cs229-linalg.pdf>
- Elements of Statistical Learning (Chapter 3.2)
<https://web.stanford.edu/~hastie/ElemStatLearn/>
- Learning from Data (Chapter 3.2)
<https://work.caltech.edu/telecourse>

Appendix

Solving the Minimum Norm problem

Consider this problem

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \|\boldsymbol{\theta}\|^2 \quad \text{subject to } \mathbf{A}\boldsymbol{\theta} = \mathbf{y}. \quad (5)$$

The Lagrangian is

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \|\boldsymbol{\theta}\|^2 + \boldsymbol{\lambda}^T (\mathbf{A}\boldsymbol{\theta} - \mathbf{y}).$$

Take derivative with respect to $\boldsymbol{\theta}$ and $\boldsymbol{\lambda}$ yields

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} = 2\boldsymbol{\theta} + \mathbf{A}^T \boldsymbol{\lambda} = 0, \quad \nabla_{\boldsymbol{\lambda}} \mathcal{L} = \mathbf{A}\boldsymbol{\theta} - \mathbf{y} = 0$$

- First equation gives us $\boldsymbol{\theta} = -\mathbf{A}^T \boldsymbol{\lambda} / 2$.
- Substitute into second equation yields $\boldsymbol{\lambda} = -2(\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{y}$.
- Therefore, $\boldsymbol{\theta} = \mathbf{A}^T (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{y}$.