

4.1 Least Squares

Background

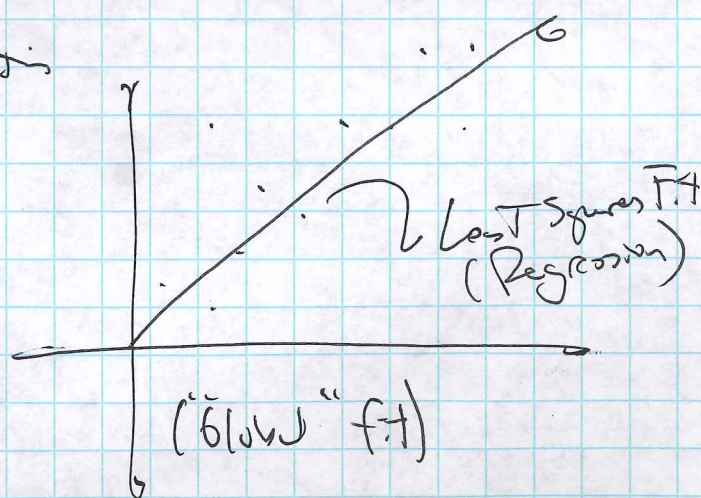
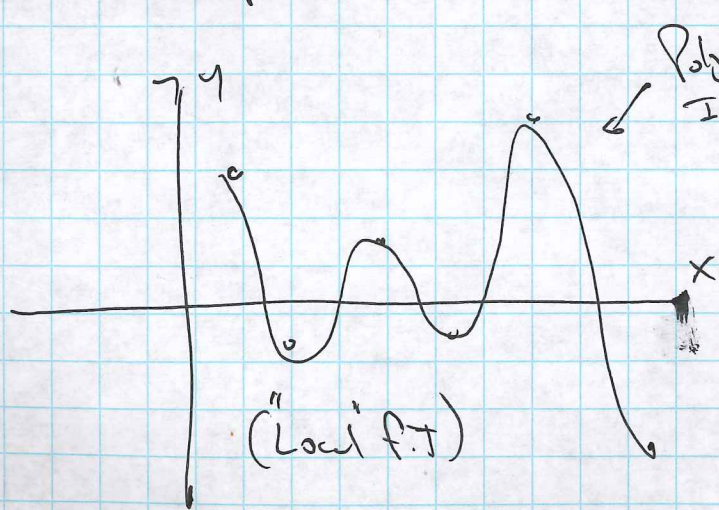
Recall that in Chapter 2 we developed a repertoire of ^{numerical} techniques for solving consistent systems of \mathbb{R} form: $A\vec{x} = \vec{b}$.

In Chapter 3 we studied techniques of polynomial & Trig function interpolation.

Note that interpolation, in general, provides a model that "locally" fits a data set - however, polynomial interpolation is sometimes subject to the Runge Phenomenon of "overfitting" issues in general.

With Least Squares we investigate an alternative (of sorts) to function interpolation by building a model that "globally" fits a data set by minimizing "residual error" (e.g. least squares).

In particular, we investigate in detail systems of \mathbb{R} form: $A\vec{x} = \vec{b}$ that are inconsistent ($n > p$), where \mathbb{R} # of equations (i.e. data points) exceeds the number of model parameters (p).



Ex. Consider the system:

$$\begin{aligned}x_1 + x_2 &= 2 \\x_1 - x_2 &= 1 \\x_1 + x_2 &= 3\end{aligned}$$

$\rightarrow A\vec{x} = \vec{b} \rightarrow \begin{bmatrix} 1 & 1 & | & 2 \\ 1 & -1 & | & 1 \\ 1 & 1 & | & 3 \end{bmatrix}$. Note that this system is inconsistent!.

i.e. it has no exact solution. Why? Most obviously, $x_1 + x_2 = 2 \neq x_1 + x_2 = 3$ cannot both hold. Furthermore, using the language of linear algebra, A is singular (i.e. non-invertible) - i.e. the column space of A does not span \mathbb{R}^3 .

Q: why study inconsistent systems? One Answer: Most problems in applied science involve data sets where the number of data points exceeds the number of model parameters used for the model ~~with which~~ to which we want to model the data (or "compress" the data, as the case may be).

So, how do we "solve" a system that is unsolvable?

Answer: Find the "solution" $\vec{x} \in \text{Col}(A)$ so that $\|A\vec{x} - \vec{b}\|$ is as small as possible!

Recall that the system above, $A\vec{x} = \vec{b}$ can be expressed in vector form: $\begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix} \rightarrow x_1 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + x_2 \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix}$.

In fact, any $m \times n$ system: $A\vec{x} = \vec{b}$ can be viewed as \mathbb{R}^m vector equation:

$$x_1\vec{v}_1 + x_2\vec{v}_2 + \dots + x_n\vec{v}_n = \vec{b} \quad \left(\text{where } A = \begin{bmatrix} | & | & \dots & | \\ \vec{v}_1 & \vec{v}_2 & \dots & \vec{v}_n \\ | & | & \dots & | \end{bmatrix} \right)$$

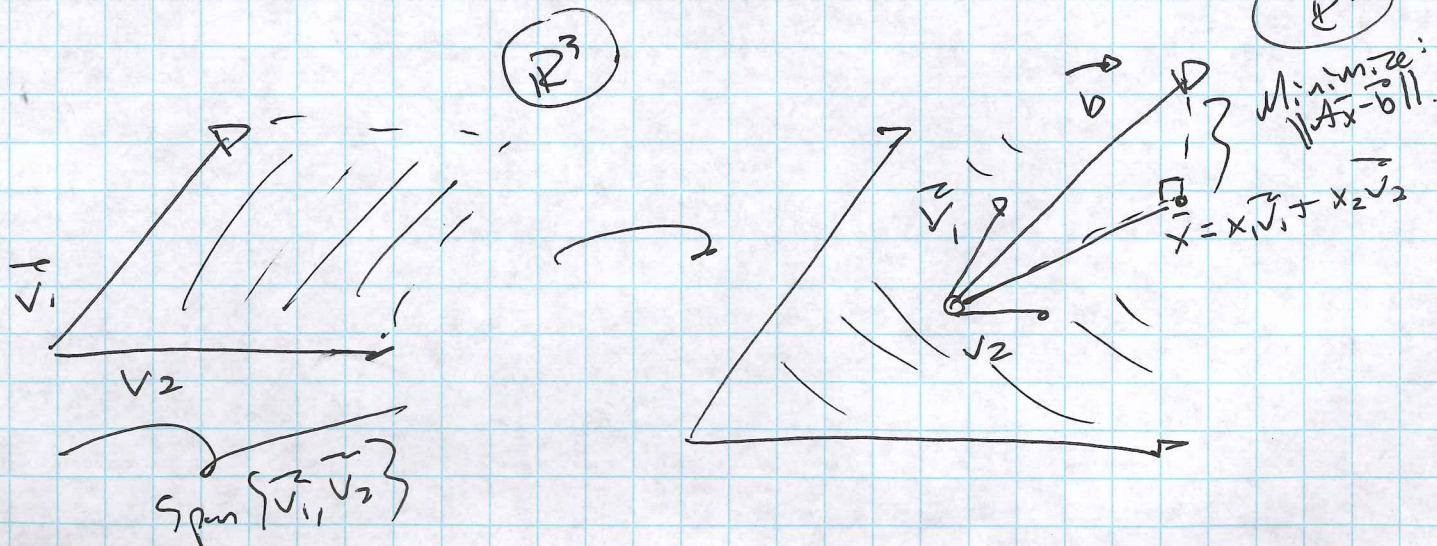
which expresses \vec{b} as a linear combination of \mathbb{R}^m columns \vec{v}_i of A , with coefficients: x_1, \dots, x_n .

In our case, we are attempting to "hit" \mathbb{R}^3 target vector \vec{b} as a linear combination of (2) (only) 3-dimensional vectors!

Since a linear combination of 2 3-dim vectors forms a plane in \mathbb{R}^3 , the vector equation above has a solution

only if \vec{b} lies in the aforementioned plane. Too many equations makes this problem "overspecified."

We proceed to find a pair: x_1, x_2 describing a vector in the plane spanned by $\{\vec{v}_1, \vec{v}_2\}$ that is closest to \vec{b} !



Recall that two vectors \vec{u}, \vec{v} are said to be orthogonal if $\vec{u} \cdot \vec{v} = 0$ (geometrically: $\vec{u} \perp \vec{v}$).

4

From the previous picture: The residual vector $(\vec{b} - A\vec{x}) \perp$ to the plane: $\{A\vec{x} \mid \vec{x} \in \mathbb{R}^n\}$ (i.e. $\text{Col}(A)$).

$$\downarrow (\vec{b} - A\vec{x}) \perp \{A\vec{x} \mid \vec{x} \in \mathbb{R}^n\}$$

Note: A dot product $(\vec{u} \cdot \vec{v})$ can be expressed as matrix multiplication: $\vec{u}^T \vec{v} = [u_1 \dots u_n] \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \vec{u} \cdot \vec{v}$.

We use this fact + the notion of orthogonality to write:

$$(\vec{b} - A\vec{x}) \perp \{A\vec{x} \mid \vec{x} \in \mathbb{R}^n\} \rightarrow (A\vec{x})^T (\vec{b} - A\vec{x}) = 0 \quad \forall \vec{x} \in \mathbb{R}^n$$

(Noting that $(AB)^T = B^T A^T$) This yields:

$$\vec{x}^T A^T (\vec{b} - A\vec{x}) = 0 \quad \forall \vec{x} \in \mathbb{R}^n$$

This means that $A^T (\vec{b} - A\vec{x}) \perp \vec{x}$ for any choice of $\vec{x} \in \mathbb{R}^n$.

The only way for this to hold is if: $A^T (\vec{b} - A\vec{x}) = 0$

Simplifying this equation produces the following:

$$A^T A \bar{x} = A^T \bar{b}$$

This system is known as the "normal equations."

The solution to the normal equations, \bar{x} , is known as the least squares solution of the system: $A\bar{x} = \bar{b}$.

Normal Equations for Least Squares

① Given the inconsistent system: $A\bar{x} = \bar{b}$

solve: $A^T A \bar{x} = A^T \bar{b}$ for the least squares solution

\bar{x} , that minimizes the Euclidean length of the residual:

$$\bar{r} = \bar{b} - A\bar{x}$$

Ex. Use the normal equations to find the least squares solution of the inconsistent system.

$$A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & 1 \end{bmatrix}, \bar{b} = \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix}$$

$$A^T A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 1 \\ 1 & 3 \end{bmatrix}$$

$$A^T \bar{b} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix} = \begin{bmatrix} 6 \\ 4 \end{bmatrix} \xrightarrow{\text{Normal Equations}} \begin{bmatrix} 3 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 4 \end{bmatrix}$$

$$\bar{x} = \langle \bar{x}_1, \bar{x}_2 \rangle = \langle 3/4, 3/4 \rangle$$

Substituting the least squares solution into the original problem yields:

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \frac{7}{4} \\ \frac{3}{4} \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \\ 2.5 \end{bmatrix} \neq \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix}$$

$$\text{The residual: } \vec{r} = \vec{b} - A\vec{x} = \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix} - \begin{bmatrix} 2.5 \\ 1 \\ 2.5 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0 \\ 0.5 \end{bmatrix}$$

Note if $\vec{r} = \vec{0}$, we have found the solution $A\vec{x} = \vec{b}$ exactly.

Here we can measure our degree of success for finding a "best" approximate solution to the inconsistent system: $A\vec{x} = \vec{b}$ in at least (3) ways.

① Using Euclidean distance

$$\|\vec{r}\|_2 = \sqrt{r_1^2 + \dots + r_m^2}$$

② The Squared Error

$$SE = r_1^2 + \dots + r_m^2$$

③ The Root Mean Squared Error

$$RMSE = \sqrt{SE/m} = \sqrt{\frac{r_1^2 + \dots + r_m^2}{m}}$$

Note:

$$RMSE = \frac{\sqrt{SE}}{\sqrt{m}} = \frac{\|\vec{r}\|_2}{\sqrt{m}}$$

⑦

Thus, finding \bar{x} that minimizes each type of error measure above, minimizes all.

For instance, $SE = (.5)^2 + 0^2 + (-.5)^2 = .5$; $\|\bar{r}\|_2 = \sqrt{.5} \approx .707$;

$RMSSE = \sqrt{\frac{.5}{3}} \approx .408$.

Ex. Solve the least squares problem.

$$\begin{bmatrix} 1 & -4 \\ 2 & 3 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 15 \\ 9 \end{bmatrix}$$

→ Normal Equations: $A^T A \bar{x} = A^T b$ →

$$\begin{bmatrix} 9 & 6 \\ 6 & 29 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 45 \\ 75 \end{bmatrix}$$

→ $\bar{x}_1 = 3.8$, $\bar{x}_2 = 1.8$.

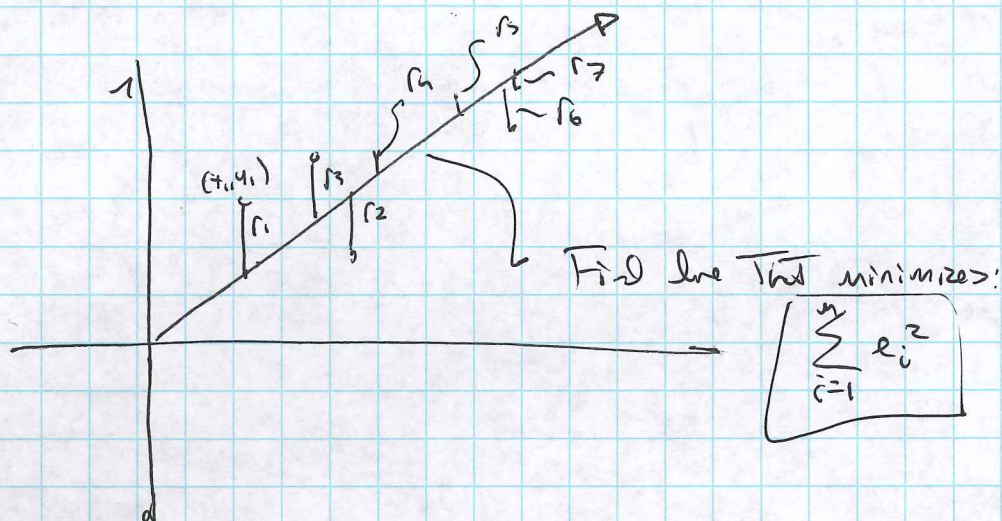
$$\bar{r} = \bar{b} - A\bar{x} = \begin{bmatrix} -3 \\ 15 \\ 9 \end{bmatrix} - \begin{bmatrix} 1 & -4 \\ 2 & 3 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 3.8 \\ 1.8 \end{bmatrix} = \begin{bmatrix} -3 \\ 15 \\ 9 \end{bmatrix} - \begin{bmatrix} -3.4 \\ 13 \\ 11.2 \end{bmatrix} = \begin{bmatrix} .4 \\ 2 \\ -2.2 \end{bmatrix}$$

$$\|\bar{r}\|_2 = \sqrt{(.4)^2 + 2^2 + (-2.2)^2} = \sqrt{3}$$

Fitting Models To Data (Regression)

Let $\{(t_1, y_1), \dots, (t_m, y_m)\}$ be a data set. Given a fixed "class of models":

$\hat{y} = c_1 + c_2 T$, find the model that "best fits" the data by minimizing the sum of residual errors squared.



Ex. Given data: $\{(1, 2), (-1, 1), (1, 3)\}$ Find the best fit linear model.

$$\hat{y} = c_1 + c_2 T$$

$$\begin{cases} c_1 + c_2(1) = 2 \\ c_1 + c_2(-1) = 1 \\ c_1 + c_2(1) = 3 \end{cases}$$

↓ Matrix Form

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix}$$

Has no pure solution!
(points aren't colinear)

From our previous result, $\bar{c}_1 = 7/4$, $\bar{c}_2 = 3/4 \rightarrow \hat{y} = \frac{7}{4} + \frac{3}{4} T$

The RMSE is $\sqrt{\frac{1}{6}}$.

(3) Step Procedure for fitting data (by least squares)

Given n data points: $\{(T_i, y_i)\}_{i=1}^n$
(1) Choose a model, e.g. $\hat{y} = c_1 + c_2 T$,

OR $\hat{y} = c_1 + c_2 T + c_3 T^2$, OR $\hat{y} = e^{c_1 + c_2 T}$ ("logistic regression"), etc.

(2) Force the model to fit the data.

Substitute points into model. This gives system: $A\bar{x} = \bar{b}$ with $A (m \times p)$

(3) Solve the normal equations: $A^T A \bar{x} = A^T \bar{b}$

Ex. Find both the best line & best parabola to fit: $\{(-1, 1), (0, 0), (1, 0), (2, -2)\}$.

(1) $\hat{y} = c_1 + c_2 T \rightarrow$

$$\left. \begin{aligned} c_1 + c_2(-1) &= 1 \\ c_1 + c_2(0) &= 0 \\ c_1 + c_2(1) &= 0 \\ c_1 + c_2(2) &= -2 \end{aligned} \right\} \rightarrow \begin{bmatrix} 1 & -1 \\ 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ -2 \end{bmatrix}$$

$A \qquad \qquad \qquad \bar{b}$

$\rightarrow A^T A \bar{x} = A^T \bar{b} \rightarrow \begin{bmatrix} 4 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} -1 \\ -5 \end{bmatrix} \rightarrow \hat{y} = 0.2 - 0.9T$

SE(Linear) = $(-1)^2 + (0)^2 + (1)^2 + (-2)^2 = 6.7$

RMSE = $\frac{\sqrt{6.7}}{\sqrt{4}} \approx 1.418$

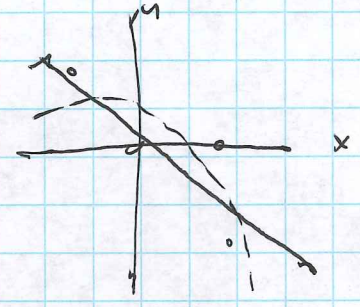
(2) $\hat{y} = c_1 + c_2 T + c_3 T^2 \rightarrow$

$$\left. \begin{aligned} c_1 + c_2(-1) + c_3(-1)^2 &= 1 \\ c_1 + c_2(0) + c_3(0)^2 &= 0 \\ c_1 + c_2(1) + c_3(1)^2 &= 0 \\ c_1 + c_2(2) + c_3(2)^2 &= -2 \end{aligned} \right\} \rightarrow \text{G/A}$$

$$\begin{bmatrix} 1 & -1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 6 \\ -2 \end{bmatrix} \rightarrow A^T A \vec{x} = A^T \vec{b} \rightarrow \begin{bmatrix} 4 & -2 & 6 \\ 2 & 6 & 8 \\ 6 & 8 & 18 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} -1 \\ -5 \\ -7 \end{bmatrix}$$

$$\hat{y} = 0.45 - 0.65x - 0.25x^2$$

$$SE(\text{quad}) = .45, \quad RMSE(\text{quad}) \approx .335$$



(Quad fit appears to be superior).

Compression Least Squares is a classic example of data compression, where we fit a data set with a minimally complex model (i.e. a model with few parameters). Often times, in practice least squares is used to replace "noisy" data with a plausible underlying model (see: signal processing).

Conditioning Note that in general, $\text{cond}(A^T A) \approx \text{cond}(A)^2$, so least squares can greatly increase the possibility of an ill-conditioned problem. Other methods such as QR factorization (4.3) compute least squares without need for $A^T A$.

In this section we explored the so-called **Jordan-Moreau Matrix** (Used in polynomial curve fitting).

$$\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & x_2 & x_2^2 & \dots & x_2^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & x_m^2 & \dots & x_m^n \end{bmatrix} \begin{bmatrix} c_1 \\ \vdots \\ c_{n+1} \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$