

Lecture 16: Least squares (class)

Admin:

Natural problems for solving linear systems:

1. Solve $Ax = b$ for x
2. Assuming $Ax = b$ has a solution,
find the shortest solution x .
3. "Least squares": If $Ax = b$ has no solution ($b \notin R(A)$), find an x that minimizes $\|Ax - b\|$.
4. "Compressed sensing": Find the sparsest x such that $Ax = b$.

Simple application of SVD:

Theorem: For any matrix A ,

- $\text{rank}(ATA) = \text{rank}(A) \checkmark$
- $R(ATA) = R(A^T) \checkmark \quad R(AA^T) = R(A)$
- $N(A^TA) = N(A) \checkmark \quad N(AA^T) = N(A^T)$

$$\text{SVD: } A = \sum_j \sigma_j \vec{v}_j \vec{u}_j^T \quad A\vec{x} = \sum_j \sigma_j \vec{v}_j (\vec{u}_j \circ x) \\ R(A) = \text{Span} \{ \vec{v}_j \mid \sigma_j > 0 \} \quad \text{for } x = \frac{\vec{v}_j}{\sigma_j}, A\vec{x} = \vec{v}_j$$

$$R(A^T) = \text{Span} \{ \vec{u}_j \mid \sigma_j > 0 \}$$

$$\text{Rank}(A) = \dim R(A) = \dim R(A^T) = \#\{j \mid \sigma_j > 0\}.$$

$$\hookrightarrow A^T A = \sum_{jk} \sigma_j \sigma_k \vec{v}_j \vec{v}_k^T \vec{u}_k \vec{u}_k^T = \sum_j \sigma_j^2 \vec{v}_j \vec{v}_j^T \quad R(A^T A) = \text{Span} \{ \vec{v}_j \mid \sigma_j^2 \neq 0 \} \\ \text{in the } \vec{v}_j \text{'s basis } [A^T A] = \begin{pmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \ddots \end{pmatrix} = \text{Span} \{ \vec{v}_j \mid \sigma_j \neq 0 \} = R(A^T)$$

$$\text{let } U = \sum_j \vec{v}_j e_j^T \Rightarrow U^T A^T A U = \sum_j \sigma_j^2 e_j e_j^T \\ \text{unitary } \quad \text{e-values} \quad = \begin{pmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \ddots \end{pmatrix}$$

↑ e-vectors

⇒ to find the SVD for A ,
you can find the spectral decomposition
of $A^T A$

$$\text{or of } AA^T = \sum_j \sigma_j^2 \vec{v}_j \vec{v}_j^T$$

Example: Let

$$A = \begin{pmatrix} -2 & 1 & & & 1 \\ 1 & -2 & 1 & & \\ & 1 & -2 & 1 & \\ & & 1 & -2 & \\ & & & \ddots & 1 \\ & & & & 1 & -2 \\ & & & & & \ddots & 1 \\ & & & & & & 1 \end{pmatrix} = -2I + \begin{pmatrix} 0 & 1 & & & \\ 0 & 0 & 1 & & \\ 0 & 0 & 0 & 1 & \\ & \ddots & & & \\ & & & & 1 \end{pmatrix} + P^T$$

What is $\|A\|$? $\approx 4?$

Why? $3.99 \leq \|A\| \leq 4$

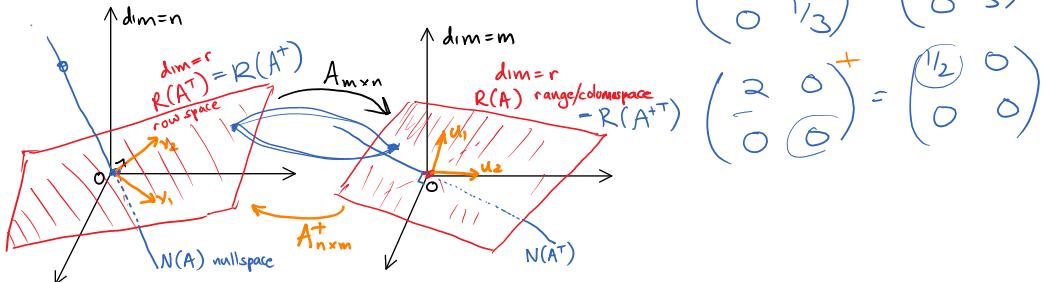
$$\|A\| \leq \| -2I \| + \|P\| + \|P^T\| = 2\|I\| + 2\|P\| = 4 \\ \rightarrow \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \\ \vdots \\ 1 \end{pmatrix} \quad \|A\| \geq \frac{\|A\vec{x}\|}{\|\vec{x}\|} \quad A\vec{x} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} = \vec{0}$$

$$\|A\| \leq \| -2I \| + \|P\| + \|T\| = 2\|I\| + \|P\|$$

$$\vec{x} = \begin{pmatrix} 1 \\ -1 \\ -1 \end{pmatrix} \quad \|A\| \geq \frac{\|A\vec{x}\|}{\|\vec{x}\|} \quad A\vec{x} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} = \vec{0}$$

Moore-Penrose

Definition: Matrix pseudoinverse



$$\begin{pmatrix} 2 & 0 \\ 0 & 1/3 \end{pmatrix}^{-1} = \begin{pmatrix} 1/2 & 0 \\ 0 & 3 \end{pmatrix}$$

$$\begin{pmatrix} 2 & 0 \\ 0 & 0 \end{pmatrix}^+ = \begin{pmatrix} 1/2 & 0 \\ 0 & 0 \end{pmatrix}$$

The pseudoinverse of A , denoted A^+ , is the unique linear operator $\mathbb{R}^m \rightarrow \mathbb{R}^n$ with

- $R(A^+) = R(A^T)$ ✓
- $R(A^{+T}) = R(A)$ ✓
- $A^+ A = P_{R(A^T)}$ ✓
- $A A^+ = P_{R(A)}$ ✓

If the SVD of A is

$$A = \sum_j \sigma_j v_j u_j^T$$

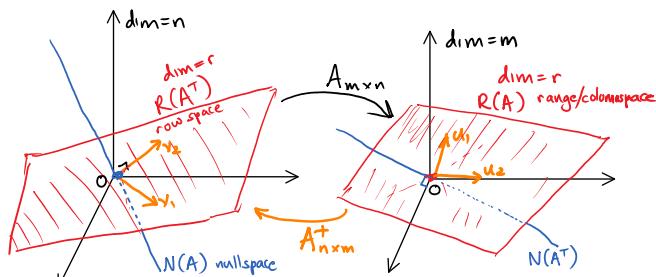
then A^+ is given by

$$A^+ = \sum_{j: \sigma_j > 0} \frac{1}{\sigma_j} u_j v_j^T$$

$$A^+ A = \sum_{j: \sigma_j > 0} \frac{1}{\sigma_j} u_j v_j^T \underbrace{\sum_k \sigma_k v_k u_k^T}_{= I} = \sum_{j: \sigma_j > 0} \frac{1}{\sigma_j} u_j v_j^T = \sum_{j: \sigma_j > 0} u_j v_j^T = P_{R(A^T)}$$

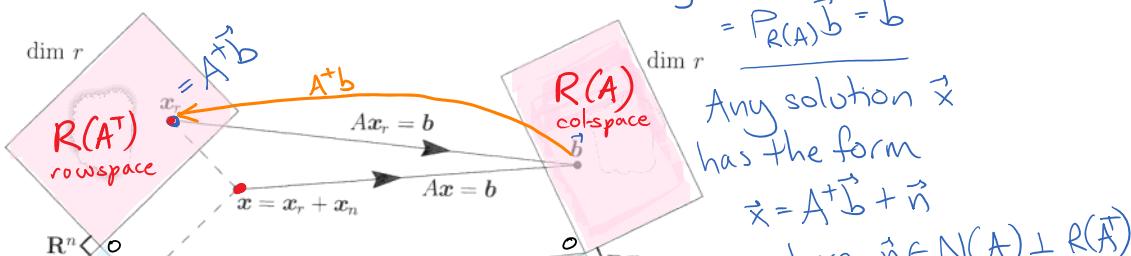
$$A A^+ = \sum_{j: \sigma_j > 0} v_j v_j^T = P_{R(A)}$$

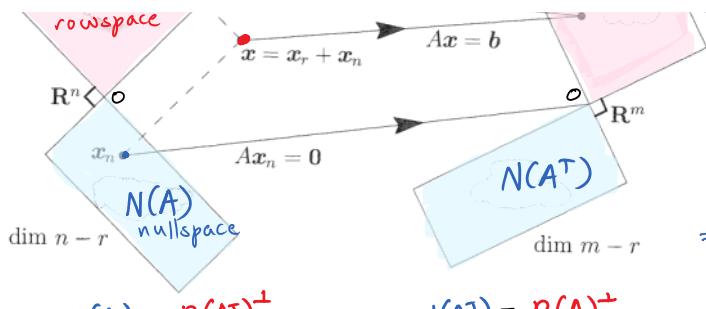
$$\Rightarrow (\text{for example}) A(A^+ A) = A P_{R(A^T)} A = A$$



- If A is invertible/nonsingular, then $A^+ = A^{-1}$.

Corollary: If $A\vec{x} = \vec{b}$ has multiple solutions, ie, $N(A) \neq \{\vec{0}\}$, then $\vec{x} = A^+ \vec{b}$ is the shortest solution.





$\vec{x} = \vec{A}^+ \vec{b} + \vec{n}$
 where $\vec{n} \in N(A) \perp R(\vec{A})$
 $\|\vec{x}\|^2 = \|\vec{A}^+ \vec{b}\|^2 + \|\vec{n}\|^2$
 \Rightarrow shortest solution
 has $\vec{n} = \vec{0}$,
 i.e. $\vec{x} = \vec{A}^+ \vec{b}$. \square

Why?

Any solution $x \in A^+ b + N(A)$

\downarrow

particular solution

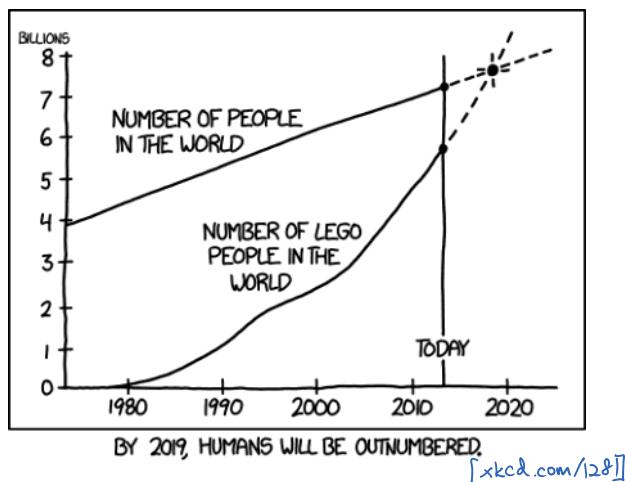
\uparrow

set of solutions to
the homogeneous equations
 $Ax = 0$

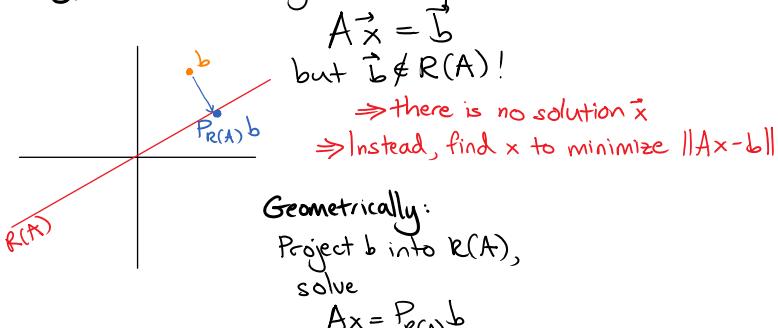
Since $A^+ b \in R(A^\top)$, which is \perp to $N(A)$,
 $\|x\|^2 = \|A^+ b\|^2 + \|\text{its component in } N(A)\|^2$
 which is minimal if the component in $N(A)$ is 0 ,
 i.e., $x = A^+ b$. \square

"LEAST SQUARES" FITTING & APPLICATIONS TO DATA ANALYSIS

(Reading:
 Meyer 4.6
 5.13-14
 Strang 3.3)



Typical problem: System of equations



Example:

$$\begin{pmatrix} 1 & -1 \\ 2 & 3 \\ 0 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 5 \\ 10 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 3 \\ 2 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 5 \\ 10 \\ 0 \end{pmatrix}$$

$\overset{\text{A}}{\parallel}$ $\overset{\text{b}}{\parallel}$

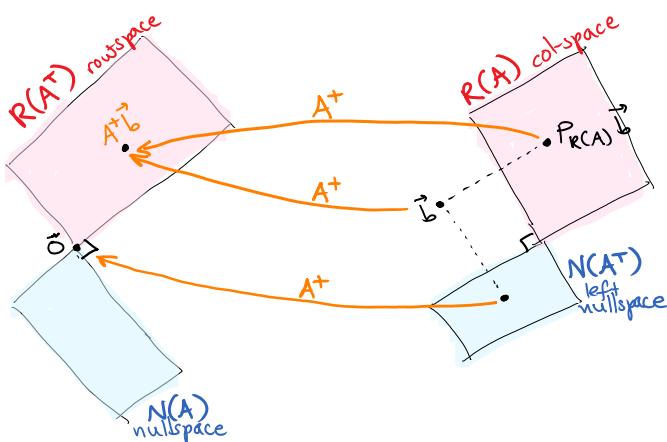
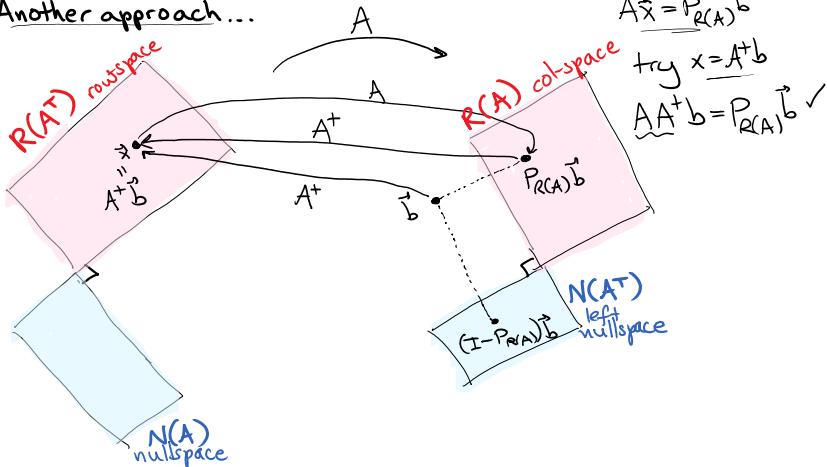
There is no solution, since $(Ax)_3 = 0$ always.

But $R(A) = \text{Span} \left\{ \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \\ 2 \end{pmatrix} \right\}$
 $= xy\text{-plane in } \mathbb{R}^3$

$$x = A^+ \vec{b} \text{ achieves } \min_{\vec{x}} \|A\vec{x} - \vec{b}\| \quad \checkmark$$

$A\vec{x} = \vec{b}$

Another approach...



$\Rightarrow x = A^+ b$
 minimizes $\|Ax - b\|$

Algebraically:

$$A = \sum_i \sigma_i \vec{u}_i \vec{v}_i^T \quad (\text{SVD})$$

$$R(A) = \text{Span}(\{\vec{u}_i | \sigma_i > 0\})$$

$$P_{R(A)} = \sum_{i: \sigma_i > 0} \vec{u}_i \vec{u}_i^T \quad (\text{since they are orthonormal})$$

We want to solve

$$A\vec{x} = P_{R(A)}\vec{b}$$

$$\sum_{i:\sigma_i > 0} \sigma_i (\mathbf{v}_i \cdot \mathbf{x}) \vec{u}_i = \sum_{i:\sigma_i > 0} (\mathbf{u}_i \cdot \mathbf{b}) \vec{u}_i$$

$$\Rightarrow \mathbf{v}_i \cdot \mathbf{x} = \begin{cases} \frac{1}{\sigma_i} (\mathbf{u}_i \cdot \mathbf{b}) & \text{if } \sigma_i > 0 \\ 0 & \text{if } \sigma_i = 0 \end{cases}$$

$$\Rightarrow \vec{x} = \sum_i (\mathbf{v}_i \cdot \mathbf{x}) \vec{v}_i$$

$$= \sum_{i:\sigma_i > 0} \frac{1}{\sigma_i} \vec{v}_i (\mathbf{u}_i \cdot \mathbf{b})$$

$$= \mathbf{A}^+ \vec{b}$$

$\Rightarrow \mathbf{x} = \mathbf{A}^+ \vec{b}$
 minimizes $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|$

Example:

>> A = [1 -1; 2 3; 0 0]

A =

$$\begin{pmatrix} 1 & -1 \\ 2 & 3 \\ 0 & 0 \end{pmatrix}$$

>> format long e
 >> [U,D,V] = svd(A)

U =

$$\begin{pmatrix} -8.980559531591714e-02 \\ 9.959593139531121e-01 \\ 0 \end{pmatrix} \begin{pmatrix} 9.959593139531120e-01 \\ 8.980559531591698e-02 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1.000000000000000e+00 \end{pmatrix}$$

left singular vectors

D =

$$\begin{pmatrix} 3.618033988749895e+00 & 0 \\ 0 & 1.381966011250106e+00 \\ 0 & 0 \end{pmatrix}$$

singular values

V =

$$\begin{pmatrix} 5.257311121191336e-01 \\ 8.506508083520399e-01 \end{pmatrix} \begin{pmatrix} 8.506508083520399e-01 \\ -5.257311121191336e-01 \end{pmatrix}$$

right singular vectors

>> Apseudoinv = V * [1/3.618033988749895 0 0; 0 1/1.381966011250106 0] * U'

Apseudoinv =

$$\begin{pmatrix} 5.999999999999995e-01 & 1.999999999999999e-01 & 0 \\ -3.99999999999997e-01 & 2.000000000000000e-01 & 0 \end{pmatrix}$$

>> Apseudoinv * [0; 5; 10]

$$= \begin{pmatrix} .6 & .2 & 0 \\ -.4 & .2 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

ans =

$$\begin{pmatrix} 9.99999999999996e-01 \\ 1.000000000000000e+00 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \checkmark$$

>> format short e
 >> Apseudoinv2 = pinv(A)

Apseudoinv2 = Matlab's built-in pseudo-inverse function

$$\begin{pmatrix} 6.0000e-01 & 2.0000e-01 & 0 \\ -4.0000e-01 & 2.0000e-01 & 0 \end{pmatrix}$$

same as we got above!

>> A * Apseudoinv

Observe:
 $\mathbf{A}\mathbf{A}^+ = \text{projection onto xy-plane}$

ans =

$$\begin{pmatrix} 1.0000e+00 & -1.3878e-16 & 0 \\ 0 & 1.0000e+00 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

>> Apseudoinv * A

$\mathbf{A}^+\mathbf{A} = \mathbf{I}$ on \mathbb{R}^2

ans =

```

>> format short e
>> Apseudoinv2 = pinv(A)
= Matlab's built-in pseudo-inverse function
Apseudoinv2 =
  6.0000e-01  2.0000e-01  0  ← same as we got above!
 -4.0000e-01  2.0000e-01  0

>> A * Apseudoinv
Observe:
ans =
  AA+ = projection onto xy-plane
  1.0000e+00 -1.3878e-16  0
  0  1.0000e+00  0
  0  0  0

>> Apseudoinv * A
ans =
  A+A = I on R2
  1.0000e+00  1.1102e-16
  3.8858e-16  1.0000e+00

```

APPLICATION : LINEAR REGRESSION

Example: Find the equation for the line $y = a + bx$ that goes through the points

(2, 5) and (4, 11).

Answer: $a + 2 \cdot b = 5$ $\Leftrightarrow \begin{pmatrix} 1 & 2 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 5 \\ 11 \end{pmatrix}$

>> [1 2; 1 4] \ [5; 11]

ans =

$$\begin{pmatrix} -1 \\ 3 \end{pmatrix} \Rightarrow y = -1 + 3x$$

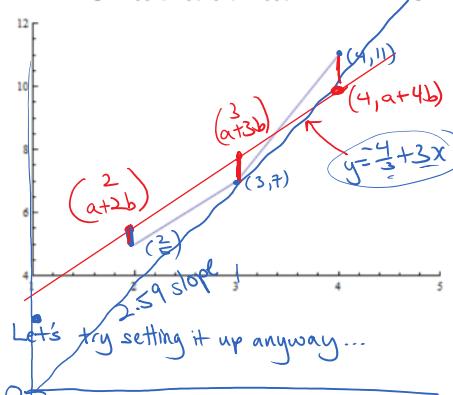
Observe: $\begin{pmatrix} 1 & 2 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 5 \\ 11 \end{pmatrix}$
 ↑ ↑ ↑
 all 1s x-values y-values
 to account for
 the affine shift a

Example: Find the equation for the line $y = a + bx$ that goes through the points

(2, 5) — $a + 2b = 5$
 (3, 7) — $a + 3b = 7$
 (4, 11) — $a + 4b = 11$

Answer: There is no such line!

ListPlot[{{2, 5}, {3, 7}, {4, 11}},
 PlotRange -> {{1, 5}, {4, 12}}, Joined -> True]



This value for \vec{x} minimizes
 $\|A\vec{x} - \vec{b}\|^2 = (a + 2b - 5)^2 + (a + 3b - 7)^2 + (a + 4b - 11)^2$
 $= \sum_{\text{data points}} (a + x_i b - y_i)^2$

$$\begin{pmatrix} 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 5 \\ 7 \\ 11 \end{pmatrix}$$

$$y = \vec{b}x$$

$$b = 2.59$$

```
>> pinv([1 2; 1 3; 1 4]) * [5;7;11]
```

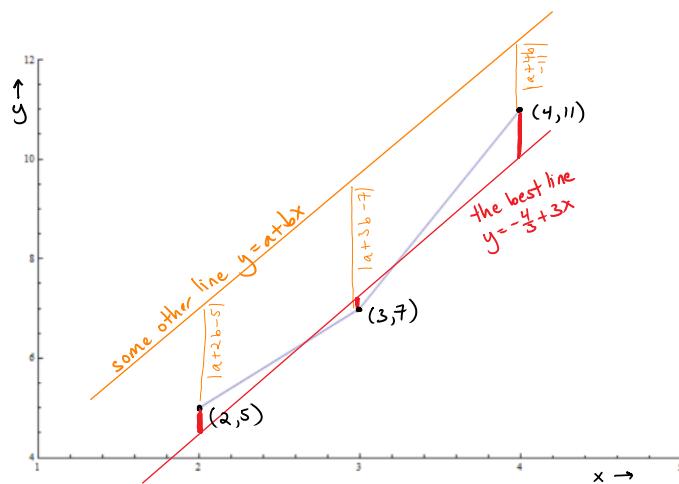
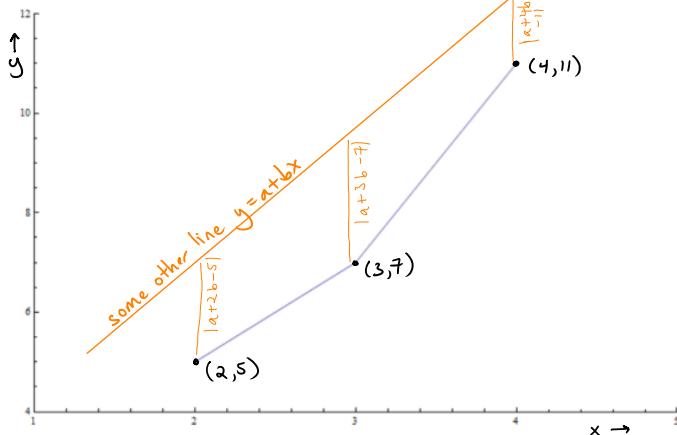
ans =

-1.3333
3.0000

$$\Rightarrow (a, b) = \left(-\frac{4}{3}, 3\right) \text{ minimizes } \left\| \begin{pmatrix} 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} - \begin{pmatrix} 5 \\ 7 \\ 11 \end{pmatrix} \right\|^2$$

$$= |a+2b-5|^2 + |a+3b-7|^2 + |a+4b-11|^2$$

the sum of the squared errors to each data point



Example:

Predict what percent of Americans 16 or older will be employed in 2050.

Answer:

- ① Get the data.

Google percent us population working

Web News Images Shopping Videos More Search tools

About 247,000,000 results (0.47 seconds)

According to the October jobs report, the seasonally adjusted employment-to-population ratio was 59.2% last month, one percentage point higher than it was a year earlier. Over that same period, the "official" unemployment rate fell from a seasonally adjusted 7.2% to 5.8%. Nov 7, 2014

Employment, unemployment and underemployment ...
www.pewresearch.org/.../employment-vs-unemployment... Pew Research Center

Feedback

→ **Employment-Population Ratio - Bureau of Labor Statistics ...**
data.bls.gov/timeseries/LNS12300000 Bureau of Labor Statistics
 Follow Us Follow BLS on Twitter | What's New | Release ... Labor Force Statistics from the Current Population Survey ... Type of data: Percent or rate, Age: 16 ...

Table A-1. Employment status of the civilian population by ...
www.bls.gov/news.release/empsit.t01.htm Bureau of Labor Statistics
 Follow Us Follow BLS on Twitter | What's New | Release Calendar ... Table A-1. Employment status of the civilian population by sex and age ... Employed: 146,941, 149,228, 148,980, 146,607, 148,795, 148,739, 148,840, 149,036, 148,800.

Employment-to-population ratio - Wikipedia, the free ...
https://en.wikipedia.org/wiki/Employment-to-population_ratio Wikipedia
 Employment-to-population ratio in the world[edit] In general, a high ratio is considered to be above 70 percent of the working-age population whereas a ratio below 50 percent is considered to be low.

UNITED STATES DEPARTMENT OF LABOR
BUREAU OF LABOR STATISTICS

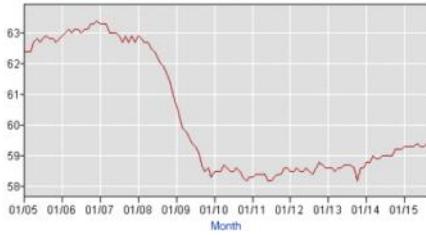
Home Subjects Data Tools Publications Economic Releases Students Beta

Databases, Tables & Calculators by Subject

Change Output Options: From: 2005 To: 2015 Go
 include graphs include annual averages
 Data extracted on: October 27, 2015 (10:44:11 AM) 1948

Labor Force Statistics from the Current Population Survey

Series Id: LNS12300000
 Seasonally Adjusted
 Series title: (Seas) Employment-Population Ratio
 Labor force status: Employment-population ratio
 Type of data: Percent or rate
 Age: 16 years and over



Download: [CSV](#) [XLS](#)

Year	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
2005	62.4	62.4	62.4	62.7	62.8	62.7	62.8	62.9	62.8	62.8	62.7	62.8
2006	62.9	63.0	63.1	63.0	63.1	63.1	63.0	63.1	63.1	63.1	63.3	63.4
2007	63.3	63.3	63.3	63.0	63.0	63.0	62.9	62.7	62.9	62.7	62.9	62.7
2008	62.9	62.8	62.7	62.7	62.5	62.4	62.2	62.0	61.9	61.7	61.4	61.0
2009	60.6	60.3	59.9	59.8	59.6	59.4	59.3	59.1	58.7	58.5	58.6	58.3
2010	58.5	58.5	58.5	58.7	58.6	58.5	58.5	58.6	58.5	58.3	58.2	58.3
2011	58.3	58.4	58.4	58.4	58.2	58.2	58.3	58.4	58.4	58.6	58.6	58.6
2012	58.5	58.5	58.6	58.5	58.5	58.6	58.5	58.4	58.6	58.8	58.7	58.6

② Load the data into Mathematica / Matlab, and set up the matrices

```
data = Import["employment-population.csv"];
plot1 = ListPlot[data[[;; -17]], PlotMarkers -> {Automatic, Small}, PlotStyle -> Blue];
plot2 = ListPlot[data[[16 ;;]], PlotMarkers -> {Automatic, Small}, PlotStyle -> Red];
plot = Show[plot1, plot2, PlotRange -> {{1948, 2015}, {55, 65}}];

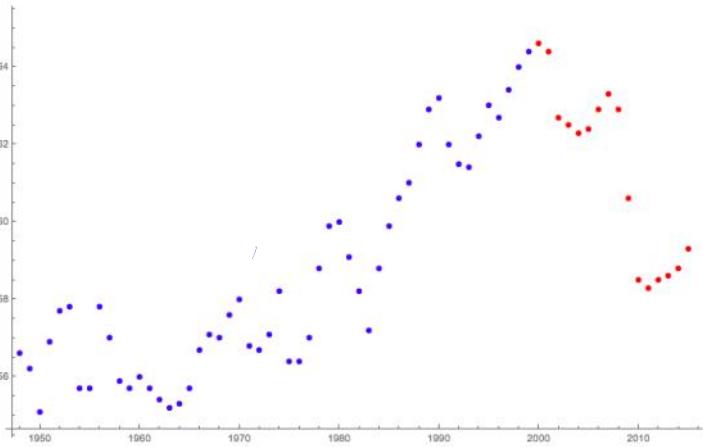
{{1948, 56.6}, {1949, 56.2}, {1950, 56.1}, {1951, 56.9}, {1952, 57.7}, {1953, 57.8}, {1954, 57.7}, {1955, 57.8}, {1956, 57.8}, {1957, 57.1}, {1958, 55.9}, {1959, 55.7}, {1960, 56.1}, {1961, 55.7}, {1962, 55.4}, {1963, 55.2}, {1964, 55.3}, {1965, 55.7}, {1966, 56.7}, {1967, 57.1}, {1968, 57.1}, {1969, 57.6}, {1970, 58.1}, {1971, 56.8}, {1972, 56.7}, {1973, 57.1}, {1974, 58.2}, {1975, 56.4}, {1976, 56.4}, {1977, 57.1}, {1978, 58.8}, {1979, 59.9}, {1980, 60.1}, {1981, 59.1}, {1982, 58.2}, {1983, 57.2}, {1984, 58.8}, {1985, 59.9}, {1986, 60.6}, {1987, 61.1}, {1988, 62.1}, {1989, 62.1}, {1990, 62.1}, {1991, 62.1}, {1992, 62.1}, {1993, 62.1}, {1994, 62.1}, {1995, 62.1}, {1996, 62.1}, {1997, 63.4}, {1998, 64.1}, {1999, 64.1}, {2000, 64.1}, {2001, 64.1}, {2002, 62.7}, {2003, 62.5}, {2004, 62.3}, {2005, 62.4}, {2006, 62.8}, {2007, 63.3}, {2008, 62.8}, {2009, 60.6}, {2010, 58.5}, {2011, 58.3}, {2012, 58.5}, {2013, 58.6}, {2014, 58.8}, {2015, 59.3}}
```



```

data = Import["employment-population.csv"]
plot1 = ListPlot[data[[;; -17]], PlotMarkers -> {Automatic, Small}, PlotStyle -> Blue];
plot2 = ListPlot[data[[-16 ;;]], PlotMarkers -> {Automatic, Small}, PlotStyle -> Red];
plot = Show[plot1, plot2, PlotRange -> {{1948, 2015}, {55, 65}}]

```



```

A = data[[-16 ;;, 1]];
A = {1, #} & /@ A ← insert 1's in each row
b = data[[-16 ;;, 2]];
{{1, 2000}, {1, 2001}, {1, 2002}, {1, 2003}, {1, 2004}, {1, 2005}, {1, 2006}, {1, 2007},
{1, 2008}, {1, 2009}, {1, 2010}, {1, 2011}, {1, 2012}, {1, 2013}, {1, 2014}, {1, 2015}}
{64.6, 64.4, 62.7, 62.5, 62.3, 62.4, 62.9, 63.3, 62.9, 60.6, 58.5, 58.3, 58.6, 58.8, 59.3}

```

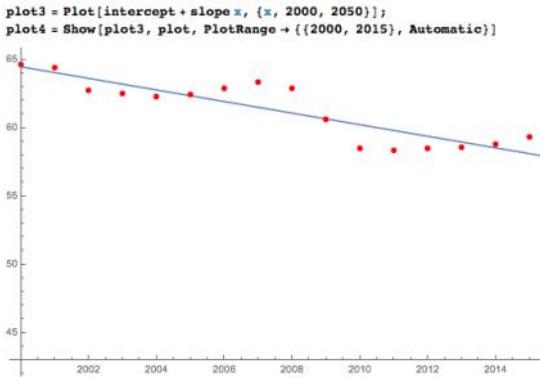
③ Compute the best-fitting line

```

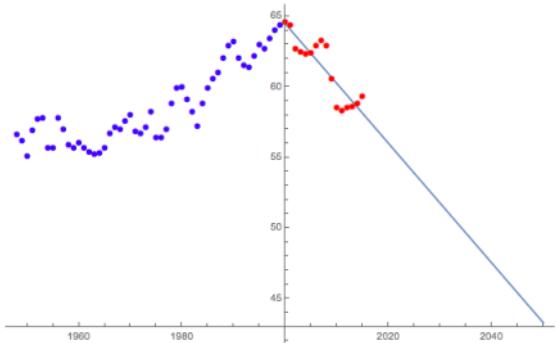
{intercept, slope} = PseudoInverse[A].b
{912.113, -0.423824}

```

④ Evaluate it, and get the answer!



```
Show[plot4, PlotRange -> {{1948, 2050}, Automatic}]
```



$$\begin{aligned}
 & \text{intercept + slope 2050} \\
 & \text{intercept + slope 1900} \\
 & \text{intercept + slope 2200} \\
 & 43.275 \text{ in 2050} \\
 & 106.849 \text{ in 1900} \\
 & -20.2985 \text{ in 2200}
 \end{aligned}$$

Extrapolations often don't work well. ☹

LET'S GENERALIZE!

Setting: m data points

Data pt. 1: $(x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_k^{(1)}, y^{(1)})$
 $(x_1^{(2)}, x_2^{(2)}, \dots, x_k^{(2)}, y^{(2)})$

Data pt. m : $(x_1^{(m)}, x_2^{(m)}, \dots, x_k^{(m)}, y^{(m)})$

components known exactly,
 e.g., date/time component
 we want to predict

Goal: Find the best linear predictor for y ,

$a_0, a_1, a_2, \dots, a_k \in \mathbb{R}$
 to minimize total squared error.

$$\sum_{j=1}^m |a_0 + a_1 x_1^{(j)} + \dots + a_k x_k^{(j)} - y^{(j)}|^2$$

Answer: linear/affine prediction

each row is one data point

$$\begin{pmatrix}
 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_k^{(1)} \\
 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_k^{(2)} \\
 1 & x_1^{(3)} & x_2^{(3)} & \dots & x_k^{(3)} \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_k^{(m)}
 \end{pmatrix}
 \begin{pmatrix}
 a_0 \\
 a_1 \\
 a_2 \\
 \vdots \\
 a_k
 \end{pmatrix} = \begin{pmatrix}
 y^{(1)} \\
 y^{(2)} \\
 y^{(3)} \\
 \vdots \\
 y^{(m)}
 \end{pmatrix}$$

\mathbf{A}

The least-squares solution is $\mathbf{A}^+ \vec{y}$.

Example:

Consider the time (T) it takes for a runner to complete a marathon (26 miles and 385 yards). Many factors such as height, weight, age, previous training, etc. can influence an athlete's performance, but experience has shown that the following three factors are particularly important:

$$x_1 = \text{Ponderal index} = \frac{\text{height (in.)}}{[\text{weight (lbs.)}]^{\frac{1}{3}}} = \frac{1}{\text{BMI}}$$

x_2 = Miles run the previous 8 weeks,

x_3 = Age (years).

A linear model hypothesizes that the time T (in minutes) is given by $T = a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + \varepsilon$, where ε is a random function accounting for all other factors and whose mean value is assumed to be zero. On the basis of the five observations given below, estimate the expected marathon time for a 43-year-old runner of height 74 in., weight 180 lbs., who has run 450 miles during the previous eight weeks. Answer: ~3:50

runner	time T (minutes)	x_1	x_2	x_3	
1	181	13.1	619	23	$a_1 < 0$
2	193	13.5	803	42	$a_2 < 0$
3	212	13.8	207	31	$a_3 > 0$
4	221	13.1	409	38	
5	248	12.51	482	45	

What is your personal predicted mean marathon time?

$$T = \frac{181 + 193 + 212 + 221 + 248}{5} = 210.2$$

$$T = a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_4$$

$$\min \left(\sum_{j=1}^m |a_0 + x_1^{(j)} a_1 + \dots + x_k^{(j)} a_k - y^{(j)}|^2 \right)$$

$$\rightarrow \text{int} \left(\begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_k^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_k^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} & \dots & x_k^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_k^{(m)} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_k \end{pmatrix} \right) = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \vdots \\ y^{(m)} \end{pmatrix}$$

⇒ if you want to minimize

$$\sum_{j=1}^m w_j |a_0 + a_1 x_1^{(j)} - y^{(j)}|^2$$

$$\sqrt{w} A(a) = \sqrt{w} y \quad \text{where } w = \begin{pmatrix} w_1 & 0 & \dots & w_m \end{pmatrix}$$

⇒ solution is

$$(\sqrt{w} A)^+ \sqrt{w} \vec{y} \neq A^+ \vec{y} \quad \text{in general}$$

$$[(AB)^- = B^- A^- \text{ if both are invertible}]$$

Answer: Start with a sanity check: $d_1 < 0, d_2 < 0$
 $d_3 > 0$

① Enter the data ② Find best-fit plane ③ Predict!

```
>> A = [
1 13.1 619 23 ;
1 13.5 803 42 ;
1 13.8 207 31 ;
1 13.1 409 38 ;
1 12.5 482 45 ;
];
x1 xL x3
>> b = [
181; 193; 212; 221; 248
];
A
times in minutes
```

$\alpha = \text{pinv}(A) * b$
 $\alpha' = [1; 74/180^{(1/3)}; 450; 43]$
 $\alpha =$
 $[492.0442$
 -23.4355
 -0.8761
 $1.8624]$
every year older
= 2 miles
more/week
= 2 minutes slower
= 1 min. faster

Example (continued): Same problem, but now fit a quadratic curve to x_2 (distance over previous 8 weeks) — running too much slows you down?

Answer: Want $T = d_0 + d_1 x_1 + d_2 x_2 + d_3 x_3 + d_4 x_2^2$

① Enter the data & find best-fit hyperplane

```
>> A2 = [A, A(:,3).^2];
alpha2 = pinv(A2) * b
>> alpha2' = [1; 74/180^{(1/3)}; 450; 43; 450.^2]
```

② Predict!

① Enter the data & find best-fit hyperplane

```
>> A2 = [A, A(:,3).^2];
alpha2 = pinv(A2) * b
alpha2 =
748.7697
-38.7075
-0.2441
1.5801
0.0002
ans =
224.3594 = 3 hours 44 minutes
```

(Of course, be careful of overfitting the data by using too many parameters.)

Fitting other curves than lines:

Example: Predict US gross domestic product (GDP) in 2050.

Answer: Google "US GDP"



An exponential should fit the data better than a straight line

$$G = e^{a+bT}$$

$$\Rightarrow \log G = a + bT$$

```
years = [1966:2014];
GDP = [3800176644539.76, 389518060653.26, 4082149751564.6, 4208696393863.11, 4343661175265.73, 4486805518455.5,
4722957883148.85, 4989480292992.24, 4963676968025.7, 4953864844037.75, 5220684465530.35, 5461284794411.99,
5765024247749.54, 5948103589527.53, 5933554752676.08, 6087499073703.77, 5971173597636.82, 6247785657762.76,
6701317851694.21, 6985369125864.67, 7230668360903.64, 748097585592.56, 7795473984765.14, 8082388278264.36,
8237519238199.96, 8231416510730, 8524075976250.96, 8758134889170.93, 9111757146664.33, 9359503617414.74,
9714779258395.62, 10150683977472.6, 10602380376634.4, 11099123060521.4, 11553318760428.6, 11666077052746.7,
11874448085025.7, 12207737238841, 12669890778469.9, 13093726000000, 13442886679117.7, 13681977860942.7,
13642078277526.4, 13263438360402.5, 13599258090661.7, 13817044044776, 14137749306899.7, 14451509918869.6,
14796640462032];
>> A = [ones(length(years), 1) years'];
>> b = log(GDP)';
>> x = pinv(A) * b;
>> x = A \ b;
```

A =

	x =	x =
1	1966	-29.0848
1	1967	0.0295
1	1968	-29.0848
1	1969	0.0295
1	1970	-29.0848
1	1971	0.0295
1	1972	-29.0848

ans =

$$4.7132e+13 = \$47\text{ trillion}$$

Also easy to fit other functions, e.g., $a + bT + cT^2 = \log G$:

```
>> B = [A (years') .^2]
>> y = B \ b
>> exp([1 2050 2050^2] * y)
B =
y =
ans =
1 1966 3865156 -588.2974
1 1967 3869089 0.5916
1 1968 3873024 -0.0001
1 1969 3876961
1 1970 3880900
1 1971 3884841
1 1972 3888794
```

Four ways to find x to minimize $\|Ax - b\|$:

② Compute $P_{R(A)} b$ (e.g., using Gram-Schmidt)
(to get on. basis for $R(A)$)

Solve $\vec{Ax} = P_{R(A)} \vec{b}$

① Compute pseudoinverse A^+ , set $\vec{x} = A^+ \vec{b}$.

Pros: Easy to remember, one line $\text{pinv}(A) * b$ in Matlab

Cons: Slow for large A , numerically unstable

/ Caution! Generalized inverses are useful in formulating theoretical statements \

Pros: Easy to remember, one line $\text{pinv}(A) * b$ in Matlab

Cons: Slow for large A , numerically unstable

Meyer p.423: Caution! Generalized inverses are useful in formulating theoretical statements such as those above, but, just as in the case of the ordinary inverse, generalized inverses are not practical computational tools. In addition to being computationally inefficient, serious numerical problems result from the fact that A^+ need not be a continuous function of the entries of A .

② Find x manually, with calculus:

$$E = \left\| \begin{pmatrix} 1 & a_1 \\ 1 & a_2 \\ 1 & a_3 \\ \vdots & \vdots \\ 1 & a_m \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix} \right\|^2$$

$$= \sum_{j=1}^m (x_1 + a_j x_2 - b_j)^2$$

$$\frac{\partial E}{\partial x_1} = 2 \sum_j (x_1 + a_j x_2 - b_j) = 0 \Rightarrow m x_1 + (\sum_j a_j) x_2 = (\sum_j b_j)$$

$$\frac{\partial E}{\partial x_2} = 2 \sum_j a_j (x_1 + a_j x_2 - b_j) = 0 \Rightarrow (\sum_j a_j) x_1 + (\sum_j a_j^2) x_2 = \sum_j a_j b_j$$

$$\text{Let } \bar{a} = \frac{1}{m} \sum_j a_j \text{ (average value of } a_j)$$

$$\bar{b} = \frac{1}{m} \sum_j b_j \text{ (average of } b_j\text{'s)}$$

First equation

$$\Rightarrow \text{intercept } x_1 = \bar{b} - \bar{a} x_2$$

(this is 0 if the data is centered so $\bar{a} = \bar{b} = 0$)

Second equation

$$\Rightarrow x_2 = \frac{\sum_j a_j b_j - m \cdot \bar{a} \bar{b}}{\sum_j a_j^2 - m (\bar{a})^2}$$

$$\left(= \frac{\text{Cov}(A, B)}{\text{Var}(A)} \quad \text{if } A \text{ and } B \text{ are random variables} \right. \\ \left. \text{with } \Pr[A, B] = (a_j, b_j)] = \frac{1}{m}. \right)$$

Pros: None (okay, it gives an easy closed-form solution)

③ Solve $A^T A x = A^T b$ exact!

Claim: This has the same solution, i.e. $A^+ b$.

in other words, in Matlab

$$\text{pinv}(A) * b = (A^T A) \setminus (A^T b)$$

Proof: Plug in $x = A^+ b$:

$$A^T (A A^+) b = A^T P_{R(A)} \overset{\uparrow}{b} = A^T b . \quad \checkmark$$

$$\underset{P_{R(A)}}{\underset{\parallel}{A^T}} \quad \underset{\underset{\parallel}{P_{R(A)}}}{\underset{\parallel}{A^T}} = A^T$$

$$\underset{\underset{\parallel}{P_{R(A)}}}{\underset{\parallel}{A^T P_{R(A)}^T}} = A^T$$

Theorem: Let A be an $m \times n$ real matrix with rank n .

- The equation $A^T A x = A^T b$ is feasible
(i.e., it has a solution x)
- The unique solution is $A^+ b$

- (ie, it has a solution x)
- The unique solution is
$$(A^T A)^{-1} A^T b = A^+ b$$
pseudo inverse

Proof: There is a solution because

- $R(A^T A) = R(A^T)$ (shown above)
- $A^T b \in R(A^T)$

$$\Rightarrow A^T b \in R(A^T A) \checkmark$$

The solution is unique because $N(A) = \{0\}$
 (because $\text{rank}(A) = \dim R(A^T) = n$, by assumption,
 and by rank-nullity $\dim R(A^T) + \dim N(A) = n$)

So why is $(A^T A)^{-1} A^T b = A^+ b$, as claimed?

There are several ways of seeing it.

Algebraically: Using the SVD of A , $A = \sum_i \lambda_i u_i v_i^T$,

$$\begin{aligned} A^+ &= \sum_{i: \lambda_i > 0} \frac{1}{\lambda_i} v_i u_i^T \\ (A^T A)^{-1} A^T &= \left[\sum_i \lambda_i^2 v_i v_i^T \right]^{-1} \left(\sum_j \lambda_j v_j u_j^T \right) \\ &= \left(\sum_i \frac{1}{\lambda_i^2} v_i v_i^T \right) \left(\sum_j \lambda_j v_j u_j^T \right) \\ &\quad \left. \begin{array}{l} \text{since } \text{rank}(A) = n, \\ \lambda_1, \dots, \lambda_n \text{ are all } > 0 \end{array} \right. \\ &= \sum_i \frac{1}{\lambda_i} v_i u_i^T \\ &= A^+ \checkmark \end{aligned}$$

More geometrically:

$$\begin{aligned} x = A^+ b \text{ solves } Ax &= P_{R(A)} b \\ \Rightarrow A^T A x &= A^T P_{R(A)} b \\ &= A^T b \quad \text{since } A^T = A^T P_{R(A)} \\ &\quad (\text{in general } A = P_{R(A)} A P_{R(A^T)}) \quad \square \end{aligned}$$

Corollary: For any $m \times n$ real matrix A , of rank n ,

$$A^+ = (A^T A)^{-1} A^T$$

Proof: Since we just showed $A^+ b = (A^T A)^{-1} A^T b$
 for any b . $A^T A \vec{x} = A^T \vec{b}$ □

Pros: Fast

Cons: condition number of $A^T A$

$$= (\text{condition number of } A)^2$$

Method

④ Solve $\begin{pmatrix} I & A \\ A^T & 0 \end{pmatrix} \begin{pmatrix} y \\ x \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix}$

Pros: Maintains sparsity of A .

[Problem 4.6.9, p. 237 of Meyer]

③ Latest research has focused on finding fast randomized approximation algorithms, based on dimension reduction
 à la Johnson-Lindenstrauss Lemma, e.g.,

[Clarkson Woodruff 2012 <http://arxiv.org/abs/1207.6365>]

[Nelson & Nguyen 2012 <http://arxiv.org/abs/1211.1002>]

$$A = \begin{pmatrix} 100 & 0 \\ 0 & 1 \end{pmatrix} \quad \|A^{-1}\| = 1$$

$$B = \begin{pmatrix} 100 & 0 \\ 0 & \frac{1}{1000} \end{pmatrix} \quad \|B^{-1}\| = 1000$$

A sing values $\sigma_1 > \sigma_2 > \dots > \sigma_n > 0$

$$F = \begin{pmatrix} 1 & 1 & -1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{pmatrix}$$

$$G = \begin{pmatrix} 1 & -2 & 0 \\ 0 & 1 & 1 \end{pmatrix}$$

$$\|A\| = \sqrt{\lambda_{\max}(A^T A)}$$

$$\tilde{v} = \frac{Fu}{\|Fu\|}, \quad \sigma = \|Fu\|$$

$$v = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$

$$G^T = \begin{pmatrix} 1 & 0 \\ -2 & 1 \\ 0 & 1 \end{pmatrix} \left| \begin{pmatrix} 1 \\ y \\ 1 \end{pmatrix} \right. \quad \frac{d}{dy} \frac{\|G^T(y)\|^2}{\|y\|^2} = 0$$

$$\left\| \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \right\| = \sqrt{2}$$

$$\left\| \begin{pmatrix} 1 & 0 \\ -0.01 & 1 \end{pmatrix} \right\| = \frac{1}{\sqrt{1+0.01}} \begin{pmatrix} 1 \\ 0.01 \end{pmatrix}$$

$$\begin{aligned} \mathbb{E}_{x \sim D} [(f(x) - p(x))^2] \\ &= \int dx q_D(x) (f(x) - p(x))^2 \\ &= \|f - p\|^2 = \langle f - p, f - p \rangle \\ \text{where } \langle a, b \rangle &= \int dx q_D(x) a(x) b(x) \end{aligned}$$

