

Pr. 1. (sol/hs121a)

- (a) Suppose
- \mathbf{A}
- is
- $p \times m$
- ,
- \mathbf{X}
- is
- $m \times n$
- , and
- \mathbf{B}
- is
- $q \times n$
- . Let us represent these matrices as

$$\mathbf{X} = [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_n], \quad \mathbf{B}^T = [\mathbf{b}_1 \quad \dots \quad \mathbf{b}_q].$$

Using these definitions, we have

$$\begin{aligned} \mathbf{A}\mathbf{X}\mathbf{B}^T &= \mathbf{A} [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_n] [\mathbf{b}_1 \quad \dots \quad \mathbf{b}_q] = \mathbf{A} \left[\left(\sum_{k=1}^n b_{1k} \mathbf{x}_k \right) \quad \dots \quad \left(\sum_{k=1}^n b_{qk} \mathbf{x}_k \right) \right] \\ &= \left[\mathbf{A} \left(\sum_{k=1}^n b_{1k} \mathbf{x}_k \right) \quad \dots \quad \mathbf{A} \left(\sum_{k=1}^n b_{qk} \mathbf{x}_k \right) \right] \Rightarrow \text{vec}(\mathbf{A}\mathbf{X}\mathbf{B}^T) = \begin{bmatrix} \mathbf{A} \sum_{k=1}^n b_{1k} \mathbf{x}_k \\ \vdots \\ \mathbf{A} \sum_{k=1}^n b_{qk} \mathbf{x}_k \end{bmatrix}. \end{aligned}$$

On the other hand,

$$(\mathbf{B} \otimes \mathbf{A}) \text{vec}(\mathbf{X}) = \begin{bmatrix} b_{11} \mathbf{A} & \dots & b_{1n} \mathbf{A} \\ \vdots & \ddots & \vdots \\ b_{q1} \mathbf{A} & \dots & b_{qn} \mathbf{A} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} = \begin{bmatrix} \sum_{k=1}^n b_{1k} \mathbf{A} \mathbf{x}_k \\ \vdots \\ \sum_{k=1}^n b_{qk} \mathbf{A} \mathbf{x}_k \end{bmatrix} = \begin{bmatrix} \mathbf{A} \sum_{k=1}^n b_{1k} \mathbf{x}_k \\ \vdots \\ \mathbf{A} \sum_{k=1}^n b_{qk} \mathbf{x}_k \end{bmatrix}.$$

This establishes the vec trick equality.

- (b) Multiplying
- $(\mathbf{A}\mathbf{X})\mathbf{B}^T$
- requires
- $N^3 + N^3 = 2N^3$
- scalar multiplications for
- $N \times N$
- dense matrices.

In contrast, computing the Kronecker product $\mathbf{Y} = \mathbf{B} \otimes \mathbf{A}$ requires N^4 multiplies, and then the product $(\mathbf{B} \otimes \mathbf{A}) \text{vec}(\mathbf{X}) = \mathbf{Y} \text{vec}(\mathbf{X})$ needs another N^4 multiplies, for a total of $2N^4$. Clearly the LHS is more efficient. But the RHS often is useful for paper analysis.

Pr. 2. (sol/hs029)

(This solution considers a more general case with some $\theta > 0$. Note that $\theta = 1$ in the problem statement.)

- (a) Method 1. An eigenvalue λ of $\mathbf{B} = \mathbf{A} + \theta \mathbf{x} \mathbf{x}'$ satisfies $\det(\mathbf{B} - \lambda \mathbf{I}) = \det(\mathbf{A} + \theta \mathbf{x} \mathbf{x}' - \lambda \mathbf{I}) = 0$. Note that when $\lambda \neq A_{ii}$ for any i then $\mathbf{A} - \lambda \mathbf{I}$ is invertible, so for now we focus on such values of λ . Thus

$$\begin{aligned} \det(\mathbf{A} + \theta \mathbf{x} \mathbf{x}' - \lambda \mathbf{I}) &= \det((\mathbf{A} - \lambda \mathbf{I}) \cdot (\mathbf{I} + (\mathbf{A} - \lambda \mathbf{I})^{-1} \theta \mathbf{x} \mathbf{x}')) = \det(\mathbf{A} - \lambda \mathbf{I}) \cdot \det(\mathbf{I} + \underbrace{\theta (\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{x} \mathbf{x}'}_{\triangleq \mathbf{w}}) \\ &= \det(\mathbf{A} - \lambda \mathbf{I}) \cdot (1 + \underbrace{\theta \mathbf{x}' (\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{x}}_{=w}). \quad (\text{From from HW 1, } \det(\mathbf{I} + \mathbf{w} \mathbf{x}') = 1 + \mathbf{x}' \mathbf{w}.) \end{aligned}$$

Thus an eigenvalue of \mathbf{B} that is not an eigenvalue of \mathbf{A} must satisfy the equation

$$\mathbf{x}' (\lambda \mathbf{I} - \mathbf{A})^{-1} \mathbf{x} = \frac{1}{\theta}.$$

Since \mathbf{A} is diagonal, a_{ii} are the eigenvalues of \mathbf{A} , and we obtain the (implicit) relationship

$$\frac{|x_1|^2}{\lambda - a_{11}} + \frac{|x_2|^2}{\lambda - a_{22}} + \dots + \frac{|x_n|^2}{\lambda - a_{nn}} = \frac{1}{\theta}. \quad (1)$$

Because $x_i \neq 0$ is given, this equation will have n solutions for λ whenever \mathbf{A} does not have repeated eigenvalues. When \mathbf{A} has repeated eigenvalues, some of the eigenvalues of \mathbf{A} and \mathbf{B} will coincide because the degree of the above polynomial will not be n . The problem statement says \mathbf{A} has distinct entries so this technicality is avoided.

Note that the solutions of (1) will *not* correspond to any of the diagonal elements of \mathbf{A} , so our use of the inverse of $\mathbf{A} - \lambda \mathbf{I}$ was legitimate.

Method 2. Alternatively, one can use properties 16, 17 in Section 1.4 of Laub (see HW1, Question 3), as follows:

$$\begin{aligned} \det(\mathbf{A} - \lambda \mathbf{I} + \theta \mathbf{x} \mathbf{x}') &= \det \left(\begin{bmatrix} \mathbf{A} - \lambda \mathbf{I} & -\theta \mathbf{x} \\ \mathbf{x}' & 1 \end{bmatrix} \right) \\ &= \det(\mathbf{A} - \lambda \mathbf{I}) \det(1 + \theta \mathbf{x}' (\mathbf{A} - \lambda \mathbf{I})^{-1} \mathbf{x}). \end{aligned}$$

Method 3. Let z be an eigenvalue of \mathbf{B} associated with the eigenvector \mathbf{v} . Then

$$\begin{aligned} \mathbf{B} \mathbf{v} &= z \mathbf{v} \Rightarrow (\mathbf{A} + \theta \mathbf{x} \mathbf{x}') \mathbf{v} = z \mathbf{v} \Rightarrow (z \mathbf{I} - \mathbf{A}) \mathbf{v} = \theta \mathbf{x} \mathbf{x}' \mathbf{v} \\ &\Rightarrow \mathbf{v} = \theta (z \mathbf{I} - \mathbf{A})^{-1} \mathbf{x} \mathbf{x}' \mathbf{v} \quad (\text{for } z \neq \lambda(\mathbf{A})) \\ &\Rightarrow (\mathbf{x}' \mathbf{v}) = \theta \mathbf{x}' (z \mathbf{I} - \mathbf{A})^{-1} \mathbf{x} (\mathbf{x}' \mathbf{v}) \\ &\Rightarrow 1 = \theta \mathbf{x}' (z \mathbf{I} - \mathbf{A})^{-1} \mathbf{x} \quad (\text{assuming that the scalar } \mathbf{x}' \mathbf{v} \neq 0). \end{aligned}$$

The assumption that $\mathbf{x}' \mathbf{v} \neq 0$ implies that $\mathbf{v} \notin \mathcal{N}(\mathbf{x})$. More work would be needed to complete this argument, although the equivalence with the previous derivations shows that this assumption is valid whenever $x_i \neq 0$ when considering the eigenvalues of \mathbf{B} that are not equal to the eigenvalues of \mathbf{A} .

Here is a proof that a_{ii} cannot be an eigenvalue of \mathbf{B} under the given conditions. When $\lambda = a_{ii}$, clearly the corresponding eigenvector of \mathbf{A} is \mathbf{e}_i . Now suppose that λ is also an eigenvalue of \mathbf{B} . Then we would have for some vector \mathbf{v} : $\mathbf{B} \mathbf{v} = (\mathbf{A} + \theta \mathbf{x} \mathbf{x}') \mathbf{v} = \mathbf{A} \mathbf{v} + \theta \mathbf{x} (\mathbf{x}' \mathbf{v}) = \lambda \mathbf{v}$. Now examine the i th element of both sides of that last equality: $a_{ii} v_i + \theta x_i (\mathbf{x}' \mathbf{v}) = a_{ii} v_i \Rightarrow x_i (\mathbf{x}' \mathbf{v}) = 0$. Because each element of \mathbf{x} is nonzero, it must be the case that $\mathbf{x}' \mathbf{v} = 0$. So we have $\mathbf{B} \mathbf{v} = \mathbf{A} \mathbf{v} = \lambda \mathbf{v} = a_{ii} \mathbf{v}$, so \mathbf{v} must be the eigenvector of \mathbf{A} corresponding to eigenvalue a_{ii} , which means $\mathbf{v} = \mathbf{e}_i$. But that would mean $\mathbf{x}' \mathbf{e}_i = 0$, which would contradict the assumption that all elements of \mathbf{x} are nonzero.

Here is a proof that \mathbf{B} can have some eigenvalues that are the same as those of \mathbf{A} when \mathbf{A} is diagonal with some repeated values. For illustration (and without loss of generality via permutations), suppose $\mathbf{A} = \begin{bmatrix} \lambda \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}$ where the elements of the diagonal matrix \mathbf{D} are distinct from λ and where $k \geq 2$. Similarly partition the vector \mathbf{x} as $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$ where \mathbf{x}_1 has length k and \mathbf{x}_2 has length $n - k$. Then $\mathbf{B} = \mathbf{A} + \mathbf{x}\mathbf{x}' = \begin{bmatrix} \lambda \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \begin{bmatrix} \mathbf{x}_1' & \mathbf{x}_2' \end{bmatrix}$. Now pick $\mathbf{v}_1 \in \mathcal{N}(\mathbf{x}_1') \neq \mathbf{0}$ which exists because $k \geq 2$ and define $\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{0} \end{bmatrix}$. Then one can see that $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ and $\mathbf{B}\mathbf{v} = \lambda\mathbf{v}$ so in this case \mathbf{A} and \mathbf{B} have a common eigenvalue (and eigenvector). This is why the problem stated that the elements of \mathbf{A} are distinct.

(b) When $\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}$ and $\mathbf{x} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$, then (1) becomes

$$\frac{1}{\lambda - 1} + \frac{1}{\lambda - 2} + \frac{1}{\lambda - 3} = 1.$$

Multiplying through, this becomes a cubic equation:

$$(z - 2)(z - 3) + (z - 1)(z - 3) + (z - 1)(z - 2) = (z - 1)(z - 2)(z - 3).$$

Using Julia's `Polynomials.jl` package, this simplifies to $0 = -17 + 23z - 9z^2 + z^3$, the roots of which are approximately $\{1.32, 2.46, 5.21\}$.

```
using Polynomials: fromroots, roots
fromroot(root::Real) = fromroots([root])
(p1, p2, p3) = fromroot.(1:3)
# solve: (z-2)(z-3) + (z-1)(z-3) + (z-1)(z-2) = (z-1)(z-2)(z-3)
p = p1*p2*p3 - p1*p2 - p2*p3 - p1*p3
roots(p) # basic version that does not generalize easily to n > 3
#q = +(1 ./ fromroot.(1:3)... , -1); roots(q) # one line version: generalizes!
```

Pr. 3. (sol/hs069)

Define $\mathbf{A} = [\mathbf{q}_1 \quad \mathbf{q}_2]$, and $\mathbf{x} = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}$.

- (a) Approach 1 (using normal equations) The columns of \mathbf{A} are orthonormal, so the Gram matrix is $\mathbf{A}'\mathbf{A} = \mathbf{I}_2$. Thus from the normal equations the unique LLS solution is

$$\hat{\mathbf{x}} = (\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}'\mathbf{b} = \mathbf{A}'\mathbf{b}.$$

This expression is unsurprising, because what we get is precisely the first two “coordinates” of the vector \mathbf{b} relative to the basis whose first two basis vectors correspond to the columns of \mathbf{A} .

Thus the final answer for the optimal linear combination is

$$\mathbf{A}\hat{\mathbf{x}} = \mathbf{A}\mathbf{A}'\mathbf{b}.$$

Approach 2 (using compact SVD)

The columns of \mathbf{A} are linearly independent, so the unique **linear least squares estimate** that minimizes $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$ is given by $\hat{\mathbf{x}} = \mathbf{A}^+\mathbf{b}$. A compact SVD of \mathbf{A} is simply:

$$\mathbf{A} = \sum_{i=1}^2 1 \mathbf{q}_i \mathbf{e}_i',$$

where \mathbf{e}_i denotes the i th unit vector. Thus:

$$\mathbf{A}^+ = \sum_{i=1}^2 1 \mathbf{e}_i \mathbf{q}_i',$$

and hence

$$\hat{\mathbf{x}} = \mathbf{A}^+\mathbf{b} = \sum_{i=1}^2 1 \mathbf{e}_i \mathbf{q}_i' \mathbf{b} = \mathbf{A}'\mathbf{b} \Rightarrow \mathbf{A}\hat{\mathbf{x}} = \mathbf{A}\mathbf{A}'\mathbf{b}.$$

- (b) Approach 1

Here we have that the residual (or error) vector:

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}} = \mathbf{b} - \mathbf{A}(\mathbf{A}'\mathbf{b}) = (\mathbf{I} - \mathbf{A}\mathbf{A}')\mathbf{b} \Rightarrow \langle \mathbf{r}, \mathbf{q}_i \rangle = \mathbf{q}_i' \mathbf{r} = \mathbf{q}_i' (\mathbf{I} - \mathbf{A}\mathbf{A}')\mathbf{b} = (\mathbf{q}_i' - \mathbf{e}_i' \mathbf{A}')\mathbf{b} = (\mathbf{q}_i' - \mathbf{q}_i')\mathbf{b} = 0.$$

Approach 2

Here we have that the residual (or error) vector:

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}} = (\mathbf{I} - \mathbf{A}\mathbf{A}')\mathbf{b} = \sum_{i=3}^n \mathbf{q}_i \mathbf{q}_i' \mathbf{b},$$

where \mathbf{q}_i for $i=3, \dots, n$ are the $n-2$ (unit norm) basis vectors, orthogonal to \mathbf{q}_1 and \mathbf{q}_2 so that $\text{span}(\{\mathbf{q}_1, \dots, \mathbf{q}_n\}) = \mathbb{R}^n$. Hence $\mathbf{q}_1' \mathbf{r} = \sum_{i=3}^n \mathbf{q}_1' \mathbf{q}_i \mathbf{q}_i' \mathbf{b} = 0$ and $\mathbf{q}_2' \mathbf{r} = \sum_{i=3}^n \mathbf{q}_2' \mathbf{q}_i \mathbf{q}_i' \mathbf{b} = 0$.

This property is related to the **projection theorem**.

Pr. 4. (sol/hs049)

- Approach 1. No SVD, just using properties of **pseudo-inverse**:

$$\begin{aligned} \mathbf{x}'(\mathbf{I} - \mathbf{A}^+ \mathbf{A})'(\mathbf{A}^+ \mathbf{b}) &= \mathbf{x}'(\mathbf{I} - (\mathbf{A}^+ \mathbf{A})')(\mathbf{A}^+ \mathbf{b}) = \mathbf{x}'(\mathbf{I} - (\mathbf{A}^+ \mathbf{A}))(\mathbf{A}^+ \mathbf{b}) \\ &= \mathbf{x}'(\mathbf{A}^+ - \mathbf{A}^+ \mathbf{A} \mathbf{A}^+) \mathbf{b} = \mathbf{x}'(\mathbf{A}^+ - \mathbf{A}^+) \mathbf{b} = \mathbf{0} \end{aligned}$$

- Approach 2. **full SVD**:

$$\begin{aligned} (\mathbf{A}^+ \mathbf{b})'(\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \mathbf{x} &= \mathbf{b}'(\mathbf{A}^+)'(\mathbf{I} - \mathbf{V} \Sigma^+ \mathbf{U}' \mathbf{U} \Sigma \mathbf{V}') \mathbf{x} = \mathbf{b}' \mathbf{U} (\Sigma^+)' \mathbf{V}' (\mathbf{V} \mathbf{V}' - \mathbf{V} \Sigma^+ \Sigma \mathbf{V}') \mathbf{x} \\ &= \mathbf{b}' \mathbf{U} ((\Sigma^+)' - (\Sigma^+)' \Sigma^+ \Sigma) \mathbf{V}' \mathbf{x} = \mathbf{0}, \end{aligned}$$

because direct multiplication verifies that $(\Sigma^+)' = (\Sigma^+)' \Sigma^+ \Sigma$.

- Approach 3. **compact SVD**: $\mathbf{A} = \mathbf{U}_r \Sigma_r \mathbf{V}_r' \Rightarrow \mathbf{A}^+ = \mathbf{V}_r \Sigma_r^{-1} \mathbf{U}_r'$ so:

$$(\mathbf{A}^+ \mathbf{b})'(\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \mathbf{x} = \mathbf{b}' \mathbf{U}_r \Sigma_r^{-1} \mathbf{V}_r' (\mathbf{I} - \mathbf{V}_r \mathbf{V}_r') \mathbf{x} = \mathbf{b}' \mathbf{U}_r \Sigma_r^{-1} (\mathbf{V}_r' - \mathbf{V}_r' \mathbf{V}_r \mathbf{V}_r') \mathbf{x} = \mathbf{0}, \quad \text{because } \mathbf{V}_r' \mathbf{V}_r = \mathbf{I}.$$

Put another way in terms of **subspaces**: $\mathbf{A}^+ \mathbf{b} \in \mathcal{R}(\mathbf{V}_r)$ whereas $(\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \mathbf{x} = (\mathbf{I} - \mathbf{V}_r \mathbf{V}_r') \mathbf{x} \in \mathcal{R}(\mathbf{V}_r)^\perp$.

Pr. 5. (sol/hs114)

- (a) We must show $\forall \alpha \in [0, 1]$ that: $\|\mathbf{A}(\alpha \mathbf{x} + (1 - \alpha) \mathbf{z}) - \mathbf{b}\|_2 \leq \alpha \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2 + (1 - \alpha) \|\mathbf{A} \mathbf{z} - \mathbf{b}\|_2$.

For any $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$:

$$\begin{aligned} f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{z}) &= \|\mathbf{A}(\alpha \mathbf{x} + (1 - \alpha) \mathbf{z}) - \mathbf{b}\|_2 \\ &= \|\mathbf{A}(\alpha \mathbf{x} + (1 - \alpha) \mathbf{z}) - \mathbf{b}(\alpha + 1 - \alpha)\|_2 \\ &= \|(\alpha \mathbf{A} \mathbf{x} - \alpha \mathbf{b}) + ((1 - \alpha) \mathbf{A} \mathbf{z} - (1 - \alpha) \mathbf{b})\|_2 \\ &\leq \alpha \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2 + \|(1 - \alpha)(\mathbf{A} \mathbf{z} - \mathbf{b})\|_2, \quad (\text{via the triangle inequality}) \\ &= |\alpha| \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2 + |1 - \alpha| \|\mathbf{A} \mathbf{z} - \mathbf{b}\|_2 \\ &= \alpha \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2 + (1 - \alpha) \|\mathbf{A} \mathbf{z} - \mathbf{b}\|_2, \quad (\text{because } \alpha \in [0, 1]) \\ &= \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{z}). \end{aligned}$$

- (b) Next we need to show that, given two matrices \mathbf{A}, \mathbf{B} , $\forall \alpha \in [0, 1]$:

$$\sigma_1(\alpha \mathbf{A} + (1 - \alpha) \mathbf{B}) \leq \alpha \sigma_1(\mathbf{A}) + (1 - \alpha) \sigma_1(\mathbf{B}).$$

From the hint:

$$\begin{aligned} \sigma_1(\alpha \mathbf{A} + (1 - \alpha) \mathbf{B}) &\triangleq \max_{\|\mathbf{u}\|_2=1} \|(\alpha \mathbf{A} + (1 - \alpha) \mathbf{B}) \mathbf{u}\|_2 \\ &\leq \max_{\|\mathbf{u}\|_2=1} (\|\alpha \mathbf{A} \mathbf{u}\|_2 + \|(1 - \alpha) \mathbf{B} \mathbf{u}\|_2) \quad (\text{by the triangle inequality}) \\ &= \max_{\|\mathbf{u}\|_2=1} (\alpha \|\mathbf{A} \mathbf{u}\|_2 + (1 - \alpha) \|\mathbf{B} \mathbf{u}\|_2) \\ &\leq \max_{\|\mathbf{u}\|_2=1} \alpha \|\mathbf{A} \mathbf{u}\|_2 + \max_{\|\mathbf{u}\|_2=1} (1 - \alpha) \|\mathbf{B} \mathbf{u}\|_2 \\ &= \alpha \sigma_1(\mathbf{A}) + (1 - \alpha) \sigma_1(\mathbf{B}). \end{aligned}$$

Thus, $\sigma_1(\cdot)$ is a convex function.

Pr. 6. (sol/hsj4s)

Here is the code for the LS solution for the coefficients.

```
using CSV
using Plots; default(markerstrokecolor=:auto)

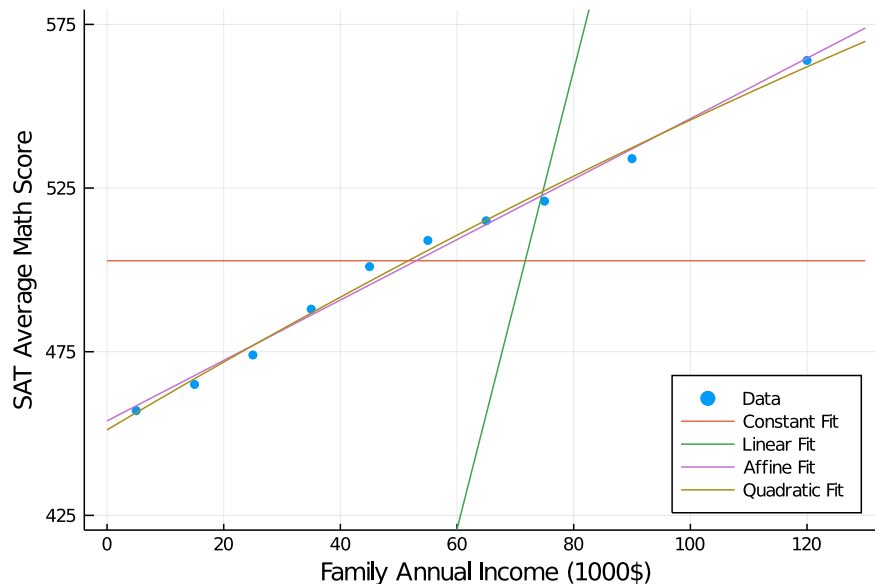
data = CSV.File("../prob/hpj4s.csv", normalizenames=true, transpose=true)
y = data.Math # math scores from the file (should be from 457 to 564)
income = [5:10:75; 90; 120]
scatter(income, y, label="Data", ylim=[425,575], ytick=425:50:575)
plot!(xlabel = "Family Annual Income (1000$)")
plot!(ylabel = "SAT Average Math Score")
plot!(legend = :bottomright)

As = [ones(length(income),1),
      reshape(income, :, 1),
      [income .^ 0 income .^ 1],
      [income .^ 0 income .^ 1 income .^ 2],
      ]
xh = Array{Any}(undef, 4)
for i=1:4
    xh[i] = As[i] \ y
end
display([xh[1]; 0; 0] [0; xh[2]; 0] [xh[3]; 0] xh[4]) # for table

inc = 0:130
x = xh[1]; plot!(inc, x[1] * inc.^0, label="Constant Fit")
x = xh[2]; plot!(inc, x[1] * inc, label="Linear Fit")
x = xh[3]; plot!(inc, x[1] + x[2] * inc, label="Affine Fit")
x = xh[4]; plot!(inc, x[1] + x[2] * inc + x[3] * inc.^2, label="Quadratic Fit")

#savefig("hsj4s.pdf")
```

(a) Scatter plot and fits:



Note that the “linear” model fits very poorly; often a “linear” model in the regression literature really means an affine model. One must be careful with the terminology: the quadratic model is quadratic in income, but it is a linear function of the coefficients $\{\beta_0, \beta_1, \beta_2\}$, so **linear regression** via a LS fit is still appropriate.

(b) Table of polynomial coefficients:

Constant Fit	Linear Fit	Affine Fit	Quadratic Fit
502.8	0	453.86	451.114
0	7.01	0.923	1.06
0	0	0	-0.0011

- (c) The β_1 coefficients of the affine and quadratic fits are both approximately 1, showing that for every \$1000 increase in family income, the SAT Math scores average about 1 point higher. Originally, SAT stood for “Scholastic Aptitude Test” but the dependence on family income raises doubts about how much it measures aptitude.

The β_2 coefficient is very small, but one must keep in mind that the $(\text{income})^2$ values can be very large, so one cannot conclude that it is unimportant solely by looking at its value. However, from the plot we see that the quadratic fit is nearly the same as the linear fit and there is no obvious evidence for a quadratic trend in this data. So we can conclude that average SAT scores have roughly an affine relationship with family income.

Tuesday October 05 2021 21:32
yuzhanji@umich.edu

Pr. 7. (sol/hs068)

Define for integer values d :

$$\mathbf{y} \triangleq \begin{bmatrix} f(t_1) \\ \vdots \\ f(t_{16}) \end{bmatrix} \in \mathbb{R}^{16}, \quad \mathbf{A} \triangleq \begin{bmatrix} 1 & t_1 & t_1^2 & \dots & t_1^d \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & t_{16} & t_{16}^2 & \dots & t_{16}^d \end{bmatrix} \in \mathbb{R}^{16 \times (d+1)}, \quad \mathbf{x} \triangleq \begin{bmatrix} x_1 \\ \vdots \\ x_{d+1} \end{bmatrix} \in \mathbb{R}^{d+1}.$$

The solution $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y}$ yields the the desired optimal least-squares estimate of the coefficients of the degree- d polynomial $p_d(t) = \sum_{i=1}^{d+1} x_i t^{i-1}$ that minimizes the residual norm $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$. The \mathbf{A} of the form above is called a **Vandermonde matrix**.

Here is Julia code:

```
using LinearAlgebra: norm
using Random: seed!
using Plots; default(markerstrokecolor=:auto)

T = LinRange(0, 2, 16)
seed!(3); e = randn(length(T))
b = 0.5 * exp.(0.8*T) # part a
y = b + e # part b
d = length(T)-1

A15 = [tt^j for tt in T, j=0:d]
A2 = A15[:,1:3]
x15b = A15 \ b # pinv also acceptable due to small problem size
x2b = A2 \ b
x15y = A15 \ y
x2y = A2 \ y
t = LinRange(0,2,500)

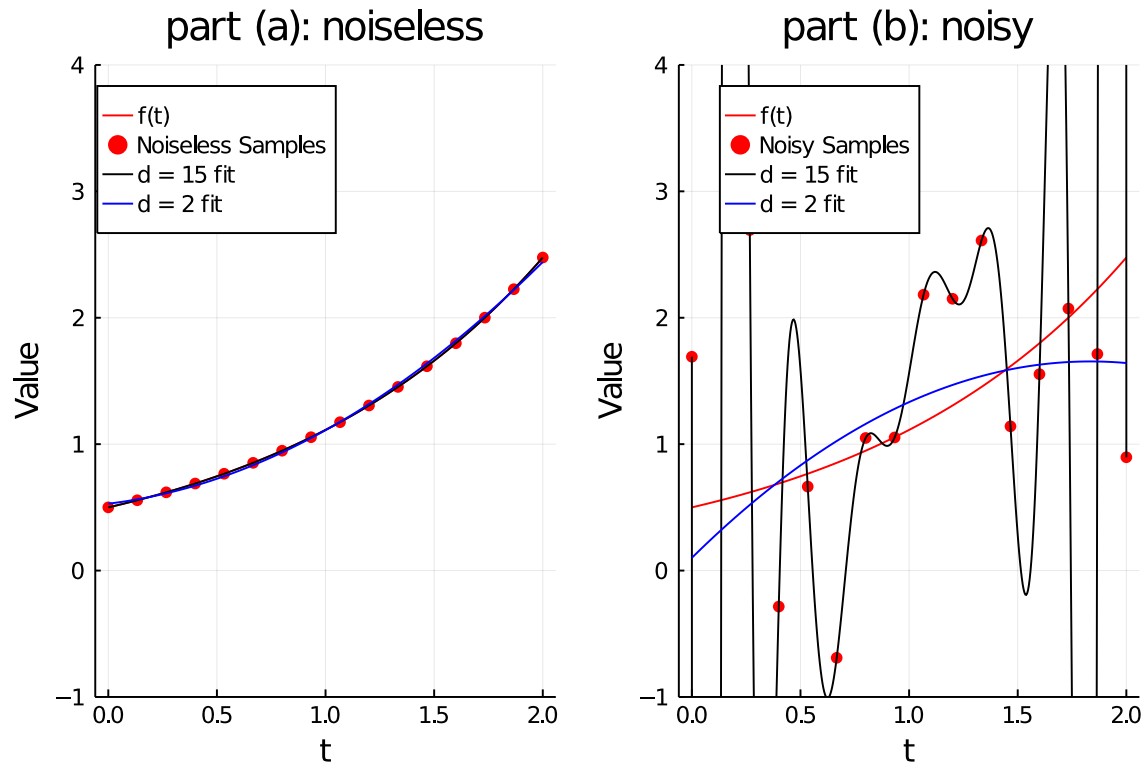
@show norm(A2 * x2b - b), norm(A15 * x15b - b)
@show norm(A2 * x2y - y), norm(A15 * x15y - y)
@show norm(A2 * x2y - b), norm(A15 * x15y - b)

a15 = [tt^j for tt in t, j=0:d]
a2 = a15[:,1:3]

plot(t, 0.5*exp.(0.8*t), color=:red, label = "f(t)", ylim=(-1,4))
scatter!(T,b, marker=:circle,:red), label = "Noiseless Samples")
plot!(t, a15*x15b, color=:black, label = "d = 15 fit", legend=:top)
plot!(t, a2*x2b, color=:blue, label = "d = 2 fit", xaxis = "t", yaxis = "Value")
plot!(title = "part (a): noiseless")
pa = current()

plot(t, 0.5*exp.(0.8*t), color=:red, label = "f(t)", ylim=(-1,4))
scatter!(T, y, marker=:circle,:red), label = "Noisy Samples")
plot!(t, a15*x15y, color=:black, label = "d = 15 fit", legend=:top)
plot!(t, a2*x2y, color=:blue, label = "d = 2 fit", xaxis = "t", yaxis = "Value")
plot!(title = "part (b): noisy")
pb = current()
plot(pa, pb)
#savefig("hs068.pdf")
```

- (a) For the noiseless case, this yielded the left figure below. Note that the error for $d = 15$ is much smaller than for $d = 2$. In fact the $d = 15$ error is exactly zero (to within numerical precision) because $\text{rank}(\mathbf{A}) = 16 = \dim(\mathbf{y})$ so the solution is exact! Comparing the least-squares coefficients obtained for $d = 2$ and $d = 15$ reveals that the first three terms are about the same.
- (b) When noise is added we get the right figure below. The polynomial for $d = 15$ passes *exactly* through the noisy samples. This property follows from the previous argument because $\text{rank}(\mathbf{A}) = 16$ so that error will be zero, relative to the samples, but large relative to the function. This example illustrates why choosing large model orders ($= d$) can make things worse by “over fitting” noise.



Use the `cond(A)` command in Julia to obtain the condition number: the ratio of the largest and smallest singular value of a matrix. The condition number is important because the measurement error manifests as a $\delta \mathbf{y}$ so that the error in the least squares solution (relative to the noise-less case) is given by exactly by $\delta \hat{\mathbf{x}} = \mathbf{A}^+(\mathbf{y} + \delta \mathbf{y}) - \mathbf{A}^+ \mathbf{y} = \mathbf{A}^+ \delta \mathbf{y}$. A large condition number means that small $\delta \mathbf{y}$ can get amplified, producing large $\delta \mathbf{x}$ and hence instabilities as manifested in the overfitting.

An estimate of the degree to which $\delta \mathbf{x}$ can change as a function of the condition number κ is given by the relation:

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \kappa \frac{\|\delta \mathbf{y}\|}{\|\mathbf{y}\|}.$$

Here $\|\delta \mathbf{y}\|$ is the norm of the noise. The condition number $\kappa = \sigma_1(\mathbf{A})/\sigma_{\min(m,n)}(\mathbf{A})$ is large when $d = 15$ (check it) compared to when $d = 2$.

polynomial degree:	$d = 2$	$d = 15$
(c) Residual norm $\ \mathbf{A}\hat{\mathbf{x}}(\mathbf{b}) - \mathbf{b}\ _2$ noiseless (a)	0.071	7.6e-16
Residual norm $\ \mathbf{A}\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{y}\ _2$ noisy (b)	4.438	8.3e-5
Fitting error $\ \mathbf{A}\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{b}\ _2$	1.271	4.641

Grader: the numerical values for $d = 15$ are unreliable; accept anything small (below $1\text{e-}4$) for the first two and anything larger than 1 for the last one.

- (d) The residual norm for degree 15 is smaller than that of degree 2 because $\mathbf{A}_{15} = [\mathbf{A}_2 \quad \mathbf{A}_{3:15}]$ so $\mathcal{R}(\mathbf{A}_2) \subseteq \mathcal{R}(\mathbf{A}_{15})$.

Non-graded problem(s) below**Pr. 8.** (sol/hs136)

- (a)
- $\mathbf{x}_{\text{opt}} = z\mathbf{v}_1$, $\mathbf{y}_{\text{opt}} = z\mathbf{u}_1$, where $|z| = 1$
 - \mathbf{x}_{opt} and \mathbf{y}_{opt} are unique to within a sign ambiguity when $N = 1$ or when $N > 1$ and $\sigma_1(\mathbf{A}) > \sigma_2(\mathbf{A})$.
 - Negating \mathbf{A} does not change the answer (aside from a potential sign flip), as the singular values are unchanged, but one of the left/right singular vectors will have a sign-flip.
 - Adding an orthogonality constraint of $\mathbf{x} \perp \text{span}(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{r-1})$ suffices. Equivalently, projecting the columns of \mathbf{A} onto the orthogonal complement of $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{r-1})$ also suffices, i.e., we replace \mathbf{A} with $(\mathbf{I} - \mathbf{P}_{\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{r-1})}) \mathbf{A}$.

We may also iterate this algorithm: given \mathbf{u}_1 through \mathbf{u}_{k-1} , we may seek \mathbf{u}_k that is orthogonal to $\text{span}(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k-1})$, or, replace \mathbf{A} with $(\mathbf{I} - \mathbf{P}_{\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{k-1})}) \mathbf{A}$.

Indeed, given $\mathbf{B} = (\mathbf{I} - \mathbf{P}_{\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{k-2})}) \mathbf{A}$, we may form $(\mathbf{I} - \mathbf{P}_{\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{k-1})}) \mathbf{A} = (\mathbf{I} - \mathbf{P}_{\mathbf{u}_{k-1}}) \mathbf{B}$.

- (b)
- Consider the vector $\mathbf{z} = \mathbf{V}'\mathbf{x}$. We know that $\|\mathbf{x}\|_2 = 1 \iff \|\mathbf{z}\|_2 = 1$. Now,

$$\mathbf{x}'\mathbf{A}\mathbf{x} = \mathbf{x}'\mathbf{V}\mathbf{\Lambda}\mathbf{V}'\mathbf{x} = \mathbf{z}'\mathbf{\Lambda}\mathbf{z} = \sum_{k=1}^N \lambda_k |z_k|^2 \leq \lambda_1 \sum_{k=1}^N |z_k|^2 = \lambda_1 \|\mathbf{z}\|_2^2 = \lambda_1.$$

This is an equality if $\mathbf{z} = z\mathbf{e}_1$ for $|z| = 1$. Hence, $\mathbf{x}_{\text{opt}} = z\mathbf{v}_1$.

If $\lambda_1 = \lambda_2$ then we also achieve equality when $\mathbf{z} = z\mathbf{e}_2$, so in general the solution is not unique.

- $\mathbf{x}'_{\text{opt}}\mathbf{A}\mathbf{x}_{\text{opt}} = z\mathbf{v}_1'\mathbf{V}\mathbf{\Lambda}\mathbf{V}'(z\mathbf{v}_1) = \lambda_1$. (This equality holds even in the non-unique case.)
 - The solution \mathbf{x}_{opt} will be unique (to within $|z| = 1$) when $\lambda_1 > \lambda_2$.
- (c)
- Use the same projection as the previous problem. If $\mathbf{x} \perp \mathbf{v}_1 \perp \dots \perp \mathbf{v}_{K-1}$, then $\mathbf{z} = \mathbf{V}'\mathbf{x}$ yields that $z_1 = z_2 = \dots = z_{K-1} = 0$. Hence, we have

$$\mathbf{x}'\mathbf{A}\mathbf{x} = \sum_{k=1}^N \lambda_k |z_k|^2 = \sum_{k=K}^N \lambda_k |z_k|^2 \leq \lambda_K \sum_{k=1}^N |z_k|^2 = \lambda_K.$$

Equality holds if $\mathbf{z} = z\mathbf{e}_K$, i.e., $\mathbf{x}_{\text{opt}} = z\mathbf{v}_K$.

- $\mathbf{x}'_{\text{opt}}\mathbf{A}\mathbf{x}_{\text{opt}} = \lambda_K$
- The solution \mathbf{x}_{opt} is unique (to within $|z| = 1$) when $\lambda_{K-1} > \lambda_K > \lambda_{K+1}$.

Pr. 9. (sol/hs050)

Let $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}'$. If $\text{rank}(\mathbf{A}) = N$, then \mathbf{A} has N linearly independent columns, $N \leq M$, and $\mathbf{A}'\mathbf{A}$ is invertible. Then:

$$\begin{aligned} (\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}' &= (\mathbf{V}\mathbf{\Sigma}'\mathbf{U}'\mathbf{U}\mathbf{\Sigma}\mathbf{V}')^{-1}\mathbf{V}\mathbf{\Sigma}'\mathbf{U}' = (\mathbf{V}\mathbf{\Sigma}'\mathbf{\Sigma}\mathbf{V}')^{-1}\mathbf{V}\mathbf{\Sigma}'\mathbf{U}' \\ &= \mathbf{V}(\mathbf{\Sigma}'\mathbf{\Sigma})^{-1}\mathbf{\Sigma}'\mathbf{U}' = \mathbf{V} \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2^2} & \dots & 0 \\ & \ddots & \ddots & \\ 0 & 0 & \dots & \frac{1}{\sigma_n^2} \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 & \dots & 0 \\ & \ddots & & & \ddots & \\ 0 & 0 & \dots & \sigma_n & \dots & 0 \end{bmatrix} \mathbf{U}' \\ &= \mathbf{V} \begin{bmatrix} \frac{1}{\sigma_1} & 0 & \dots & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2} & \dots & 0 & \dots & 0 \\ & \ddots & & & \ddots & \\ 0 & 0 & \dots & \frac{1}{\sigma_n} & \dots & 0 \end{bmatrix} \mathbf{U}' = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}' = \mathbf{A}^\dagger. \end{aligned}$$

Pr. 10. (sol/hs067)

- (a) The problem of finding the line $y = \alpha x + b$ that best fits the points $(1, 2)$, $(2, 1)$ and $(3, 3)$ is equivalent to the problem: Find β that minimizes $\|z - A\beta\|_2$, with $z \triangleq \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix}$, $A \triangleq \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{bmatrix}$, $\beta = \begin{bmatrix} \alpha \\ b \end{bmatrix}$.

Thus the optimal LLS solution is $\hat{\beta} = A^+ z = V \Sigma^+ U' z = (A' A)^{-1} A' z$.

Here $(A' A)^{-1} = \begin{bmatrix} 14 & 6 \\ 6 & 3 \end{bmatrix}^{-1} = \frac{1}{14 \cdot 3 - 6 \cdot 6} \begin{bmatrix} 3 & -6 \\ -6 & 14 \end{bmatrix} = \frac{1}{6} \begin{bmatrix} 3 & -6 \\ -6 & 14 \end{bmatrix}$, and $A' z = \begin{bmatrix} 13 \\ 6 \end{bmatrix}$,
 so $\hat{\beta} = \frac{1}{6} \begin{bmatrix} 3 & -6 \\ -6 & 14 \end{bmatrix} \begin{bmatrix} 13 \\ 6 \end{bmatrix} = \frac{1}{6} \begin{bmatrix} 3 \\ 6 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix}$.

We repeat this calculation using the SVD. Let $A = U \Sigma V'$ denote a SVD of A . Numerically computing an SVD yields:

$$U = \begin{bmatrix} -0.3231 & 0.8538 & 0.4082 \\ -0.5475 & 0.1832 & -0.8165 \\ -0.7719 & -0.4873 & 0.4082 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 4.0791 & 0 \\ 0 & 0.6005 \\ 0 & 0 \end{bmatrix}, \quad V = \begin{bmatrix} -0.9153 & -0.4027 \\ -0.4027 & 0.9153 \end{bmatrix},$$

$$\text{so } \hat{\beta} = V \Sigma^+ U' = V \begin{bmatrix} 0.2451 & 0 & 0 \\ 0 & 1.6653 & 0 \end{bmatrix} U' z = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix}.$$

Unsurprisingly we get the same answer since $A^+ = (A' A)^{-1} A'$ when A has full column rank.

- (b) Here we are asked to find the line $y = \alpha x + b$ that best fits the points $(2, 1)$, $(1, 2)$ and $(3, 3)$. We could repeat the above computation, or we can recognize that all that has changed is that the first two points are swapped between (a) and (b). Thus the solution β of the two problems will be identical because the least-squares criterion is the same for any ordering of the (x_i, y_i) points. A more formal way of seeing is to first note that the new A matrix for (b) is exactly same as the A matrix in (a) except for the first two rows swapping places. Denote the

“ A ” matrix in part (a) by A_a and the “ A ” matrix in part (b) by A_b . Then $A_b = \underbrace{\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\triangleq P} A_a$.

Note that P is a orthogonal matrix because $PP' = I$. Such a P is called a permutation matrix. We also have that

the new $z_b = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = P z_a$. Thus the new problem is Find β_b that minimizes $\|z_b - A_b \beta\|_2 = \|P z_a - P A_a \beta\|_2 = \|z_a - A_a \beta\|_2$, because P is an orthogonal matrix. Thus we get the same solution β .