

Linear regression

CS 446 / ECE 449

2022-02-19 22:48:15 -0600 (383ae58)

Plan for today

- ▶ Linear regression setup revisited.
- ▶ Normal equations, SVD, and pseudoinverse.
- ▶ Example (if time).

“pytorch meta-algorithm”

1. Clean/augment data.
2. Pick model/architecture.
3. Pick a loss function measuring model fit to data.
4. Run a gradient descent variant to fit model to data.
5. Tweak 1-4 until training error is small.
6. Tweak 1-5, possibly reducing model complexity, until testing error is small.

Is that all of ML?

No, but these days it's much of it!

Linear regression — basic setup

1. Start from **training data** $((\mathbf{x}_i, y_i))_{i=1}^n$, with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$.
2. **Model** is a linear predictor: pick $\mathbf{w} \in \mathbb{R}^d$ with

$$\mathbf{x}_i \mapsto \mathbf{w}^\top \mathbf{x}_i =: \hat{y}_i \approx y_i.$$

3. **Loss function** ℓ is squared loss ℓ_{sq} (standard regression loss):

$$\ell_{\text{sq}}(\mathbf{w}^\top \mathbf{x}_i, y_i) = \frac{1}{2}(\mathbf{w}^\top \mathbf{x}_i - y_i)^2.$$

We will minimize the **empirical risk** (average loss over training examples):

$$\hat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{sq}}(\mathbf{w}^\top \mathbf{x}_i, y_i) = \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \quad \text{where } \mathbf{X} := \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix}.$$

4. Basic method: **gradient descent**. Set $\mathbf{w}_0 = 0$, and thereafter

$$\mathbf{w}_{i+1} := \mathbf{w}_i - \eta \nabla \hat{\mathcal{R}}(\mathbf{w}_i) = \mathbf{w}_i - \frac{\eta}{n} \mathbf{X}^\top (\mathbf{X}\mathbf{w}_i - \mathbf{y}),$$

where η is a **learning rate** (step size).

2. **Model** is a **linear predictor**: pick $\mathbf{w} \in \mathbb{R}^d$ with

$$\mathbf{x}_i \mapsto \mathbf{w}^\top \mathbf{x}_i \approx y_i.$$

- ▶ Our **model/architecture/function class** is $\{\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} : \mathbf{w} \in \mathbb{R}^d\}$.
For each $\mathbf{w} \in \mathbb{R}^d$, we have another predictor.
- ▶ This is a simple model; we'll build off of it to get more powerful ones!
- ▶ This model is insufficient for complicated tasks, but often does well, and forms a good baseline.

3. **Loss function** is squared loss (standard regression loss):

$$\ell(\mathbf{w}^\top \mathbf{x}_i, y_i) = \frac{1}{2}(\mathbf{w}^\top \mathbf{x}_i - y_i)^2.$$

We will minimize the **empirical risk**:

$$\widehat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}^\top \mathbf{x}_i, y_i) = \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

► **Regression towards the mean**: if $\mathbf{x}_i = 1 \in \mathbb{R}^1$ for all i , then

$$\arg \min_{\mathbf{w} \in \mathbb{R}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{n} \sum_{i=1}^n y_i.$$

Seems a reasonable notion of loss/error.

► There are many choices for ℓ . Next lecture we'll use **logistic loss** ℓ_{logistic}

$$\ell_{\text{logistic}}(\hat{y}, y) = \ln(1 + \exp(-\hat{y}y)).$$

This and squared loss are the most common.

4. Basic method: **gradient descent**. Set $\mathbf{w}_0 = 0$, and thereafter

$$\mathbf{w}_{i+1} := \mathbf{w}_i - \eta \nabla \widehat{\mathcal{R}}(\mathbf{w}_i) = \mathbf{w}_i - \frac{\eta}{n} \mathbf{X}^\top (\mathbf{X} \mathbf{w}_i - \mathbf{y}),$$

where η is a **learning rate (step size)**.

- ▶ In a few lectures, we'll see that this **globally minimizes $\widehat{\mathcal{R}}$** .
- ▶ We'll spend most of this lecture on other solutions via SVD.

Normal equations and SVD.

We want to find $\hat{\mathbf{w}}$ so that

$$2n\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}\|^2 = \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Normal equations and SVD.

We want to find $\hat{\mathbf{w}}$ so that

$$2n\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}\|^2 = \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Idea from calculus: set gradient to zero and solve:

$$0 = \nabla_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = 2\mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}),$$

meaning we want $\hat{\mathbf{w}}$ so that

$$\mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

These are called the **normal equations**.

The **normal equations** are the system of linear equalities

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

Proposition. $\hat{\mathbf{w}}$ satisfies $\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$ iff $\hat{\mathbf{w}}$ satisfies the normal equations.

The **normal equations** are the system of linear equalities

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

Proposition. $\hat{\mathbf{w}}$ satisfies $\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$ iff $\hat{\mathbf{w}}$ satisfies the normal equations.

Proof (one direction). Consider \mathbf{w} with $\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}$, and any \mathbf{w}' ; then

$$\begin{aligned} \|\mathbf{X} \mathbf{w}' - \mathbf{y}\|^2 &= \|\mathbf{X} \mathbf{w}' - \mathbf{X} \mathbf{w} + \mathbf{X} \mathbf{w} - \mathbf{y}\|^2 \\ &= \|\mathbf{X} \mathbf{w}' - \mathbf{X} \mathbf{w}\|^2 + 2(\mathbf{X} \mathbf{w}' - \mathbf{X} \mathbf{w})^\top (\mathbf{X} \mathbf{w} - \mathbf{y}) + \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2. \end{aligned}$$

Since

$$(\mathbf{X} \mathbf{w}' - \mathbf{X} \mathbf{w})^\top (\mathbf{X} \mathbf{w} - \mathbf{y}) = (\mathbf{w}' - \mathbf{w})^\top (\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y}) = 0,$$

then

$$2n\hat{\mathcal{R}}(\mathbf{w}') = \|\mathbf{X} \mathbf{w}' - \mathbf{y}\|^2 = \|\mathbf{X} \mathbf{w}' - \mathbf{X} \mathbf{w}\|^2 + \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2 \geq \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2 = 2n\hat{\mathcal{R}}(\mathbf{w}).$$

□

Later we'll get a general version by convexity, but it's nice that we can check this directly so easily!

The **normal equations** are the system of linear equalities

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

Proposition. $\hat{\mathbf{w}}$ satisfies $\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$ iff $\hat{\mathbf{w}}$ satisfies the normal equations.

The **normal equations** are the system of linear equalities

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}.$$

Proposition. $\hat{\mathbf{w}}$ satisfies $\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$ iff $\hat{\mathbf{w}}$ satisfies the normal equations.

How do we solve for $\hat{\mathbf{w}}$?

- ▶ If $\mathbf{X}^T \mathbf{X}$ is invertible, we can use $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.
- ▶ In general, we will use the **SVD**.

The SVD (Singular Value Decomposition).

Let $M \in \mathbb{R}^{n \times d}$ be given. $((s_i, \mathbf{u}_i, \mathbf{v}_i))_{i=1}^r$ is an **SVD of M** if:

- ▶ M has rank r ;
- ▶ $s_1 \geq s_2 \cdots \geq s_r > 0$;
- ▶ $(\mathbf{u}_i)_{i=1}^r$ are orthonormal (orthogonal and unit length), and span the column space of M ;
- ▶ $(\mathbf{v}_i)_{i=1}^r$ are orthonormal, and span the row space of M .
- ▶ $M = \sum_i s_i \mathbf{u}_i \mathbf{v}_i^\top$.

The SVD (Singular Value Decomposition).

Let $M \in \mathbb{R}^{n \times d}$ be given. $((s_i, \mathbf{u}_i, \mathbf{v}_i))_{i=1}^r$ is an **SVD of M** if:

- ▶ M has rank r ;
- ▶ $s_1 \geq s_2 \geq \dots \geq s_r > 0$;
- ▶ $(\mathbf{u}_i)_{i=1}^r$ are orthonormal (orthogonal and unit length), and span the column space of M ;
- ▶ $(\mathbf{v}_i)_{i=1}^r$ are orthonormal, and span the row space of M .
- ▶ $M = \sum_i s_i \mathbf{u}_i \mathbf{v}_i^T$.

- ▶ The SVD always exists, and is real-valued.
(When do real eigendecompositions not exist?)
- ▶ The ordered tuple (s_1, \dots, s_r) is unique, but the SVD is in general not unique (why not?).
- ▶ For $k < r$, the **low rank approximation** $\sum_{i=1}^k s_i \mathbf{u}_i \mathbf{v}_i^T \approx M$ has many applications (wait for the PCA lecture).

Pseudoinverse.

Given SVD $M = \sum_i s_i \mathbf{u}_i \mathbf{v}_i^\top$, the **pseudoinverse** is

$$M^+ := \sum_{i=1}^r \frac{1}{s_i} \mathbf{v}_i \mathbf{u}_i^\top.$$

Pseudoinverse.

Given SVD $M = \sum_i s_i \mathbf{u}_i \mathbf{v}_i^\top$, the **pseudoinverse** is

$$M^+ := \sum_{i=1}^r \frac{1}{s_i} \mathbf{v}_i \mathbf{u}_i^\top.$$

- ▶ The SVD may fail to be unique, but M^+ is unique.
- ▶ $MM^+ = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^\top$ and $M^+M = \sum_{i=1}^r \mathbf{v}_i \mathbf{v}_i^\top$; in general, neither is an identity matrix. (Consider the case $M = \mathbf{e}_1 \mathbf{e}_1^\top$.)
- ▶ On the other hand,

$$MM^+M =$$

$$M^+MM^+ =$$

- ▶ If M^{-1} exists, then $M^+ = M^{-1}$.
- ▶ If $M = 0$, then $r = 0$ and $M^+ = 0$.

Pseudoinverse.

Given SVD $M = \sum_i s_i \mathbf{u}_i \mathbf{v}_i^\top$, the **pseudoinverse** is

$$M^+ := \sum_{i=1}^r \frac{1}{s_i} \mathbf{v}_i \mathbf{u}_i^\top.$$

- ▶ The SVD may fail to be unique, but M^+ is unique.
- ▶ $MM^+ = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^\top$ and $M^+M = \sum_{i=1}^r \mathbf{v}_i \mathbf{v}_i^\top$; in general, neither is an identity matrix. (Consider the case $M = \mathbf{e}_1 \mathbf{e}_1^\top$.)
- ▶ On the other hand,

$$MM^+M = \left(\sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^\top \right) \left(\sum_{j=1}^r \frac{1}{s_j} \mathbf{v}_j \mathbf{u}_j^\top \right) \left(\sum_{k=1}^r s_k \mathbf{u}_k \mathbf{v}_k^\top \right) = M,$$

$$M^+MM^+ = \left(\sum_{i=1}^r \frac{1}{s_i} \mathbf{v}_i \mathbf{u}_i^\top \right) \left(\sum_{j=1}^r s_j \mathbf{u}_j \mathbf{v}_j^\top \right) \left(\sum_{k=1}^r \frac{1}{s_k} \mathbf{v}_k \mathbf{u}_k^\top \right) = M^+.$$

- ▶ If M^{-1} exists, then $M^+ = M^{-1}$.
- ▶ If $M = 0$, then $r = 0$ and $M^+ = 0$.

OLS (Ordinary Least Squares) solution via SVD.

Given a least squares problem $\hat{\mathcal{R}}(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2/(2n)$, the **OLS solution**

$$\hat{\mathbf{w}}_{\text{ols}} = \mathbf{X}^+ \mathbf{y}$$

satisfies the normal equations (whereby $\hat{\mathcal{R}}(\hat{\mathbf{w}}_{\text{ols}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$).

OLS (Ordinary Least Squares) solution via SVD.

Given a least squares problem $\hat{\mathcal{R}}(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2/(2n)$, the **OLS solution**

$$\hat{\mathbf{w}}_{\text{ols}} = \mathbf{X}^+ \mathbf{y}$$

satisfies the normal equations (whereby $\hat{\mathcal{R}}(\hat{\mathbf{w}}_{\text{ols}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$).

Easy to check: writing $\mathbf{X} = \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^T$,

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}}_{\text{ols}} &= \mathbf{X}^T \mathbf{X} \mathbf{X}^+ \mathbf{y} \\ &= \left(\sum_{i=1}^r s_i \mathbf{v}_i \mathbf{u}_i^T \right) \left(\sum_{j=1}^r s_j \mathbf{u}_j \mathbf{v}_j^T \right) \left(\sum_{k=1}^r \frac{1}{s_k} \mathbf{v}_k \mathbf{u}_k^T \right) \mathbf{y} \\ &= \mathbf{X}^T \mathbf{y}. \end{aligned}$$

SVD $M = \sum_i s_i \mathbf{u}_i \mathbf{v}_i^T$ and orthonormal bases.

We can extend $(\mathbf{u}_i)_{i=1}^r$ and $(\mathbf{v}_i)_{i=1}^r$ to full orthonormal bases for \mathbb{R}^n and \mathbb{R}^d respectively: write $M \in \mathbb{R}^{n \times d}$ as

$$\left[\begin{array}{ccc|ccc} \uparrow & & \uparrow & \uparrow & & \uparrow \\ \mathbf{u}_1 & \cdots & \mathbf{u}_r & \mathbf{u}_{r+1} & \cdots & \mathbf{u}_n \\ \downarrow & & \downarrow & \downarrow & & \downarrow \end{array} \right] \cdot \left[\begin{array}{ccc|ccc} s_1 & & & & & 0 \\ & \ddots & & & & 0 \\ & & & & & 0 \\ \hline 0 & & & s_r & & 0 \\ & & & & & 0 \end{array} \right] \cdot \left[\begin{array}{ccc|ccc} \uparrow & & \uparrow & \uparrow & & \uparrow \\ \mathbf{v}_1 & \cdots & \mathbf{v}_r & \mathbf{v}_{r+1} & \cdots & \mathbf{v}_d \\ \downarrow & & \downarrow & \downarrow & & \downarrow \end{array} \right]^T .$$

The old parts span the column and row spaces of M ;
the new vectors span the left and right nullspaces.
Some call this a “full” SVD.

SVD and relationship to eigenvalues.

SVD and relationship to eigenvalues.

Note

$$MM^T = \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^T \sum_{j=1}^r s_j \mathbf{v}_j \mathbf{u}_j^T = \sum_{i=1}^r s_i^2 \mathbf{u}_i \mathbf{u}_i^T,$$

thus left singular vectors $(\mathbf{u})_{i=1}^r$ are top eigenvectors of MM^T , with eigenvalues $s_1^2 \geq \dots \geq s_r^2$.

Similarly,

$$M^T M = \sum_{i=1}^r s_i \mathbf{v}_i \mathbf{u}_i^T \sum_{j=1}^r s_j \mathbf{u}_j \mathbf{v}_j^T = \sum_{i=1}^r s_i^2 \mathbf{v}_i \mathbf{v}_i^T,$$

obtaining right singular vectors from $M^T M$.

Summary on least squares solutions

We want to approximately solve the **empirical risk minimization problem**

$$\min_{\mathbf{w} \in \mathbb{R}^d} \widehat{\mathcal{R}}(\mathbf{w}) = \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Three approaches:

1. **Gradient descent**: $\mathbf{w}_0 := 0$, thereafter $\mathbf{w}_{i+1} := \mathbf{w} - \eta \nabla \widehat{\mathcal{R}}(\mathbf{w}_i)$.
2. Pick any $\hat{\mathbf{w}}$ satisfying the **normal equations**

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

3. Use the **ordinary least squares (OLS)** solution $\hat{\mathbf{w}}_{\text{ols}} = \mathbf{X}^+ \mathbf{y}$.

Summary on least squares solutions

We want to approximately solve the empirical risk minimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} \widehat{\mathcal{R}}(\mathbf{w}) = \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Three approaches:

1. **Gradient descent:** $\mathbf{w}_0 := 0$, thereafter $\mathbf{w}_{i+1} := \mathbf{w} - \eta \nabla \widehat{\mathcal{R}}(\mathbf{w}_i)$.
2. Pick any $\hat{\mathbf{w}}$ satisfying the normal equations

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

3. Use the **ordinary least squares (OLS)** solution $\hat{\mathbf{w}}_{\text{ols}} = \mathbf{X}^+ \mathbf{y}$.

(Side note: are these different?...))

Example: Old Faithful geyser (Yellowstone)

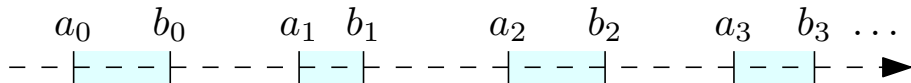
Example: Old Faithful geyser (Yellowstone)



Task: Predict time of next eruption.

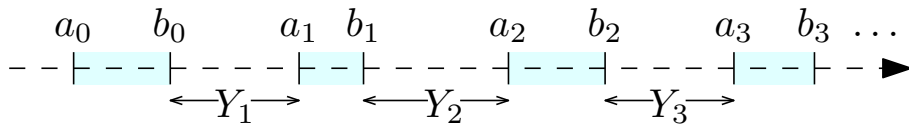
Time between eruptions

Source data: start and end times (a_i, b_i) of $n = 136$ eruptions.



Time between eruptions

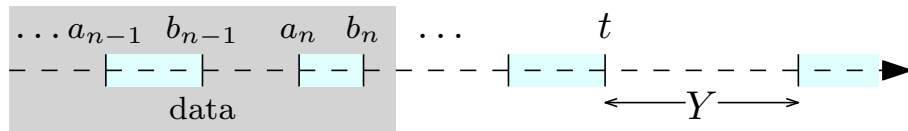
Source data: start and end times (a_i, b_i) of $n = 136$ eruptions.



Let's pre-process: form time between eruptions $y_i := a_{i+1} - b_i$.

Time between eruptions

Source data: start and end times (a_i, b_i) of $n = 136$ eruptions.



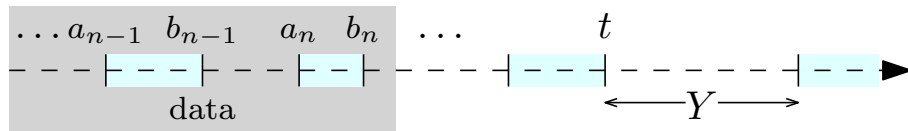
Let's pre-process: form time between eruptions $y_i := a_{i+1} - b_i$.

Reformulated task:

to estimate next eruption, find last end time t , compute \hat{y} , and output $t + \hat{y}$.

Time between eruptions

Source data: start and end times (a_i, b_i) of $n = 136$ eruptions.



Let's pre-process: form time between eruptions $y_i := a_{i+1} - b_i$.

Reformulated task:

to estimate next eruption, find last end time t , compute \hat{y} , and output $t + \hat{y}$.

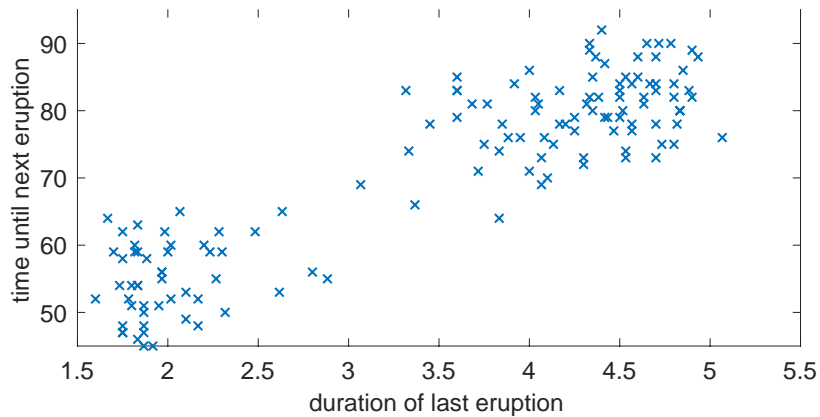
Let's use linear regression.

- ▶ Set $x_i = 1$ and the OLS solution is the mean:

$$\hat{y} = \frac{1}{136} \sum_{i=1}^{136} y_i = 70.7941.$$

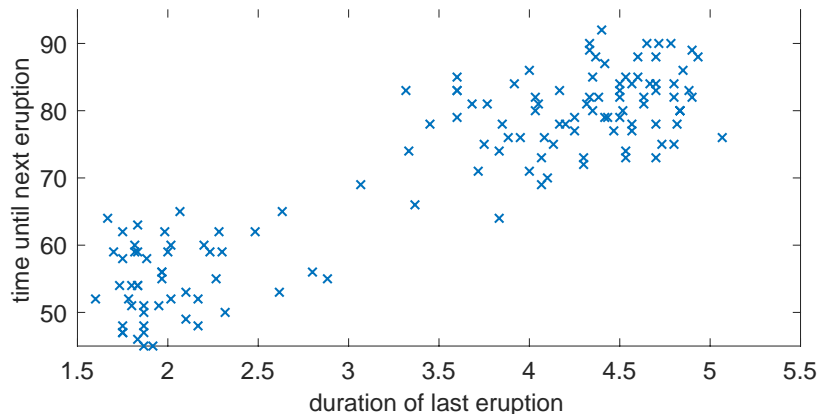
- ▶ Can we do better with another x_i ?

Eruption length and time to eruption are correlated.



Let choose $\mathbf{x}_i := \begin{bmatrix} b_i - a_i \\ 1 \end{bmatrix}$.

Eruption length and time to eruption are correlated.



Let choose $\mathbf{x}_i := \begin{bmatrix} b_i - a_i \\ 1 \end{bmatrix}$.

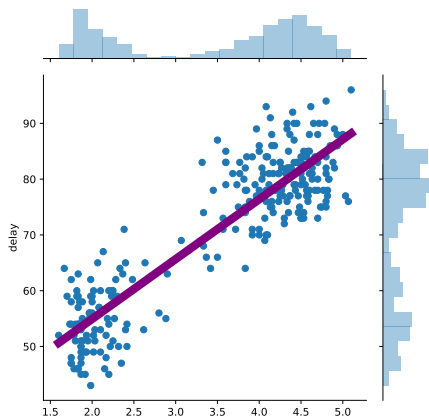
(Side note: the extra "1" will be discussed extensively later.)

1. Form pairs $\mathbf{x}_i := \begin{bmatrix} b_i - a_i \\ 1 \end{bmatrix}$, and matrix

$$\mathbf{X} := \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} b_1 - a_1 & 1 \\ & \vdots \\ b_n - a_n & 1 \end{bmatrix} \in \mathbb{R}^{n \times 2}.$$

Form labels $\mathbf{y} \in \mathbb{R}^n$, $y_i := a_{i+1} - b_i$.

2. Choose OLS solution $\hat{\mathbf{w}}_{\text{ols}} := \mathbf{X}^+ \mathbf{y}$.
3. Given a new eruption (a, b) , estimate next eruption time $b + \mathbf{w}^\top \begin{bmatrix} b - a \\ 1 \end{bmatrix}$.



“pytorch meta-algorithm” on Old Faithful data

1. Clean/augment data.

From (a_i, b_i) , form $\mathbf{x}'_i = (1,)$ or $\mathbf{x}_i = (b_i - a_i, 1)$, and $y_i = a_i - b_{i-1}$.

2. Pick model/architecture (anything from lectures 2-13).

Linear predictor.

3. Pick a loss function measuring model fit to data.

Squared loss.

4. Run a gradient descent variant to fit model to data.

5. Tweak 1-4 until training error is small.

\mathbf{x}'_i was bad, so we added a feature and got \mathbf{x}_i .

6. Tweak 1-5, possibly reducing model complexity, until testing error is small.

We didn't try this!

Summary for today

- ▶ Linear regression setup revisited.
- ▶ Normal equations, SVD, and pseudoinverse.
- ▶ Example (if time).

(Appendix.)

The **normal equations** are the system of linear equalities

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

Proposition. $\hat{\mathbf{w}}$ satisfies $\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$ iff $\hat{\mathbf{w}}$ satisfies the normal equations.

The **normal equations** are the system of linear equalities

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

Proposition. $\hat{\mathbf{w}}$ satisfies $\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$ iff $\hat{\mathbf{w}}$ satisfies the normal equations.

Proof (other direction).

Suppose \mathbf{w} is optimal; since $\hat{\mathbf{w}}_{\text{ols}}$ satisfies the normal equations, then expanding the square as in the proof of the other direction gives

$$\|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2 = \|\mathbf{X} \mathbf{w} - \mathbf{X} \hat{\mathbf{w}}_{\text{ols}}\|^2 + \|\mathbf{X} \hat{\mathbf{w}}_{\text{ols}} - \mathbf{y}\|^2.$$

Since \mathbf{w} and $\hat{\mathbf{w}}_{\text{ols}}$ are optimal, then $\hat{\mathcal{R}}(\mathbf{w}) = \hat{\mathcal{R}}(\hat{\mathbf{w}}_{\text{ols}})$, so the preceding implies

The **normal equations** are the system of linear equalities

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}.$$

Proposition. $\hat{\mathbf{w}}$ satisfies $\hat{\mathcal{R}}(\hat{\mathbf{w}}) = \min_{\mathbf{w}} \hat{\mathcal{R}}(\mathbf{w})$ iff $\hat{\mathbf{w}}$ satisfies the normal equations.

Proof (other direction).

Suppose \mathbf{w} is optimal; since $\hat{\mathbf{w}}_{\text{ols}}$ satisfies the normal equations, then expanding the square as in the proof of the other direction gives

$$\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{X}\hat{\mathbf{w}}_{\text{ols}}\|^2 + \|\mathbf{X}\hat{\mathbf{w}}_{\text{ols}} - \mathbf{y}\|^2.$$

Since \mathbf{w} and $\hat{\mathbf{w}}_{\text{ols}}$ are optimal, then $\hat{\mathcal{R}}(\mathbf{w}) = \hat{\mathcal{R}}(\hat{\mathbf{w}}_{\text{ols}})$, so the preceding implies

$$0 = \|\mathbf{X}\mathbf{w} - \mathbf{X}\hat{\mathbf{w}}_{\text{ols}}\|^2 = \|\mathbf{X}(\mathbf{w} - \hat{\mathbf{w}}_{\text{ols}})\|^2,$$

therefore $\mathbf{X}(\mathbf{w} - \hat{\mathbf{w}}_{\text{ols}}) = \mathbf{0}$ and $\mathbf{X}\mathbf{w} = \mathbf{X}\hat{\mathbf{w}}_{\text{ols}}$, which by the normal equations for $\hat{\mathbf{w}}_{\text{ols}}$ means

$$\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \hat{\mathbf{w}}_{\text{ols}} = \mathbf{X}^\top \mathbf{X} \mathbf{w},$$

thus \mathbf{w} satisfies the normal equations. □

Geometric interpretation of least squares ERM

Let $\mathbf{a}_j \in \mathbb{R}^n$ be the j -th column (not row!) of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, so

$$\mathbf{X} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{a}_1 & \cdots & \mathbf{a}_d \\ \downarrow & & \downarrow \end{bmatrix}.$$

Geometric interpretation of least squares ERM

Let $\mathbf{a}_j \in \mathbb{R}^n$ be the j -th **column** (not row!) of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, so

$$\mathbf{X} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{a}_1 & \cdots & \mathbf{a}_d \\ \downarrow & & \downarrow \end{bmatrix}.$$

Minimizing $\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2$ means finding $\hat{\mathbf{b}} \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_d)$ closest to \mathbf{b} .

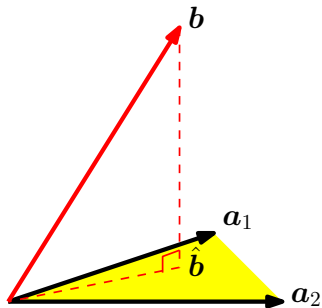
Geometric interpretation of least squares ERM

Let $\mathbf{a}_j \in \mathbb{R}^n$ be the j -th **column** (not row!) of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, so

$$\mathbf{X} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{a}_1 & \cdots & \mathbf{a}_d \\ \downarrow & & \downarrow \end{bmatrix}.$$

Minimizing $\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2$ means finding $\hat{\mathbf{b}} \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_d)$ closest to \mathbf{b} .

Solution $\hat{\mathbf{b}}$ is orthogonal projection of \mathbf{b} onto $\text{range}(\mathbf{A}) = \{\mathbf{A}\mathbf{w} : \mathbf{w} \in \mathbb{R}^d\}$.



Geometric interpretation of least squares ERM

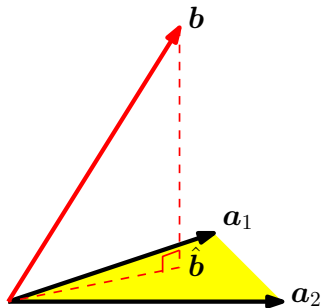
Let $\mathbf{a}_j \in \mathbb{R}^n$ be the j -th **column** (not row!) of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, so

$$\mathbf{X} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{a}_1 & \cdots & \mathbf{a}_d \\ \downarrow & & \downarrow \end{bmatrix}.$$

Minimizing $\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2$ means finding $\hat{\mathbf{b}} \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_d)$ closest to \mathbf{b} .

Solution $\hat{\mathbf{b}}$ is orthogonal projection of \mathbf{b} onto $\text{range}(\mathbf{A}) = \{\mathbf{A}\mathbf{w} : \mathbf{w} \in \mathbb{R}^d\}$.

- ▶ $\hat{\mathbf{b}}$ is uniquely determined; indeed,
 $\hat{\mathbf{b}} = \mathbf{A}\mathbf{A}^+ \mathbf{b} = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^\top \mathbf{b}$.



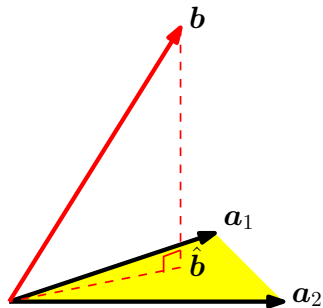
Geometric interpretation of least squares ERM

Let $\mathbf{a}_j \in \mathbb{R}^n$ be the j -th **column** (not row!) of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, so

$$\mathbf{X} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{a}_1 & \cdots & \mathbf{a}_d \\ \downarrow & & \downarrow \end{bmatrix}.$$

Minimizing $\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2$ means finding $\hat{\mathbf{b}} \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_d)$ closest to \mathbf{b} .

Solution $\hat{\mathbf{b}}$ is orthogonal projection of \mathbf{b} onto $\text{range}(\mathbf{A}) = \{\mathbf{A}\mathbf{w} : \mathbf{w} \in \mathbb{R}^d\}$.



- ▶ $\hat{\mathbf{b}}$ is uniquely determined; indeed,
 $\hat{\mathbf{b}} = \mathbf{A}\mathbf{A}^+ \mathbf{b} = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^\top \mathbf{b}$.
- ▶ If $r = \text{rank}(\mathbf{A}) < d$, then >1 way to write $\hat{\mathbf{b}}$ as linear combination of $\mathbf{a}_1, \dots, \mathbf{a}_d$.

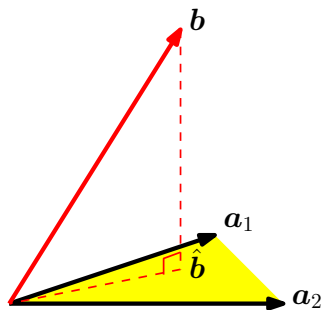
Geometric interpretation of least squares ERM

Let $\mathbf{a}_j \in \mathbb{R}^n$ be the j -th **column** (not row!) of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, so

$$\mathbf{X} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{a}_1 & \cdots & \mathbf{a}_d \\ \downarrow & & \downarrow \end{bmatrix}.$$

Minimizing $\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2$ means finding $\hat{\mathbf{b}} \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_d)$ closest to \mathbf{b} .

Solution $\hat{\mathbf{b}}$ is orthogonal projection of \mathbf{b} onto $\text{range}(\mathbf{A}) = \{\mathbf{A}\mathbf{w} : \mathbf{w} \in \mathbb{R}^d\}$.



- ▶ $\hat{\mathbf{b}}$ is uniquely determined; indeed,
 $\hat{\mathbf{b}} = \mathbf{A}\mathbf{A}^+ \mathbf{b} = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^\top \mathbf{b}$.
- ▶ If $r = \text{rank}(\mathbf{A}) < d$, then >1 way to write $\hat{\mathbf{b}}$ as linear combination of $\mathbf{a}_1, \dots, \mathbf{a}_d$.

If $\text{rank}(\mathbf{A}) < d$, then **ERM solution is not unique**.

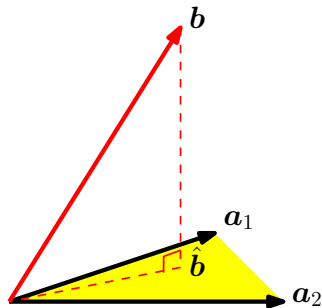
Geometric interpretation of least squares ERM

Let $\mathbf{a}_j \in \mathbb{R}^n$ be the j -th **column** (not row!) of matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, so

$$\mathbf{X} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{a}_1 & \cdots & \mathbf{a}_d \\ \downarrow & & \downarrow \end{bmatrix}.$$

Minimizing $\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2$ means finding $\hat{\mathbf{b}} \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_d)$ closest to \mathbf{b} .

Solution $\hat{\mathbf{b}}$ is orthogonal projection of \mathbf{b} onto $\text{range}(\mathbf{A}) = \{\mathbf{A}\mathbf{w} : \mathbf{w} \in \mathbb{R}^d\}$.



- ▶ $\hat{\mathbf{b}}$ is uniquely determined; indeed,
 $\hat{\mathbf{b}} = \mathbf{A}\mathbf{A}^+ \mathbf{b} = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^\top \mathbf{b}$.
- ▶ If $r = \text{rank}(\mathbf{A}) < d$, then >1 way to write $\hat{\mathbf{b}}$ as linear combination of $\mathbf{a}_1, \dots, \mathbf{a}_d$.

If $\text{rank}(\mathbf{A}) < d$, then **ERM solution is not unique**.

To get \mathbf{w} from $\hat{\mathbf{b}}$:
solve system of linear equations $\mathbf{A}\mathbf{w} = \hat{\mathbf{b}}$.

Computing the SVD

Typical solver is an iterative, greedy method.

For more information, see the excellent data science book by Blum, Hopcroft, Kannan.

Why GD?

Why include GD, since pseudoinverse seems sufficient?

- ▶ GD is easy to implement, pseudoinverse more painful.
- ▶ Pseudoinverse after all implemented as an iterative solver.
- ▶ GD generalizes to other cases of squared loss (e.g., deep network training with squared loss).

Why $1/2$ in squared loss?

Note: this material is beyond the scope of this course.

1. (Sanity check: doesn't affect methods.) Note that multiplying the risk by a constant just scales the whole surface, and does not change the sign of $\hat{\mathcal{R}}(\mathbf{v}) - \hat{\mathcal{R}}(\mathbf{w})$ for any two choices (\mathbf{v}, \mathbf{w}) , and therefore to some extent it doesn't matter (and can be treated as just rescaling the step size in gradient descent).
2. (Superficial reasons: many calculations seem nicer.) The usual answer is certain calculations are nicer: the gradient for linear regression becomes $\mathbf{X}^\top(\mathbf{X}\mathbf{w} - \mathbf{y})$, the Hessian becomes $\mathbf{X}^\top\mathbf{X}$, and if $\|\mathbf{x}_i\| \leq 1$ for all examples then this has 1 as its largest eigenvalue and we can use a step size of 1 with gradient descent.
3. (Deeper, but related reason.) The squared loss is convex, and in fact the related function $g(\mathbf{v}) = \|\mathbf{v}\|_2^2/2$ (in $\hat{\mathcal{R}}$) satisfies a very interesting property. In convex analysis, a key concept for any convex function f is its conjugate f^* ; interestingly, g is the unique convex function satisfying $g = g^*$. To some extent this is equivalent to the previous point, since conjugacy has connections to gradients, and moreover is also a bit arbitrary, since conjugacy relies upon a choice of primal and dual spaces and we could have obtained a different unique pair.

Overall, I think it is nice to use the $1/2$, and less controversial than the choice of placing 2π in the definition of the Fourier transform, which has become sort of standardized but does have various tradeoffs.

- ▶ Shalev-Shwartz/Ben-David: chapters 9, 24.
- ▶ SVD: see the related chapter in the “Foundations of Data Science” book by Blum, Hopcroft, Kannan; they use the same convention $\mathbf{M} = \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^T$ as we use here.