# Linear regression

**CS 446 / ECE 449**

2022-01-17 04:49:48 -0600 (d748537)

# Plan for today

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

# "pytorch meta-algorithm"

1. Clean/augment data. $\longrightarrow$ *after this step, input is vector* $x \in \mathbb{R}^d$.
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 $((\boldsymbol{x}_i, y_i))_{i=1}^n$, with $\boldsymbol{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$.

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

$$\boldsymbol{x}_i \mapsto \boldsymbol{w}^\mathsf{T}\boldsymbol{x}_i =: \hat{y}_i \approx y_i.$$

3. Loss function $\ell$ is squared loss $\ell_{\mathrm{sq}}$ (standard regression loss):

$$\ell_{\mathrm{sq}}(\boldsymbol{w}^\mathsf{T}\boldsymbol{x}_i, y_i) = \frac{1}{2}(\boldsymbol{w}^\mathsf{T}\boldsymbol{x}_i - y_i)^2.$$

*makes it easy to differentiate.*

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

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

$$\frac{\partial \|Xw - y\|^2}{\partial X} = 2X^\mathsf{T}(Xw - y)$$

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

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

where $\eta$ is a learning rate (step size).

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

$$\boldsymbol{x}_i \mapsto \boldsymbol{w}^\mathsf{T}\boldsymbol{x}_i \approx y_i.$$

▶ Our model/architecture/function class is $\{\boldsymbol{x} \mapsto \boldsymbol{w}^\mathsf{T}\boldsymbol{x} : \boldsymbol{w} \in \mathbb{R}^d\}$.
For each $\boldsymbol{w} \in \mathbb{R}^d$, 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):

$$\ell(\boldsymbol{w}^\mathsf{T}\boldsymbol{x}_i, y_i) = \frac{1}{2}(\boldsymbol{w}^\mathsf{T}\boldsymbol{x}_i - y_i)^2.$$

We will minimize the empirical risk:

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

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

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

Seems a reasonable notion of loss/error.

▶ There are many choices for $\ell$. Next lecture we'll use logistic loss $\ell_{\text{logistic}}$

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

This and squared loss are the most common.

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

$$\boldsymbol{w}_{i+1} := \boldsymbol{w}_i - \eta\nabla\widehat{\mathcal{R}}(\boldsymbol{w}_i) = \boldsymbol{w}_i - \frac{\eta}{n}\boldsymbol{X}^\mathsf{T}\left(\boldsymbol{X}\boldsymbol{w}_i - \boldsymbol{y}\right),$$

where $\eta$ 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 $\hat{\boldsymbol{w}}$ so that

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

Idea from calculus: set gradient to zero and solve:

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

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

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

These are called the normal equations.

The normal equations are the system of linear equalities

$$X^\top X w = X^\top y.$$

**Proposition.** $\hat{w}$ satisfies $\widehat{\mathcal{R}}(\hat{w}) = \min_w \widehat{\mathcal{R}}(w)$ iff $\hat{w}$ satisfies the normal equations.

*（手写）Optimal training solution ⟹ global minimum.*

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

$$\|X w' - y\|^2 = \|X w' - X w + X w - y\|^2$$
$$= \|X w' - X w\|^2 + 2(X w' - X w)^\top (X w - y) + \|X w - y\|^2.$$

Since

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

then

$$\|X w' - y\|^2 = \|X w' - X w\|^2 + \|X w - y\|^2 \geq \|X w - y\|^2.$$

$\square$

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

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

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

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

▶ If $\boldsymbol{X}^\top \boldsymbol{X}$ is invertible, we can use $(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \boldsymbol{y}$.

▶ In general, we will use the SVD.

# The SVD (Singular Value Decomposition).

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

*Right. $d \times 1$*

*Left $n \times 1$*

▶ $M$ has rank $r$;

▶ $s_1 \geq s_2 \cdots \geq s_r > 0$;

▶ $(\boldsymbol{u}_i)_{i=1}^r$ are orthonormal (orthogonal and unit length), and span the column space of $M$;

▶ $(\boldsymbol{v}_i)_{i=1}^r$ are orthonormal, and span the row space of $M$.

▶ $M = \sum_i s_i \boldsymbol{u}_i \boldsymbol{v}_i^\mathsf{T}$. *"decomposition."*

▶ The SVD always exists, and is real-valued.
(When do real eigendecompositions not exist?)

▶ The ordered tuple $(s_1, \ldots, s_r)$ is unique,
but the SVD is in general not unique (why not?).

▶ For $k < r$, the low rank approximation $\sum_{i=1}^k s_i \boldsymbol{u}_i \boldsymbol{v}_i^\mathsf{T} \approx M$ has many applications (wait for the PCA lecture).

# Pseudoinverse.

Given SVD $M = \sum_i s_i u_i v_i^\mathsf{T}$, the pseudoinverse is

$$M^+ := \sum_{i=1}^{r} \frac{1}{s_i} v_i u_i^\mathsf{T}.$$

▶ The SVD may fail to be unique, but $M^+$ is unique.
▶ $MM^+ = \sum_{i=1}^{r} u_i u_i^\mathsf{T}$ and $M^+M = \sum_{i=1}^{r} v_i v_i^\mathsf{T}$; in general, neither is an identity matrix. (Consider the case $M = e_1 e_1^\mathsf{T}$.)

$$MM^+ = \begin{bmatrix} I_{r\times r} & 0 \\ 0 & 0 \end{bmatrix}_{n\times n} \qquad M^+M = \begin{bmatrix} I_{r\times r} & 0 \\ 0 & 0 \end{bmatrix}_{d\times d}$$

▶ On the other hand,

$$MM^+M = \left( \sum_{i=1}^{r} s_i u_i v_i^\mathsf{T} \right) \left( \sum_{j=1}^{r} \frac{1}{s_j} v_j u_j^\mathsf{T} \right) \left( \sum_{k=1}^{r} s_k u_k v_k^\mathsf{T} \right) = M,$$

$$M^+MM^+ = \left( \sum_{i=1}^{r} \frac{1}{s_i} v_i u_i^\mathsf{T} \right) \left( \sum_{j=1}^{r} s_j u_j v_j^\mathsf{T} \right) \left( \sum_{k=1}^{r} \frac{1}{s_k} v_k u_k^\mathsf{T} \right) = M^+.$$

▶ 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 $\widehat{\mathcal{R}}(\boldsymbol{w}) = \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|^2/(2n)$, the OLS solution

$$\hat{\boldsymbol{w}}_{\mathsf{ols}} = \boldsymbol{X}^+\boldsymbol{y}$$

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

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

$$\boldsymbol{X}^{\mathsf{T}}\boldsymbol{X}\hat{\boldsymbol{w}}_{\mathsf{ols}} = \boldsymbol{X}^{\mathsf{T}}\boldsymbol{X}\boldsymbol{X}^+\boldsymbol{y}$$

$$= \left(\sum_{i=1}^{r} s_i \boldsymbol{v}_i \boldsymbol{u}_i^{\mathsf{T}}\right) \left(\sum_{j=1}^{r} s_j \boldsymbol{u}_j \boldsymbol{v}_j^{\mathsf{T}}\right) \left(\sum_{k=1}^{r} \frac{1}{s_k} \boldsymbol{v}_k \boldsymbol{u}_k^{\mathsf{T}}\right) \boldsymbol{y}$$

$$= \boldsymbol{X}^{\mathsf{T}}\boldsymbol{y}.$$

# SVD $M = \sum_i s_i u_i v_i^\mathsf{T}$ and orthonormal bases.

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

$$
\begin{bmatrix} \uparrow & & \uparrow & \uparrow & & \uparrow \\ u_1 & \cdots & u_r & u_{r+1} & \cdots & u_n \\ \downarrow & & \downarrow & \downarrow & & \downarrow \end{bmatrix} \cdot \begin{bmatrix} s_1 & & 0 & \\ & \ddots & & 0 \\ 0 & & s_r & \\ \hline & 0 & & 0 \end{bmatrix} \cdot \begin{bmatrix} \uparrow & & \uparrow & \uparrow & & \uparrow \\ v_1 & \cdots & v_r & v_{r+1} & \cdots & v_d \\ \downarrow & & \downarrow & \downarrow & & \downarrow \end{bmatrix}^\top .
$$

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.

Note

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

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

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

obtaining right singular vectors from $M^\mathsf{T}M$.

# Summary on least squares solutions

We want to approximately solve the empirical risk minimization problem

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

<u>Three approaches:</u>

1. Gradient descent: $\boldsymbol{w}_0 := 0$, thereafter $\boldsymbol{w}_{i+1} := \boldsymbol{w} - \eta \nabla \widehat{\mathcal{R}}(\boldsymbol{w}_i)$.

2. Pick any $\hat{\boldsymbol{w}}$ satsifying the normal equations

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

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

*GD vs. OLS.*

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

*both solved by iterative method , faster ? depends on.*

*OLS usually quicker*

# 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).