"Linear regression" (or Just least squares...)

Let
$$x_i \in \mathbb{R}^n, \quad i = 1, ..., m$$
$$y_i \in \mathbb{R}, \quad i = 1, ..., m$$

$$\overline{X} := \begin{pmatrix} - & x_1 & - \\ & \vdots & \\ - & x_m & - \end{pmatrix} \in \mathbb{R}^{m \times n}, \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$$

Want to solve

$$\overline{X}\hat{w} = y \qquad (*)$$

## Case I  $m < n$.

• $\overline{X}$ is fat

• few eqns, many unknowns

• $N(\overline{X}) \neq \{0\}$

  └ obvious way to see this, what if $\overline{X} = e_1^T$.

To make life easier, assume for now that $\text{rank}(\overline{X}) = m$. Then $\dim \text{col } \overline{X} = m$, so $\overline{X}$ is surjective & $y \in \text{col } \overline{X}$. Since $N(\overline{X}) \neq \{0\}$, ∃ many solutions to $(*)$.

Aside  ~~dim col $\overline{X}$ = dim row $\overline{X}$ follows from rank-nullity. By construction, we have the lin. indep rows~~

How can we narrow down solutions?

  options:

    • Pick $\hat{w}$ w/ small $\|\hat{w}\|_2$

      └ Jives w/ notion of regularization & "low complexity" solution
      └ also, can work w/ $\|\hat{w}\|_2^2$, and use properties of Hilbert spaces, which should make analysis _really_ simple.

    • Pick "sparsest" w. small $\|\hat{w}\|_0$ norm ~ compressed sensing. Revisit sometime

We will go w/ option 1 for now.

If $\hat{w}, w$ both solve (*), then $\hat{w} - w \in N(\underline{X})$.

If $\hat{w}$ lies in the rowspace of $\underline{X}$, then $\hat{w} \perp (\hat{w} - w)$,

Since $\text{row}(\underline{X}) \perp N(\underline{X})$. Suppose $w \neq \hat{w}$. In that case,

$$\|\hat{w}\|^2 = \langle w + (\hat{w} - w), w + (\hat{w} - w) \rangle$$

$$= \|w\|^2 + \underbrace{\|\hat{w} - w\|^2}_{> 0} > \|w\|^2.$$

Hence, If $\exists$ a soln in the rowspace, it has minimum norm and is unique.

To find such a solution, we need a map that brings $\text{col } \underline{X}$ back to $\text{row } \underline{X}$. In general, $\underline{X}$ is bijective from $\text{row } \underline{X}$ to $\text{col } \underline{X}$ & the pseudo inverse provides the reverse mapping. For now, we'll cheat a little rather than use the general pseudo inverse.

claim

$$\text{rank } \underline{X} = \text{rank } \underline{X}^T \underline{X} = \text{rank } \underline{X} \underline{X}^T$$

pf

1. claim $\text{rank } \underline{X} = \text{rank } \underline{X}^T$. don't want to show.
2. Use rank-nullity. (ie show $X$ & $X^T X$ have same Nullspace.)

   - Suppose $\underline{X} w = 0$. Then $\underline{X}^T \underline{X} w = 0$. So, $N(\underline{X}) \subset N(\underline{X}^T \underline{X})$
   - Suppose $\underline{X}^T \underline{X} w = 0$, Then $w^T X^T X w = 0 \Rightarrow |\underline{X} w|^2 = 0$
     $\Rightarrow \underline{X} w = 0 \Rightarrow N(\underline{X}^T \underline{X}) \subset N(\underline{X})$.

Claim follows by rank-nullity.

$$X : m \times n$$

Since we assumed rank $X = m$ ($X$ fat, linearly ind. rows), we get rank $X X^T = m \iff (X X^T)^{-1}$ exists.

Want to solve

$$X \hat{w} = y$$

Anything we can set $\hat{w}$ to to make this work?

Let $\hat{w} = X^T (X X^T)^{-1} y$

$$X \hat{w} = X X^T (X X^T)^{-1} y = y.$$

Note that $\hat{w} = X^T \underbrace{(X X^T)^{-1} y}_{=: z}$. This means $\hat{w} \in \text{row}(X^T)$

Hence, $\hat{w}$ is our unique soln.

Note:

- $X^+ := X^T (X X^T)^{-1}$ is pseudo-inverse when rows are lin indep.

- If $X$ were rank deficient, what would it mean? Most importantly, $\dim \text{col } X < m$. (since $\dim \text{row } X = \dim \text{col } X = \text{rank } X$.)

- Analysis would go through if we restricted to $y \in \text{col}(X)$ and had a more general pseudo-inv. taking $\text{col } X$ back to row $X$.

<u>Case II</u>: M = n

    • dont care right now since we're looking at full row/col
      rank matrices to isolate effects of each asymmetry
      separately first.

Case III    m > n

    • $X$ is tall

    ~~• linearly indep. cols~~

    • $\dim \text{col } X \leq n < m$

        └ so, $X$ not surjective

    • What about $N(X)$?

        └ may or may not be empty

    • If we assume rank $X = n$, then

        – $X$ has lin indep cols

        – $\dim \text{col } X = n$

        – $N(X) = \{0\}$

           └ <u>rank – nullity</u>:

             $A \in \mathbb{R}^{m \times n}$    $\underbrace{\text{rank } A}_{} + \dim N(A) = n.$

                              $(= \dim \text{col } A$
                               $= \dim \text{row } A)$

Assume rank $X = n$.

    • Buys you: $N(X) = \{0\}$.

Suppose $\hat{y} \in \text{col } X$. <u>claim</u>: $\exists$ a unique $\hat{w}$ solvs $X\hat{w} = \hat{y}$ $(**)$

<u>pf</u>    rank $X^T X = n$ $\Rightarrow$ $(X^T X)^{-1}$ exists. Also, since $\hat{y} \in \text{col } X$,

a $\hat{w}$ solvs $(**)$ $\exists$. observe

$$X\hat{w} = \hat{y} \iff (X^T X)^{-1} X^T X \hat{w} = (X^T X)^{-1} X^T y$$

$$\iff \hat{w} = (X^T X)^{-1} X^T y$$

Also, $(\mathbb{X}^T\mathbb{X})^{-1}\mathbb{X}^T$ has nullspace $= \{0\}$.

   └ why? TODO.

This proves claim.

Suppose $y$ possibly outside col $\mathbb{X}$. Want $\hat{w}$ s.t.

$$\hat{w} \in \underset{w}{\arg\min} \|\mathbb{X}w - y\|_2$$

But we're in a Hilbert space, so Soln to this is just projection.

  i.e., $\|\hat{y} - y\|_2^2 \leq \|y' - y\|^2 \quad \forall y' \in \text{col } \mathbb{X}$

for $\hat{y}$ in col $\mathbb{X}$ if $\hat{y}$ is the unique vector satisfying

$$(\hat{y} - y) \perp \text{col } \mathbb{X}.$$

   <u>ToDo:</u> Formally recall Hilbert proj thrm.

Hence, the optimal $\hat{w}$ satisfies

$$(\mathbb{X}\hat{w} - y) \perp \text{each col of } \mathbb{X}$$

$$\Updownarrow$$

$$(\mathbb{X}\hat{w} - y)^T\mathbb{X} = \vec{0}$$

$\iff$ $\hat{w}^T\mathbb{X}^T\mathbb{X} = y^T\mathbb{X}$ $\iff$ $\boxed{\mathbb{X}^T\mathbb{X}\,\hat{w} = \mathbb{X}^Ty.}$

know $(\mathbb{X}^T\mathbb{X})^{-1}$ exists.

$$\hat{w} = (\mathbb{X}^T\mathbb{X})^{-1}\mathbb{X}^Ty$$

## Observations

- $X^+ := (X^T X)^{-1} X^T$ is pseudo-inv. when cols are lin indep.

- If $X$ were rank deficient, what would it mean? $N(X) \neq 0$. So, many $\hat{w}$ satisfy $X\hat{w} = \hat{y}$ where $\hat{y} = \text{proj}_{col(x)}(y)$. So, many least squares solns. Which to pick?

## Computational consideration:

Have to invert $(X^T X) \in \mathbb{R}^{n \times n}$. What if $n$ (dim of feature space) is big?

Can solve optimization problem directly.

$$\min_{w} \underbrace{\| Xw - y \|^2}_{=: J(w)}$$

$$\nabla J(w) = (Xw - y)^T X$$

← double check. Can confirm by computing $\frac{\partial J}{\partial w_i}$ and see if this is $i$-th component of $(Xw-y)^T X$.

Use GD:

$$w_{t+1} = w_t - (Xw - y)^T X$$

- GD: $O(mn)$ flops per iteration
- Normal eq. / pseudo-inverse: $\approx O(n^3)$
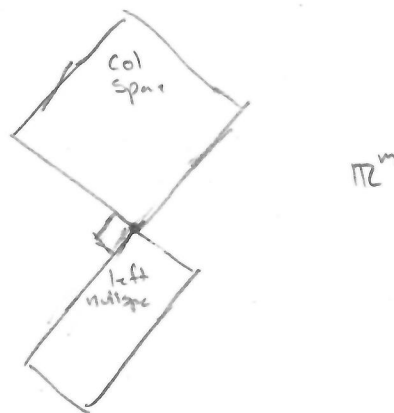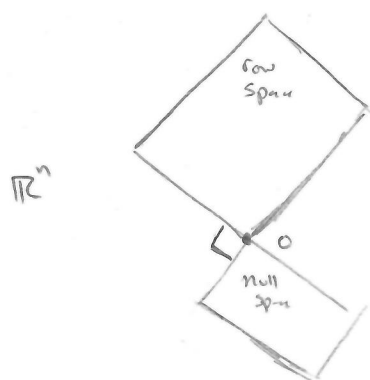- If $m \ll n$, GD seems probably better.

Omissions / take on faith so far:
- rank $A$ = rank $A^T$
- rank-nullity thm
- Hilbert proj. thm.
- form of $\nabla J$

# SVD

Diagonalization makes life incredibly easy, when possible. Life is particularly easy, when eigenvectors are orthogonal. The SVD provides a ___an orthonormal___ diagonal decomposition for arbitrary (non-square) matrices. The catch: The basis vectors used for domain / range can be different. So, SVD not necessarily suitable for applications w/ powers of matrices since you ___don't___ get

$$A^2 = S \Delta S^T S \Delta S^T = S \Delta^2 S^T \quad \text{But useful}$$

when this doesn't matter.

$$A \in \mathbb{R}^{m \times n}$$



Let $r = \text{rank } A$

We want an orthonormal basis $v_1, ..., v_r$ & $u_1, ..., u_r$ for col space s.t.

$$A v_i = \sigma_i u_i$$

What's a good candidate for $(v_i)_{i=1}^r$? How about eigenvectors of

$\boxed{A^T A \in \mathbb{R}^{n \times n}}$ $\quad AA^T \in \mathbb{R}^{m \times m}$
$\quad\quad$ ↳ basis for rowspace

Note: $A^T A$ & $AA^T$ are positive semi def. Why?
↳ b/c $x^T A^T A x = \|Ax\|^2 \geq 0 \ \forall x.$

Have, $A^T A v_i = \sigma_i^2 v_i \implies v_i^T A^T A v_i = \sigma_i^2 v_i^T v_i \implies \|A v_i\| = \sigma_i$

~~Let $u_i = A v_i$~~

Let $\hat{u}_i = A v_i$. Lets see what happens.

Have

$$A^T A v_i = \sigma_i^2 v_i \iff A A^T A v_i = \sigma_i^2 A v_i$$

$$\iff A A^T \hat{u}_i = \sigma_i^2 \hat{u}_i \quad \text{so,} \quad u_i := \frac{1}{\sigma_i} A v_i \text{ is a unit eigenvector}$$
of $A A^T$

Have

- $(u_i)_{i=1}^r$ ᴸorthonormal basis for col space
- $(v_i)_{i=1}^r$ ᴸorthonormal basis for row space
- $A v_i = \sigma_i u_i$
- $\sigma_i > 0$, $\sigma_i = \sqrt{\lambda_i}$, when $\lambda_i \in \sigma(A^T A) \setminus \{0\}$

Complete each basis w/ orthonormal basis for nullspace / left nullspace to
get $(v_1, \ldots, v_r, \ldots, v_n)$, ..., $(u_1, \ldots, u_r, \ldots u_m)$

- $V$ is $n \times n$ matrix w/ cols $(v_i)$
- $U$ is $m \times m$ w/ cols $(u_i)$
- $\Sigma$ is ? w/ diagonal entries $\sigma_i$
  ↳ $n \times n$ = dim $A$. See ↘

$$A v_i = \sigma u_i \iff A V = U \Sigma \iff \boxed{A = U \Sigma V^T}$$
$$\underbrace{\quad}_{SVD}$$

Pseudo - inverse

We know $A$ is an injective map from row span to col. space.
Want to construct an inverse that will bring any $y \in col(A)$
back to unique $x$ s.t. $Ax = y$. This is what pseudo-inv. does.
This all becomes easy if we can map $\sigma_i u_i$ back to $v_i$,
ie., want $A^+$ s.t.

$$A^+(\sigma_i u_i) = v_i$$

$$\Updownarrow$$

$$A^+ u_i = \frac{1}{\sigma_i} v_i$$

$$\Updownarrow$$

$$A^+ U = V \Sigma^+, \text{ where } \Sigma^+ \text{ is same as } \Sigma \text{ but w/}$$
$$\frac{1}{\sigma_i} \text{ on diag, where } \underbrace{\sigma_i \neq 0}$$
$$\text{holds by def.}$$
$$\text{of } \sigma_i, \text{ since only}$$
$$\text{defined up to } i = rank(A)$$

$$\Updownarrow$$

$$\boxed{A^+ := V \Sigma^+ U^T}$$

As a sanity check, say $Ax = y$ w/ $x \in row(A)$, $x = \sum_{i=1}^{r} \alpha_i v_i$.

$y = \sum_{i=1}^{r} \alpha_i \sigma_i u_i \Rightarrow A^+ y = \sum_{i=1}^{r} \alpha_i A^+ \sigma_i u_i = \sum_{i=1}^{r} \alpha_i v_i = x$.

So,

• $A^+ A x = x$, for $x \in row(x)$.

other implications:

more generally,

- If $x \in \mathbb{R}^n$, $\quad A^+ A x = \text{proj}_{\text{row}(A)}(x)$

- If $y \in \mathbb{R}^m$, $\quad A^+ y = A^+ (\hat{y} + z)$ for $\hat{y} = \text{proj}_{\text{col}(A)}(y)$,

and any $z \in N(A^T)$. Continuing,

$$A y = \hat{x}$$

where $\hat{x}$ is unique element of row$(A)$ solving $A\hat{x} = \hat{y}$,

w/ $\hat{y} = \text{proj}_{\text{col}(A)}(y)$.