# Today: Connections

- Maximum Likelihood estimation
- Maximum a-posteriori
- Tikhonov regression

## Recall: Ridge Regression

$$\min_{\vec{x}} \| A\vec{x} - \vec{b} \|^2 + \lambda^2 \| \vec{x} \|^2$$

$\underbrace{\qquad}_{\text{regularizer}}$

Note! → the solution to this $\neq$ solution to the least squares problem

$$\boxed{\vec{x}^* = (A^T A + \lambda I)^{-1} A^T \vec{b}}$$

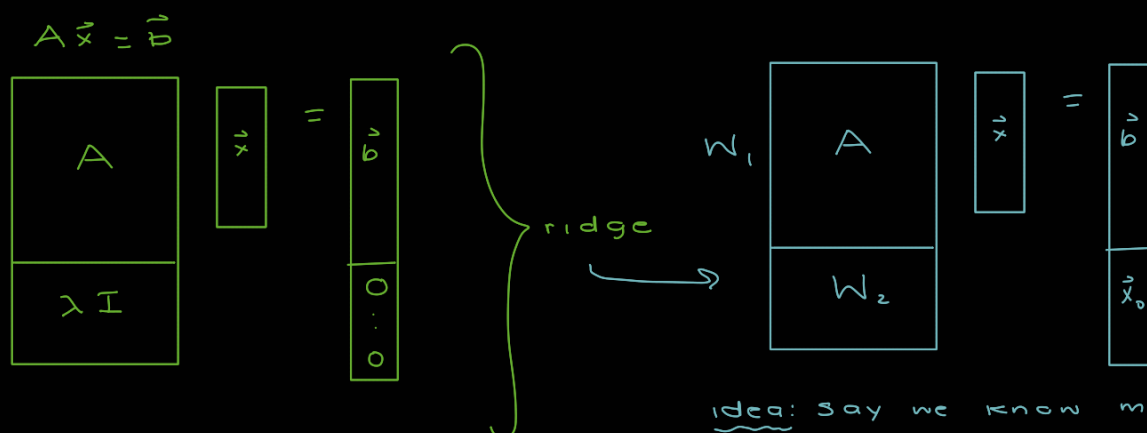↳ Multiple interpretations of this optimization problem

① Wanted to be robust to perturbations ⇒ reduce sensitivity towards them

↳ wanted to make sure $\sigma_i \{A\}$ not too large ⇒ shifted them away from 1

→ w/o this ridge coefficient/lambda term, our predicted coefficients were really large

↳ said we knew that $\vec{x}$ wasn't very large

② Ghost data: Add measurements to least squares setup saying how confident we were about closeness to $\vec{0}$

## Tikhonov Regularization

↳ generalization of ridge

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



ridge

idea: Say we know more about the structure of the matrices; can we weight them by some values to tell us more about our data?

### Tikhonov regularization

$$\min_{\vec{x}} \| W_1 (A\vec{x} - \vec{b}) \|_2^2 + \| W_2 (\vec{x} - \vec{x}_0) \|_2^2$$

## Probabilistic Perspective

- to figure out what side information to incorporate

$$y_i = g(x_i) + z_i$$

$\uparrow$ data points

- say we have a linear model:

$$g(\vec{x}_i) = \vec{x}_i^T \vec{\omega}$$

↳ $\vec{\omega}$ is our model

→ can rewrite this as

$$y_i = g(\vec{x}_i^T \vec{\omega}) + z_i$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \text{---} \ \vec{x}_1^T \ \text{---} \\ \vdots \\ \text{---} \ \vec{x}_n^T \ \text{---} \end{bmatrix} \begin{bmatrix} \vdots \\ \vec{\omega} \\ \vdots \end{bmatrix} + \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix}$$

- say noise is drawn from a Normal distribution: $z_i \sim N(0, \sigma_i^2)$

⌈ Recall:
The density fcn for a normal distribution:

$$F(z_i) = \frac{(e)^{-z_i^2 / 2\sigma_i^2}}{\sqrt{2\pi}\,\sigma_i}$$

→ assuming these are all iid (independent, identically distributed)

$$\vec{y} \approx X\vec{w}$$

$\hookrightarrow$ Note: we can solve this with least squares but that doesn't tell us everything about our measurements; don't take the distribution of our noise into account

solution!?
MLE

## Maximum Likelihood estimation: $\vec{w}$ that makes observed

$$\underset{\vec{w}_0}{\text{argmax}} \; f(Y_1 = y_1, Y_2 = y_2, \cdots, Y_n = y_n \mid \vec{w} = \vec{w}_0)$$

$\hookrightarrow$ we're assuming we know nothing about $\vec{w}$ (no prior)

$\hookrightarrow$ ask us: what is the unknown that makes our data most likely?

$\rightarrow$ trying to maximize the density in the case of continuous RV's

- Since $z_i$'s are iid, we can write $y_i$'s as a product (because the only thing they're contingent on is the $\vec{w}$ but we're conditioning that out):

$$\underset{\vec{w}_0}{\text{argmax}} \; f(Y_1 = y_1, Y_2 = y_2, \cdots, Y_n = y_n \mid \vec{w} = \vec{w}_0)$$

$$= \underset{\vec{w}_0}{\text{argmax}} \; \prod_{i=1}^{n} \boxed{f(Y_i = y_i \mid \vec{w} = \vec{w}_0)}$$

$\hookrightarrow$ want to understand this conditional density

### Ayan's attempt

$$f(Y_i = y_i \mid \vec{w} = \vec{w}_0)$$

$$f(z_i) = \frac{(e)^{-z_i^2/2\sigma_i^2}}{\sqrt{2\pi}\,\sigma_i}$$

$\hookrightarrow$ use Baye's rule? idk

### Ranade

$$f(Y_i = y_i \mid \vec{w} = \vec{w}_0) = f(\vec{x}_i^T \vec{w}_0 + z_i = y_i \mid \vec{w} = \vec{w}_0)$$

$$= f(z_i = y_i - \vec{x}_i^T \vec{w}_0 \mid \vec{w} = \vec{w}_0)$$

$$= \frac{c^{-\frac{(y_i - \vec{x}_i^T \vec{w}_0)^2}{2\sigma_i^2}}}{\sqrt{2\pi}\cdot\sigma_i}$$

can now rewrite this as

$$\underset{\vec{w}_0}{\text{argmax}} \; \prod_{i=1}^{n} \frac{e^{-\frac{(y_i - \vec{x}_i^T \vec{w}_0)^2}{2\sigma_i^2}}}{\sqrt{2\pi}\cdot\sigma_i}$$

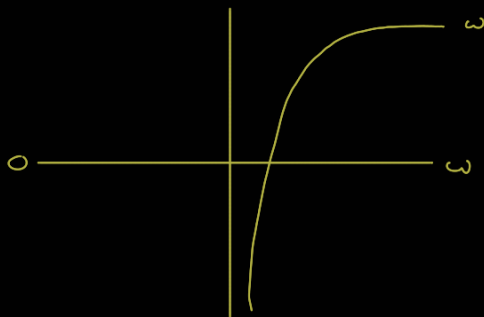ⓦ since denominator not contingent on $\vec{\omega}_0$, can rewrite as:

$$\underset{\vec{\omega}_0}{\text{argmax}} \left(\frac{1}{\sqrt{2\pi}\prod\limits_{i=1}^{n}\sigma_i^2}\right)^n \prod_{i=1}^{n} e^{-\frac{(y_i - \vec{x}_i^T\vec{\omega}_0)^2}{2\sigma_i^2}}$$

$$\prod_{i=1}^{n} e^{x_i} = e^{x_1}e^{x_2}\cdots e^{x_n} = e^{\sum_i^n x_i} = \exp\left\{\sum_{i=1}^{n}x_i\right\}$$

⌐ Recall:

• $\underset{\vec{x}}{\text{argmax}} \|A\vec{x} - \vec{b}\|_2 = \underset{\vec{x}}{\text{argmax}} \|A\vec{x} - \vec{b}\|_2^2$

⋮

• $\underset{\omega}{\text{argmax}} F(\omega) = \underset{\omega}{\text{argmax}} \log F(\omega)$



because it's monotonically increasing. This is the case for any increasing fcn. ⌐

⌐ IMPORTANT

$$\underset{\vec{x}}{\max} \|A\vec{x} - \vec{b}\|_2 \neq \underset{\vec{x}}{\max} \|A\vec{x} - \vec{b}\|_2^2$$

lets us rewrite as

$$\underset{\vec{\omega}_0}{\text{argmax}} \frac{1}{(\sqrt{2\pi})^n} \frac{1}{\prod\limits_{i=n}^{}\sigma_i} \exp\left\{-\sum \frac{(y_i - \vec{x}_i^T\vec{\omega}_0)^2}{2\sigma_i^2}\right\}$$

taking logs ⸴ dropping constants

$$= \underset{\vec{\omega}_0}{\text{argmax}} \left\{-\sum_{i=1}^{n} \frac{(\vec{y}_i - \vec{x}_i^T\vec{\omega}_0)^2}{2\sigma_i^2}\right\}$$

$$= \underset{\vec{\omega}_0}{\text{argmin}} \sum_{i=1}^{n} \frac{(\vec{y}_i - \vec{x}_i^T\vec{\omega}_0)^2}{\boxed{2\sigma_i^2}}$$

⌐ Recall: our least-squares cost is

$$\sum_{i=1}^{n}(x_i^T\vec{\omega} - y_i)^2$$

this is the same as our least squares cost except it's ⌐weighted⌐ by some value

$$= \underset{\vec{\omega}_0}{\text{argmin}} \|S(X\vec{\omega}_0 - \vec{y})\|_2^2$$

$$\hookrightarrow S = \begin{bmatrix} \sqrt{\frac{1}{2\sigma_1^2}} & & & O \\ & \sqrt{\frac{1}{2\sigma_2^2}} & & \\ & & \ddots & \\ O & & & \sqrt{\frac{1}{2\sigma_n^2}} \end{bmatrix}$$

# MAP (Maximum a Posteriori)

↳ used if we have a prior on $\vec{w}$

$$y_i = \vec{x}_i^T \vec{w} + z_i \qquad z_i \sim N(0, \sigma_i^2) \qquad w_i \sim N(\mu_i, \rho_i^2)$$

$$\vec{w} \sim N(\vec{\mu}, \Sigma_w)$$

$$\Sigma_w = \begin{bmatrix} \rho_1^2 & & & O \\ & \rho_2^2 & & \\ & & \ddots & \\ O & & & \rho_n^2 \end{bmatrix}$$

↑ covariance matrix

↳ variances along diagonal

→ 0 covariance

( $w_i$ not correlated w/ $w_j$ if $i \neq j$ )

Recall from IGA, imaging:

$$w_1 \; w_2 \; w_3$$
$$\vdots$$

↳ if we have

What is most likely data given $y_1, y_2, \dots, y_n$?

$$\underset{\vec{w}}{\text{argmax}} \; F(\vec{w} \mid Y = \vec{y})$$

↳ in MLE, these were swapped

→ using Baye's rule, rewrite as:

$$f(\vec{w} \mid Y = \vec{y}) = \frac{f(Y = \vec{y} \mid \vec{w}) f(\vec{w})}{F(\vec{y})}$$

→ $F(\vec{y})$ this is a constant; your data is what it is (can ignore it as a constant in our argmax)

$$= \underset{\vec{w}}{\text{argmax}} \; F(Y = \vec{y} \mid \vec{w}) f(\vec{w})$$

$$= \underset{\vec{w}}{\text{argmax}} \; \prod_{i=1}^{n} F(Y = y_i \mid \vec{w}) \cdot f(\vec{w})$$

↳ already computed this previously

Multivariate normal
$$w \sim \frac{e^{-\frac{(x-\mu)^2}{2\sigma_i^2}}}{\sqrt{2\pi}\,\sigma^2}$$

Just the multivariate Gaussian density

$$= \underset{\vec{w}}{\text{argmax}} \; \prod_{i=1}^{n} \exp\left\{ \frac{-(y_i - \vec{x}_i^T \vec{w})^2}{2\sigma_i^2} \right\} \cdot \frac{\exp\left\{ -(\vec{w} - \vec{\mu})^T \Sigma_w^{-1} (\vec{w} - \vec{\mu}) \right\}}{(\sqrt{2\pi})^n \left(\prod_{i=1}^{n} \rho_i\right)}$$

⟵ $\sqrt{2\pi}\,\sigma_i$ & ⟶ neglible constants

$$= \underset{\vec{w}}{\text{argmax}} \; \exp\left\{ -\sum_{i=1}^{n} \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2\sigma_i^2} + -(\vec{w} - \vec{\mu})^T \Sigma_w^{-1} (\vec{w} - \vec{\mu}) \right\}$$

$$= \underset{\vec{w}}{\text{argmin}} \; \sum_{i=1}^{n} \frac{(y_i - \vec{x}_i^T \vec{w})^2}{2\sigma_i^2} + (\vec{w} - \vec{\mu})^T \Sigma_w^{-1} (\vec{w} - \vec{\mu})$$

Because $\Sigma_w$ is symmetric can write as $\Sigma_w^{-1} = \sqrt{\Sigma_w^{-1}}\sqrt{\Sigma_w^{-1}}$

$$= \underset{\vec{w}}{\text{argmin}} \; \| S(X\vec{w} - \vec{y}) \|_2^2 + \left\| \sqrt{\Sigma_w^{-1}} (\vec{w} - \vec{\mu}) \right\|_2^2$$

↳ if $\vec{\mu} = \vec{0}$ we'd get ridge regression