

Generative Models (Just-so-stories)

Suppose we know $y = w_*^T x + \text{noise}$

w_* "true" model, noise Gaussian w/ mean zero and variance σ^2 .

In vector form $\vec{y} = \bar{X} w_* + \vec{n}$

(this is never really true, but is instructive)

~~What is~~

PREDICTION ERROR: Let $y_i^{\text{true}} = w_*^T x_i$

For a learned model \hat{w} , $y_i^{\text{pred}} = \hat{w}^T x_i$

prediction error: $\frac{1}{n} \sum_{i=1}^n (y_i^{\text{pred}} - y_i^{\text{true}})^2 = \frac{1}{n} \|\bar{X}(\hat{w} - \vec{w}_*)\|^2$

OLS: $\frac{1}{n} \sum_{i=1}^n (y_i^{\text{pred}} - y_i^{\text{true}})^2 = \frac{1}{n} \|\bar{X}(\vec{w}_{\text{OLS}} - \vec{w}_*)\|^2$

$\vec{w}_{\text{OLS}} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \vec{y}$ (assume $d \leq n$ and \bar{X} full rank)

$$\begin{aligned} \bar{X}(\vec{w}_{\text{OLS}} - \vec{w}_*) &= \bar{X}((\bar{X}^T \bar{X})^{-1} \bar{X}^T \vec{y} - w_*) \\ &= \bar{X}((\bar{X}^T \bar{X})^{-1} \bar{X}^T (\bar{X} \vec{w}_* + \vec{n}) - w_*) \\ &= \bar{X}(\bar{X}^T \bar{X})^{-1} \bar{X}^T \vec{n} \end{aligned}$$

$$\bar{X} = \bar{U} \bar{S} \bar{V}^T$$

$$\begin{aligned} \bar{X} (\bar{X}^T \bar{X})^{-1} X^T &= (\bar{U} \bar{S} \bar{V}^T) (\bar{V} \bar{S}^T \bar{S} \bar{V}^T)^{-1} (\bar{V} \bar{S}^T \bar{U}^T) \\ &= (\bar{U} \bar{S} \bar{V}^T) (\bar{V} (\bar{S}^T \bar{S})^{-1} \bar{V}^T) (\bar{V} \bar{S}^T \bar{U}^T) \\ &= \bar{U} \bar{S} (\bar{S}^T \bar{S})^{-1} \bar{S}^T \bar{U}^T \\ &= \bar{U} \begin{bmatrix} \bar{I}_d & \bar{O} \\ \bar{O} & \bar{O} \end{bmatrix} \bar{U}^T \end{aligned}$$

$$\begin{aligned} \mathbb{E} \left[\frac{1}{n} \| \bar{X} (\bar{w}_{OLS} - \bar{w}_*) \|^2 \right] &= \frac{1}{n} \mathbb{E} \left[\left\| \bar{U} \begin{bmatrix} \bar{I}_d & \bar{O} \\ \bar{O} & \bar{O} \end{bmatrix} \bar{U}^T \bar{n} \right\|^2 \right] \\ &= \frac{1}{n} \mathbb{E} \left[\bar{n}^T \bar{U} \begin{bmatrix} \bar{I}_d & \bar{O} \\ \bar{O} & \bar{O} \end{bmatrix} \bar{U}^T \bar{n} \right] \\ &= \frac{1}{n} \mathbb{E} \text{Trace} \left[\begin{bmatrix} \bar{I}_d & \bar{O} \\ \bar{O} & \bar{O} \end{bmatrix} \bar{U}^T \bar{n} \bar{n}^T \bar{U} \right] \\ &= \frac{1}{n} \text{Trace} \left[\begin{bmatrix} \bar{I}_d & \bar{O} \\ \bar{O} & \bar{O} \end{bmatrix} \bar{U}^T \mathbb{E} [\bar{n} \bar{n}^T] \bar{U} \right] \\ &= \frac{1}{n} \text{Trace} \left[\begin{bmatrix} \bar{I}_d & \bar{O} \\ \bar{O} & \bar{O} \end{bmatrix} \bar{U}^T (v^2 \mathbf{I}_n) \bar{U} \right] \\ &= v^2 \frac{d}{n} \end{aligned}$$

PIN V:

assume

$w_{\star} \in \text{range}(\bar{X})$

$$\vec{w}_{\text{PINV}} = \bar{X}^{\dagger} \vec{y}, \quad \text{rank}(\bar{X}) = r$$

$$\begin{aligned} \bar{X}(\vec{w}_{\text{PINV}} - \vec{w}_{\star}) &= \bar{X} \left[(\bar{X}^{\top} \bar{X})^{-1} \bar{X}^{\top} (\bar{X} \vec{w}_{\star} + \vec{n}) - \vec{w}_{\star} \right] \\ &= \bar{X} (\bar{X}^{\top} \bar{X})^{-1} \bar{X}^{\top} \vec{n} \end{aligned}$$

BUT

$$\begin{aligned} \mathbb{E} \left[\frac{1}{n} \|\bar{X}(\vec{w}_{\text{OLS}} - \vec{w}_{\star})\|^2 \right] &= \frac{1}{n} \text{Tr} \left[\begin{bmatrix} \bar{I}_r & \vec{0} \\ \vec{0} & \delta \end{bmatrix} \bar{U}^{\top} (n^{-1} \mathbf{I}) \bar{U} \right] \\ &= \frac{\sigma^2 r}{n} \end{aligned}$$

RIDGE attempts to approximate PINV by
choice of δ

Random vectors

$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$p(\vec{x})$ is probability density for \vec{x} .

$$p(\vec{x}) = \frac{\partial^d}{\partial x_1 \dots \partial x_d} \text{Pr}[x_1 \leq a_1, x_2 \leq a_2, \dots, x_d \leq a_d]$$

$$\int p(\vec{x}) = 1$$

$$p(x_1) = \int p(x) dx_2 dx_3 \dots dx_d$$

Mean: $\mathbb{E}[\vec{x}] = \vec{\mu}$

Covariance: $\bar{\Sigma} = \mathbb{E}[(\vec{x} - \vec{\mu})(\vec{x} - \vec{\mu})^T]$
 $= \mathbb{E}[\vec{x} \vec{x}^T] - \vec{\mu} \vec{\mu}^T$

Σ_{ii} ~~variance~~ = variance of x_i

Σ_{ij} = covariance of x_i and x_j

If \vec{x} and \vec{z} are random vectors and

$$\vec{u} = \bar{A} \vec{x} + \bar{B} \vec{z} + \vec{c}, \text{ then}$$

$$\mu_u = \bar{A} \vec{\mu}_x + \bar{B} \vec{\mu}_z + \vec{c}$$

$$\Sigma_u = \bar{A} \bar{\Sigma}_x \bar{A}^T + \bar{A} \bar{\Sigma}_{xz} \bar{B}^T + \bar{B} \bar{\Sigma}_{zx} \bar{A}^T + \bar{B} \bar{\Sigma}_z \bar{B}^T$$

If $z = \vec{v}^T \bar{x}$ for non random \vec{v} then

$$\bar{\Sigma}_z = \sigma_z^2 = \vec{v}^T \bar{\Sigma}_x \vec{v} \geq 0$$

$\bar{\Sigma}_x$ is positive semidefinite (symmetric, all eigenvalues non-negative)

Least-squares revisited

Suppose (\vec{x}, y) has mean $(0, 0)$

Consider $\mathbb{E}[(\vec{w}^\top \vec{x} - y)^2] = \text{var}(\vec{w}^\top \vec{x} - y)$

$$= \vec{w}^\top \Sigma_x \vec{w} - 2 \vec{w}^\top \Sigma_{xy} + \sigma_y^2$$

Then ~~min~~ the minimizer of expected squared loss is

$$\vec{w}_{LS} = \Sigma_x^{-1} \Sigma_{xy} \quad (\text{compare to } (\bar{X}^\top \bar{X})^{-1} (\bar{X}^\top \bar{y}))$$

$$\min_{\vec{w}} \mathbb{E}[(\vec{w}^\top \vec{x} - y)^2] = \sigma_y^2 - \vec{\Sigma}_{yx} \vec{\Sigma}_x^{-1} \vec{\Sigma}_{xy}$$

($\vec{w}_{LS}^\top \vec{x} - y$ is a random variable w/ mean 0 and variance $\sigma_y^2 - \vec{\Sigma}_{yx} \vec{\Sigma}_x^{-1} \vec{\Sigma}_{xy}$)

~~min~~

Multivariate Gaussian: Everyone's favorite

(i) only need to know first and second moments

(ii) CLT implies most data resembles Gaussians when n is large

$$p(\vec{x}) = \frac{1}{(2\pi)^d |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right) \\ = \mathcal{N}(\vec{\mu}, \Sigma)$$

$$\int p(\vec{x}) dx_1, \dots, dx_d = \mathcal{N}(\mu_1, \sigma_1^2)$$

Diagonal Case:

$$p(\vec{x}) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2} \frac{(x_i - \mu_i)^2}{\sigma_i^2}\right)$$

For Gaussians, covariance = 0 \Rightarrow independence.

Isocontours: $p(\vec{x}) = p_0 \iff (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) = c_0$

Level sets of probability are ellipses.

Center is at $\vec{\mu}_x$, Σ^{-1} determine the axes.

$$\Sigma = \sum_{i=1}^d \lambda_i \vec{V}_i \vec{V}_i^T, \quad \vec{V}_i \text{ determine major axes of ellipse} \\ \lambda_i \text{ determine lengths}$$

$$\underset{\vec{w}}{\text{minimize}} \quad \vec{w}^T \overline{Q} \vec{w} - 2 \vec{p}^T \vec{w} + r$$

complete the square

$$(\vec{w} - \overline{Q}^{-1} \vec{p})^T \overline{Q} (\vec{w} - \overline{Q}^{-1} \vec{p}) + \vec{p}^T \overline{Q}^{-1} \vec{p} + r$$

$$\vec{w} = \overline{Q}^{-1} \vec{p}$$

$$\text{cost} = r - \vec{p}^T \overline{Q}^{-1} \vec{p}$$