

Notes for CS 181: Probabilistic Regression

Last time: linear regression

$$\hat{y} = w^T x$$

Chose loss: least squares: $\mathcal{L}_D(w) = \frac{1}{2} \sum_n (y_n - w^T x_n)^2$

Interpretation 1: Projection View

$$\underbrace{Y^T}_{\substack{\text{N-dim} \\ \text{vect.}}} \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = W^T X \underbrace{\begin{bmatrix} x_{11} & \dots & x_{1N} \\ \vdots & & \vdots \\ x_{D1} & \dots & x_{DN} \end{bmatrix}}_{\substack{\text{D} \\ \text{N-dim} \\ \text{vecs}}} \quad \text{3D diagram}$$

if $D < N$, then $x_1 \dots x_D$ define some linear subspace in \mathbb{R}^N

recall: compute projection: $\sum_d \underbrace{x_d}_{\substack{\text{N-dim} \\ \text{vec.}}} \cdot \frac{\langle x_d, y \rangle}{\langle x_d, x_d \rangle} \Rightarrow \underbrace{X^T (X X^T)^{-1} X}_{\substack{\text{projection} \\ \text{operator}}} Y$

note: linear algebra applies to random variables also.

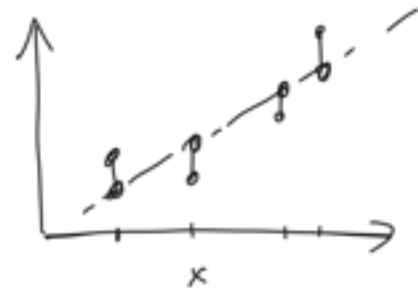
Probabilistic View: (idea: We're going to make a "story" for how our variables were created. "generative model")

ex: $y_n = w^T x_n + \epsilon_n$, $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$

we've made explicit:

→ form of noise (Gaussian)

→ indep of noise across samples



"graphical model"

$$\Pr(\text{data} | \text{model}) = ?$$

"likelihood" of model



Quick Review of Gaussian distributions:

$$r(z; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(z-\mu)^2\right\}$$



$$= \prod_n \Pr(y_n | x_n, w, \sigma^2) \quad \left[\begin{array}{l} \text{because noise} \\ \text{indep, } y_n \text{ only depends} \\ \text{on } x_n, w, \sigma^2 \end{array} \right]$$

First, let's take logs (monotonic):

$$\begin{aligned} \log \Pr(\text{data} | \text{model}) &= \sum_n \log \Pr(y_n | x_n, w, \sigma^2) \quad y_n \sim \mathcal{N}(w^T x_n, \sigma^2) \\ &= \sum_n \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_n - w^T x_n)^2 \right\} \\ &= \underbrace{N \log \frac{1}{\sqrt{2\pi\sigma^2}}}_{\text{no dep on } w} + \underbrace{\sum_n \left(-\frac{1}{2\sigma^2} \right) (y_n - w^T x_n)^2}_{\text{if we min } \sum_n (y_n - w^T x_n)^2, \text{ we max } \Pr(\text{data} | \text{model})} \end{aligned}$$

for practice, we can do the same w/ matrices...

$$\Pr(Y | X, w, \sigma^2)$$

$$\Rightarrow \text{multivariate Gaussian: } \underline{\mu} = w^T X, \quad \underline{\Sigma} = \begin{bmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{bmatrix} = \mathbf{I}_N \sigma^2$$

$$Y \sim \mathcal{N}(w^T X, \underline{\Sigma})$$



$$\text{formula for multivariate Gaussian: } \frac{1}{\sqrt{2\pi|\underline{\Sigma}|}} \exp \left\{ -\frac{1}{2} (z - \underline{\mu})^T \underline{\Sigma}^{-1} (z - \underline{\mu}) \right\}$$

$$\begin{aligned} \log \frac{1}{\sqrt{2\pi|\underline{\Sigma}|}} &+ \underbrace{\left(-\frac{1}{2} \right) (Y - w^T X)^T \underline{\Sigma}^{-1} (Y - w^T X)}_{-\frac{1}{2\sigma^2} (Y - w^T X)^T (Y - w^T X)} \\ &\underbrace{-\frac{N}{2} \log 2\pi\sigma^2} \end{aligned}$$

we solve (max of w) \Rightarrow "maximum likelihood estimation"

we can do more: estimate σ_{ML}^2
 \sim max likelihood

$$\frac{\partial}{\partial \sigma^2} : -\frac{N}{2} \cdot \frac{1}{\sigma^2} - \frac{1}{2} (Y - w^T X)^T (Y - w^T X) \left(\frac{1}{\sigma^2} \right)^2 = 0$$

$$\sigma_{ML}^2 = \frac{1}{N} (Y - w^T X)^T (Y - w^T X) \quad \left. \vphantom{\sigma_{ML}^2} \right\} \begin{array}{l} \text{empirical} \\ \text{variance} \end{array}$$

Going further: we can add to story: what if w is a random variable?





write down as ~~the~~ generative model:

$$w \sim p(w)$$

for n in $1 \dots N$

$$e_n \sim \mathcal{N}(0, \sigma_n^2)$$

$$y_n = w^T x_n + e_n$$

and now, we can write joint prob: $\frac{1}{Z} p(\{y_n\}, w | X, \sigma_n^2)$

$$= \underbrace{p(\{y_n\} | \{x_n\}, w, \sigma_n^2)}_{\text{likelihood from before}} \underbrace{p(w | X, \sigma_n^2)}_{p(w)} = p(\{y_n\}, w | X, \sigma_n^2) \propto p(w | X, \sigma_n^2, \{y_n\})$$

taking logs:

$$\log \text{JOINT} = \underbrace{\log \text{LIKELIHOOD}} + \log \text{PRIOR}$$

ex. $w_d \sim \mathcal{N}(0, \sigma_w^2)$

$$N \log \frac{1}{\sqrt{2\pi}\sigma_n^2} + \sum_n \left(-\frac{1}{2\sigma_n^2} \right) (y_n - w^T x_n)^2 + D \log \frac{1}{\sqrt{2\pi}\sigma_w^2} + \sum_d \left(-\frac{1}{2\sigma_w^2} \right) (w_d)^2$$

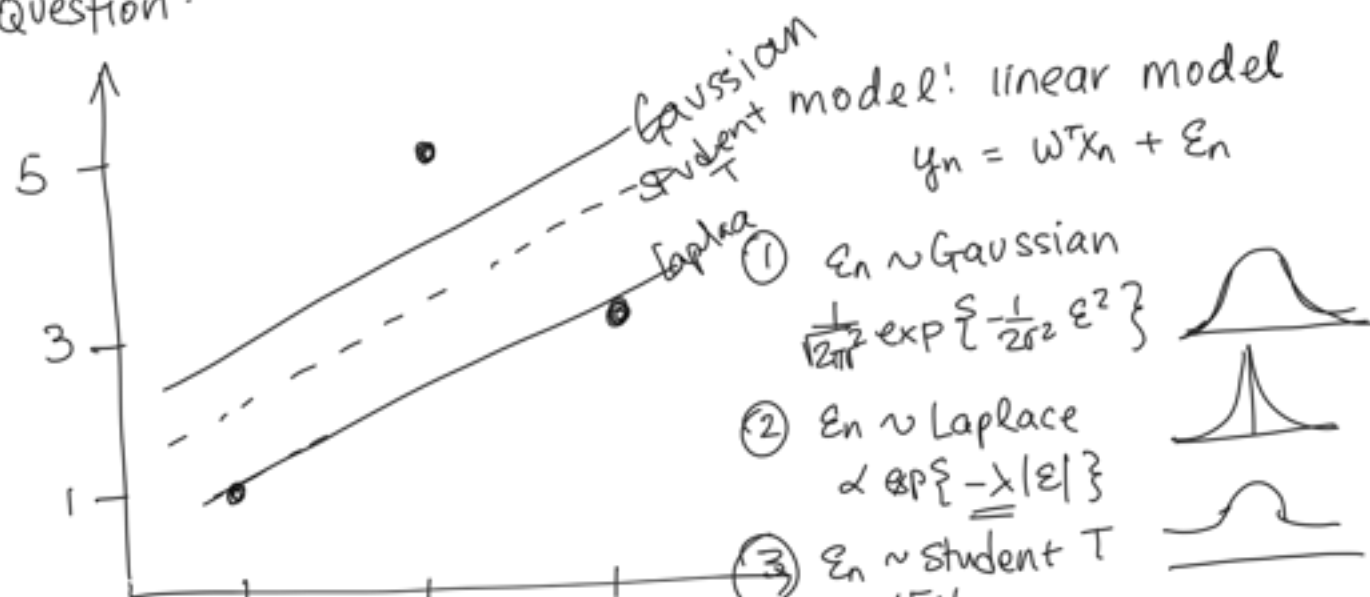
We can ask same kinds of questions:

what is the w that maximizes JOINT prob? (best σ_n^2 , best σ_w^2)

before "maximum likelihood" solution (ML)

now: "maximum a posteriori" solution (MAP)

Concept Question:



$$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 2(1+\epsilon^2)^{-1} \end{matrix}$$

1. Convert all three options into log losses to minimize
2. Think about what effect diff losses will have on fitted w ?

$$\begin{aligned} \text{Gaussian: } \epsilon^2 &\rightarrow \sum_n (y_n - w^T x_n)^2 \\ \text{Laplace: } &\rightarrow \sum_n |y_n - w^T x_n| \\ \text{Student T: } &\rightarrow \sum_n \log(1 + (y_n - w^T x_n)^2) \end{aligned}$$