

CS181: Probabilistic Regression

Housekeeping HW1: due Thurs
 TI-Q4.py
 • Notes, Notation

Supervised Learning ($x \rightarrow y$)

continuous - y			
discrete y			
	non prob	prob	

Last time:
 $\hat{y} = \underline{w}^T \underline{x}$ model
 $J(\underline{w}) = \frac{1}{2} \sum_n (y_n - \hat{y}_n)^2$ loss
 $\underline{w}^* = (\underbrace{X^T X}_{\text{looks like var}})^{-1} (\underbrace{X^T Y}_{\text{looks like cov}})$

This time:

1) provide geometric view of linear regression.

$\underline{y}^T \approx \underline{w}^T X^T$
 $\underbrace{N\text{-dim}}_{\underline{y}^T} \quad \underbrace{D}_{\underline{w}^T} \quad \underbrace{D \times N}_{X^T}$ } D vectors of length N

$\sum_d x_d = \frac{\langle x_d, y \rangle}{\langle x_d, x_d \rangle}$
 N-dim vector that is all elements of dth dim

Book 2.6.3 $\underline{w} = X(X^T X)^{-1} X^T Y$
 our weight vector

2) Probabilistic Methods.

Key Idea: Generative Model.
 "Story for how the data came to be"

Story / Generative Model:

- we some input \underline{x}_n
- multiply \underline{x}_n by \underline{w} : $\underline{w}^T \underline{x}_n$
- output is the product + noise, $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$

Book 2.6.2 $y_n = \underline{w}^T \underline{x}_n + \epsilon_n$

Probability of the data given $\underline{w}, \sigma^2 \rightarrow$ likelihood of the model

$$\Pr(\text{data} | \text{model}) = \prod_n \Pr(y_n | \underline{x}_n, \underline{w}, \sigma^2)$$

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

$$\sum_n \log \Pr(y_n | \underline{x}_n, \underline{w}, \sigma^2)$$

recall: $\mathcal{N}(z; \mu, \sigma^2) =$

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

$$\sum_n \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(y_n - \underline{w}^T \underline{x}_n)^2\right\}$$

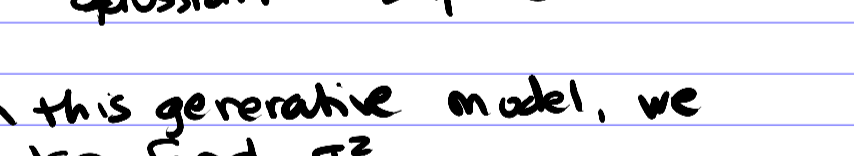
$$= \sum_n \underbrace{\log \frac{1}{\sqrt{2\pi}\sigma}}_{\text{no dep. on } \underline{w}} - \frac{1}{2\sigma^2} (y_n - \underline{w}^T \underline{x}_n)^2$$

has dep on \underline{w} , looks just like our least squares loss!!

So now, if we maximize prob of data, we will get the same solution as least squares!

• solution is known as the "maximum likelihood solution" (for this generative model)

Note: we could use other noise models.



Given this generative model, we can also find $\hat{\sigma}_{MLE}^2$.

First, just for practice, let's write our probability of data w/ vectors

$$\log \frac{1}{\sqrt{2\pi|\Sigma|}} \exp\left\{-\frac{1}{2}(y - \underline{w}^T \underline{x})^T \Sigma^{-1} (y - \underline{w}^T \underline{x})\right\}$$

$\Sigma = \sigma^2 I$ (vectors of size N)

$$= -\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y - \underline{w}^T \underline{x})^T (y - \underline{w}^T \underline{x})$$

Take gradient w.r.t. σ^2 (note: grad w.r.t. σ^2 not σ)

$$\frac{\partial}{\partial \sigma^2} \text{ is: } -\frac{N}{2} \cdot \frac{1}{\sigma^2} + \frac{1}{2} (y - \underline{w}^T \underline{x})^T (y - \underline{w}^T \underline{x}) \left(\frac{1}{\sigma^2}\right)^2$$

Solve by multiplying through by σ^4

$$0 = -N\sigma^2 + (y - \underline{w}^T \underline{x})^T (y - \underline{w}^T \underline{x})$$

$$\hat{\sigma}_{MLE}^2 = \frac{(y - \underline{w}^T \underline{x})^T (y - \underline{w}^T \underline{x})}{N}$$

vec. of length N that contains diff btw true y and our predicted y \rightarrow product: $\sum_n (\text{err})^2$

$$\hat{\sigma}_{MLE}^2 = \frac{\sum_n (\text{err})^2}{N}$$

What we did: prob approach to regression) 1) Generative model 2) Max likelihood of model w.r.t params.

Note: in a few lectures, we will return to this model, put priors over e.g. \underline{w}

$\underline{w} \sim p(\underline{w})$, then $y = \underline{w}^T \underline{x} + \epsilon$ today