

03/14/18

$y = A(x) + e$  (e and y, random)  
 $P(y|x)$  likelihood function

- frequentist:  $\hat{x}(y)$  estimator of  $x$ , random
- Bayesian:  $P(x|y)$  posterior  
 $P(x|y) \sim P(y|x)P(x)$  ← prior

Def MSE =  $E \|\hat{x} - x\|^2$

Prop. =  $E \|\hat{x} - E\hat{x}\|^2 + \|E\hat{x} - x\|^2$   
Var  $\hat{x}$  Bias( $\hat{x}, x$ )<sup>2</sup>

Reminder:  $E[f(y)] = \int f(y) p(y|x) dy$ .

Def: Fisher score of  $y$ : (random)

$V = \frac{\partial}{\partial x} \ln p(y|x) = \frac{\frac{\partial}{\partial x} p(y|x)}{p(y|x)}$

Prop.:  $E V = \int \frac{\partial}{\partial x} p(y|x) \cdot \frac{1}{p(y|x)} p(y|x) dy$   
 $= \frac{\partial}{\partial x} \int p(y|x) dy = \frac{\partial}{\partial x} 1 = 0$ .

Def: Fisher information of  $y$ . (number)  
 ("how much information  $y$  contains about  $x$ ").

$J(x) = E \left[ \left( \frac{\partial}{\partial x} \ln p(y|x) \right)^2 \right] (= EV^2 = \text{Var } V)$

ex.  $y = (y_1, \dots, y_n)$  iid, then  
 $P(y|x) = \prod_i P(y_i|x) \Rightarrow V_y = \sum_{i=1}^n V_{y_i}$   
 $J_y(x) = n J_{y_1}(x)$

(Additive, and relative to what we want to learn, here  $x$ )

Thm (Gomén-Rao lower bound,  $x \in \mathbb{R}$ )

Let  $\hat{x}$  be unbiased:  $E\hat{x} = x$ .

Then,  $\text{Var } \hat{x} \geq \frac{1}{J(x)}$

Pf Recall  $J_f(x) = EV^2$  ( $EV=0$ )

$\text{Var } \hat{x} = E(\hat{x}-x)^2$  ( $E\hat{x}=x$ )

Cauchy-Schwarz  $(\int fg)^2 \leq (\int f^2)(\int g^2)$

$$\left( E[V(\hat{x}-x)] \right)^2 \leq \underbrace{(EV^2)}_{J(x)} \underbrace{(E(\hat{x}-x)^2)}_{\text{Var } \hat{x}}$$

$$\stackrel{||}{=} E[V\hat{x}] = \int \frac{\frac{\partial}{\partial x} p(y|x)}{p(y|x)} \hat{x}(y) p(y|x) dy$$

$$= \frac{\partial}{\partial x} \int \hat{x}(y) p(y|x) dy$$

$$= \frac{\partial}{\partial x} E\hat{x} = \frac{\partial}{\partial x} x = 1$$

□

Thm Biased case:

$$\text{Var } \hat{x} \geq \frac{\left( 1 + \frac{\partial}{\partial x} \overbrace{\text{Bias}(\hat{x}, x)}^{E\hat{x}-x} \right)^2}{J(x)}$$

Thm Multivariate, unbiased case:

$$\Sigma_{ij} = \text{Cov}(\hat{x}_i, \hat{x}_j) = E[\hat{x}_i \hat{x}_j] \quad (\Sigma_{ii} = \text{Var } \hat{x}_i)$$

$$J_{ij}(x) = E \left[ \frac{\partial}{\partial x_i} \ln p(y|x) \frac{\partial}{\partial x_j} \ln p(y|x) \right]$$

(Fisher information matrix)

then  $\Sigma \geq J^{-1}(\alpha)$  meaning  $\Sigma - J^{-1}(\alpha) \succeq 0$ .

in particular  $\text{Var } \hat{x}_i \geq (J^{-1})_{ii}$

MLE:  $\hat{x} = \text{argmax}_x p(y|x)$ .

Properties: ( $x \in \mathcal{R}$ ) (y)  
cons. in probability

- consistency:  $\hat{x} \xrightarrow{P} x$  as  $n \rightarrow \infty$  ( $y \in \mathcal{R}^n$ )
- asymptotic efficiency:

$$\frac{(\text{Var } \hat{x}) J(x)}{(1 + \text{Bias}(\hat{x}, x))^2} \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

reaches Cramer-Rao  
D (with or without denom) as  $n \rightarrow \infty$

- asymptotic normality:

$$\sqrt{n} (\hat{x} - x) \xrightarrow{d} N(0, J^{-1}(x))$$

convergence in distribution

in the vector case as well

suggests  $\text{Bias}(\hat{x}, x) \sim 1/\sqrt{n}$   
although:

- asymptotic unbiasedness:  
 $\text{Bias}(\hat{x}, x) \sim 1/\sqrt{n}$

Rule:  $J_{ij} = -E \left[ \frac{\partial^2}{\partial x_i \partial x_j} \ln p(y|x) \right]$ .

Ex.  $y \in \mathbb{R}^m$ ,  $x \in \mathbb{R}^n$ ,  $e \sim N(0, \Sigma)$

$$y = Ax + e$$

$$y|x \sim N(Ax, \Sigma)$$

$$P(y|x) = C \exp\left[-\frac{1}{2}(y - Ax)^T \Sigma^{-1} (y - Ax)\right]$$

MLE:  $\max_x P(y|x) \Leftrightarrow \min_x -\ln P(y|x)$

$$\Leftrightarrow \min_x \frac{1}{2} (y - Ax)^T \Sigma^{-1} (y - Ax)$$

$$\|y - Ax\|_{\Sigma^{-1}}^2$$

$$\nabla_x = A^T \Sigma^{-1} (Ax - y) = 0$$

$$\Rightarrow A^T \Sigma^{-1} A x = A^T \Sigma^{-1} y$$

$$\hat{x} = (A^T \Sigma^{-1} A)^{-1} A^T \Sigma^{-1} y$$

(provided invertible (overdet).)

$\hat{x}|x$  is Normal with expectation

$$(A^T \Sigma^{-1} A)^{-1} A^T \Sigma^{-1} A x = x \rightarrow E \hat{x} = x, \text{ unbiased.}$$

Lemma:  $y \sim N(0, \Sigma)$   
 $\Rightarrow My \sim N(0, M \Sigma M^T)$   
 here  $M = (A^T \Sigma^{-1} A)^{-1} A^T \Sigma^{-1}$

$$\text{and Cov}(\hat{x}) = (A^T \Sigma^{-1} A)^{-1} A^T \Sigma^{-1} \Sigma \Sigma^{-1} A (A^T \Sigma^{-1} A)^{-1}$$

$$= (A^T \Sigma^{-1} A)^{-1}$$

suitably decreases as rows are added to A.

$$J(\alpha) = E \left[ \frac{\partial}{\partial x_i} \ln p(y|\alpha) \frac{\partial}{\partial x_j} \ln p(y|\alpha) \right]$$

$$\frac{\partial}{\partial x_i} \ln p(y|\alpha) = - \left[ A^T \Sigma^{-1} (Ax - y) \right]_i$$

$$J(\alpha) = E \left[ A^T \Sigma^{-1} (Ax - y) (Ax - y)^T \Sigma^{-1} A \right]$$

$$= A^T \Sigma^{-1} E(y y^T) \Sigma^{-1} A$$

$$= A^T \Sigma^{-1} \Sigma \Sigma^{-1} A$$

$$J^{-1}(\alpha) = \frac{A^T \Sigma^{-1} A}{A^T \Sigma^{-1} A} = \text{Cov}(\hat{x})$$

→ efficient (rare occurrence!).

Rule:  $A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \Rightarrow J(\alpha) = A_1^T \Sigma^{-1} A_1 + A_2^T \Sigma^{-1} A_2$   
additive, also known as  
"precision matrix" =  $\text{Cov}(\hat{x})^{-1}$

→ build up precision by adding data.

Rule:  $f(x) = \frac{1}{2} (y - Ax)^T \Sigma^{-1} (y - Ax)$

$\nabla^2 f(x) = A^T \Sigma^{-1} A$  Hessian

→ another interpretation of  $J(\alpha)$ .

ex.  $y = x + e$

$x \in \mathbb{R}$

$A = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

"ones(m, 1)"

$e \sim N(0, \sigma^2 I)$

$$A^T \Sigma^{-1} A = m \sigma^{-2}$$

$$\hat{x} = (A^T \Sigma^{-1} A)^{-1} A^T \Sigma^{-1} y = \frac{\sigma^2}{m} \sigma^{-2} [1 \dots 1] y = \frac{1}{m} \sum y_i$$

$$E \hat{x} = x$$

$$\text{Var} \hat{x} = (A^T \Sigma^{-1} A)^{-1} = \frac{\sigma^2}{m}$$