

Lecture 22

Wednesday, April 13, 2022 8:31 AM

Admin

- ① HW4 due Friday.
- ② Class cancelled 04/18
- ③ Projects due 05/04.

Parameter Covariance and Akaike's Information Criterion

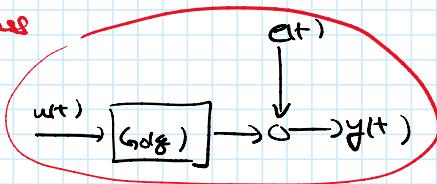
Variance of parameters.

Last time

$$\text{var}(\hat{\theta}(e^{iw})) = \frac{1}{N} \frac{\mathbb{E}_{\theta}(w)}{\mathbb{E}_{\theta}(w)^2}$$

parameters
data points

evaluating model goodness



- True for a linear model + Gaussian noise.
- Central Limit theorem: Average tends to Gaussian for general distribution.

Asymptotic Covariance of the parameter estimate

- True system model: $y(t) = \varphi^T(t) \theta_0 + e(t)$ ← Gaussian noise.
{ parameters of $G(\varphi)$ }

$$\mathbb{E}[e(t)e(t)\mathbf{T}] = \begin{cases} 1, \tau=0 \\ 0, \tau \neq 0 \end{cases} \quad \hat{y}(t, \theta) - y(t)$$

Then $\hat{\theta}_N$ minimizes $V_N(\theta, \hat{y}) = \frac{1}{N} \sum_{t=1}^N \hat{y}(t, \theta)^2$

↑
input-output
data

- Can show:

$$\text{cov}(\hat{\theta}_N) = \underbrace{\frac{1}{N} \left[\mathbb{E}[\varphi(t) \varphi^T(t)] \right]^{-1}}_{\text{Auto-covariance}} \quad \left. \begin{array}{l} \text{for linear models} \\ \frac{1}{N} \sum_{t=1}^N \varphi(t) \varphi^T(t) \end{array} \right\}$$

Asymptotic variance

- The asymptotic variance is
 - Inversely proportional to # of samples
 - proportional to variance in noise
 - Inversely related to parameter sensitivity.
→ For a linear model inversely related to input auto-covariance
- For a general non-linear model,

$$\text{cov}(\hat{\theta}_N) \sim \frac{1}{N} \left[\mathbb{E}[\psi(t, \theta_0) \psi^T(t, \theta_0)] \right]^{-1} \leftarrow$$

$$\text{cov}(\hat{\theta}_N) \sim \frac{1}{N} \left[E[\gamma(t, \theta_0) \gamma^T(t, \theta_0)] \right]^{-1} \leftarrow$$

$$\begin{aligned} \gamma(t, \theta_0) &= \frac{\partial}{\partial \theta} [\varepsilon(t, \theta)] \Big|_{\theta=\theta_0} = \frac{1}{\partial \theta} [\hat{y}(t, \theta) - y(t)] \Big|_{\theta=\theta_0} \\ &= \frac{\partial}{\partial \theta} [\hat{y}(t, \theta)] \Big|_{\theta=\theta_0} \end{aligned}$$

- For a linear model,

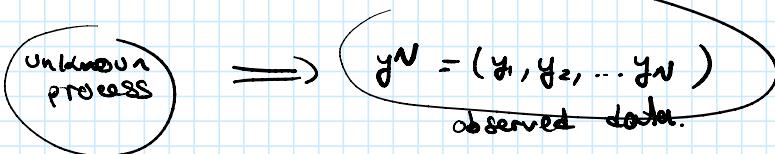
$$\begin{aligned} \hat{y}(t, \theta) &= \varphi^T(t) \theta \\ \Rightarrow \gamma(t, \theta_0) &= \frac{\partial \hat{y}}{\partial \theta} = \varphi(t) \end{aligned}$$

- For nonlinear model, θ_0 isn't known. \rightarrow asymptotic variance can't be determined.
- In practice, use empirical estimates of $\hat{\theta}_N, \lambda_N$ for large N .

maximum likelihood Estimation (MLE)

- Methods of parameter estimation.
 - Least squares.
 - \rightarrow Find a set of parameters to minimize output error.
 - \rightarrow Requires minimal info. about distribution of data.
 - Artificial Neural Nets.
 - \rightarrow useful for nonlinear modeling.
 - \rightarrow little understanding of how each parameter relates to output.
 - Maximum likelihood Estimates.
 - \rightarrow Relies on probability density functions (PDFs); distribution info.
 - \rightarrow we can estimate parameters. \Rightarrow Akaike's Information Criterion (AIC).

Principle of maximum likelihood



- Assume y_i is generated by a stochastic process having a PDF

$$f(\mu, \lambda, x) = \frac{1}{\sqrt{2\pi\lambda}} e^{-(x-\mu)^2/2\lambda} \quad \leftarrow x \text{ is random variable associated w/ } y_i.$$

μ = mean

associate w/ y_i .

$\mu = \text{mean}$

$\lambda = \text{variance} = \sigma^2$

$$\mu = \frac{1}{N} \sum_{i=1}^N y_i \quad ; \quad \lambda = \frac{1}{N} \sum_{i=1}^N (y_i - \mu)^2 \leftarrow \text{WANT THESE USING MLE.}$$

- look at distribution of data from the MLE perspective.
- Assume N observations $y^N = y_1, \dots, y_N$ are independent.
→ write joint probability

$$f(\underbrace{\mu, \lambda}_{\substack{| \\ \uparrow}} | x_1, x_2, \dots, x_N) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\lambda}} e^{-\frac{(x_i - \mu)^2}{2\lambda}} \leftarrow$$

- Given data that want to find μ, λ that most likely generated data.

$$\max f(\underbrace{\mu, \lambda | y_1, y_2, \dots, y_N}_{\substack{| \\ \text{known}}})$$

$$\hat{\theta} = \underset{\Theta}{\operatorname{argmax}} f(\mu, \lambda | y^N)$$

$$\begin{bmatrix} \mu \\ \lambda \end{bmatrix} = \underset{\Theta}{\operatorname{argmax}} \log f(\mu, \lambda | y^N)$$
$$= \underset{\Theta}{\operatorname{argmin}} \sum_{i=1}^N \left(\log \frac{1}{\sqrt{2\pi\lambda}} - \frac{(y_i - \mu)^2}{2\lambda} \right)$$

$$\Rightarrow \frac{\partial \log f}{\partial \mu} = \sum_{i=1}^N \frac{1}{\lambda} (y_i - \mu) = 0 = \mu - \frac{1}{N} \sum_{i=1}^N y_i$$

$$\frac{\partial \log f}{\partial \lambda} = -\frac{N}{2} \cancel{\frac{d \log \lambda}{d \lambda}} - \frac{1}{2} \lambda' \sum_{i=1}^N \frac{(y_i - \mu)^2}{2} \rightarrow \lambda = \frac{1}{N} \sum_{i=1}^N (y_i - \mu)^2$$

- MLE gives a probability distribution approach to thinking about optimal parameter estimation.

Information Theory in Parameter Estimation

- maximum likelihood estimate (MLE) :

$$\hat{\theta}^{ML}(z^N) = \underset{\Theta}{\operatorname{argmax}} \log L(\Theta)$$

- likelihood $L(\Theta) = \prod_{i=1}^N f(\theta; y_1, \dots, y_N)$

$$= \prod_{t=1}^N f_{\epsilon}(\hat{y}(t) - \hat{y}(t|\theta), t, \theta)$$

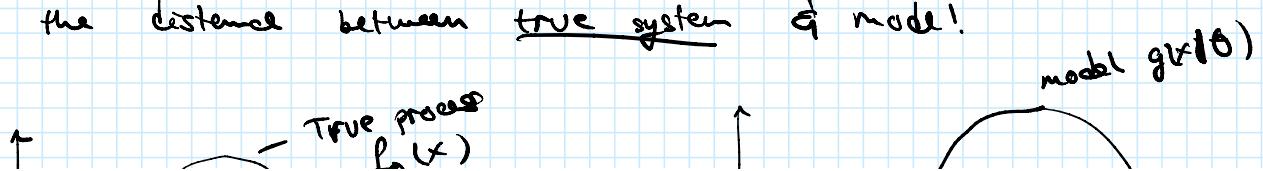
$$\hat{\theta}^{MI}(z^N) = \underset{\Theta}{\operatorname{argmax}} \underbrace{\frac{1}{N} \sum_{t=1}^N \log f_{\epsilon}(\epsilon, t; \theta)}_{\text{related entropy}}$$

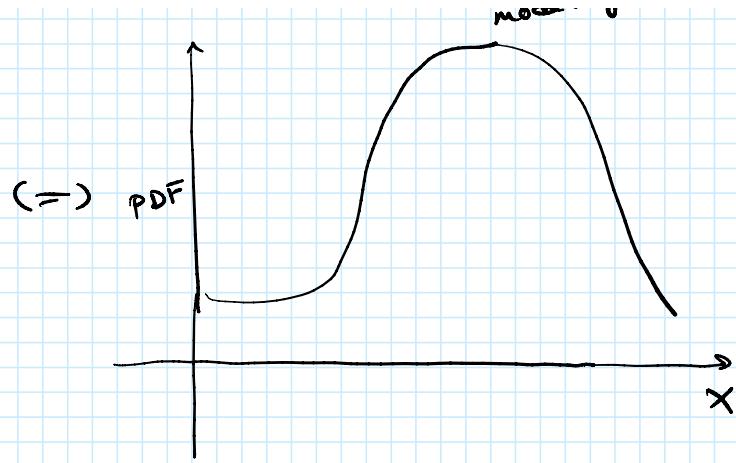
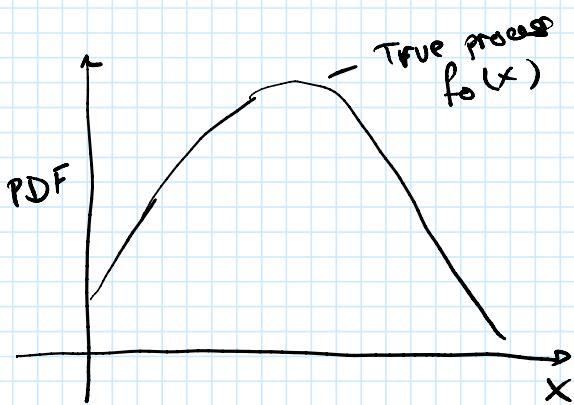
Information Theoretic Approach.

- Model - data agreement is quantified in terms of the amount of information captured by the model.
- A "good" model fully exploits all of the information in the data.
 - Least squares estimate
 - prediction error ϵ input data are uncorrelated with each other!
 - Maximum likelihood estimate.
 - The log joint probability is maximized
 - The degree of randomness in prediction error is maximized.
 - The entropy is maximized.
 - We need a unified measure to enable us to compare model structures on the same basis, providing an objective to select an optimal model.
 - => based on a trade-off between bias in estimates and variance in output error.
 - > Accuracy vs. reliability \leftarrow Liang.

Kullback - Leibler Information Distance

- measure the distance between true system & model





→ If $f_0(x) = g(x|\theta)$ for all x , then distance must be zero.

→ To evaluate the distance between true process & model,
consider

$$I(f_0, g) = \int_x f_0(x) \log \underbrace{\frac{f_0(x)}{g(x|\theta)}}_{\equiv (\log f_0(x)) - \log g(x|\theta)} dx$$

\leftarrow distance between funcs.

Goal is to minimize,