

System Identification

6.435

SET 12

- Identification in Practice
- Error Filtering
- Order Estimation
- Model Structure Validation
- Examples

Munther A. Dahleh

“Practical” Identification

- Given: $Z^N = \{y(t), u(t); t \leq N\}$
- Want
 - 1) a model for the plant
 - 2) a model for the noise
 - 3) an estimate of the accuracy
- choice of the model structure

flexibility

parsimony

- What do we know?

We know methods for identifying “models” inside a “priori” given model structures.

- How can we use this knowledge to provide a model for the plant, the process noise, with reasonable accuracy.

Considerations

- Pre-treatment of data
 - Remove the bias (may not be due to inputs)
 - Filter the high frequency noise
 - Outliers
- Introduce filtered errors. Emphasize certain frequency range. (The filter depends on the criterion).

- Pick a model structure (or model structures)
 - Which one is better?
 - How can you decide which one reflects the real system?
 - Is there any advantage from picking a model with a large number of parameters, if the input is “exciting” only a smaller number of frequency points?
- What are the important quantities that can be computed directly from the data (inputs & outputs), that are important to identification?

Pre-treatment of Data

- Removing the bias

$$Ay(t) = Bu(t) + v(t)$$

- If $Ev(t) = 0$, then the relation between the static input $u(t) = \bar{u}$ and output $y(t) = \bar{y}$ is given by

$$A(1)\bar{y} = B(1)\bar{u}$$

- The static component of $y(t) = \bar{y}$ may not be entirely due to \bar{u} , i.e. the noise might be biased.

- Method I: Subtract the means:

Define $\bar{y} = \frac{1}{N} \sum_{t=1}^N y^m(t)$ $[y^m = \text{meas. data}]$

$$\bar{u} = \frac{1}{N} \sum_{t=1}^N u^m(t)$$

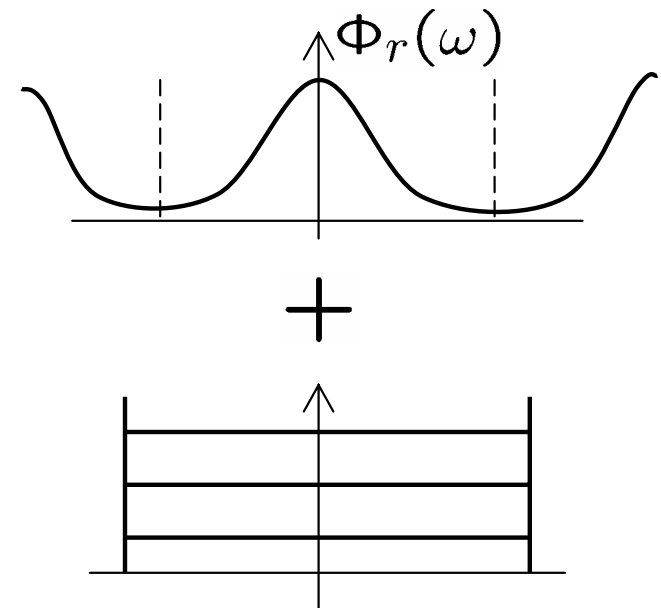
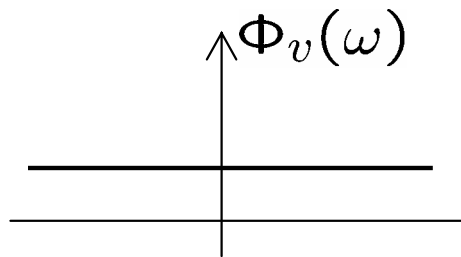
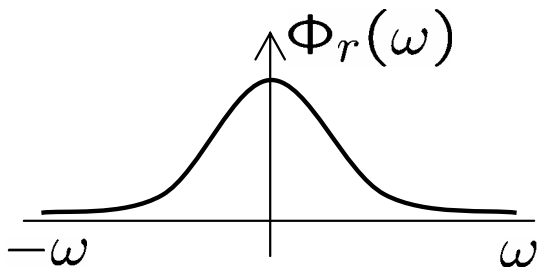
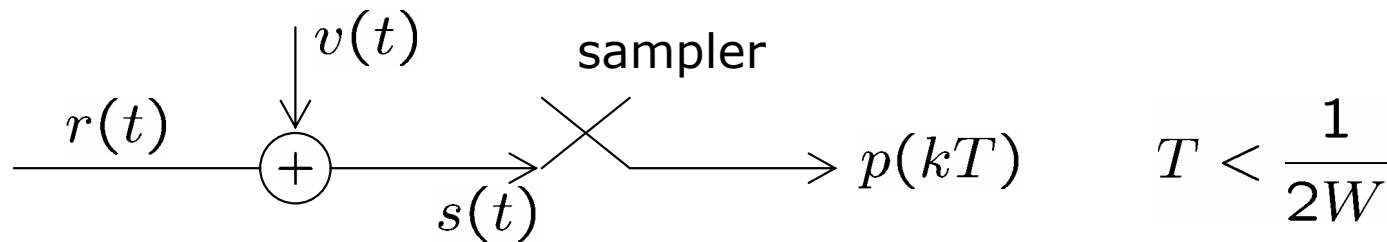
New Data: $y(t) = y^m(t) - \bar{y}$

$$u(t) = u^m(t) - \bar{u}$$

- Method II: Model the offset by an unknown constant α , and estimate it.

$$m : Ay(t) = Bu(t) + \alpha + v(t)$$

- High Frequency disturbances in the data record.
 - “High” means above the frequency of interest.
 - Related to the choice of the sampling period.



- Without an anti-aliasing filter, high frequency noise is folded to low frequency.



- high frequency noise depends on:
 - a) high frequency noise due to $v(t)$
 - b) aliasing.
- Problem occurs at both the inputs and outputs.

$$y_F = Ly \quad \mathbf{L} - \text{LTI LP filter.}$$

$$u_F = Lu$$

- m : $A(t)y_F = B(t)u_F + v(t) \quad [v = He]$

equivalently $A(t)y = B(t)u + \frac{1}{L}v(t)$

i.e. multiply the noise filter by $\frac{1}{L}$

- Outliers, Bursts

- Either erroneous or high-disturbed data point.
- Could have a very bad effect on the estimate.
- Solution:
 - a) Good choice of a criterion (Robust to changes)
 - b) Observe the residual spectrum. Sometimes it is possible to determine bad data.
 - c) Remove by hand!!! Messing up with real data.
 - d) Failure-detection using hypothesis testing or statistical methods. (Need to define a threshold).

Role of Filters: Affecting the Biase Distribution

- $m_1 : Gu + H_1 e$

$$y_F = Ly$$

$$u_F = Lu$$

$$\Rightarrow m : y(t) = Gu + He \quad H = \frac{1}{L}H_1$$

- Frequency domain interpretation of parameter estimation:

$$\theta^* = \operatorname{argmin} \int_{-\pi}^{\pi} \frac{\Phi_{ER}(\omega, \theta)}{|H(e^{i\omega}, \theta)|^2} d\omega$$

$$\Phi_{ER} = |G_o - G(e^{i\omega}, \theta)|^2 \Phi_u(\omega) + \Phi_v(\omega)$$

If $\theta = \begin{bmatrix} \rho \\ \eta \end{bmatrix}$: independently parametrized model structure

$$\rho^* = \underset{\rho}{\operatorname{argmin}} \int_{-\pi}^{\pi} |G_o(e^{i\omega}) - G(e^{i\omega}, \theta)|^2 Q(\omega, \eta^*) d\omega$$

$$Q(\omega, \eta^*) = \frac{\Phi_u}{|H(e^{i\omega}, \eta^*)|^2}$$

$$\eta^* = \underset{\eta}{\operatorname{argmin}} \int_{-\pi}^{\pi} \left| \frac{1}{N(e^{i\omega}, \rho^*)} - \frac{1}{H(e^{i\omega}, \eta)} \right|^2 d\omega$$

$$Q_{ER}(\omega, \rho^*) = \lambda^* N N^*$$

- Heuristically, θ is chosen as a compromise between minimizing the integral of $|G_o - G(e^{i\omega}, \theta)|^2 Q(\omega, \theta^*)$ and matching the error spectrum Φ_{ER} .
- $Q(\omega, \theta^*) = \frac{\Phi_u}{|H(e^{i\omega}, \theta^*)|^2}$ weighting function
 - Input spectrum
 - Noise spectrum
- With a pre-filter: $Q(\omega, \theta^*) \rightarrow |L(e^{i\omega})|^2 Q(\omega, \theta)$
- Can view the pre-filters as weighting functions to emphasize certain frequency ranges. This interpretation may not coincide with “getting rid of high frequency components of the data”.
- Depending on the criterion, the choice of L can be different.

OE Model Structures

- $m : Gu(t) + e(t) \quad H = 1$

$$\rho^* = \operatorname{argmin} \int_{-\pi}^{\pi} |G_o(e^{i\omega}) - G(e^{i\omega}, \theta)|^2 \Phi_u(\omega) d\omega$$

- If G_o rolls off, then as long as $G(e^{i\omega}, \theta)$ is small around $\omega \simeq \pi$ the contribution of the criterion $|G_o(e^{i\omega}) - G(e^{i\omega}, \theta)|$ will be very small.
- If $\Phi_u = 1$, we expect $G(e^{i\omega}, \theta)$ to match G_o much better at low-frequency.

- Example (book)

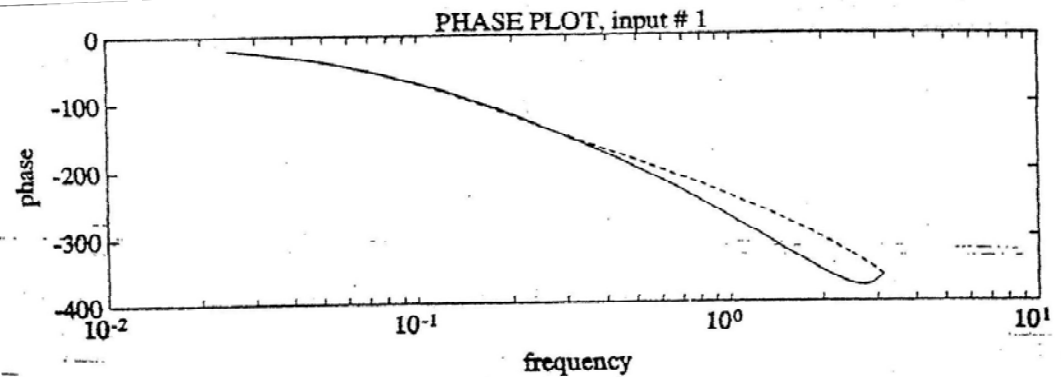
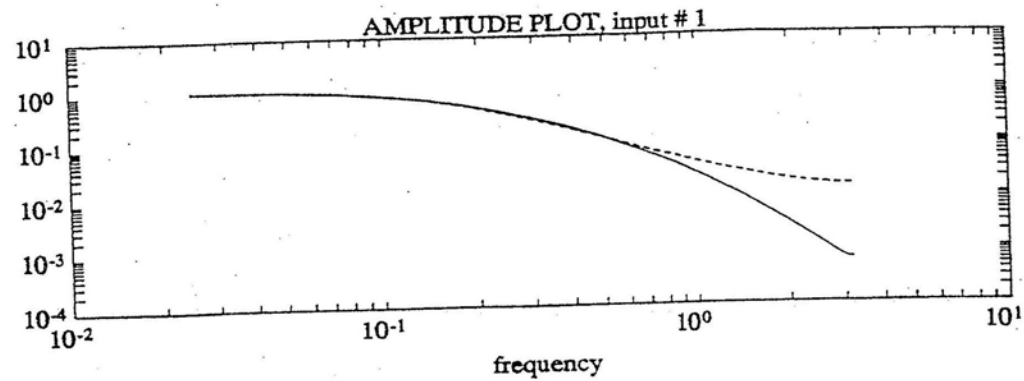
$$\delta : G_o(q) = \frac{0.001q^{-2}(10 + 7.4q^{-1} + 0.924q^{-2} + 0.1764q^{-3})}{1 - 2.14q^{-1} + 1.553q^{-2} - 0.4387q^{-3} + 0.042q^{-4}}$$

No noise.

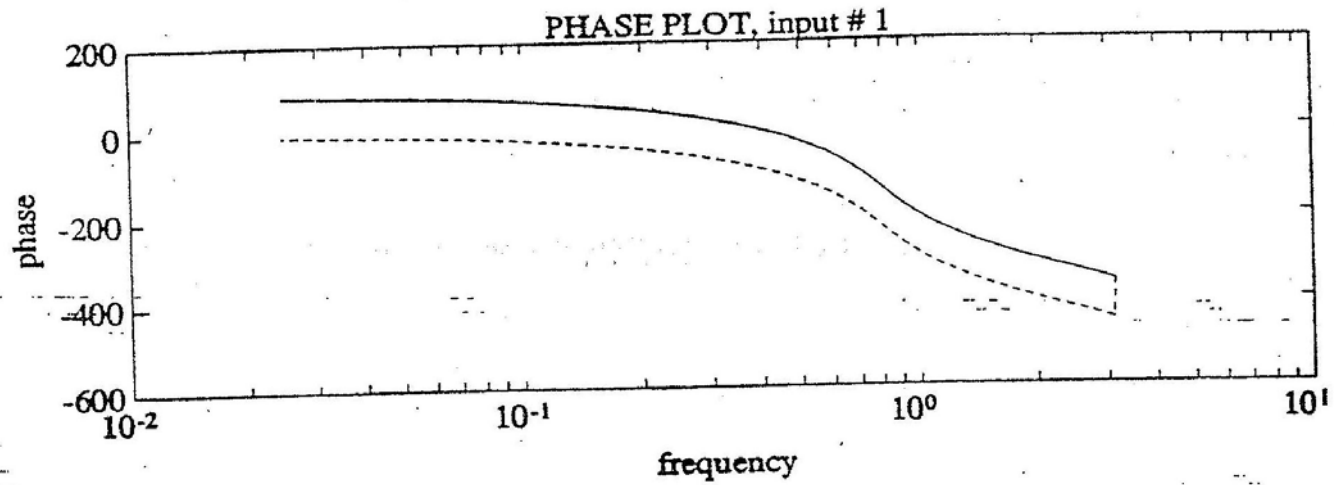
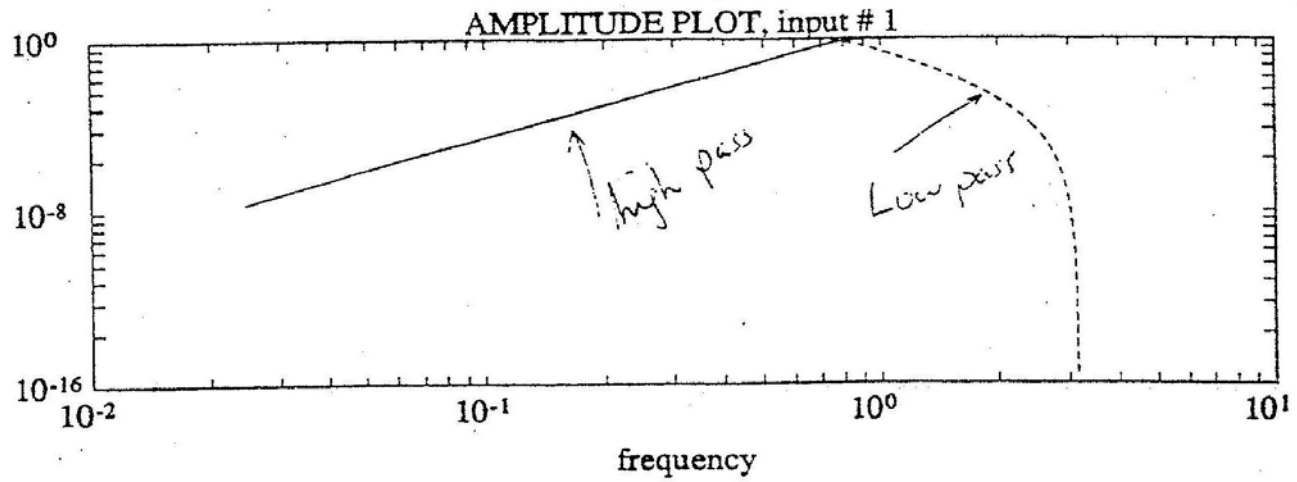
$$\Phi_u = 1 \quad \text{PSRB}$$

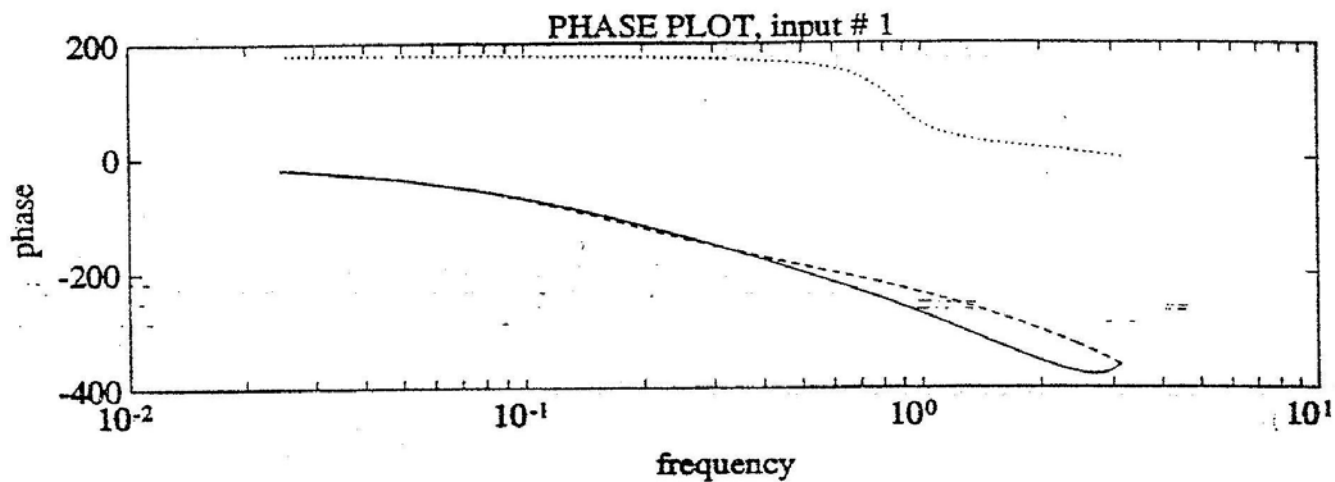
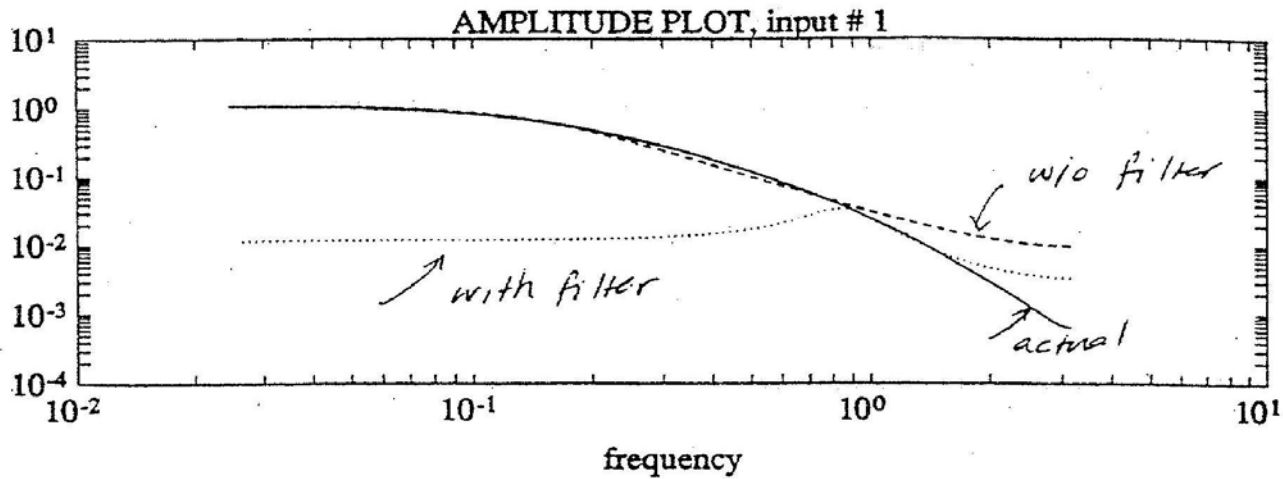
OE: $\hat{y} = G(e^{i\omega}, \theta) u$

$$G = \frac{b_1q^{-1} + b_2q^{-2}}{1 + f_1q^{-1} + f_2q^{-2}}$$



- Good match at low frequency.
Not as good at high frequency.
- Introduce a high-pass filter (5th order Butterworth filter, cut-off freq = 0.5 rad/sec.





ARX Model Structure

- $G = \frac{B}{A}$ $H = \frac{1}{A}$ not independently-parametrized.

$$\theta^* = \underset{\sim}{\operatorname{argmin}} \int_{-\pi}^{\pi} |G_o(e^{i\omega}) - G(e^{i\omega}, \theta)|^2 |A(e^{i\omega}, \theta)|^2 \Phi_u + \cancel{\Phi_v} \overset{=0}{|A|^2}$$

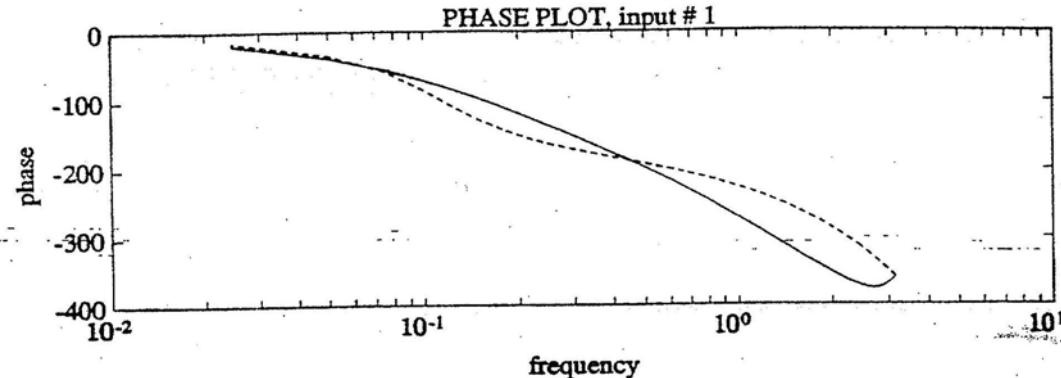
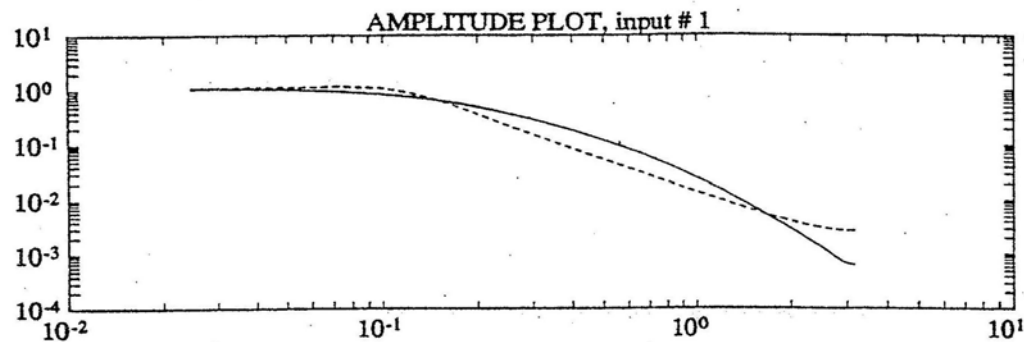
- If G_o rolls off, $B_o/A_o = G_o$ and A_o is large at high frequency.
If $A(e^{i\omega}, \theta)$ looks like A_o , then it will emphasize the high frequency part of the criterion.

- Conclusions are not as transparent in the noisy case.
However, it is in general true for large (SNR).

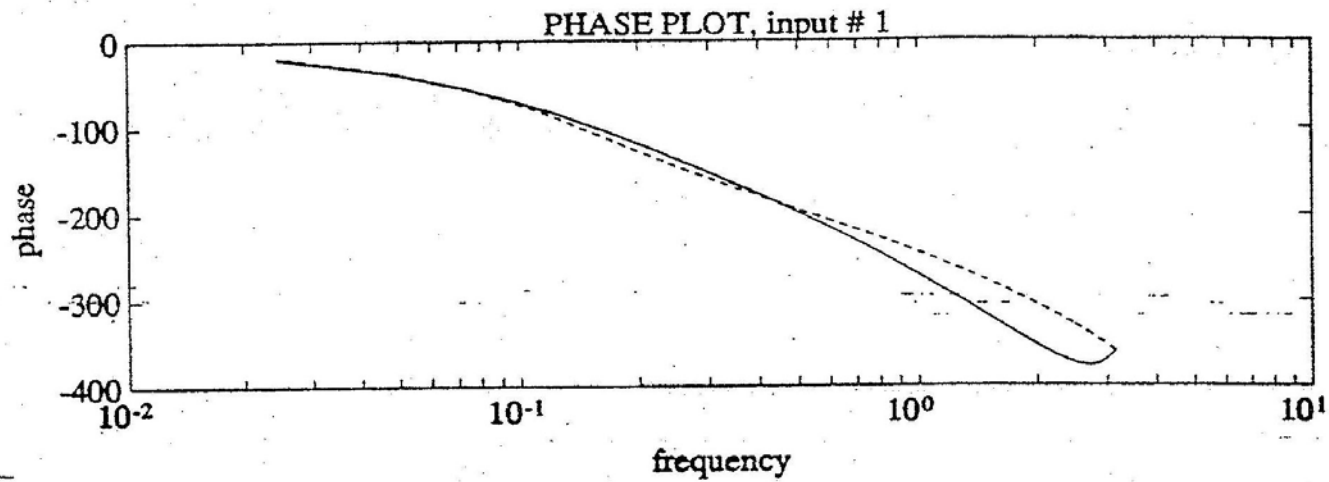
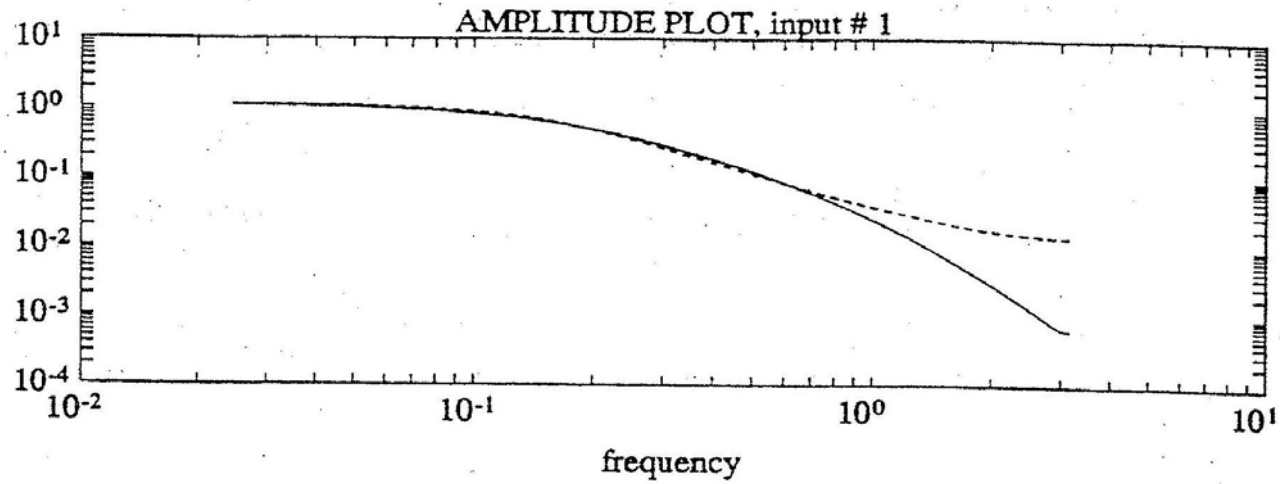
- Same example

$$m : \quad Ay = Bu + e$$

and the frequency response of $G(e^{i\omega}, \theta^*)$ is:



- Not a very good match at low frequency.
- Better than OE at high frequency.
- Can change this through a pre-filter. (5th order Butterworth, lowpass with cut-off frequency = 0.5)



- Better low frequency fit.
- Another interpretation

$$y = \frac{B}{A}u + \left(\frac{1}{A}e\right) \quad \left|\frac{1}{A}\right| \text{ small at high frequency.}$$

low frequency

Filters:

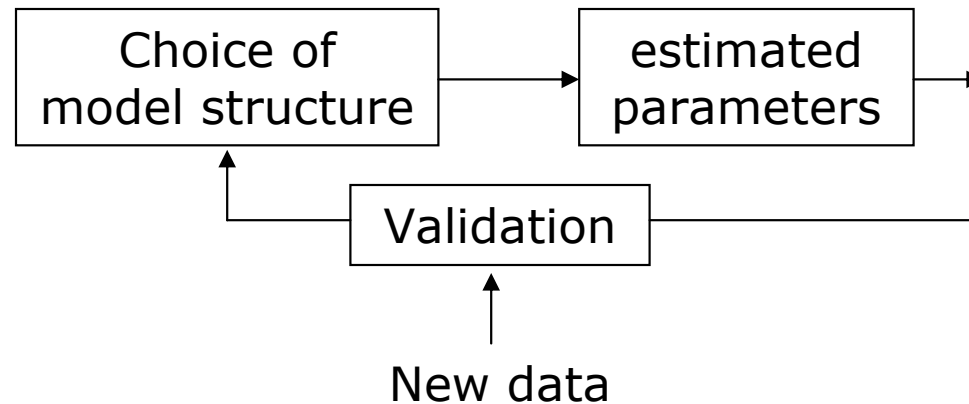
$$y = \frac{B}{A}u + \left(\frac{L}{LA}e\right) \leftarrow \text{high frequency if } L - \text{low pass.}$$

Conclusions

- Pre filters can be viewed as “design” parameters as well as the “standard” interpretation for noise reduction.
- Pre-treatment of the “data” is quite valuable, however should be done with “care”.
- Sampling can be quite tricky. Need to estimate the bandwidth of the system.

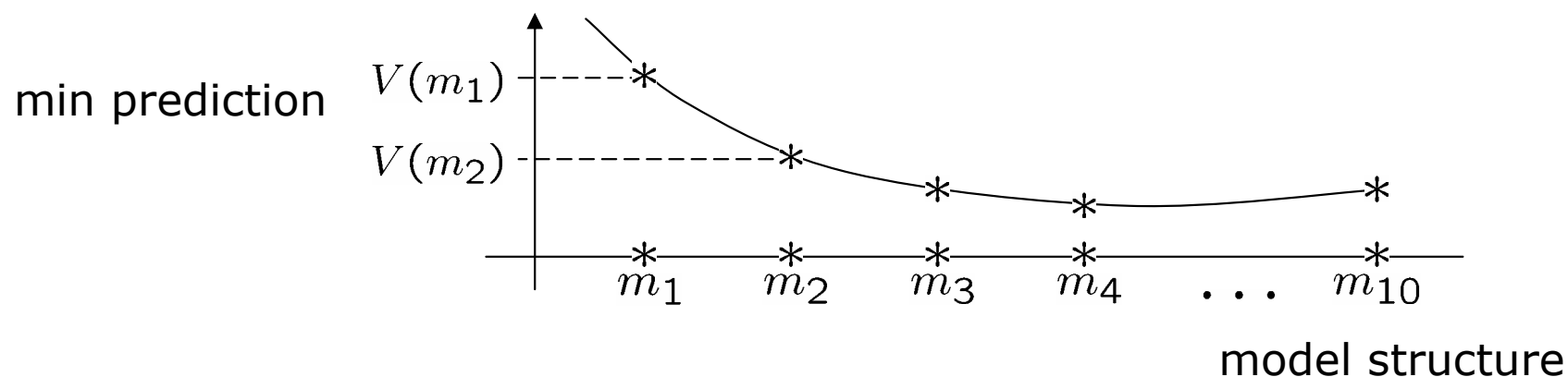
Model Structure Determination

- Flexible vs Parsimony
- Lots of trial and error. Usually more than one “experiment” is available.



- Theme: Fit the data with the least “complex” model structure. Avoid over-fitting which amounts to fitting noise.
- Better to compare similar model structures, although it is necessary to compare different structures at the end.

- Available data dictates the possible model structures and their dimensions.
- Can noise help identify the parameters?
- How “bad” is the effect of noise? Consistency in general is guaranteed for any SNR. What does that mean?
- Is there a rigorous way of comparing different model structures?



- Akaike’s Information Theoretic Criterion (AIC).
- Akaike’s final prediction error criterion (FPE).

Order Estimation

- Qualitative
 - Spectral estimate
 - Step response if available. Otherwise, step response of spectral estimate.
- Quantitative
 - Covariance matrices
 - Information matrix $E \left(\Psi(t, \theta) \Psi^T(t, \theta) \right)$
 - Residual-input correlation
 - Residual whiteness
 - \vdots
- All methods are limited by the input used.

Covariance Matrix

$$m : \overset{\curvearrowright n}{A}y = \overset{\curvearrowright n}{B}u + v \quad \phi_s(t) = (-y(t-1), \dots, -y(t-s), u(t-1), \dots, u(t-s))$$

- Two basic results

- $v = 0 \quad \bar{E}\phi_n(t)\phi_n^T(t) > 0 \Leftrightarrow u$ is p.e of order $2n$ and A, B are coprime

- $v \neq 0 \quad \bar{E}\phi_n(t)\phi_n^T(t) > 0 \Leftrightarrow u$ is p.e of order n
white or persistent

- To determine n , obtain estimates of $\bar{E}\phi_s(t)\phi_s^T(t)$

Case 1: u is WN , v is WN

Increase s until $\bar{E}\phi_s\phi_s^T(t)$ is "singular".

$$\bar{E}\phi_s(t)\phi_s^T(t) \simeq \frac{1}{N} \sum_{t=1}^N \phi_s(t)\phi_s^T(t) = \hat{R}_s$$

Use "SVD", robust rank tests.

$\det(\hat{R}_s) \stackrel{?}{=} 0$, observe a sudden drop in the rank.

Case 2: u is p.e of order n_u , Noise is white.

If $s > n_u \Rightarrow \bar{E}\phi_s(t)\phi_s^T(t)$ is singular.

you really cannot estimate the order of the system if it is larger than n_u .

Case 3: u is p.e of order n_u , Noise free case.

$n > \frac{n_u}{2}$ cannot be determined. Not a likely hypothesized model structure.

$$\hat{R}_s = \frac{1}{N} \sum_{t=1}^N \phi_s(t) \phi_s^T(t)$$

is a bad estimate of $\bar{E} \phi_s \phi_s^T(t)$. Of course N is fixed (data length).

“Enhanced criterion”

$$\hat{\hat{R}}_s = \hat{R}_s - \underbrace{\hat{\sigma} R_v}_{\text{estimated noise contribution.}}$$

- If noise level is high, use an instrumental variable

$$\xi_s(t) = [u(t-1), \dots, u(t-2s)].$$

test: rank $\bar{E}\xi_s(t)\phi_s^T(t)$

- If u is p.e , then generically

$$\bar{E}\xi_s\phi_s^T > 0 \quad \Leftrightarrow \quad s \leq n \quad , \quad n \text{ order of } A, B$$

- Other tests:

Estimates of

$$\bar{E}u(t)\varepsilon(t, \theta) \simeq 0$$

$$\bar{E}\varepsilon(t, \theta)\varepsilon(t-k, \theta) \quad k \neq 0 \quad \simeq 0$$

$$\bar{E}\Psi(t, \theta)\Psi^T(t, \theta) \quad \text{non singular}$$

Examples

- δ unknown. Study possible conclusions for different experiments and different SNR.

- Model structure

$$m : \quad Ay = Bu + v$$

- Inputs

- WN

- $\cos \omega_0 \quad \omega_0 = 2\pi/1000$

- $\cos(2\pi/1000) \quad , \quad \cos 20 \cdot (2\pi/1000)$

Can determine from the spectrum or u (or simply FFT).

- SNR:

$$\lambda = 1$$

$$\lambda = 0.1$$

$$\lambda = 0.01$$

- All examples, you can access both the inputs and outputs.

First Experiment

$$u = WN$$

$$\lambda = 1 \quad SNR \simeq 1$$

Test for model order:

$$\bar{E}\phi_s\phi_s^T(t) \cong \frac{1}{N} \sum_{t=1}^N \phi_s(t)\phi_s^T(t) = \hat{R}_s$$

From data \hat{R}_s is singular for $s \geq 2$

$$\det \hat{R}_s \simeq 0 \quad (\text{noticed a sudden drop})$$

→ Estimated system

$$a_1 = -1.44 \quad a_2 = 0.498$$

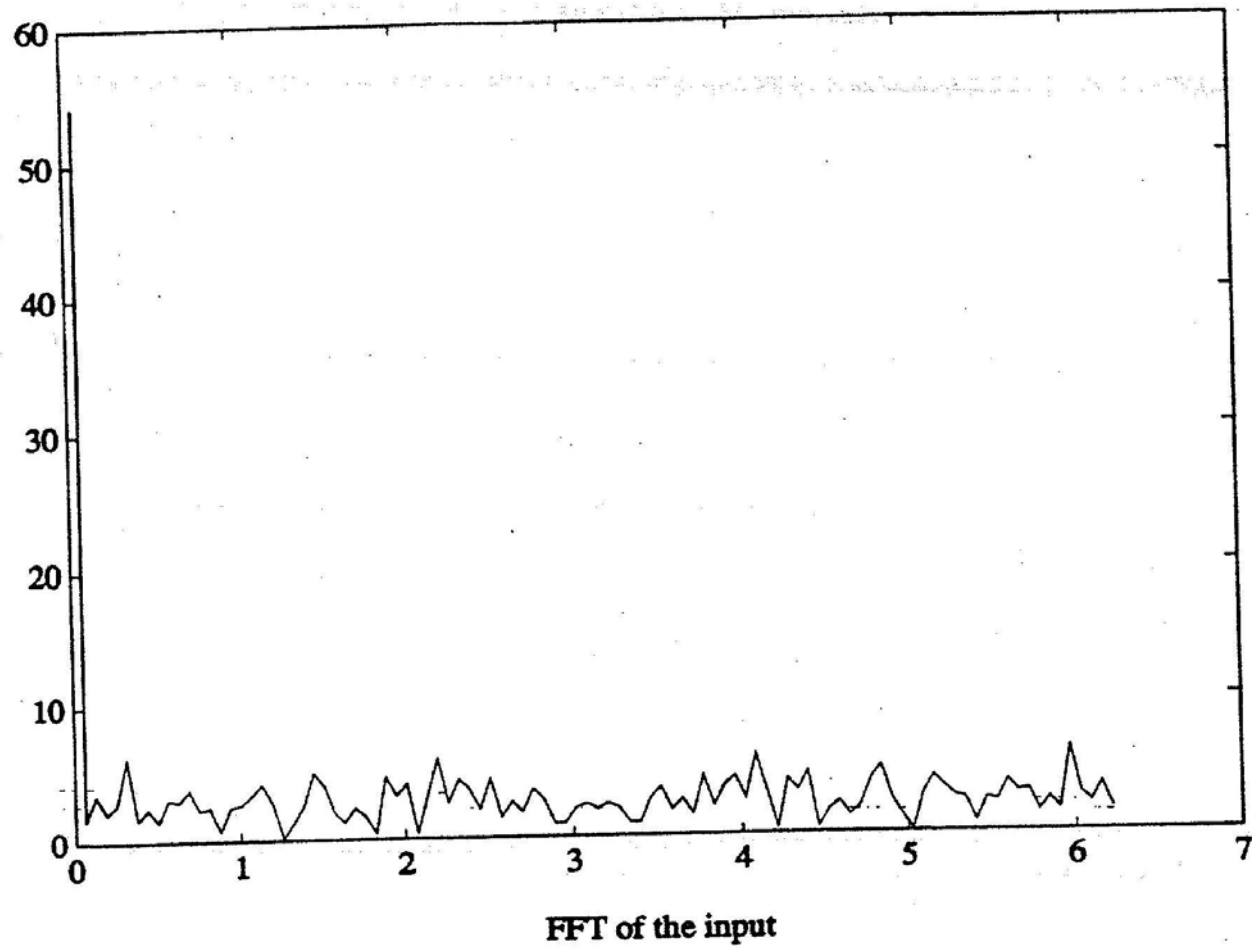
$$b_1 = 1.0221 \quad b_2 = 0.5239$$

→ COV $\hat{\theta}_N$ small

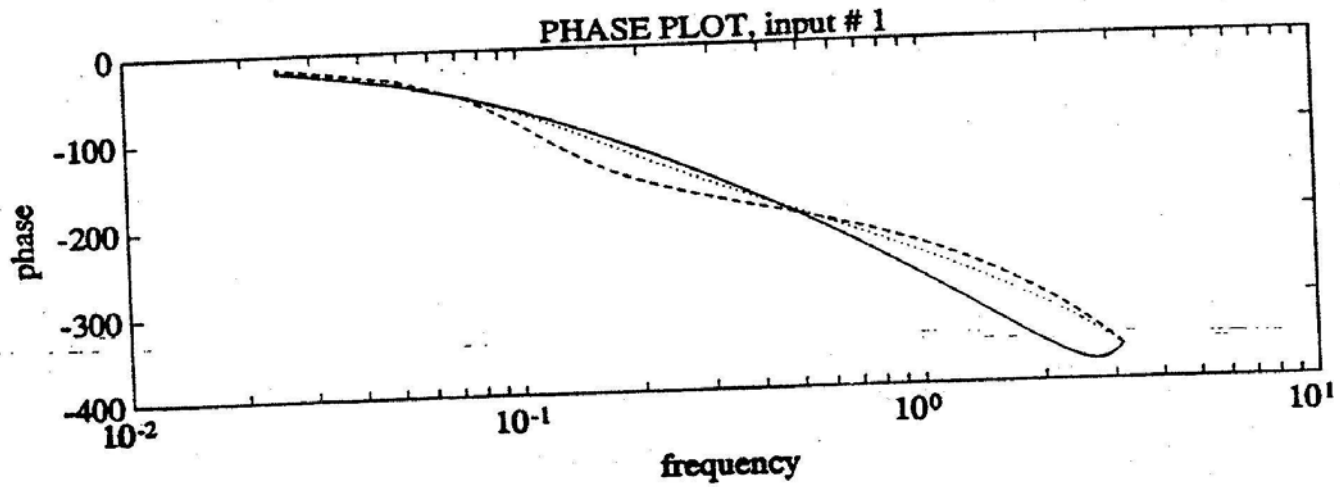
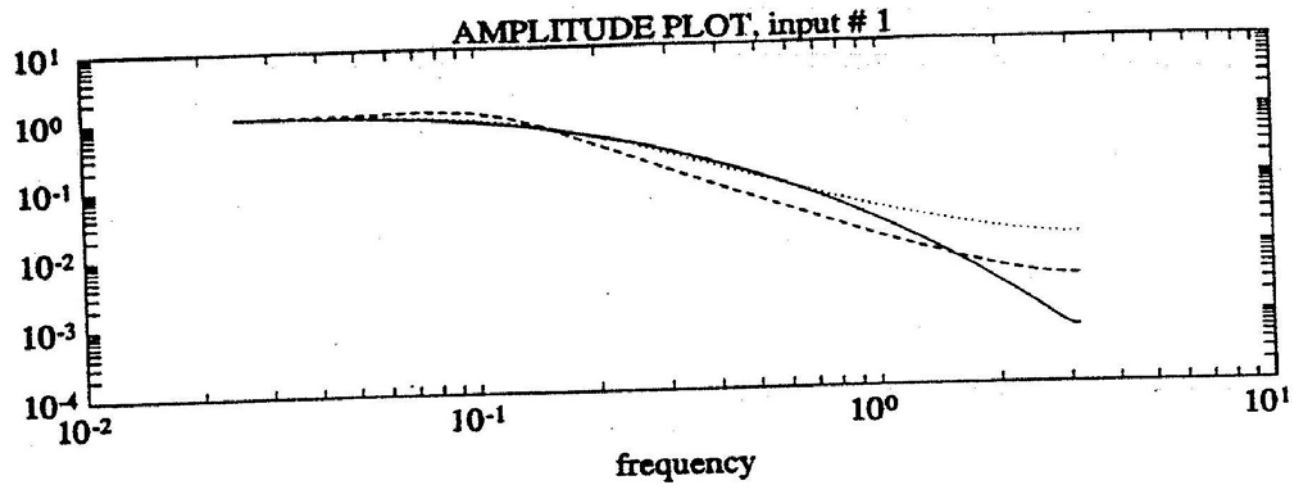
(ARX)

Comment: If m is the correct structure γ , the above are “good” estimates of the parameters.

Plot shows different SNR.



random



ARX (ARX)

Second Experiment

$$u = \cos w_0 t$$

u is p.e of order 2.

SNR	a_1	a_2	b_1	b_2	cov $\hat{\theta}_N$	det \hat{R}
1	-0.0014	0.0005	-0.0520	0.054	high	2.4×10^{-6}
0.1	-1.38	0.479	1.255	0.356	high	3.03×10^{-7}
0.01	-1.398	0.4885	0.9631	0.5437	low	2.67×10^{-7}

Theoretical Analysis:

Data is informative $\bar{E}\phi_2\phi_2^T > 0$ (although det is small)

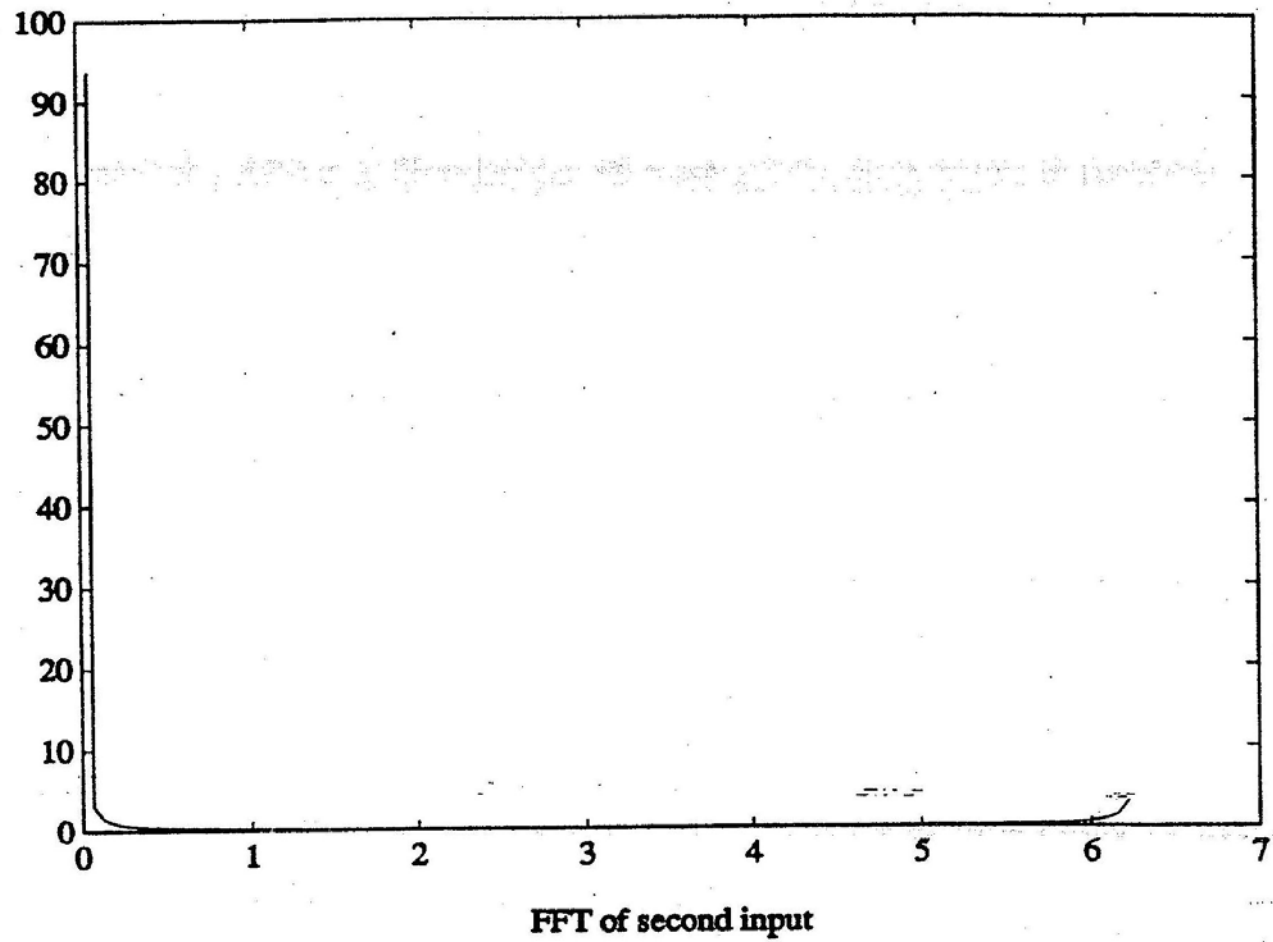
For $\cos w_0 t$, as $N \rightarrow \infty$

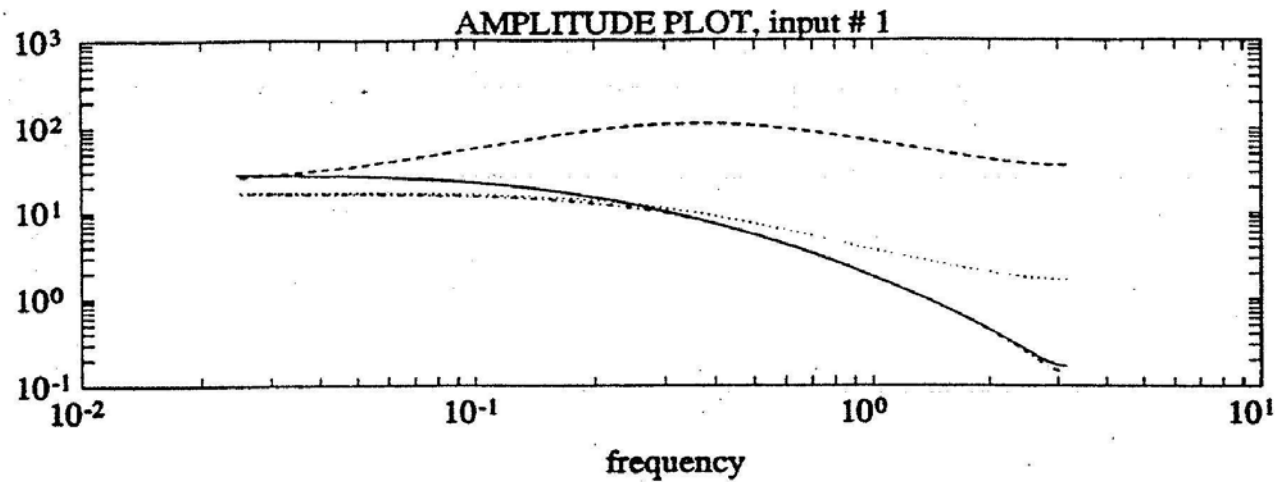
$$\hat{\theta}_N \rightarrow \hat{\theta}_o^* \text{ regardless of } \lambda$$

$N = 100$, the estimates were quite bad for $\lambda = 1$ in comparison to WN inputs.

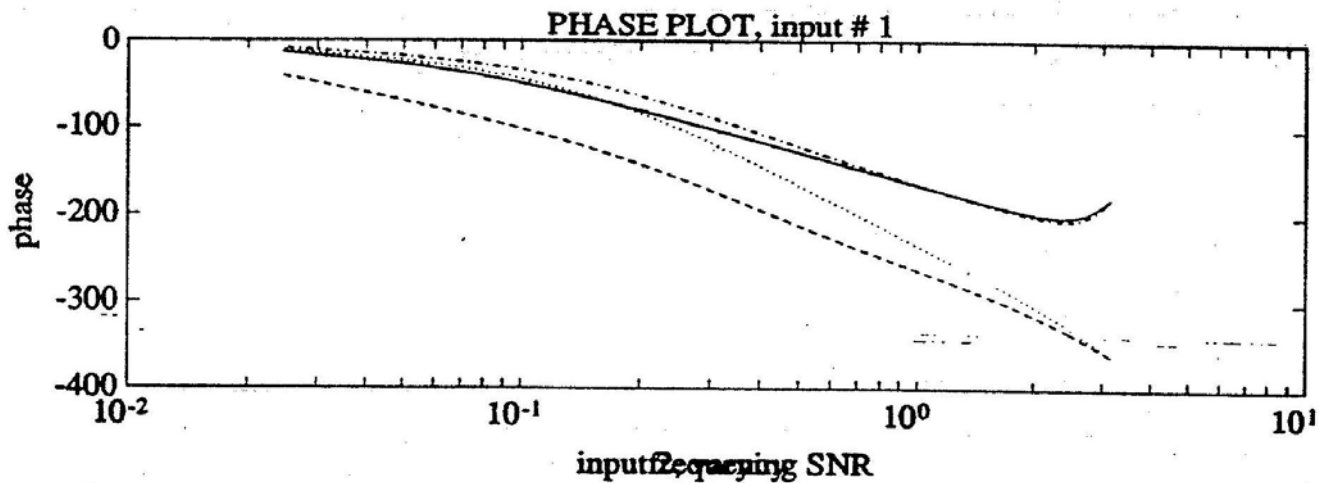
Even though noise "helps" in obtaining asymptotic convergence (through providing excitation), it is not helpful for finite-data records, since its effect cannot be averaged. Accuracy depends on

$$\lambda \left(\bar{E}\Psi\Psi^T \right)^{-1}$$





— $u = \text{rand}$ $\lambda = 1$
 - - - $u = \text{COS } w_{ot}$ $\lambda = 1$
 - . - $\lambda = 0$
 $\lambda = 0$



ARX (with ARX plant)
 2nd input

Third Experiment

$$u = \cos w_0 t + \cos 20w_0 t$$

u is p.e of order 4.

SNR	a_1	a_2	b_1	b_2	cov $\hat{\theta}_N$	det \hat{R}
1	-1.475	0.5317	2.38	-1.134	high	0.25
0.1	-1.39	0.478	1.186	0.333	smaller	0.07
0.01	-1.395	0.4862	1.022	0.487	low	0.069

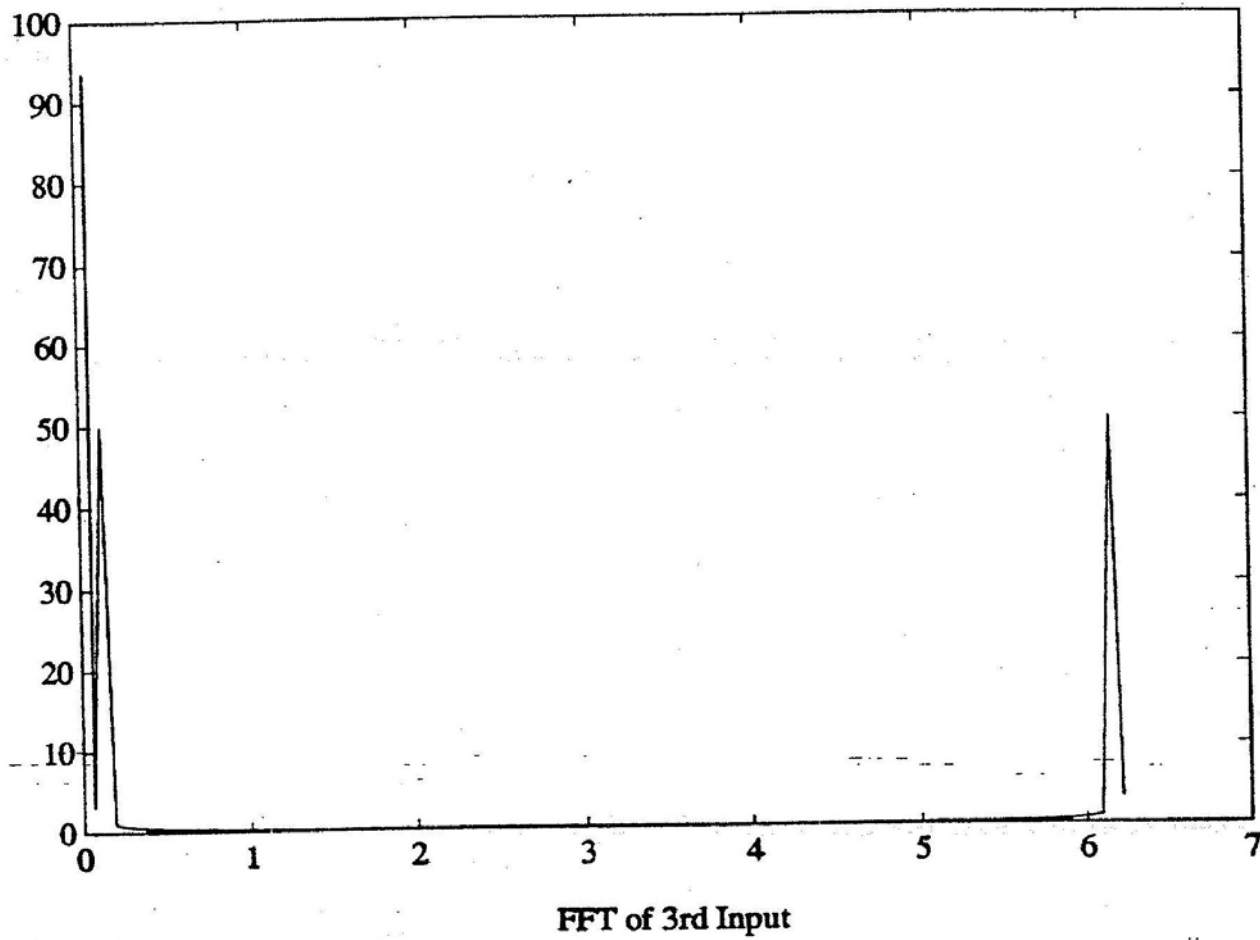
Theoretical Analysis:

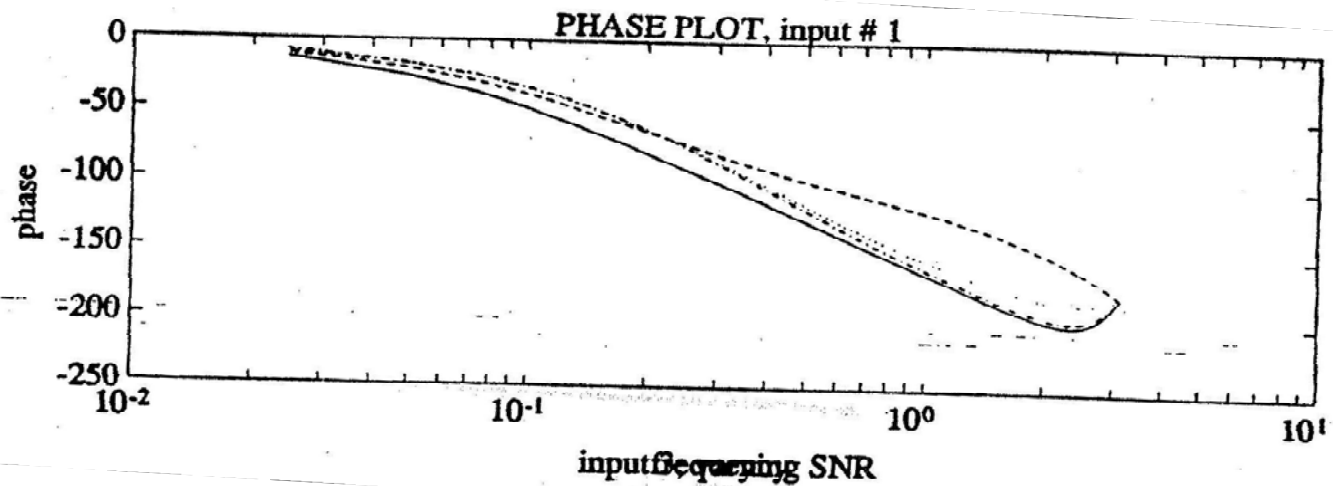
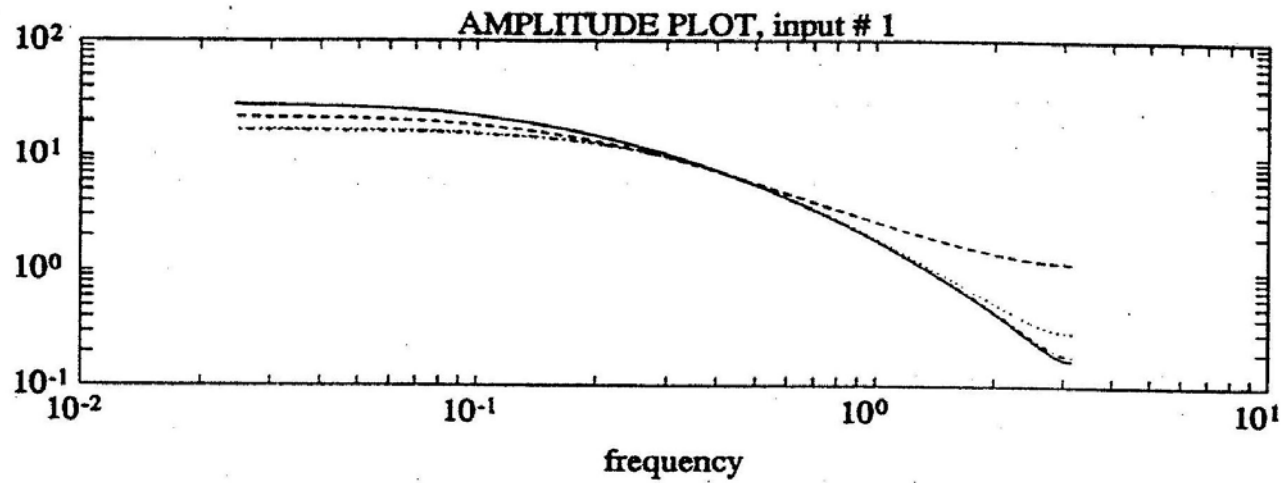
n : use $\det \hat{R}_s$ test $\Rightarrow n \leq 2$

Data is informative w. r. to $Ay = Bu + e$

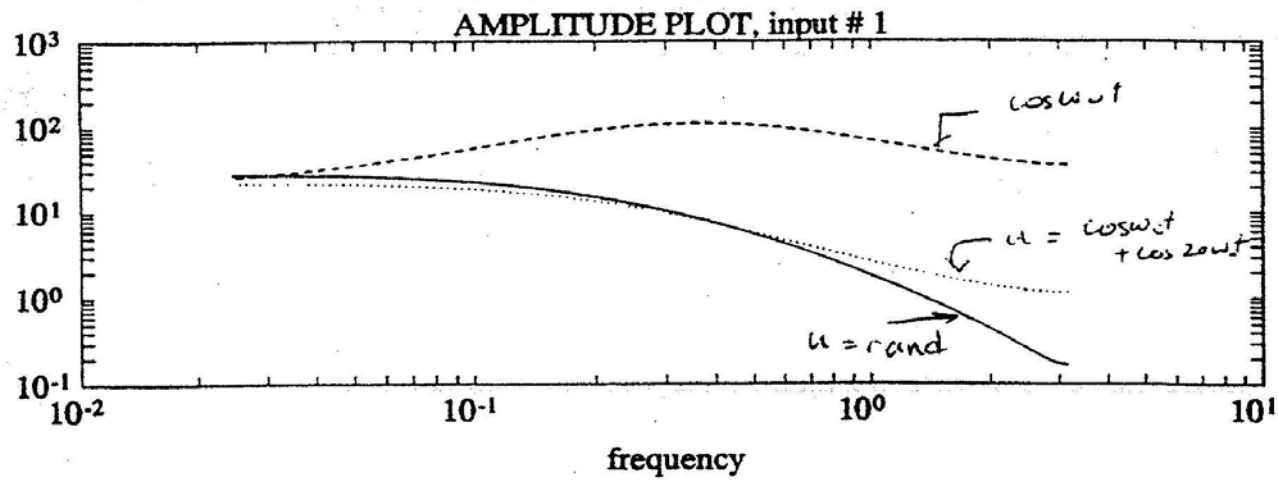
$N \rightarrow \infty \quad \hat{\theta}_N \rightarrow \theta^*$

Results for $SNR = 1$ are better in this case than $(\cos w_0 t)$.

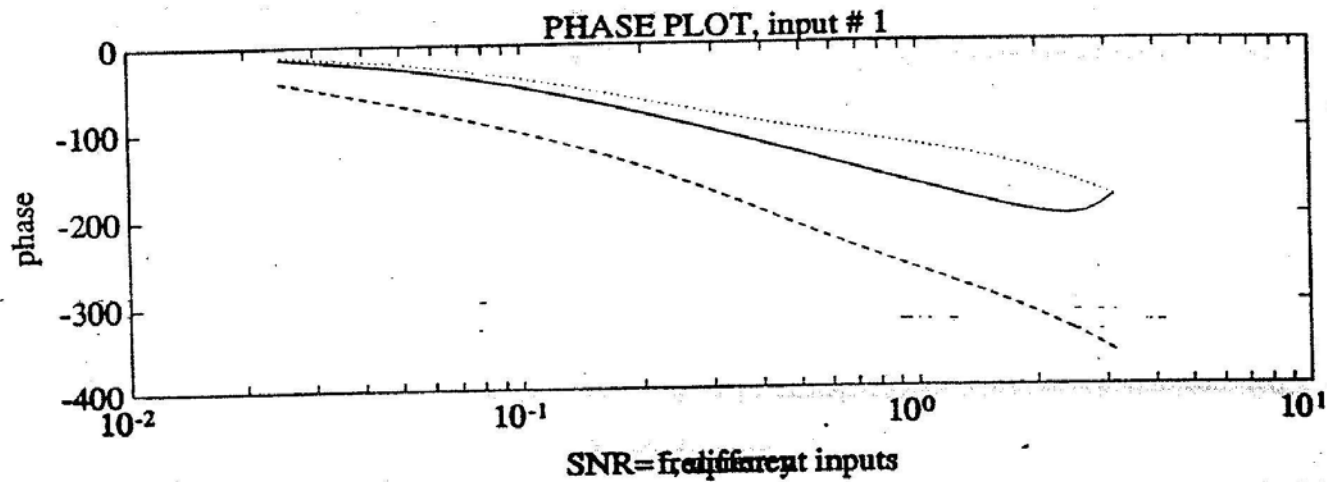




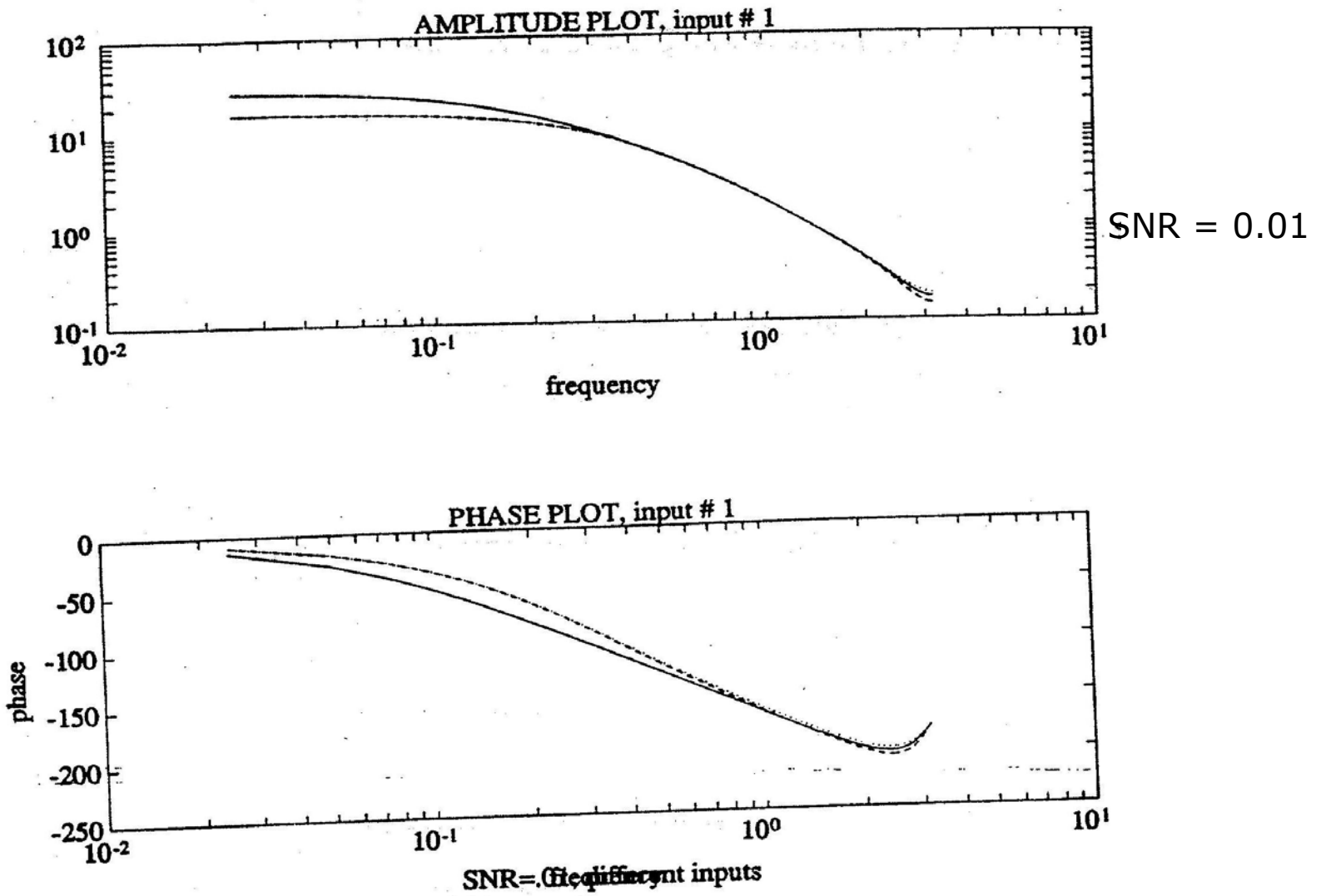
ARX (of an ARX plant)



SNR = 1



ARX (of an ARX plant)



ARX (of an ARX plant)

Conclusions

- Accuracy of estimates depend on $\lambda \left(\bar{E} \Psi(t, \theta^*) \Psi^T(t, \theta^*) \right)^{-1}$

If ARX, $\phi = \Psi$

- Estimate of accuracy $\simeq \hat{\lambda}_N \left(\sum_{t=1}^N \bar{E} \phi(t) \phi^T(t) \right)^{-1} = \hat{P}_\theta$

λ Large $\Rightarrow \hat{P}_\theta$ is large.

u is not rich $\Rightarrow \sum \bar{E} \phi(t) \phi^T(t)$ is singular.

$\Rightarrow \hat{P}_\theta$ is large.

u is rich (but close to not rich) $\Rightarrow \hat{P}_\theta$ is large.

λ small \Rightarrow better accuracy.

- Explanation (Proof of HW#3 problem 2)

$$\bar{E}\phi\phi^T(t) = \bar{E} \begin{bmatrix} -x(t-1) & \dots & -x(t-n) & u(t-1) & \dots & u(t-n) \end{bmatrix}^T \begin{bmatrix} \dots \end{bmatrix}$$

$$+ \bar{E} \begin{bmatrix} v(t-1) \\ \vdots \\ v(t-n) \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} v(t-1) & \dots & v(t-n) & 0 & \dots & 0 \end{bmatrix}$$

$$x = \frac{B}{A}u \quad v = \frac{1}{A}e$$

$$= \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} + \begin{bmatrix} \hat{A} & 0 \\ 0 & 0 \end{bmatrix}$$

$$A_{22} = \bar{E} \begin{pmatrix} u(t-1) \\ \vdots \\ u(t-n) \end{pmatrix} \begin{pmatrix} u(t-1) & \dots & u(t-n) \end{pmatrix}$$

A_{22} is near singular $\Rightarrow \det \bar{E} \phi \phi^T \simeq 0$

\Rightarrow Low accuracy for estimates. This can be countered by a small λ .

- How about other structures?

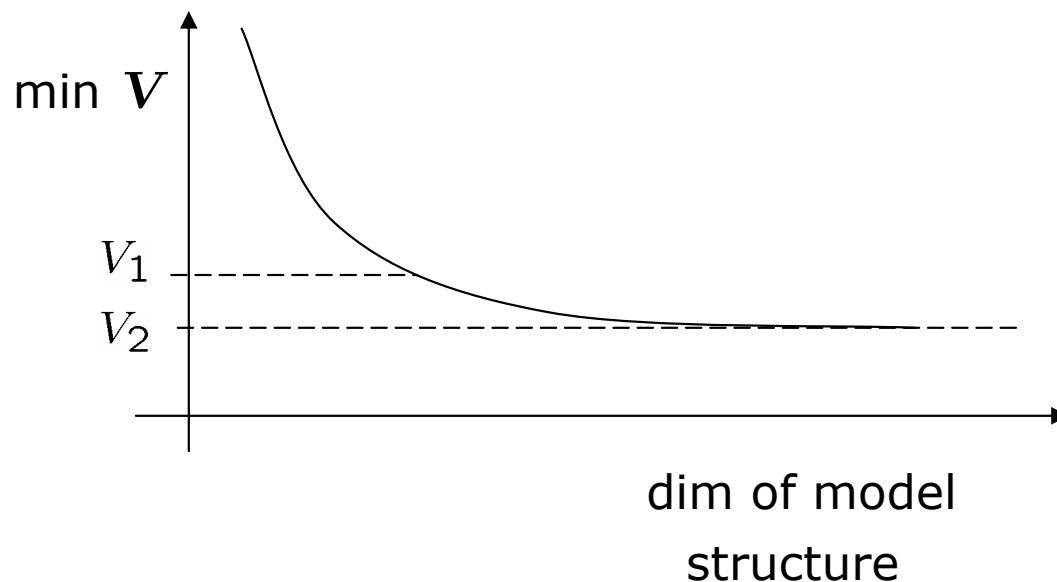
$$Ay = Bu + v$$

“same conclusions as long as v is persistent.”

OE, ARMAX,

Comparisons of Different Model Structures

- There is a trade off between the model structure complexity and the min. error



- Akaike's Final Prediction Error criterion (FPE)

model structure
↓

$$J(m) = \frac{1 + dm/N}{1 - dm/N} V_N(\hat{\theta}_N, Z^N)$$

$$V_N(\hat{\theta}_N, Z^N) = \min_{\hat{\theta}_N} \frac{1}{N} \sum_{t=1}^N \varepsilon^T \varepsilon(t|\theta)$$

- Based on minimum max likelihood. Tradeoff dm V_N
vs

- $\hat{\theta}_N^k = \underset{\theta \in \mathcal{U}_{m_R'}}{\operatorname{argmin}} V_N(\theta, Z^N)$

$$\bar{V}(\theta) = \lim_{N \rightarrow \infty} V_N(\theta, Z^N)$$

- A natural way to evaluate a model structure m_k is by

$$E\bar{V}(\hat{\theta}_N^k) = J(m_k)$$

$\hat{\theta}_N^k$ is a random variable \sim as $N(0, P_\theta)$

- Obtain estimates of both $\bar{V}(\theta)$ and \mathbf{J}

- Result

$$\begin{aligned} J(m) &= E\bar{V}(\hat{\theta}_N) \simeq EV_N(\hat{\theta}_N, Z^N) + \operatorname{tr} \bar{V}''(\theta^*)P_N \\ &\simeq V_N(\hat{\theta}_N, Z^N) + \frac{1}{N} \operatorname{tr} \bar{V}''(\theta^*)P_\theta \end{aligned}$$

$$\theta^* = \operatorname{argmin} \bar{V}(\theta)$$

Proof: expand $\bar{V}(\theta)$ around θ^*

$$\bar{V}(\hat{\theta}_N) = \bar{V}(\theta^*) + \frac{1}{2} (\hat{\theta}_N - \theta^*)^T \bar{V}''(\xi_N) (\hat{\theta}_N - \theta^*)$$

also

$$V_N(\hat{\theta}_N, Z^N) = \bar{V}(\theta^*, Z^N) - \frac{1}{2} (\hat{\theta}_N - \theta^*)^T \bar{V}_N''(\bar{\xi}_N, Z^N) (\hat{\theta}_N - \theta^*)$$

Notice:

$$E \frac{1}{2} (\hat{\theta}_N - \theta^*)^T \bar{V}''(\xi_N) (\hat{\theta}_N - \theta^*) = \frac{1}{2} E \text{ trace } \left\{ \bar{V}''(\xi_N) (\hat{\theta}_N - \theta^*) (\hat{\theta}_N - \theta^*)^T \right\}$$

$$\simeq \frac{1}{2} \text{ trace } (\bar{V}''(\theta_o^*) P_N) \simeq \frac{1}{2N} \text{ trace } \bar{V}''(\theta^*) P_\theta$$

and

$$EV_N(\theta^*, Z^N) \simeq \bar{V}(\theta^*)$$

$$\Rightarrow E\bar{V}(\hat{\theta}_N) \cong \bar{V}(\theta^*) + \frac{1}{2} \text{trace } \bar{V}''(\theta^*)P_N$$

$$EV_N(\hat{\theta}_N, Z^N) \simeq \bar{V}(\theta^*) - \frac{1}{2} \text{trace } \bar{V}''(\theta^*)P_N$$

$$\begin{aligned} J(m) = E\bar{V}(\hat{\theta}_N) &\simeq EV_N(\hat{\theta}_N, Z^N) + \text{trace } \bar{V}''(\theta^*)P_N \\ &\simeq V_N(\hat{\theta}_N, Z^N) + \frac{1}{N} \text{trace } \bar{V}'' P_\theta \end{aligned}$$

Akaike's Information Theoretic Criterion

- Let $V_N(\hat{\theta}_N, Z^N) = -\frac{1}{N}$ (Log likelihood function)
- Assume $\theta^* = \theta_o =$ true system
- The matrix $\bar{V}''(\theta_o)$ is invertible (identifiability).
- $\text{Cov } \hat{\theta}_N \simeq \frac{1}{N} [\bar{V}''(\theta_o)]^{-1} = (EL_N''(\theta_o))^{-1}$
- $J(m) = -\frac{1}{N} L_N(\theta, Z^N) + \text{trace } \bar{V}''(\theta_o) P_N$
 $= -\frac{1}{N} L_N(\theta, Z^N) + \frac{\dim D_m = d\mu}{N}$

- Model structure determination problem

$$\{\hat{\theta}_N^m, m\} = \underset{\substack{m \in M \\ \hat{\theta}^n \in D_m}}{\operatorname{argmin}} \frac{1}{N} \left[-L_N(\theta, Z^N) + \frac{d\mu}{N} \right]$$

- For every fixed m , $\frac{d\mu}{N}$ does not affect the min.

Example

- Assume that the innovations $(\varepsilon(t, \theta))$ are Gaussian with unknown variance.

- $$L_N(\theta, Z^N) = - \sum_{t=1}^N \frac{\varepsilon^2(t, \theta)}{2\lambda} - \frac{N}{2} \log \lambda - \frac{N}{2} \log 2\pi$$

- $$(\hat{\theta}_N, \hat{\lambda}_N) = \operatorname{argmin}_{\theta \in D_m} \left(-\frac{1}{N} L_N(\theta, Z^N) \right)$$

$$\hat{\lambda}_N = \frac{1}{N} \sum_{t=1}^N \varepsilon^2(t, \hat{\theta}_N)$$

$$\hat{\theta}_N = \operatorname{argmin}_{\theta} \sum_{t=1}^N \varepsilon^2(t, \theta)$$

- $L_N(\theta, Z^N) = -\frac{N}{2} - \frac{N}{2} \log \hat{\lambda}_N - \frac{N}{2} \log 2\pi$

- $\hat{m} = \operatorname{argmin} \left[\frac{1}{2} + \frac{1}{2} \log \hat{\lambda}_N^m + \frac{1}{2} \log 2\pi + \frac{dm}{N} \right]$

- Approximately minimize

$$\log \left[\frac{1}{N} \sum_{t=1}^N \varepsilon^2(t, \hat{\theta}_N) \right] + \frac{2dm}{N}$$

Akaike's Final Prediction Error Criterion

- $J(m) \simeq V_N(\hat{\theta}_N, Z^N) + \frac{1}{N} \text{trace } \bar{V}''(\theta^*) P_\theta$
- Let $V_N(\hat{\theta}_N, Z^N) = \frac{1}{2N} \sum_{t=1}^N \varepsilon^2(t, \theta) \quad \theta_o \in m$

$$\therefore P_\theta = \lambda_o (\bar{V}''(\theta_o))^{-1} \quad P_N \simeq \frac{1}{N} P_\theta$$

$$\lambda_o = 2\bar{V}(\theta_o) = (E e_o^2(t) = E \varepsilon^2(t, \theta_o))$$

$$\begin{aligned} \therefore J(m) &= V_N(\hat{\theta}_N, Z^N) + \text{trace } \frac{1}{N} (\bar{V}''(\theta_o) P_\theta) \\ &= V_N(\hat{\theta}_N, Z^N) + \lambda_o \frac{dm}{N} \end{aligned}$$

- $V_N(\hat{\theta}_N, Z^N) \simeq \bar{V}(\theta_o) - \frac{1}{2} \text{trace } \bar{V}''(\theta_o) P_N$

$$\simeq \frac{\lambda_o}{2} - \frac{1}{2} \frac{dm}{N} \lambda_o$$

estimate
of λ_o

$$\hat{\lambda}_N = \left[\frac{2V_N(\hat{\theta}, Z^N)}{N - dm} \right] \frac{N}{1}$$

$$J(m) = \frac{1 + dm/N}{1 - dm/N} V_N(\hat{\theta}_N, Z^N) \quad (\text{FPE})$$

$$= \frac{1 + dm/N}{1 - dm/N} \frac{1}{N} \sum_{t=1}^N \frac{1}{2} \varepsilon^2(t, \hat{\theta}_N)$$

Example 2

- δ is unknown.

- 3 experiments

$$u = PRBS$$

$$u = \cos \omega_o t \quad \omega_o = 2\pi/1000$$

$$u = \cos 2\omega_o t + \cos \omega_o t.$$

- Consider 2 - model structures

$$m_1 : \quad Ay = Bu + e$$

$$m_2 : y = \frac{B}{F}u$$

Experiment 1

- u = rand sequence

- $\hat{R}_s \cong \bar{E} \phi_s \phi_s^T(t)$

$$= \frac{1}{N} \sum_{t=1}^N \phi_s(t) \phi_s^T(t)$$

$$s = 2 \quad \det = 0.062$$

$$s = 3 \quad \det = 3 \times 10^{-7} \simeq 0$$

\Rightarrow system has dim = 2

- Estimated parameters

	a_1	a_2	b_1	b_2	$\bar{E}_{\varepsilon\varepsilon}^T$
ARX	-1.4	0.49	1	0.5	0
OE	-1.4	0.49	1	0.5	0

- Try different structures (n_a, n_b, n_k)

$$(1 \quad 1 \quad 1) \rightarrow 0.1022$$

$$(2 \quad 2 \quad 1) \rightarrow 0$$

$$(3 \quad 3 \quad 2) \rightarrow 0.0710$$

$$(3 \quad 4 \quad 1) \rightarrow 0$$

AIC or FPE $\Rightarrow (2 \quad 2 \quad 1)$ is the best choice.

Experiment 2

- $u = \cos w_0 t$ u is p.e of order 2

- $\hat{R}_s \simeq \frac{1}{N} \sum_{t=1}^N \phi_s(t) \phi_s^T$

$$\det(\hat{R}_2) = 2.66 \times 10^{-7}$$

$$\det(\hat{R}_3) = 3.7 \times 10^{-33}$$

$$\det(\hat{R}_1) = 8.77$$

ARX: $n \geq 1$ it may be $n = 2$ or 3

OE: $n \geq 1$ " " " "

Structure	ARX Parameters	V_N	OE Parameters	V_N
(1, 1, 1)	(-0.8627, 2.3182)	(0.034)	(-0.85, 2.51)	0.214
(2, 2, 1)	(-1.4, 0.49, 1, 0.5)	0	(-1.4, 0.49, 1, 0.5)	0
(3, 3, 1)	(- * * * *)	170	(* * * *)	0.3317
⋮				

Loss of identifiability.

Clearly (2 2 1) is the preferred model structure

Remark: Both experiments were generated from the model

$$\delta : y(t) = G_o u \quad G_o = \frac{q^{-1} (1 + 0.5q^{-1})}{1 - 1.4q^{-1} + 0.49q^{-2}}$$

Experiment 3

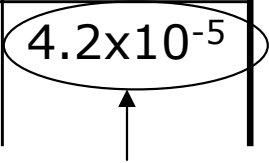
- $u = \cos w_0 t$
- $\det(\hat{R}_2) = 2.65 \times 10^{-7} \quad n \geq 1$

Structure	ARX Parameters	V_N	OE Parameters	V_N
(2 2 1)	(-1.4004, 0.4903, 0.98, 0.51)	3.4×10^{-5}	(-1.40, 0.49, 0.95, 0.54)	9.2×10^{-6}
(3 3 1)	singular		singular	

OE model is preferred.

- $u = \cos 20\omega_0 t + \cos \omega_0 t$
- $\det(\hat{R}_2) = 0.0692$ $\det(\hat{R}_3) = 3.9 \times 10^{-11}$

Structure	ARX Parameters	V_N	OE Parameters	V_N
(2 2 1)	(-1.4, 0.49, 1.0031, 0.4944)	1.6×10^{-5}	(1.399, 0.489, 0.99, 0.50)	1.3×10^{-5}
(3 3 1)	(* * *)	1.3×10^{-5}	(* * *)	4.2×10^{-5}



 Num. errors

AIC for ARX \Rightarrow (3 3 1) {In fact (4 4 1) does better}.

AIC for ARX and OE \Rightarrow (2 2 1) OE structure.

Remark: Data generated by

$$\delta : \quad y(t) = G_0 u + 0.01 e$$

Validation

- Use different sets of data to validate the model structure and the estimated model.
- You can obtain different estimates using the data and then average them. OR you can construct new input-output pairs and re-estimate.

Conclusions

- Criterion contains a penalty function for the dimension of the system.
- “AIC” is one way of doing that. (FPE) is an estimate of the AIC with a quadratic objective.
- “AIC” has connections with information theory observation

$$\begin{array}{l} Z^t \longrightarrow \text{assumed PDF } f_m(t, Z^t) \\ \text{true PDF } f_o(t, Z^t) \end{array}$$

$$\text{Entropy of } f_o \text{ w. r. to } f_m = \delta(f_o, f_m) = -I(f_o, f_m)$$

$$I(f_o, f_m) = \int f_o(t, x^t) \log \frac{f_o(t, x^t)}{f_m(t, x^t)} dx^t.$$

over the

observation

= information distance

$$\hat{\theta}_N = \underset{\theta}{\operatorname{argmin}} I(f_o, f_m(\theta, Z^N, N))$$

- "AIC" is the average information distance

$$E_{\hat{\theta}_N} I(f_o, f)$$

after some simplification

$$\hat{\theta}_{AIC} = \underset{\theta}{\operatorname{argmin}} \left\{ \frac{1}{N} \sum_{t=1}^N l(\varepsilon(t, \theta), t, \theta) + \frac{\dim \theta}{N} \right\}$$

- Careful about Numerical errors!?