

Introduction to Research

Arun K. Tangirala

Modelling Skills



Objectives

To learn the following:

- ▶ What is a model?
- ▶ First-principles (mechanistic) vs. empirical models.
- ▶ Systematic procedure for building models from data.
- ▶ Few critical aspects of data-driven (empirical) modelling.

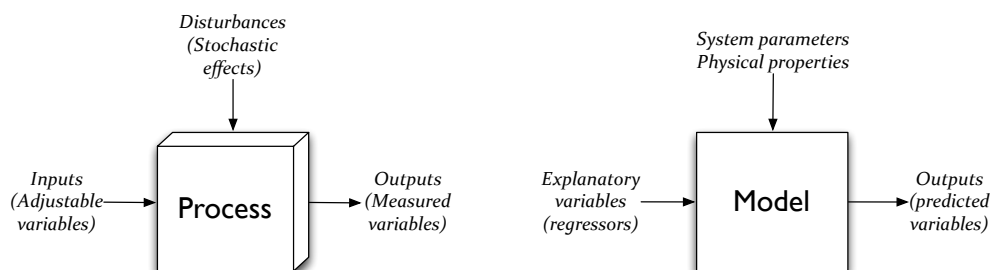
What is a model?

A model is a mathematical (or descriptive) abstraction of a process that emulates its behaviour or characteristics

For the purposes of

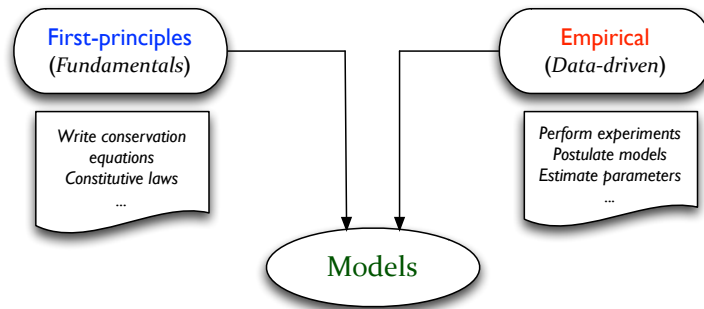
- ▶ Prediction (inferring the unknowns)
- ▶ Classification (pattern recognition)
- ▶ Fault detection
- ▶ Process simulation
- ▶ Design and optimization
- ▶ ...

Model vs. process

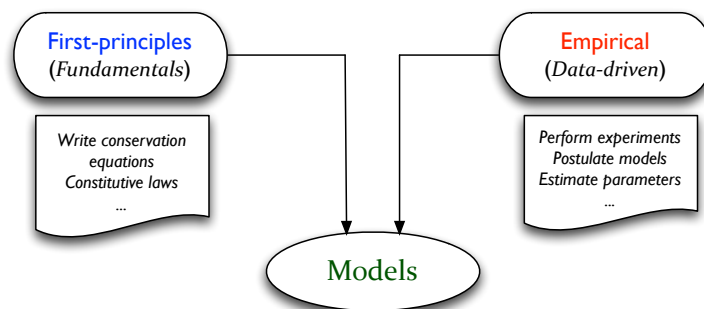


A model emulates the process given the operating conditions, parameters and physical properties. However, in general, its architecture is different from the process!

Two broad approaches to modelling



Two broad approaches to modelling



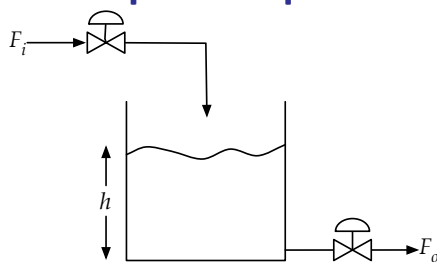
Test drive of a vehicle

Taking a vehicle for a test drive is a simple example of how experiments are a natural way of understanding processes. As the vehicle is subjected to different test (road) conditions, its response is used by the test-driver to develop a "mental model" of the vehicle.

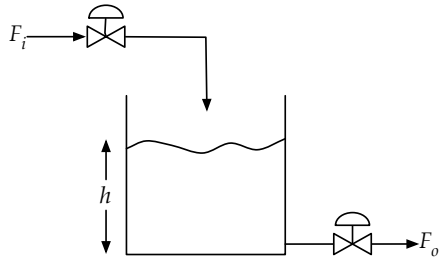
First-principles vs. Empirical modelling

First-principles	Empirical
Causal, continuous, non-linear differential-algebraic equations	Models are usually black-box and discrete-time
Model structures are quite rigid - the structure is decided by the equations describing fundamental laws.	Model structures are extremely flexible - implies they have to be assumed/known a priori.
Can be quite challenging and time-consuming to develop	Relatively much easier and are less time-consuming to develop
Require good numerical ODE and algebraic solvers.	Require good estimators (plenty of them available)
Very effective and reliable models - can be used for wide range of operating conditions	Model quality is strongly dependent on data quality - usually have poor extrapolation capabilities
Transparent and can be easily related to physical parameters/characteristics of the process	Difficult to relate to physical parameters of the process (black-box) and the underlying phenomena

Example: Liquid level system



Example: Liquid level system

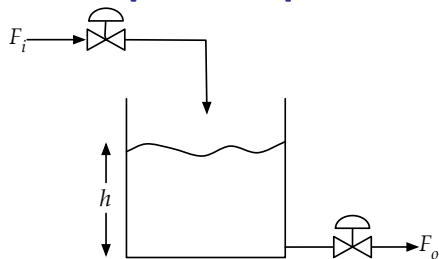


Case 1: Steady-state model
between F_o and h :

$$F_o = C_v \sqrt{h}$$

(C_v is estimated from data)

Example: Liquid level system



Case 1: Steady-state model
between F_o and h :

$$F_o = C_v \sqrt{h}$$

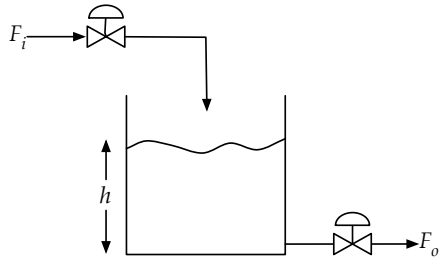
(C_v is estimated from data)

Case 2: Dynamic, non-linear model
of $h(t)$ for changes in $F_i(t)$.

$$A_c \frac{dh}{dt} = F_i - C_v \sqrt{h}$$

(requires numerical solvers)

Example: Liquid level system



Case 1: Steady-state model
between F_o and h :

$$F_o = C_v \sqrt{h}$$

(C_v is estimated from data)

Case 2: Dynamic, non-linear model
of $h(t)$ for changes in $F_i(t)$.

$$A_c \frac{dh}{dt} = F_i - C_v \sqrt{h}$$

(requires numerical solvers)

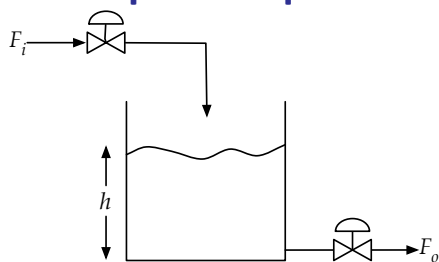
Case 3: Approximate linear, dynamic
model about an operating point

$$A_c \frac{d\tilde{h}}{dt} = \tilde{F}_i - \beta \tilde{h}, \quad \beta = \frac{C_v}{2\sqrt{h_s}}$$

$$\tilde{F}_i = F_i - F_s, \quad \tilde{h}_i = h_i - h_s$$

Typically, h_s, F_s are steady-state values.

Example: Liquid level system



Case 1: Steady-state model
between F_o and h :

$$F_o = C_v \sqrt{h}$$

(C_v is estimated from data)

Case 2: Dynamic, non-linear model
of $h(t)$ for changes in $F_i(t)$.

$$A_c \frac{dh}{dt} = F_i - C_v \sqrt{h}$$

(requires numerical solvers)

Case 3: Approximate linear, dynamic
model about an operating point

$$A_c \frac{d\tilde{h}}{dt} = \tilde{F}_i - \beta \tilde{h}, \quad \beta = \frac{C_v}{2\sqrt{h_s}}$$

$$\tilde{F}_i = F_i - F_s, \quad \tilde{h}_i = h_i - h_s$$

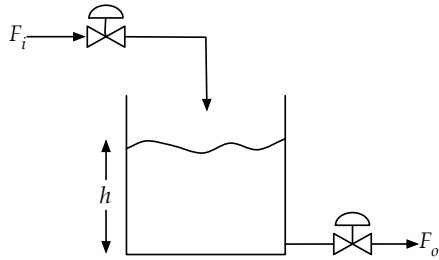
Typically, h_s, F_s are steady-state values.

Case 4: Approximate empirical, linear, grey-box, dynamic, discrete-time model,

$$h[k] = a_1 h[k-1] + b_1 F_i[k-1] + \varepsilon[k]$$

(Parameters a_1, b_1 estimated from experimental data)

Example: Liquid level system



Case 1: Steady-state model between F_o and h :

$$F_o = C_v \sqrt{h}$$

(C_v is estimated from data)

Case 2: Dynamic, non-linear model of $h(t)$ for changes in $F_i(t)$.

$$A_c \frac{dh}{dt} = F_i - C_v \sqrt{h}$$

(requires numerical solvers)

Case 3: Approximate linear, dynamic model about an operating point

$$A_c \frac{d\tilde{h}}{dt} = \tilde{F}_i - \beta \tilde{h}, \quad \beta = \frac{C_v}{2\sqrt{h_s}}$$

$$\tilde{F}_i = F_i - F_s, \quad \tilde{h}_i = h_i - h_s$$

Typically, h_s, F_s are steady-state values.

Case 4: Approximate empirical, linear, grey-box, dynamic, discrete-time model,

$$h[k] = a_1 h[k-1] + b_1 F_i[k-1] + \varepsilon[k]$$

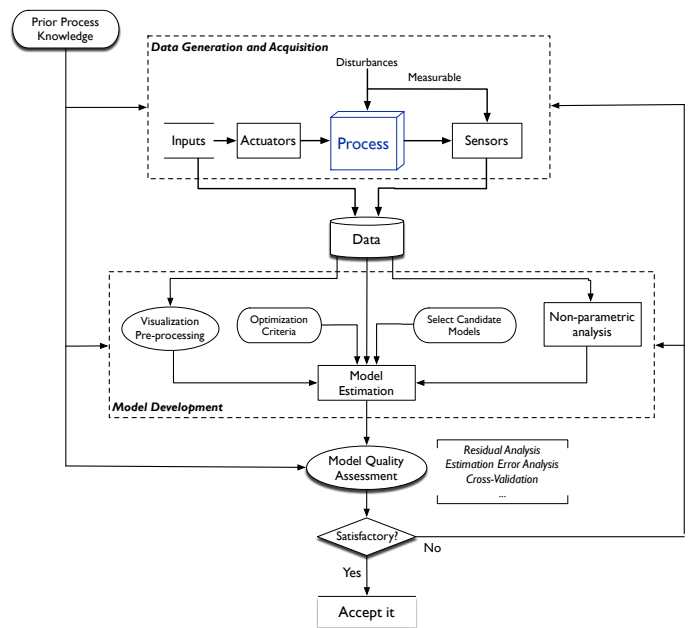
(Parameters a_1, b_1 estimated from experimental data)

Case 5: Black-box, dynamic, discrete-time model

1. Model structure based on ease of estimation and end-use
2. Model may not be physically interpretable, but designed to give good predictions.

Building models from data: systematic procedure

- Primarily one finds three stages, which again contain sub-stages
- All models should be subjected to a critical validation test
- Final model quality largely depends on the quality of data
- Generating rich (informative) data is therefore essential



Two broad classes of models

Time-series models

- ▶ Suited for modelling stochastic or random processes (e.g., stock market index, rainfall, sensor noise)
- ▶ Causes are either unknown or also random
- ▶ Usually dynamic models (in a few applications, steady-state as well)
- ▶ Challenges: choosing model structure, making the right assumptions on process characteristics, non-linearities, etc.

Two broad classes of models

... contd.

Input-output or causal models

- ▶ Suited for modelling relationship between a variable (or more) and other **explanatory** variables (a.k.a. **regressors** or **factors**) (e.g., power and current in a wire, temperature and coolant flow)
- ▶ Regressors may be known accurately or with some error
- ▶ Models could be steady-state or dynamic.
- ▶ Challenges: sufficient variability in factors, selection of regressors, measurements and regressors available at different sampling rates, choosing the order of dynamics, etc.

Two broad classes of models

... contd.

Input-output or causal models

- ▶ Suited for modelling relationship between a variable (or more) and other **explanatory** variables (a.k.a. **regressors** or **factors**) (e.g., power and current in a wire, temperature and coolant flow)
- ▶ Regressors may be known accurately or with some error
- ▶ Models could be steady-state or dynamic.
- ▶ Challenges: sufficient variability in factors, selection of regressors, measurements and regressors available at different sampling rates, choosing the order of dynamics, etc.

Remember

In practice, all modelling exercises demand a careful handling of uncertainties both at the experimental and modelling stages!

Excite the process sufficiently!

Example

Consider a steady-state model:

$$y[k] = b_0 + b_1 u[k] + b_2 u^2[k]$$

- ▶ Three unknowns
- ▶ Therefore, data corresponding to three steady-states is required
- ▶ A more general statement: The regressor matrix

$$\mathbf{U} = \begin{bmatrix} 1 & u[k_1] & u^2[k_1] \\ 1 & u[k_2] & u^2[k_2] \\ 1 & u[k_3] & u^2[k_3] \end{bmatrix} \quad (1)$$

should be non-singular, i.e., of full rank.

For dynamic systems

Example

Consider a (deterministic) process:

$$y[k] = b_1 u[k-1] + b_2 u[k-2] + b_3 u[k-3]$$

Suppose that the process is excited with $u[k] = \sin(\omega_0 k)$ (sine of single frequency).

For dynamic systems

Example

Consider a (deterministic) process:

$$y[k] = b_1 u[k-1] + b_2 u[k-2] + b_3 u[k-3]$$

Suppose that the process is excited with $u[k] = \sin(\omega_0 k)$ (sine of single frequency). With this sine wave input,

$$\begin{aligned} y[k] &= b_1 \sin(\omega_0 k - \phi) + b_2 \sin(\omega_0 k - 2\phi) + b_3 \sin(\omega_0 k - 3\phi) \\ &= \left(b_1 + \frac{b_2}{2 \cos \omega_0} \right) \sin(\omega_0 k - \phi) + \left(b_3 + \frac{b_2}{2 \cos \omega_0} \right) \sin(\omega_0 k - 3\phi) \\ &= \left(b_1 + \frac{b_2}{2 \cos \omega_0} \right) u[k-1] + \left(b_3 + \frac{b_2}{2 \cos \omega_0} \right) u[k-3] \end{aligned}$$

Only two of the three parameters can be identified! Why?

Effects of randomness (noise) in data

The second challenging aspect of empirical modelling is the **presence of random or stochastic effects**.

Randomness (uncertainties) in the process influences identification in several ways:

- ▶ Accuracy of predictions
- ▶ Errors in parameter estimates
- ▶ Goodness of the deterministic model

Signal-to-Noise Ratio (SNR)

SNR

A key measure that quantifies the effects of noise is the **Signal-to-Noise Ratio (SNR)**, which is defined as the ratio of variance of signal to the variance of noise in a measurement.

$$\text{SNR} = \frac{\sigma_{\text{signal}}^2}{\sigma_{\text{noise}}^2}$$

SNR can be interpreted as a measure of the degree of certainty (deterministic portion) to uncertainty

Example: Effect of SNR

Estimation of linear model from measured data

$$\text{Process : } x[k] = b_1 u[k - 1] + b_0; b_1 = 5; b_0 = 2$$

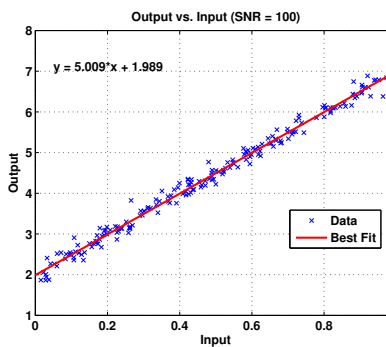
Only $y[k] = x[k] + v[k]$ (measurement) is available. Goal is to estimate b_1, b_0 .

Example: Effect of SNR

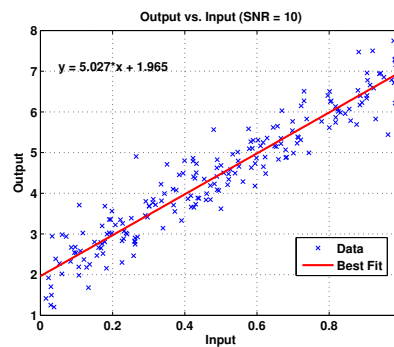
Estimation of linear model from measured data

$$\text{Process : } x[k] = b_1 u[k - 1] + b_0; b_1 = 5; b_0 = 2$$

Only $y[k] = x[k] + v[k]$ (measurement) is available. Goal is to estimate b_1, b_0 .



$$\sigma_{\hat{b}_1} = 0.036, \sigma_{\hat{b}_0} = 0.02$$



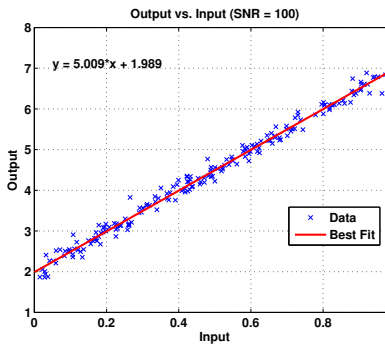
$$\sigma_{\hat{b}_1} = 0.114, \sigma_{\hat{b}_0} = 0.064$$

Example: Effect of SNR

Estimation of linear model from measured data

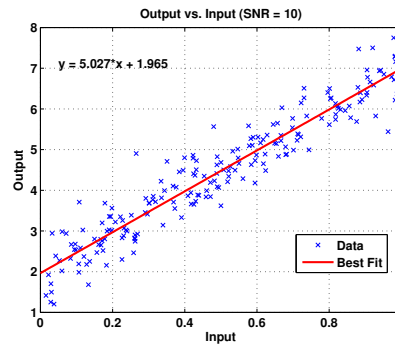
$$\text{Process : } x[k] = b_1 u[k - 1] + b_0; b_1 = 5; b_0 = 2$$

Only $y[k] = x[k] + v[k]$ (measurement) is available. Goal is to estimate b_1, b_0 .



$$\sigma_{\hat{b}_1} = 0.036, \sigma_{\hat{b}_0} = 0.02$$

Arun K. Tangirala, IIT Madras



$$\sigma_{\hat{b}_1} = 0.114, \sigma_{\hat{b}_0} = 0.064$$

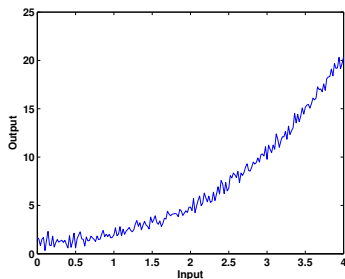
Introduction to Research

Decrease in SNR increases the error in parameter estimates (proportional to $\sqrt{1/\text{SNR}}$)

Example: Overfitting

$$\text{Process : } x[k] = 1.2 + 0.4u[k] + 0.3u^2[k] + 0.2u^3[k]$$

Only $y[k] = x[k] + v[k]$ (measurement) is available. Goal is to fit a suitable model.



Input-output data

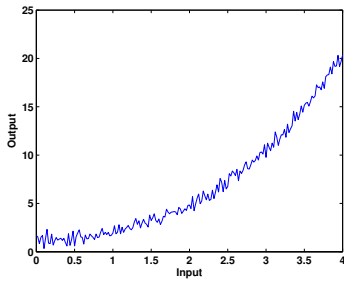
Arun K. Tangirala, IIT Madras

Introduction to Research

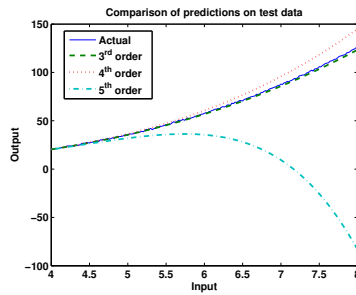
Example: Overfitting

Process : $x[k] = 1.2 + 0.4u[k] + 0.3u^2[k] + 0.2u^3[k]$

Only $y[k] = x[k] + v[k]$ (measurement) is available. Goal is to fit a suitable model.



Input-output data

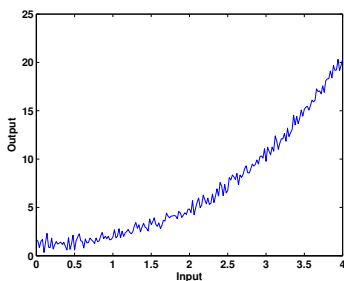


Cross-validation

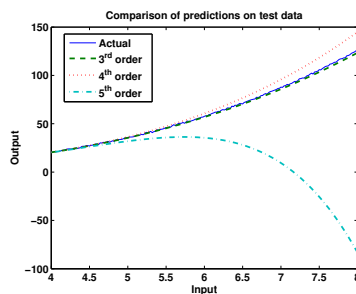
Example: Overfitting

Process : $x[k] = 1.2 + 0.4u[k] + 0.3u^2[k] + 0.2u^3[k]$

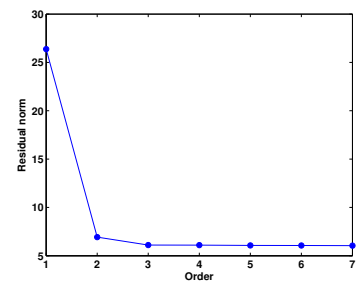
Only $y[k] = x[k] + v[k]$ (measurement) is available. Goal is to fit a suitable model.



Input-output data



Cross-validation

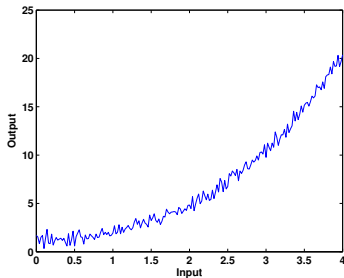


Selecting the order

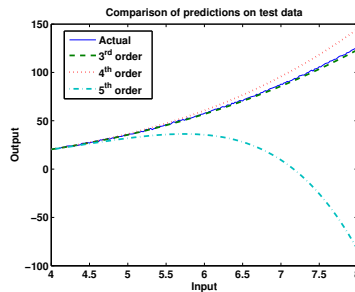
Example: Overfitting

Process : $x[k] = 1.2 + 0.4u[k] + 0.3u^2[k] + 0.2u^3[k]$

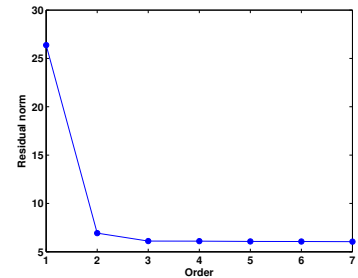
Only $y[k] = x[k] + v[k]$ (measurement) is available. Goal is to fit a suitable model.



Input-output data



Cross-validation



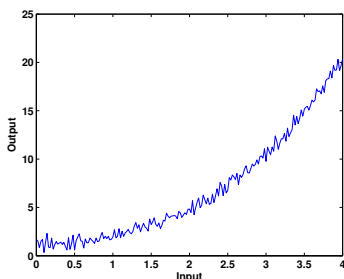
Selecting the order

Overfitting occurs whenever the local chance variations are treated as global characteristics

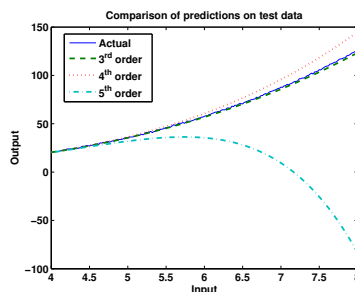
Example: Overfitting

Process : $x[k] = 1.2 + 0.4u[k] + 0.3u^2[k] + 0.2u^3[k]$

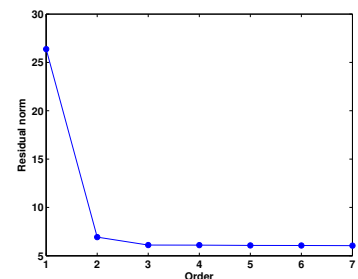
Only $y[k] = x[k] + v[k]$ (measurement) is available. Goal is to fit a suitable model.



Input-output data



Cross-validation



Selecting the order






Overfitting occurs whenever the local chance variations are treated as global characteristics

3rd order fit: $\hat{y}[k] = 1.17 + 0.35 u[k] + 0.36 u^2[k] + 0.19 u^3[k]$
 (± 0.05) (± 0.1) (± 0.06) (± 0.01)

Questions for reflection

- ▶ What type of models are possible? Which one(s) to choose?
- ▶ How do we “fit” a model that “explains” the variations observed in experimental data?
- ▶ How to “correctly” account for the deterministic and stochastic effects?
- ▶ Will the experiment influence the model that we fit? If yes, in what way?
- ▶ How do we set up and solve the problem of estimating the unknown model parameters?
- ▶ What kind of experiments should we design to obtain a good quality model?
- ▶ How much data do we collect? (what should be the sample size?)

Bibliography I

-  Bendat, J. S. and A. G. Piersol (2010). *Random Data: Analysis and Measurement Procedures*. 4th edition. New York, USA: John Wiley & Sons, Inc.
-  Johnson, R. A. (2011). *Miller and Freund's: Probability and Statistics for Engineers*. Upper Saddle River, NJ, USA: Prentice Hall.
-  Montgomery, D. C. and G. C. Runger (2011). *Applied Statistics and Probability for Engineers*. 5th edition. New York, USA: John Wiley & Sons, Inc.
-  Ogunnaike, B. A. (2010). *Random Phenomena: Fundamentals of Probability and Statistics for Engineers*. Boca Raton, FL, USA: CRC Press, Taylor & Francis Group.
-  Tangirala, A. K. (2014). *Principles of System Identification: Theory and Practice*. CRC Press, Taylor & Francis Group.