

# Parametric Statistical Modeling

ECE 275A – Statistical Parameter Estimation

Ken Kreutz-Delgado  
ECE Department, UC San Diego

# Why Parametric Statistical Models?

- They **include deterministic models** as a special case.
- They **succinctly** capture & encode properties of the perceived world.
  - Allow for data compression
  - Enable efficient explanation of past measurements
  - Enable efficient prediction of future measurements
- **Statistical** models acknowledge that **uncertainty, error, and chance exist** in our understanding of the world.
- They provide “quality of fit” measures.
  - Model-mismatch measures
  - Parameter estimate quality measures

# Contrast with Nonparametric Models

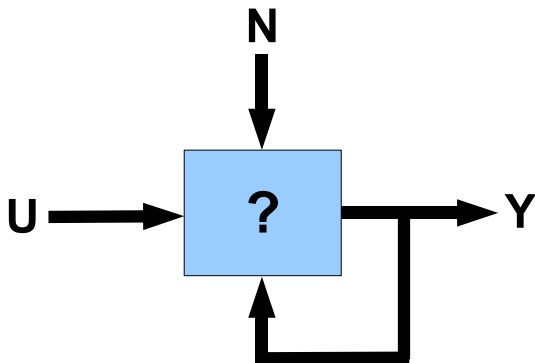
With Nonparametric Models:

- No *a priori* structural or modeling information is utilized
  - Difficult to model dynamic (nonstationary) processes.
  - Difficult to gain *insight* into physical and other processes.
- Often, all data must be kept regardless of dimensionality or amount.
  - Data processing is expensive, particularly if data is collected in an on-going, on-line manner.
- Probability density function (pdf) approximations are constructed via “binning” of data to directly form empirical density functions
  - As data is collected in an on-line manner, density-related estimates must often be recomputed via “batch processing”

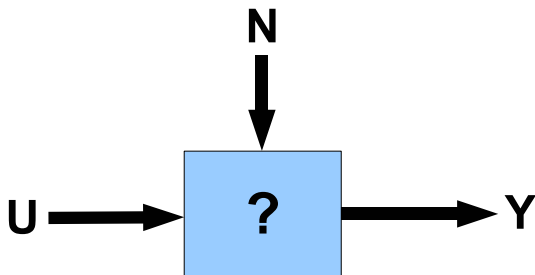
# Parametric Models = Generative Models

- Parametric Probabilistic Models are posited to explain observed phenomenon.
  - For this reason they are often referred to as *Generative Models*; models that are presumed to be have generated observed data.
  - They are also known as *forward models*, as one imagines processing inputs, noise, past observations, and *parameters* in a “forward direction” to produce observed data.
  - The problem of estimating unknown model parameters given observed inputs and past observations is known as the *inverse problem*.
- In its fullness then, the problem of parameter estimation involves:
  - 1 Proposing and constructing a candidate generative model to explain some phenomenon of interest.
  - 2 Collecting data corresponding to inputs and outputs of the model.
  - 3 Solving the statistical inverse problem of estimating the unknown parameters of the model
  - 4 Validating the model. If the statistics of the outputs of the model do not match the statistics of our observed data, and/or the estimated model yields poor predictive capabilities, we must refine and improve our posited model.
- **In this course we are primarily concerned with issues 1 and 3.**

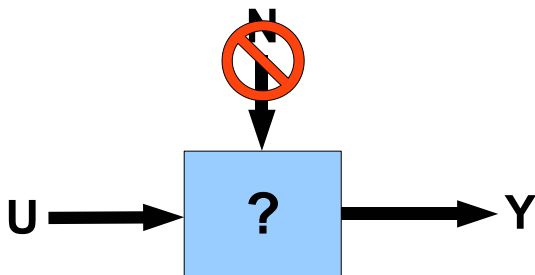
# Generative Model of World or System or ...



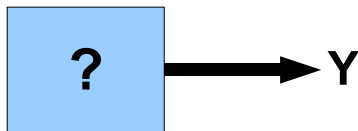
# Generative Model of World or System or ...



# Generative Model of World or System or ...



# Generative Model of World or System or ...





## Example: “Moving Average” (MA) Linear Gaussian Model

For unknown parameters  $\theta_i$ ,  $i = 1, \dots, M$ , consider the *Moving Average* (of  $f(u(t))$ ) “Linear” (in the parameters!) Time-Series Model

$$y[t] = \theta_1 f(u[t]) + \dots + \theta_M f(u[t - M + 1]) + n[t] \quad \text{with} \quad n(t) \sim \mathcal{N}(0, \sigma^2)$$

The sequence of inputs  $u[t]$  is assumed known, as is the general function  $f(\cdot)$ . The noise  $n(t)$  is considered to be iid with  $\sigma^2$  known.

Some examples are  $f(x) = x$ ,  $f(x) = \cos(x)$ ,  $f(x) = \exp(x)$ , etc.

Set  $\theta \triangleq (\theta_1, \dots, \theta_M)^T$  and  $a(t) = (f(u[t]), \dots, f(u[t - M + 1]))^T$  then

$$y[t] = a^T(t)\theta + n(t)$$

with  $y[t] \sim \mathcal{N}(a^T(t)\theta, \sigma^2)$ ,

$$P_{y[t];\theta}(y[t]) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y[t] - a^T(t)\theta)^2 \right\}$$

showing that the MA model for  $y[t]$  is entirely equivalent to a probabilistic model parameterized by  $\theta$ .

## “Moving Average” (MA) Linear Gaussian Model – Cont.

Now consider collecting a “batch” of  $N > M$  samples of  $y[t]$  ( $t = 1, \dots, N$ ) and set  $Y \triangleq (y[1], \dots, y[N])^T$ .

The MA model is entirely equivalent to the vector-matrix “batch data” parametric probabilistic model

$$Y \sim N(A\theta, C), \quad C = \text{diag}(\sigma^2, \dots, \sigma^2) = \sigma^2 I$$

$$P_\theta(Y) \triangleq P_{Y;\theta}(Y) = \frac{1}{\sqrt{(2\pi)^N \det C}} \exp \left\{ -\frac{1}{2} \|Y - A\theta\|_{C^{-1}}^2 \right\}$$

where the *data matrix*  $A$  is a known  $N \times M$  matrix whose  $N$  rows are comprised of the  $M$ -dimensional row vectors  $a^T(t)$  and

$$\|Y - A\theta\|_{C^{-1}}^2 \triangleq (Y - A\theta)^T C^{-1} (Y - A\theta)$$

with  $C$  a known (diagonal, in this example)  $N \times N$  matrix.

# Model Fitting by Likelihood Maximization

The function

$$\ell_Y(\theta) \triangleq P_\theta(Y)$$

is **the likelihood of  $\theta$**  (i.e. of the model  $P_\theta(\cdot)$ ) given the measured data  $Y$ .

The principle of maximum likelihood estimation says to find that model (parameter  $\theta$ ) for which the likelihood function takes its maximum value, given the measured data  $Y$ .

Maximizing the likelihood function is equivalent to minimizing the negative logarithm of the likelihood function (the “negative log-likelihood”). For the important *linear Gaussian model with known covariance matrix* example considered above, this corresponds to finding the parameter vector  $\theta$  that minimizes the weighted least squares loss function

$$\|Y - A\theta\|_{C^{-1}}^2 \triangleq (Y - A\theta)^T C^{-1} (Y - A\theta)$$

Note that when  $C = I$  this reduces to a (unweighted) least squares problem. In either case the problem is one of solving a linear inverse problem

$$Y \approx A\theta$$

in an appropriate least squares sense.

## Linear Gaussian Model

$$y = Ax + v, \quad v \sim N(0, C), \quad C \text{ is symmetric and full rank.}$$

Equivalent to

## Parametric Probability Model

$$y \sim N(Ax, C), \quad P_x(y) = \frac{1}{\sqrt{(2\pi)^m \det C}} \exp \left\{ -\frac{1}{2} \|y - Ax\|_{C^{-1}}^2 \right\}$$

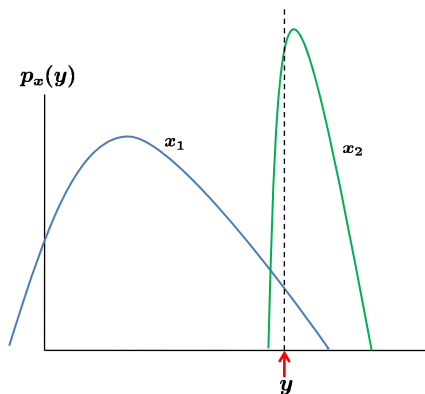
where

$$\|y - Ax\|_{C^{-1}}^2 \triangleq (y - Ax)^T C^{-1} (y - Ax)$$

# The Likelihood Function

**Likelihood** of  $x$  given  $y$  :  $\ell_y(x) \triangleq p_x(y)$

In the figure below, note that it is rational to prefer probability model  $p_{x_2}(\cdot)$  over model  $p_{x_1}(\cdot)$  given the observed value of  $y$ :



$$\ell_y(x_2) = p_{x_2}(y) > p_{x_1}(y) = \ell_y(x_1)$$

# Model Fitting by Likelihood Maximization

**Maximum Likelihood Estimate** of  $x$  given  $y$

$$\hat{x} = \arg \max_x \ell_y(x) = \arg \max_x p_x(y)$$

For the Linear Gaussian model this is equivalent to

$$\hat{x} = \arg \min_x \|y - Ax\|_{C^{-1}}^2$$

Which corresponds to solving a **Linear Inverse Problem**

$$y \approx Ax$$

in an appropriate **Minimum Norm** sense, where  $y$  and  $x$  are **Vectors**,  $A$  is a (matrix representation of) a **Linear Operator**, and  $\|\cdot\|_{C^{-1}}$  is a (weighted) **Norm**.

# Need for Vector Space Theory

- What are **Vectors** and **Linear Vector Spaces**?
- What are **Norms** and **Normed Linear Vector Spaces**?  
(Theory of **Banach Spaces**.)
- What are **Inner Products** and **Inner Product Vector Spaces**?  
(Theory of **Hilbert Spaces**.)
- What are **Linear Operators** and the **Geometry Induced by Linear Operators**.  
(The '**Four Fundamental Subspaces**' associated with a linear operator.)
- What is a **Linear Inverse Problem**?  
(**Well-Posed** and **Ill-Posed** Inverse Problems.)
- How does one solve a linear inverse problem?
  - **Minimum Norm Solution** and **Weighted Least Squares Solution**.
  - **Projection Theorem** in Hilbert Spaces. (**Orthogonality Principle**.)
  - **Generalized Inverses**. (**Pseudo-Inverse, QR-factorization, SVD**.)