

System Identification, Lecture 4

Kristiaan Pelckmans (IT/UU, 2338)

Course code: 1RT880, Report code: 61800 - Spring 2014
F, FRI Uppsala University, Information Technology

2 March 2015

Lecture 4

- Stochastic Setup.
- Interpretation.
- Maximum Likelihood.
- Least Squares Revisited.
- Instrumental Variables.

Stochastic Setup

Why:

- Analysis (abstract away).
- Constructive.

Basics:

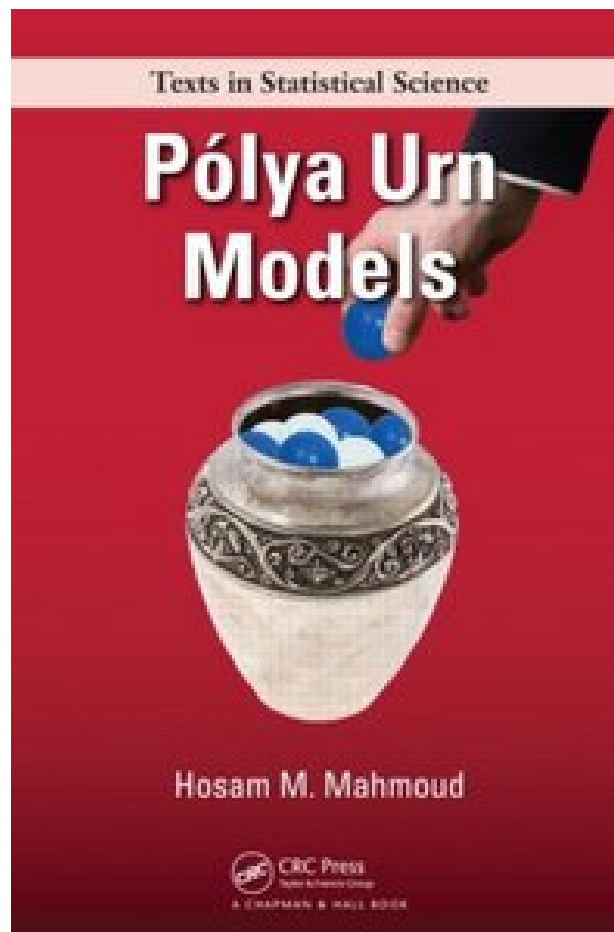
- Samples ω
- Sample space $\Omega = \{\omega\}$
- Event $A \subset \Omega$.

• Rules of probability $\Pr : \{\Omega\} \rightarrow [0, 1]$

1. $\Pr(\Omega) = 1,$

2. $\Pr(\{\}) = 0,$

3. $\Pr(A_1 \cup A_2) = \Pr(A_1) + \Pr(A_2)$ for $A_1 \cap A_2 = \{\}$

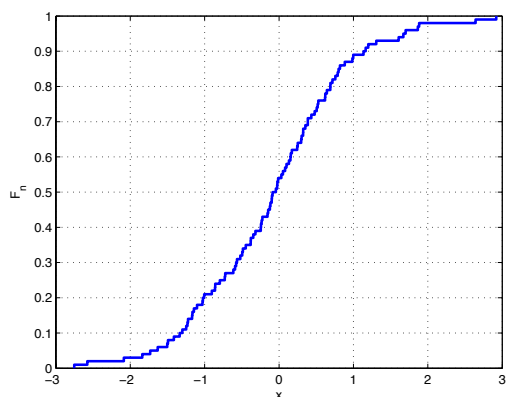


Random variables is a function $X : \Omega \rightarrow \mathbb{R}$. Examples:

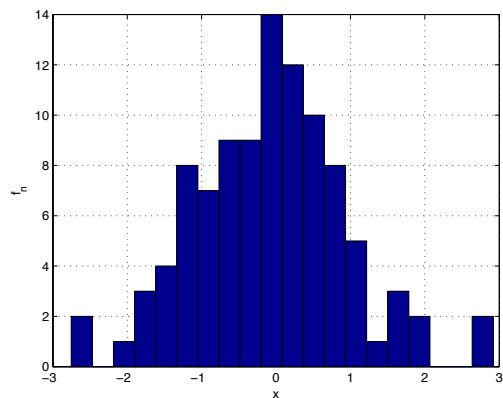
- Urn. Sample= Ball. Random variable = Color ball.
- Sample = Image. Random variable = Color/black-white.
- Sample = Speech. Random variable = Words.
- Sample = text. Random variable = Length optimal compression.
- Sample = weather. Random variable = Temperature.
- Sample = External force. Random variable = Influence.

- $\Pr(\{\omega : X(\omega) = x\}) := P(X = x)$.

- CDF $P(x) := \Pr(X \leq x)$.



- PDF $p(x) := \frac{dP(x)}{dx}$.



- Realizations x vs. random variables X .

- Expectation:

$$\mathbb{E}[X] := \int_{\Omega} x \Pr(dx) = \int_{\Omega} x dP(x) = \int_{\Omega} x p(x) dx$$

- Expectation of $g : \mathbb{R} \rightarrow \mathbb{R}$.

$$\mathbb{E}[g(X)] = \int_{\Omega} g(x) dP(x) = \int_{\Omega} g(x) p(x) dx$$

- Independence:

$$\mathbb{E}[X_1 X_2] = \mathbb{E}[X_1] \mathbb{E}[X_2]$$

Law of Large numbers: if $\{X_i\}_{i=1}^n$ I.I.D.:

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mathbb{E}[X].$$

Gaussian PDF:

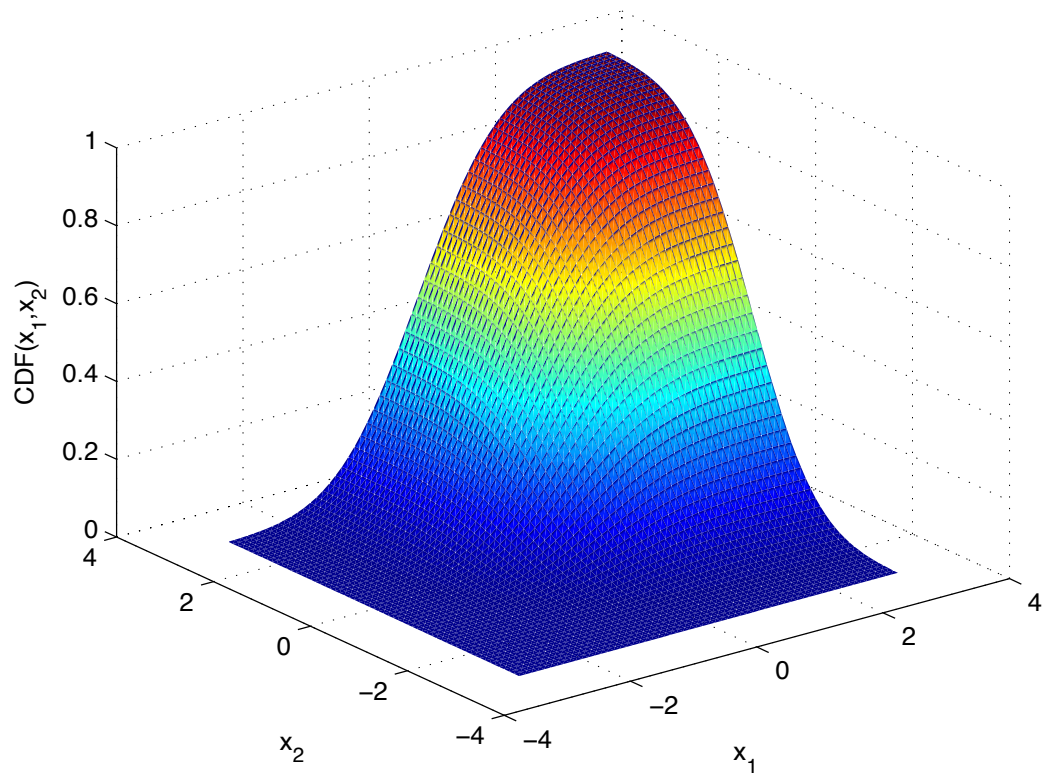
- Gaussian PDF determined by 2 moment: mean μ and variance σ^2

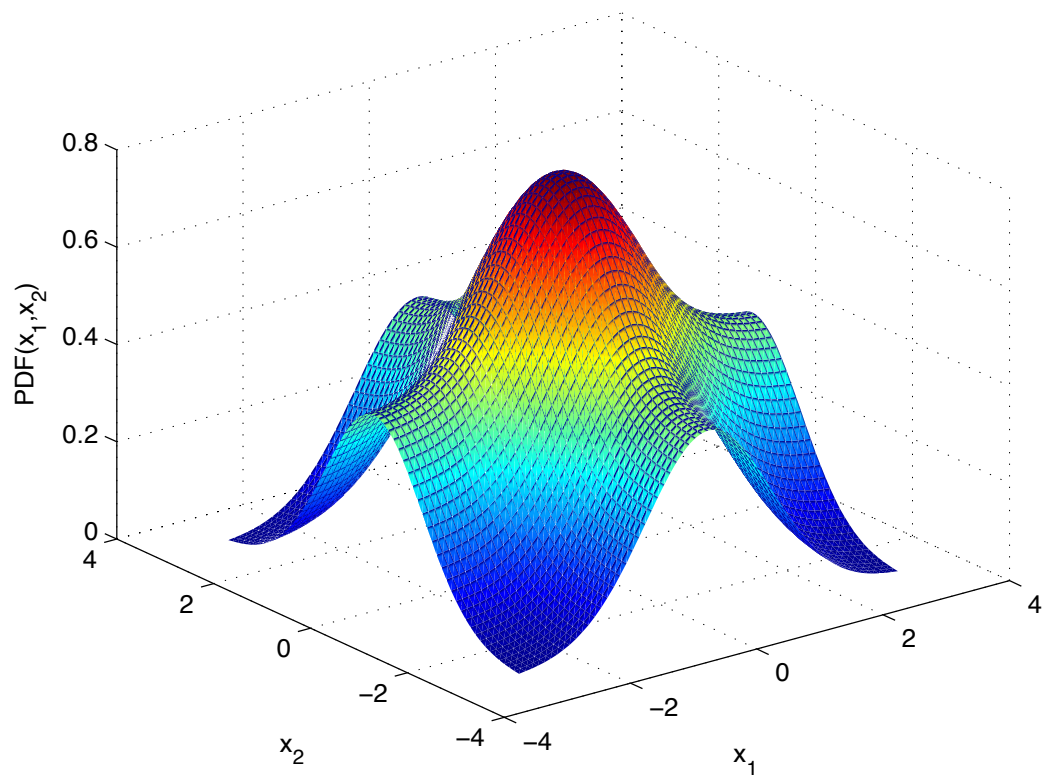
$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) = \mathcal{N}(x; \mu, \sigma)$$

(or standard deviation σ).

- Closed by convolution, conditioning and product.
- Multivariate Normal PDF for $\mathbf{x} \in \mathbb{R}^2$:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu, \Sigma) = \frac{1}{2\pi|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right)$$





- Let bivariate Gaussian with mean μ and variance matrix Σ as

$$p\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \mathcal{N}\left(\mathbf{x}; \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

Then conditional bivariate Gaussian $p(X_1|X_2 = x_2) = \mathcal{N}(\bar{\mu}, \bar{\Sigma})$ where

$$\begin{cases} \bar{\mu} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2) \\ \bar{\Sigma} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \end{cases}$$

- Central Limit Theorem (CLT): if $X_i \sim \mathcal{N}(\mathbb{E}[X], \sigma)$

$$\frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mathbb{E}[X], \frac{\sigma}{\sqrt{n}}\right)$$

Interpretations

"The metallurgist told his friend the statistician how he planned to test the effect of heat on the strength of a metal bar by sawing the bar into six pieces. The first two would go into the hot oven, the next two into the medium oven and the last two into the cool oven. The statistician, horrified, explained how he should randomize in order to avoid the effect of a possible gradient of strength in the metal bar. The method of randomization was applied, and it turned out that the randomized experiment called for putting the first two into the hot oven, the next two into the medium oven and the last two into the cool oven. "Obviously, we can't do that," said the metallurgist. "On the contrary, you have to do that," said the statistician."

Stochastic Processes

- Sequence of random variables indexed by time.
- Joint distribution function.
- Conditional PDF.
- Stationarity.
- Quasi-Stationary.
- Ergodic: different realizations vs. time averages.

Maximum Likelihood

- If the *true* PDF of X were f , then the probability of occurrence of a realization x of X were $p(x)$.
- Consider a class of such PDFs $\{p_\theta : \theta \in \Theta\}$
- Likelihood function: this PDF as a function of an unknown parameter θ

$$L(\theta; x) = p_\theta(x)$$

- Idea: given an observation y_i , find a model under which this sample was most likely to occur.

$$\theta_n = \operatorname{argmax}_{\theta \in \Theta} L(y_i, \theta)$$

- Equivalent to

$$\theta_n = \operatorname{argmax}_{\theta \in \Theta} \log L(y_i; \theta)$$

Properties of θ_n

- Assume that x_n contains n independent samples.
- Consistent, i.e. $\theta_n \rightarrow \theta_0$ with rate $\sqrt{1/n}$.
- Asymptotic normal: if n large

$$f(\sqrt{n}(\theta_n - \theta_0)) = \mathcal{N}(0, 1)$$

- Sufficient.
- Efficient.
- Regularity conditions: identifiability:

$$\exists x : L(x, \theta) \neq L(x, \theta') \Leftrightarrow \theta \neq \theta'$$

Least Squares Revisited

- Observations (realizations) $\{y_i\}_{i=1}^n$ of $\{Y_i\}_{i=1}^n$.

- Model

$$Y_i = \mathbf{x}_i^T \theta + V_i, \quad V_i \sim \mathcal{N}(0, 1)$$

and observations (realizations)

$$y_i = \mathbf{x}_i^T \theta + e_i$$

with $\{e_i\}$ realizations of $\{V_i\}_i$ *i.i.d.* .

- $\{\mathbf{x}_i\}_{i=1}^n$ and θ deterministic.

- Equivalently

$$Y_i \sim \mathcal{N}(\mathbf{x}_i^T \theta, \sigma) \text{ or } p(y_i) = \mathcal{N}(y_i; \mathbf{x}_i^T \theta, \sigma)$$

- Since $\{V_i\}$ independent

$$Y_1, \dots, Y_n \sim \prod_{i=1}^n \mathcal{N}(\mathbf{x}_i^T \theta, \sigma)$$

- or

$$p(y_1, \dots, y_n) = \prod_{i=1}^n \mathcal{N}(y_i; \mathbf{x}_i^T \theta, \sigma)$$

- Maximum Likelihood

$$\theta_n = \operatorname{argmax}_{\theta} \log L(\theta, y_1, \dots, y_n) = \log \prod_{i=1}^n \mathcal{N}(y_i; \mathbf{x}_i^T \theta, \sigma)$$

- = Ordinary Least Squares (OLS).

$$\theta_n = \operatorname{argmin}_{\theta} c \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta)^2$$

- Also when zero mean, uncorrelated errors with equal, bounded variance (Gauss-Markov).
- If noise not independent, but

$$\begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix} \sim \mathcal{N}(0_n, \Sigma)$$

Then ML leads to BLUE (Best Linear Unbiased Estimator)

$$\theta_n = \operatorname{argmin}_{\theta} \mathbf{e}^T \Sigma^{-1} \mathbf{e} = (\mathbf{y} - X\theta)^T \Sigma^{-1} (\mathbf{y} - X\theta)$$

Analysis OLS

Model with $\{V_i\}$ zero mean white noise with $\mathbb{E}[V_i^2] = \lambda^2$ and $\{\theta, \mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$

$$Y_i = \mathbf{x}_i \theta_0 + V_i.$$

$$\text{OLS: } \theta_n = \operatorname{argmin}_{\theta} \|\mathbf{y} - \mathbf{X}\theta\|_2^2.$$

$$\text{Normal Equations: } (\mathbf{X}^T \mathbf{X})\theta = \mathbf{X}^T \mathbf{y}.$$

- Unbiased:

$$\mathbb{E}[\theta_n] = \theta_0.$$

- Covariance:

$$\mathbb{E}[(\theta_n - \theta_0)^T (\theta_n - \theta_0)] = \lambda^2 (\mathbf{X}^T \mathbf{X})^{-1} = \frac{\lambda^2}{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right)^{-1}.$$

- Objective:

$$\mathbb{E} \|\mathbf{y} - \mathbf{X}\theta_n\|_2^2 = \lambda^2 (n - d).$$

- Efficient. $\sqrt{n}(\theta_n - \theta_0) \rightarrow \mathcal{N}(0, I^{-1})$ for $n \rightarrow \infty$.

Cramér-Rao Bound

- Lowerbound to the variance of any unbiased estimator θ_n to θ .
- Discrete setting: D is set of samples, with $\{D_n\}$ denoting the observed datasets.
- Given PDFs $p_\theta > 0$ over all possible datasets D_n such that

$$\sum_{D_n} p_\theta(D_n) = 1$$

- Given estimator $\theta_n(D_n)$ of $\theta \in \Theta$ such that $\forall \theta \in \Theta$

$$\mathbb{E}[\theta_n(D_n)] = \sum_{D_n} \theta_n(D_n) p_\theta(D_n) = \theta$$

- Taking derivative w.r.t. θ gives

$$\begin{cases} \sum_{D_n} \frac{dp_\theta(D_n)}{d\theta} = 0 \\ \sum_{D_n} \theta_n(D_n) \frac{dp_\theta(D_n)}{d\theta} = 1 \end{cases}$$

and hence

$$1 = \sum_{D_n} (\theta - \theta_n(D_n)) \frac{dp_\theta(D_n)}{d\theta}$$

- Combining:

$$1 = \sum_{D_n} (\theta - \theta_n(D_n)) \left(\frac{p_\theta(D_n)}{p_\theta(D_n)} \right) \left(\frac{dp_\theta(D_n)}{d\theta} \right)$$

$$1 = \sum_{D_n} (\theta - \theta_n(D_n)) \sqrt{p_\theta(D_n)} \left(\frac{\sqrt{p_\theta(D_n)} dp_\theta(D_n)}{p_\theta(D_n) d\theta} \right)$$

- Cauchy-Schwarz ($\mathbf{a}^T \mathbf{b} \leq \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$) gives

$$1 \leq \sum_{D_n} (\theta - \theta_n(D_n))^2 p_\theta(D_n) \sum_{D_n} \left(\frac{\sqrt{p_\theta(D_n)} dp_\theta(D_n)}{p_\theta(D_n) d\theta} \right)^2$$

Or

$$\mathbb{E}_\theta[\theta - \theta_n(D_n)]^2 \geq \frac{1}{I(\theta)}$$

with

$$I(\theta) = \mathbb{E}_\theta \left[\frac{dp_\theta(D_n)}{d\theta} \frac{1}{p_\theta(D_n)} \right]^2$$

Dynamical System

Describe data as coming from a ARX(1,1) model

$$Y_t + aY_{t-1} = bu_{t-1} + V_t, \quad \forall t.$$

where $\{V_t\}$ zero mean and unit variance white noise. Then application of OLS gives

- Normal equations of $\theta = (-a, b)^T$

$$\left(\begin{bmatrix} \sum_{t=1}^n Y_{t-1}^2 & \sum_{t=1}^n Y_{t-1}u_{t-1} \\ \sum_{t=1}^n u_{t-1}Y_{t-1} & \sum_{t=1}^n u_{t-1}^2 \end{bmatrix} \right) \theta_n = \begin{bmatrix} \sum_{t=1}^n Y_{t-1}y_t \\ \sum_{t=1}^n u_{t-1}Y_t \end{bmatrix}.$$

- Taking inverse.
- Taking Expectations

$$\mathbb{E}[\theta_n] = \mathbb{E} \left[\begin{bmatrix} \sum_{t=1}^n Y_{t-1}^2 & \sum_{t=1}^n Y_{t-1}u_{t-1} \\ \sum_{t=1}^n u_{t-1}Y_{t-1} & \sum_{t=1}^n u_{t-1}^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum_{t=1}^n Y_{t-1}Y_t \\ \sum_{t=1}^n u_{t-1}Y_t \end{bmatrix} \right]$$

- Approximatively

$$\mathbb{E}[\theta_n] \approx \begin{bmatrix} \mathbb{E}[Y_t^2] & \mathbb{E}[Y_t u_t] \\ \mathbb{E}[u_t Y_t] & \mathbb{E}[u_t^2] \end{bmatrix}^{-1} \begin{bmatrix} \mathbb{E}[Y_t y_{t-1}] \\ \mathbb{E}[Y_t u_{t-1}] \end{bmatrix} = \mathbf{R}^{-1} \mathbf{r}.$$

- Ill-conditioning.

Instrumental Variables

- LS: $\min \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \theta)^2$
- Normal equations

$$\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \theta = \left(\sum_{i=1}^n \mathbf{x}_i Y_i \right)$$

- Or

$$\sum_{i=1}^n \mathbf{x}_i (Y_i - \mathbf{x}_i^T \theta) = 0_d$$

- Interpretation as orthogonal projection.
- Idea:

$$\sum_{i=1}^n Z_i (Y_i - \mathbf{x}_i^T \theta) = 0_d$$

- Choose $\{Z_t\}$ taking values in \mathbb{R}^d such that
 - $\{Z_t\}$ independent from $\{V_t\}$
 - \mathbf{R} full rank.
- Example. Past input signals.

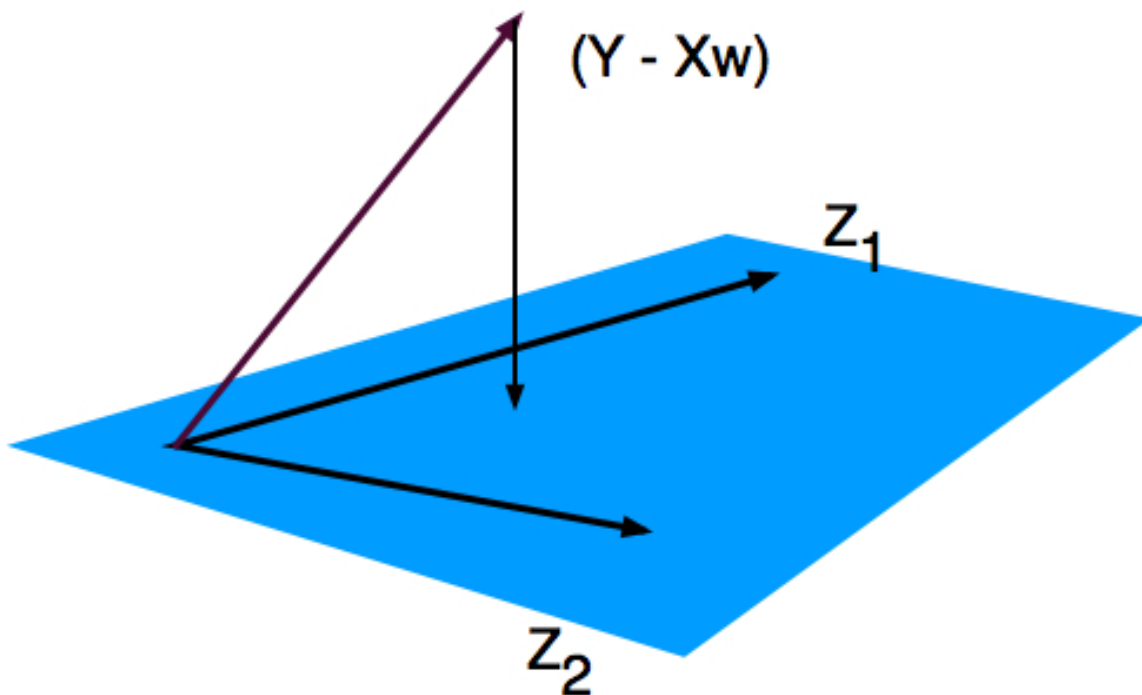


Figure 1: Instrumental Variable as Modified Projection.

Conclusions

- Stochastic (Theory) and Statistical (Data).
- Maximum Likelihood.
- Least Squares.
- Instrumental Variables.