# Astrostatistics: Tue 05 Feb 2017

https://github.com/CambridgeAstroStat/PartIII-Astrostatistics

- Examples Classes (sheets provided ~1 week prior)

  - Fri Feb 16, Fri Mar 2, Wed Mar 14 (1pm, RoomTBD)

  - One more + Revision Class in Easter Term

- Fitting Statistical Models to Astronomical Data

  - Generative / Latent Variable Modeling / Bayes

  - Hogg, Bovy & Lang. "Data analysis recipes: Fitting a model to data". https://arxiv.org/abs/1008.4686

# Statistical Modelling Wisdom

- Have an objective function [e.g. Likelihood or posterior] that you optimise or sample to fit the data - not just a procedure/recipe

- Objective function helps you evaluate relative fits of data with under different parameter values / models

- Derive your objective function from your modelling assumptions (physical or statistical)

- Write down your assumptions!

- First question: what is the likelihood $L(\theta)$ ? Derive it from the assumptions underlying your sampling distribution $P(D \mid \theta)$!

- Second question: what is your prior $P(\theta)$ ? (if Bayesian)

- Third question: How do I optimise/sample objective function to fit the data?
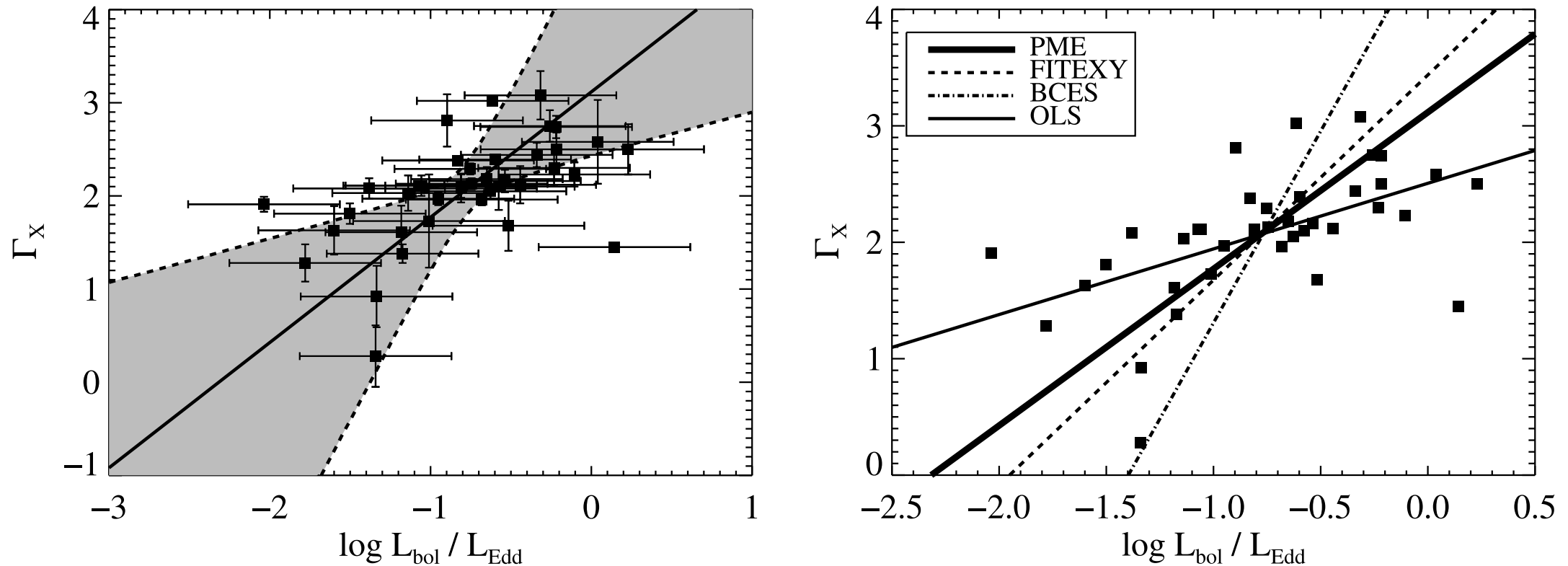
# Fitting Models to Astro Data



FIG. 10.—X-ray photon index $\Gamma_X$ as a function of $\log L_{bol}/L_{Edd}$ for 39 $z \lesssim 0.8$ radio-quiet quasars. In both plots, the thick solid line shows the posterior median estimate (PME) of the regression line. In the left panel, the shaded region denotes the 95% (2 $\sigma$) pointwise confidence intervals on the regression line. In the right panel, the thin solid line shows the OLS estimate, the dashed line shows the FITEXY estimate, and the dot-dashed line shows the BCES($Y|X$) estimate; the error bars have been omitted for clarity. A significant positive trend is implied by the data.

Modelling heteroskedastic, correlated measurement errors in both y and x, intrinsic scatter, nondetections, selection effects

B. Kelly et al. 2007, "Some Aspects of Measurement Error in Linear Regression of Astronomical Data." ApJ, 665, 1489

# Ad-hoc "$\chi^2$" approaches vs. Likelihood formulation

# FITEXY Estimator

- Press et al.(1992, *Numerical Recipes*) define an 'effective $\chi^2$' statistic:

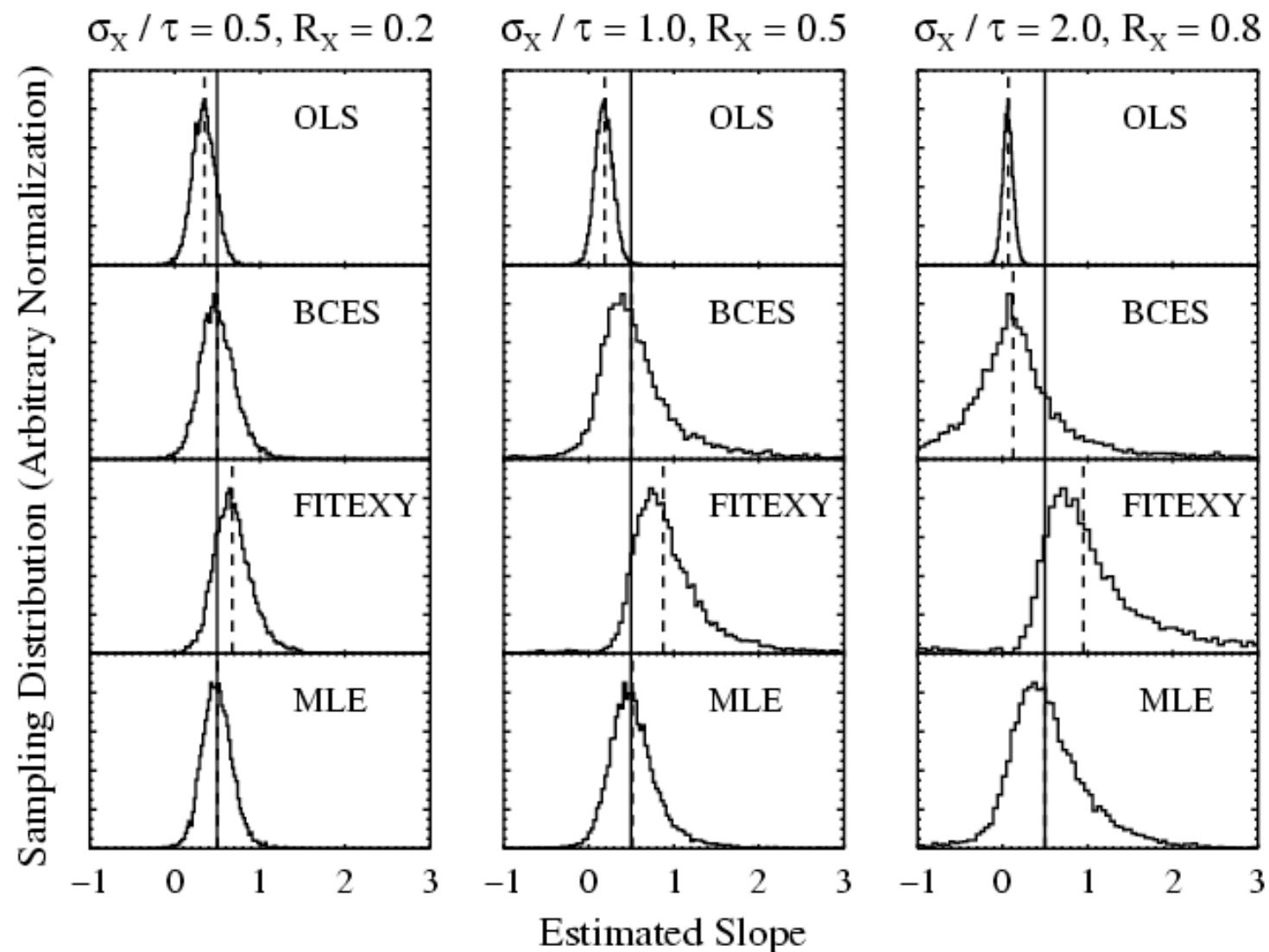$$\chi^2_{EXY} = \sum_{i=1}^{n} \frac{(y_i - \alpha - \beta x_i)^2}{\sigma^2_{y,i} + \beta^2 \sigma^2_{x,i}}$$

- Choose values of $\alpha$ and $\beta$ that minimize $\chi^2_{EXY}$

- Modified by Tremaine et al.(2002, ApJ, 574, 740), to account for intrinsic scatter:

$$\chi^2_{EXY} = \sum_{i=1}^{n} \frac{(y_i - \alpha - \beta x_i)^2}{\sigma^2 + \sigma^2_{y,i} + \beta^2 \sigma^2_{x,i}}$$

http://astrostatistics.psu.edu/su07/kelley_measerr07.pdf

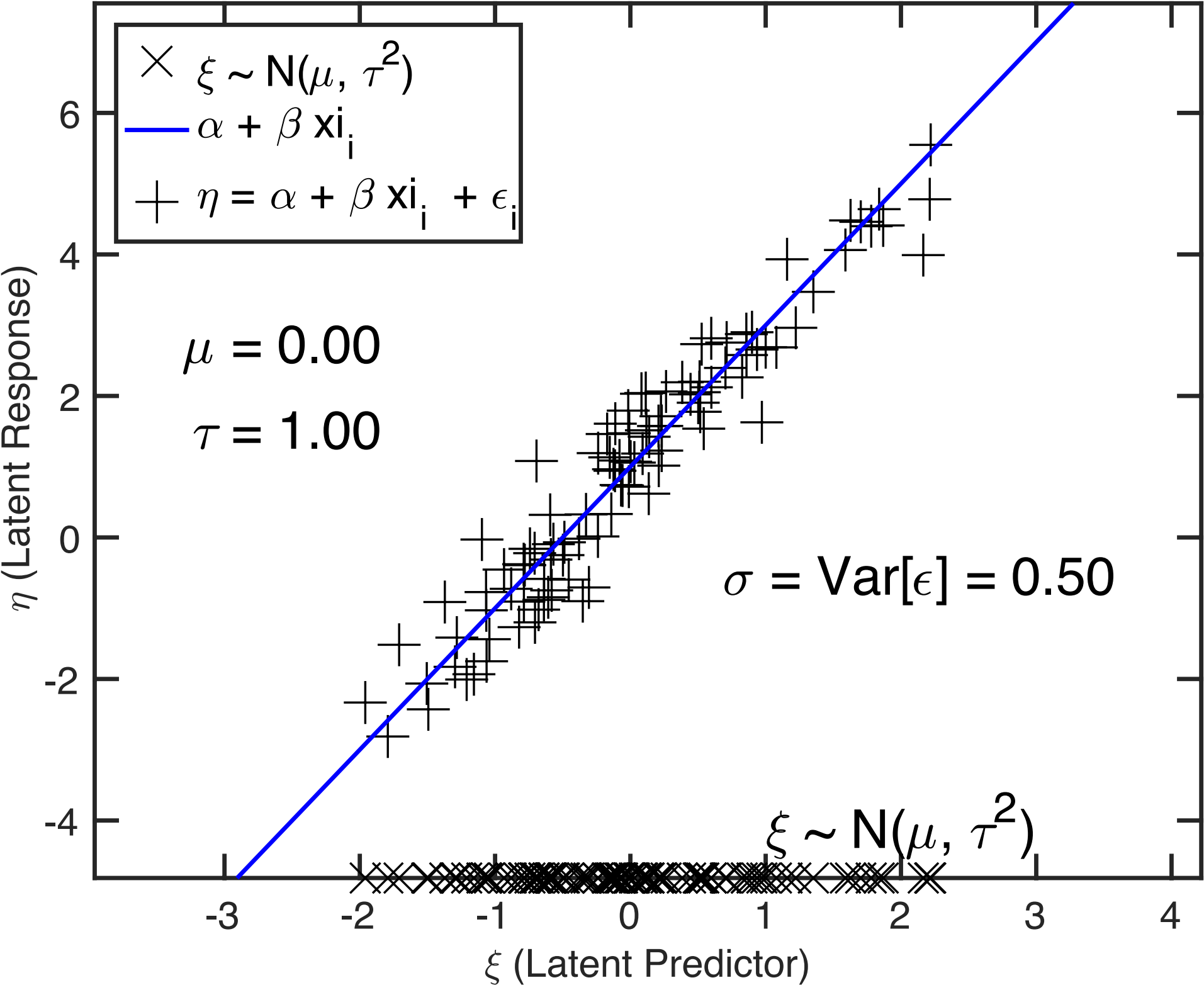# Kelly et al. 2017, Latent Variable Likelihood approach vs. Bad

## Simulation Study: Slope



Dashed lines mark the median value of the estimator, solid lines mark the true value of the slope. Each simulated data set had 50 data points, and y-measurement errors of $\sigma_y \sim \sigma$.

http://astrostatistics.psu.edu/su07/kelley_measerr07.pdf

**Latent Variable Model : N = 100**

$\times$   $\xi \sim N(\mu, \tau^2)$

——   $\alpha + \beta \, xi_i$

$+$   $\eta = \alpha + \beta \, xi_i + \epsilon_i$

$\mu = 0.00$

$\tau = 1.00$

$\sigma = \mathrm{Var}[\epsilon] = 0.50$

$\xi \sim N(\mu, \tau^2)$

$\eta$ (Latent Response)
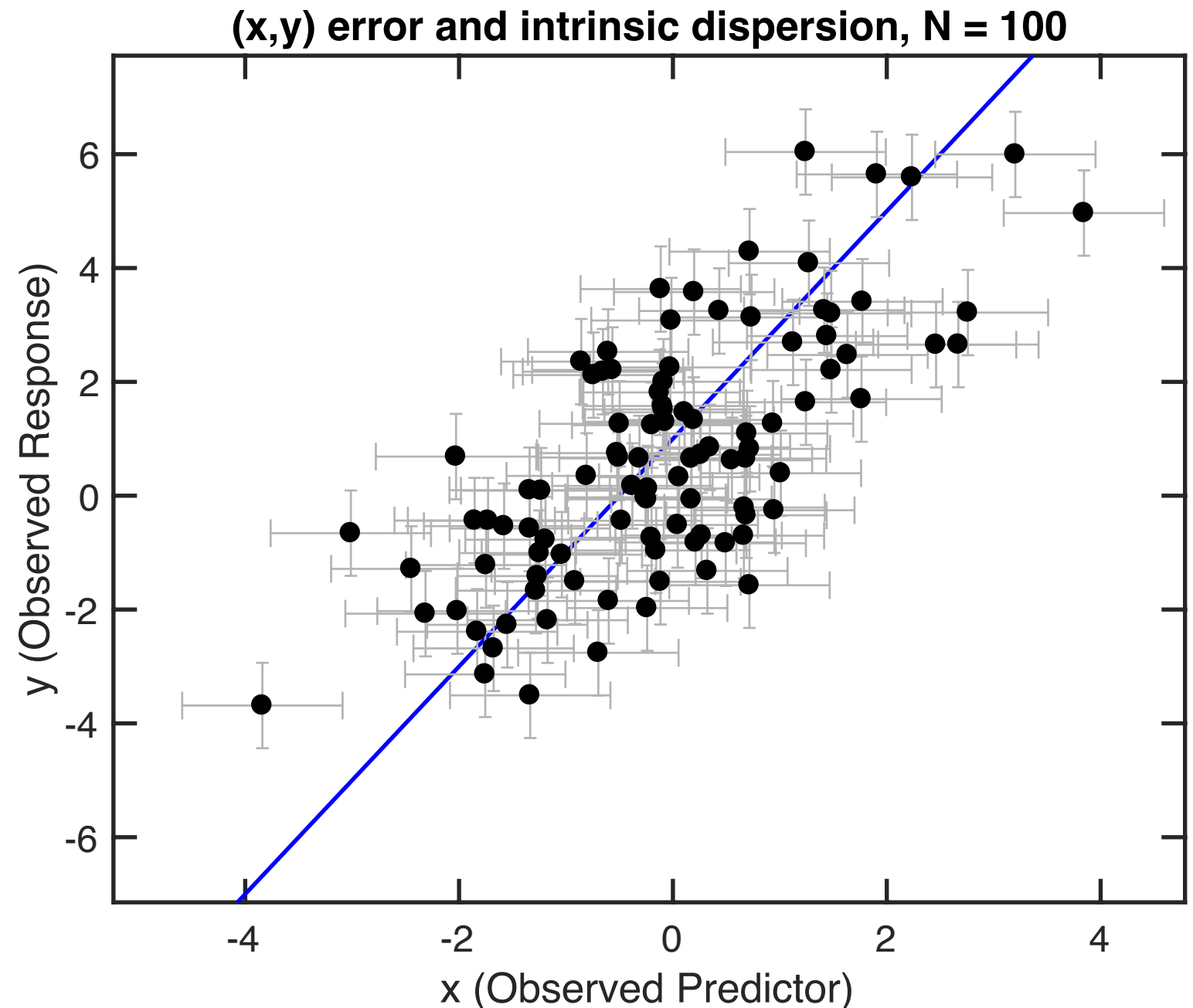
$\xi$ (Latent Predictor)

# Probabilistic Generative Modelling

- Forward Model comprises series of probabilistic steps describing conceptually how the observed data was generated from the parameters of interest

- Can introduce intermediate parameters / unobserved latent variables $\boldsymbol{\alpha}$ (e.g. true values corresponding to the observed data).

- From Forward model, derive the sampling distribution, e.g.
  $P(D \mid \theta) = \int P(D \mid \boldsymbol{\alpha}) \, P(\boldsymbol{\alpha} \mid \theta) \, d\boldsymbol{\alpha}$

- Using observed data D, draw inference from Likelihood function:
  $L(\theta) = P(D \mid \theta)$

- Or if Bayesian with prior $P(\theta)$: sample posterior:
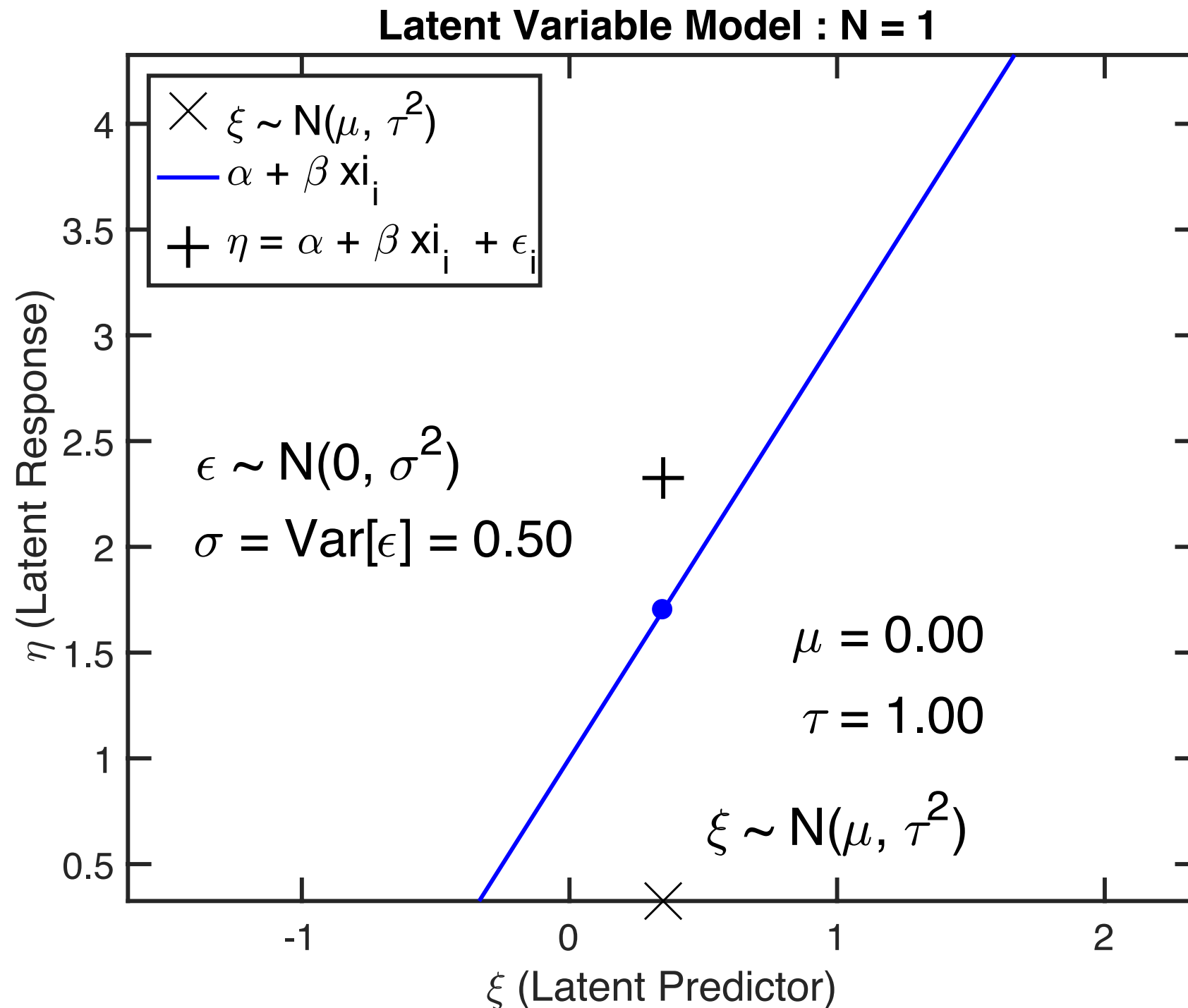  $P(\theta \mid D) = P(D \mid \theta) \, P(\theta)$

# Example: Structural Model for Linear Regression
## (B. Kelly et al. 2007, "Some Aspects of Measurement Error in Linear Regression of Astronomical Data." ApJ, 665, 1489)

- Observed data has x and y meas. errors and intrinsic dispersion

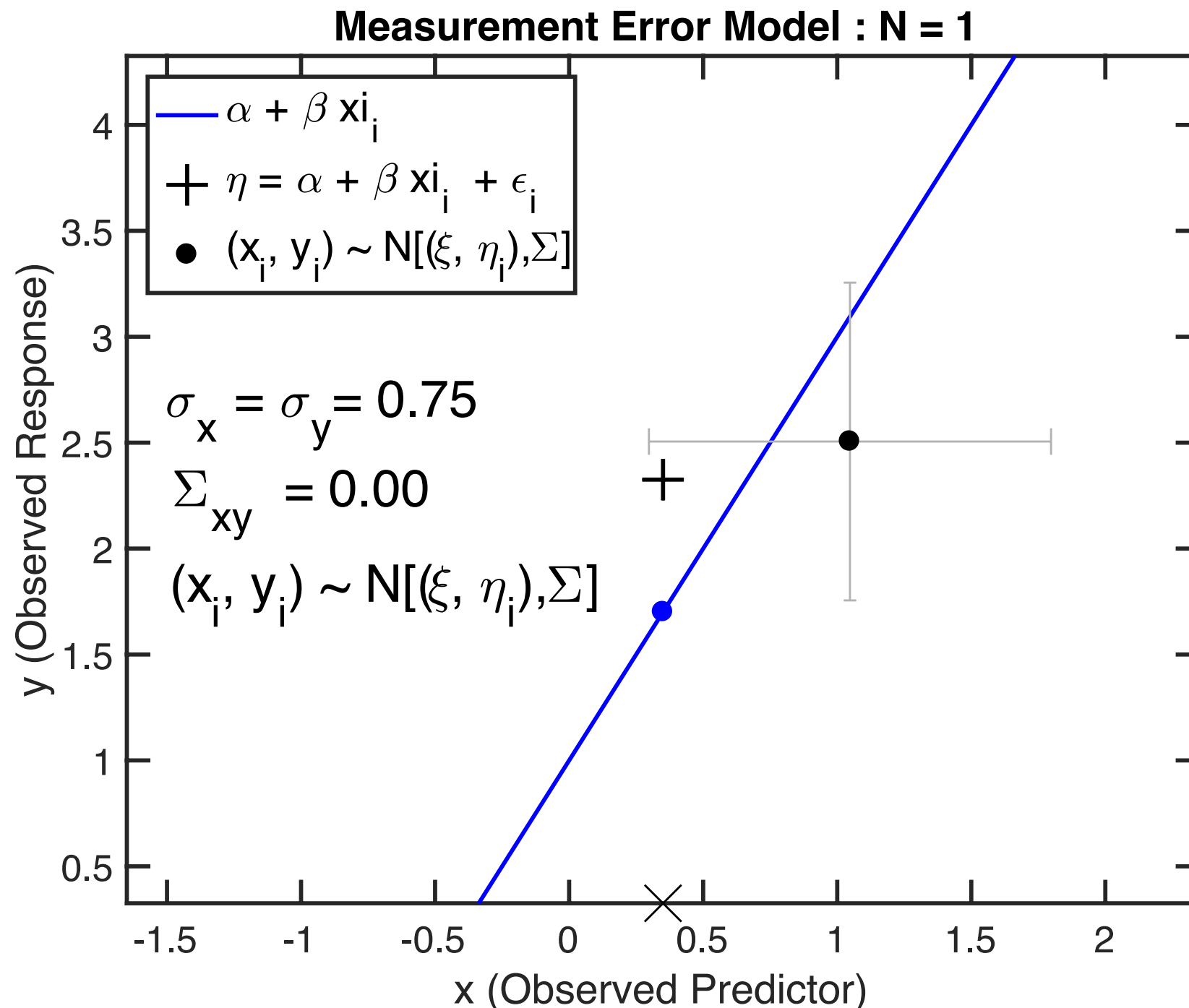- Estimate the true slope (and other parameters)

**(x,y) error and intrinsic dispersion, N = 100**

# Step 1:
# Generating Latent Variables from Parameters:

$$P(\eta_i, \xi_i | \alpha, \beta, \sigma, \mu, \tau) = P(\eta_i | \xi_i, \alpha, \beta, \sigma) \times P(\xi_i | \mu, \tau)$$

**Latent Variable Model : N = 1**



Legend:
- $\times$   $\xi \sim N(\mu, \tau^2)$
- ——   $\alpha + \beta\, \text{xi}_i$
- $+$   $\eta = \alpha + \beta\, \text{xi}_i + \epsilon_i$

$\epsilon \sim N(0, \sigma^2)$

$\sigma = \text{Var}[\epsilon] = 0.50$

$\mu = 0.00$

$\tau = 1.00$

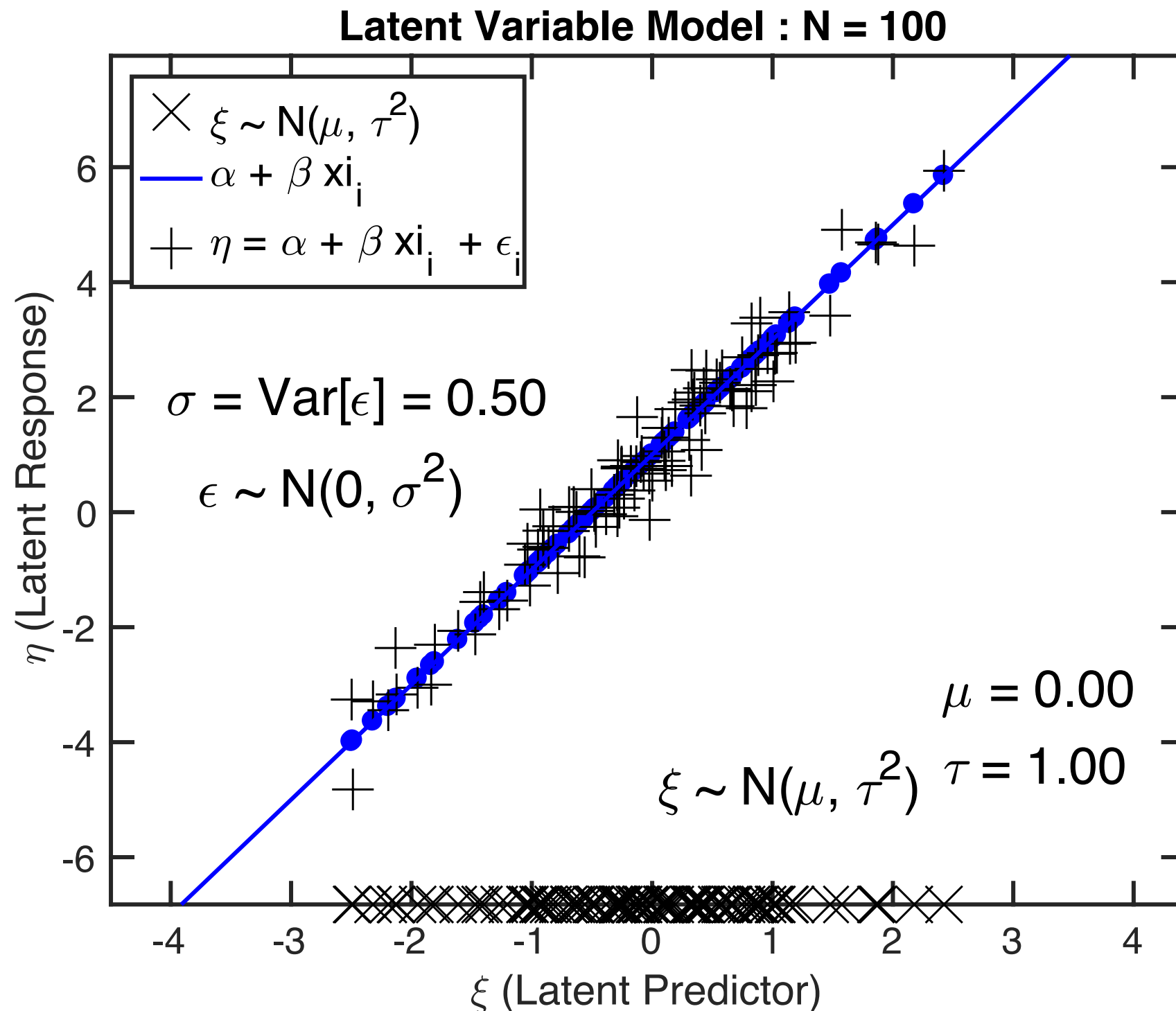$\xi \sim N(\mu, \tau^2)$

x-axis: $\xi$ (Latent Predictor)

y-axis: $\eta$ (Latent Response)

# Step 2: Generating Observed Data from Latent Variables

$$P([x_i, y_i]|\eta_i, \xi_i) = N([x_i, y_i]|[\eta_i, \xi_i], \mathbf{\Sigma})$$



**Measurement Error Model : N = 1**

Legend:
- $\alpha + \beta\, \text{xi}_i$
- $+$ : $\eta = \alpha + \beta\, \text{xi}_i + \epsilon_i$
- $\bullet$ : $(x_i, y_i) \sim N[(\xi, \eta_i), \Sigma]$

$\sigma_x = \sigma_y = 0.75$

$\Sigma_{xy} = 0.00$

$(x_i, y_i) \sim N[(\xi, \eta_i), \Sigma]$

y (Observed Response)

x (Observed Predictor)

# Now repeat for N=100 objects



**Latent Variable Model : N = 100**

# Now repeat for N=100 objects



**Measurement Error Model : N = 100**

Legend:
- $\alpha + \beta \, \mathsf{xi}_i$
- $\eta = \alpha + \beta \, \mathsf{xi}_i + \epsilon_i$
- $(x_i, y_i) \sim N[(\xi, \eta_i), \Sigma]$

$\sigma_x = \sigma_y = 0.75$

$\Sigma_{xy} = 0.00$

$(x_i, y_i) \sim N[(\xi, \eta_i), \Sigma]$

y (Observed Response)

x (Observed Predictor)

# Knowns and Unknowns

Regression Parameters

$$\boldsymbol{\theta} = (\alpha, \beta, \sigma^2)$$

Independent Variable
Population Distribution
"Hyperparameters"

$$\boldsymbol{\psi} = (\mu, \tau)$$

Latent (true) Variables

$$(\xi_i, \eta_i)$$

Observed Data

$$(x_i, y_i)$$

# Generative Model

Population
Distribution

$$\xi \sim N(\mu | \tau^2)$$

Regression

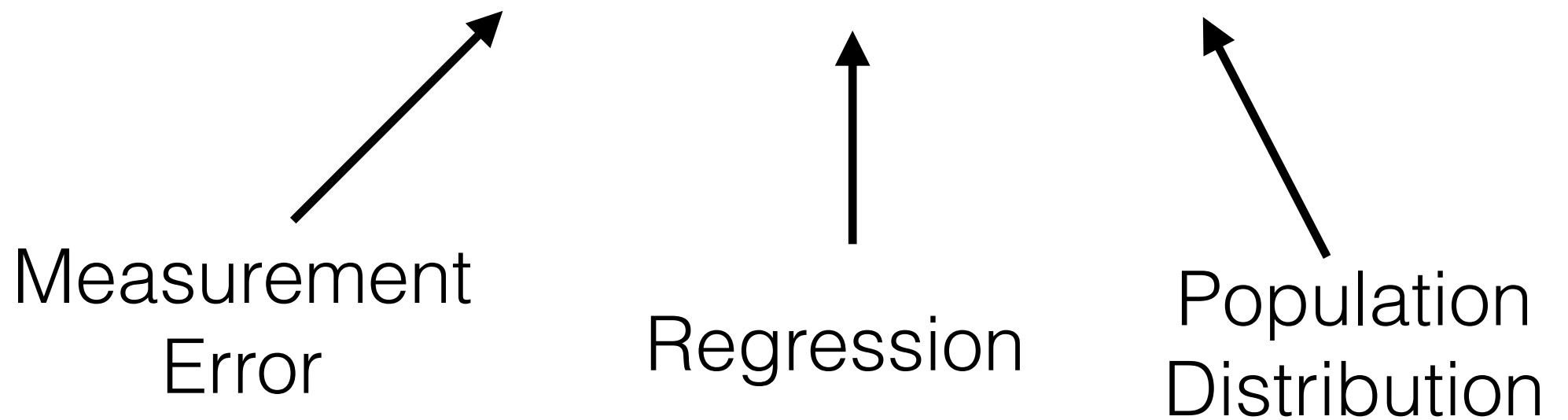$$\eta_i | \xi_i \sim N(\alpha + \beta x_i, \sigma^2)$$

Measurement
Error

$$[x_i, y_i] | \xi_i, \eta_i \sim N([\xi_i, \eta_i], \Sigma])$$

# Formulating Likelihood Function:
## Marginalising (integrating out) latent variables

$$P(x_i, y_i | \boldsymbol{\theta}, \boldsymbol{\psi}) = \int \int P(x_i, y_i, \xi_i, \eta_i | \boldsymbol{\theta}, \boldsymbol{\psi}) d\xi_i \, d\eta,$$

$$P(x_i, y_i | \boldsymbol{\theta}, \boldsymbol{\psi}) = \int \int P(x_i, y_i | \xi_i, \eta_i) P(\eta_i | \xi_i, \boldsymbol{\theta}) P(\xi_i | \boldsymbol{\psi}) \, d\xi_i \, d\eta$$

Measurement Error

Regression

Population Distribution
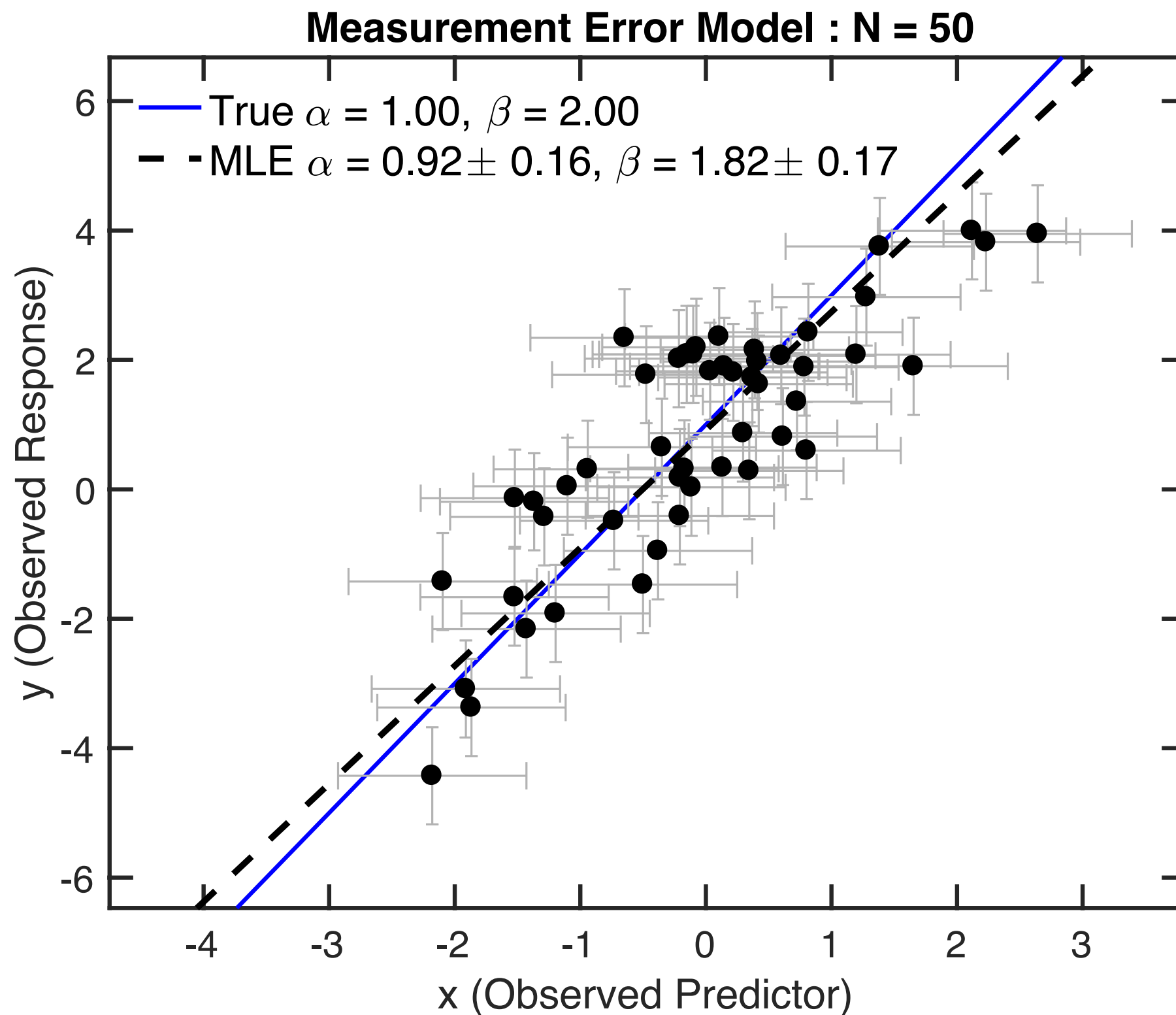
# Solution: (Kelly 2007, Eqs. 16-23)

Gave More General Solution when P(ξ|Ψ)
is a Mixture of Gaussians
(set K=1, $\pi_1 = 1$ for us)

$$p(\boldsymbol{x}, \boldsymbol{y} | \boldsymbol{\theta}, \psi) = \prod_{i=1}^{n} \sum_{k=1}^{K} \frac{\pi_k}{2\pi |\mathbf{V}_{k,i}|^{1/2}}$$

$$\times \exp\left[-\frac{1}{2}(\boldsymbol{z}_i - \boldsymbol{\zeta}_k)^T \mathbf{V}_{k,i}^{-1}(\boldsymbol{z}_i - \boldsymbol{\zeta}_k)\right], \quad (16)$$

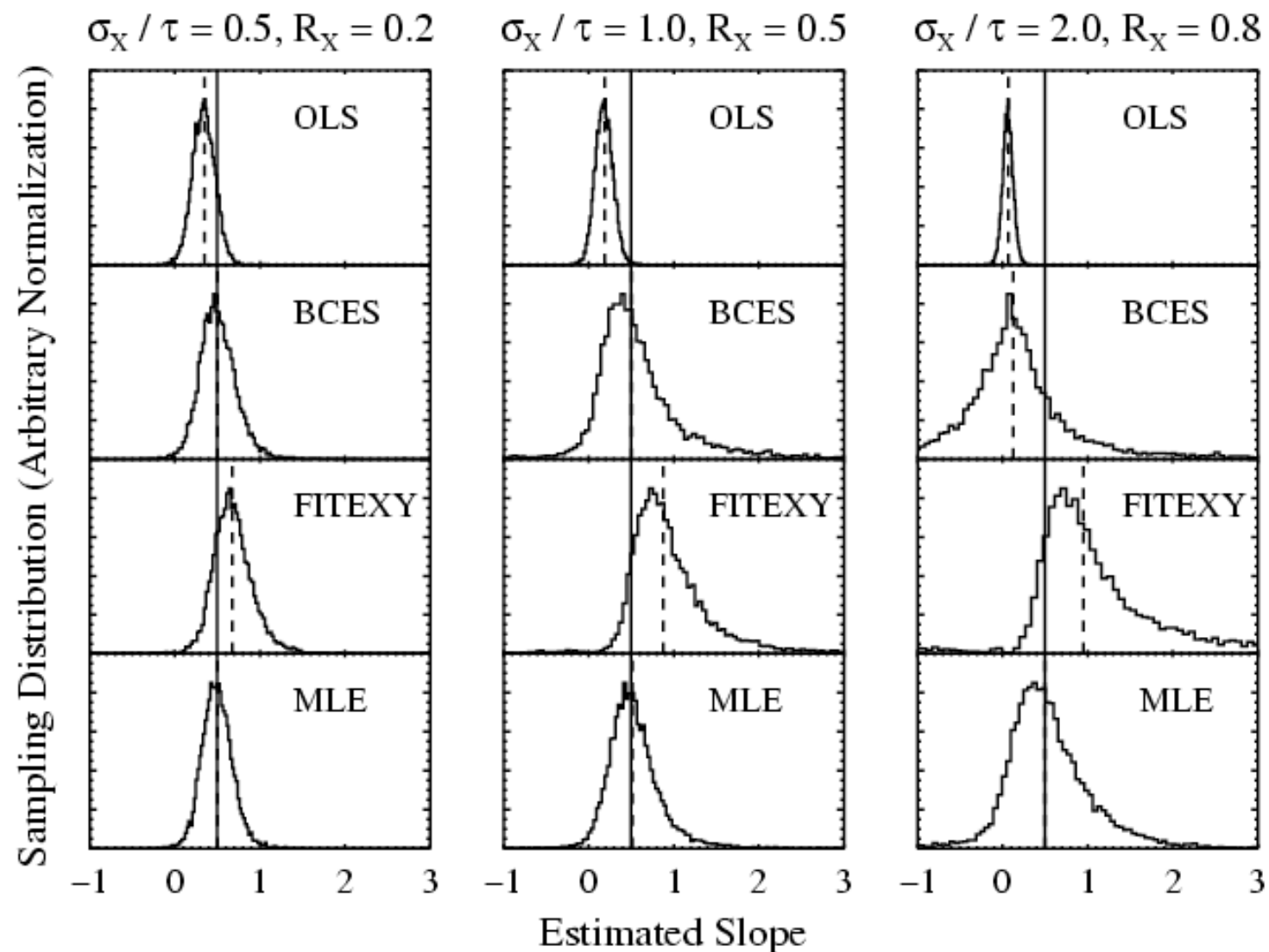$$\boldsymbol{\zeta}_k = (\alpha + \beta\mu_k, \mu_k), \quad (17)$$

$$\mathbf{V}_{k,i} = \begin{pmatrix} \beta^2\tau_k^2 + \sigma^2 + \sigma_{y,i}^2 & \beta\tau_k^2 + \sigma_{xy,i} \\ \beta\tau_k^2 + \sigma_{xy,i} & \tau_k^2 + \sigma_{x,i}^2 \end{pmatrix}, \quad (18)$$

# Example



**Measurement Error Model : N = 50**

True $\alpha$ = 1.00, $\beta$ = 2.00

MLE $\alpha$ = 0.92$\pm$ 0.16, $\beta$ = 1.82$\pm$ 0.17

x (Observed Predictor)

y (Observed Response)

# Kelly et al. 2017, Latent Variable Likelihood approach vs. Bad

## Simulation Study: Slope



Dashed lines mark the median value of the estimator, solid lines mark the true value of the slope. Each simulated data set had 50 data points, and y-measurement errors of $\sigma_y \sim \sigma$.

http://astrostatistics.psu.edu/su07/kelley_measerr07.pdf