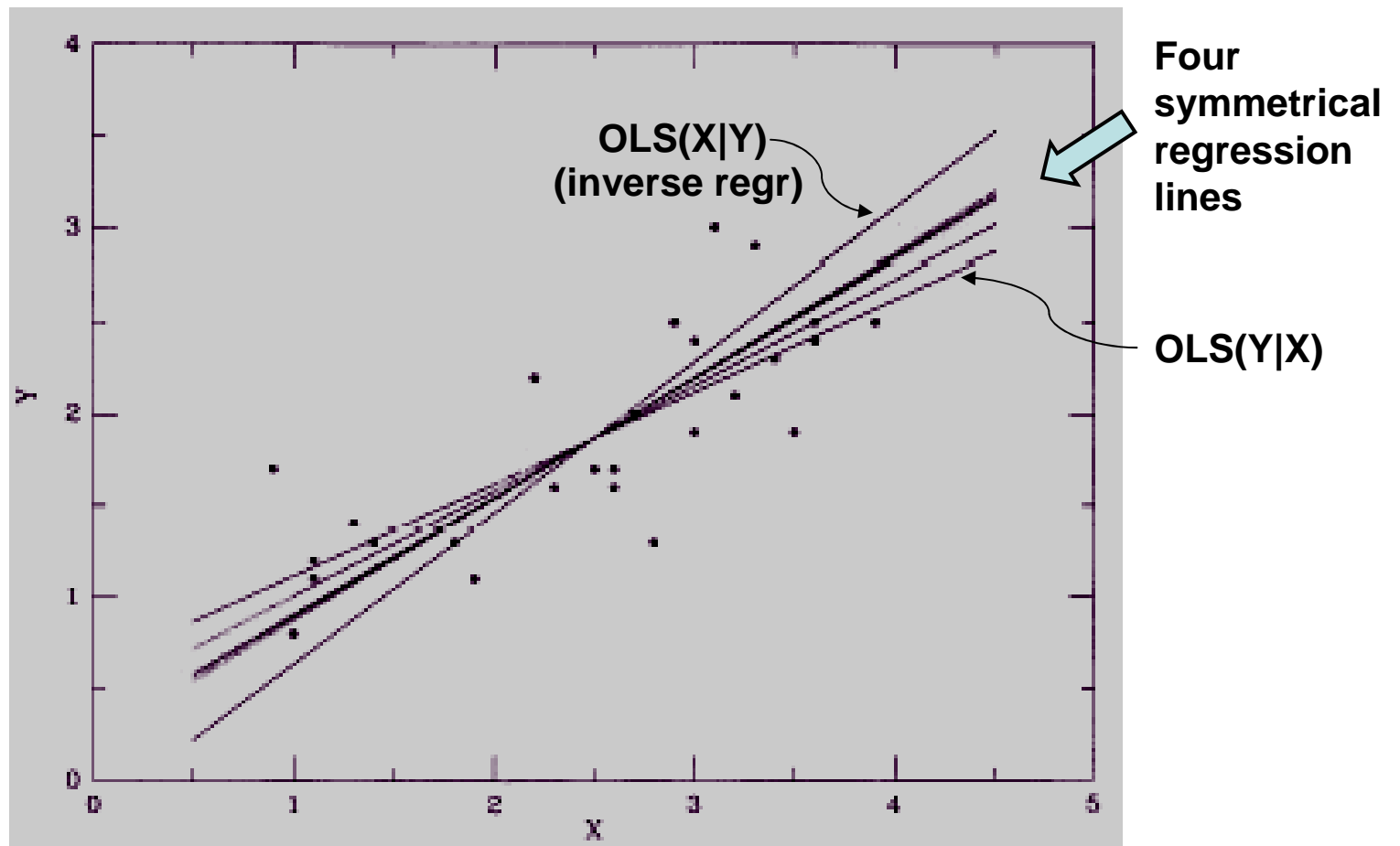# Linear regression issues in astronomy

## G. Jogesh Babu
Center for Astrostatistics

# Structural regression
## Seeking the intrinsic relationship between two properties



Four symmetrical regression lines

OLS(X|Y) (inverse regr)

OLS(Y|X)

$$E(Y)=\alpha+\beta E(X)$$

# Analytical formulae for slopes of the 6 OLS lines

## TABLE 1
### LINEAR REGRESSION FORMULAE FOR SLOPES

| Method | Expression for Slope | Estimate of the Variance of the Slope $\widehat{\text{Var}}\,(\beta_i)$ |
|---|---|---|
| OLS($X\mid Y$) .................. | $\beta_1 = \dfrac{S_{xy}}{S_{xx}}$ | $\dfrac{1}{S_{xx}^2}\left[\displaystyle\sum_{i=1}^{n}(x_i-\bar{x})^2(y_i-\beta_1 x_i-\bar{y}+\beta_1\bar{x})^2\right]$ |
| OLS($Y\mid X$) .................. | $\beta_2 = \dfrac{S_{yy}}{S_{xy}}$ | $\dfrac{1}{S_{xy}^2}\left[\displaystyle\sum_{i=1}^{n}(y_i-\bar{y})^2(y_i-\beta_2 x_i-\bar{y}+\beta_2\bar{x})^2\right]$ |
| OLS bisector .................. | $\beta_3 = (\beta_1+\beta_2)^{-1}[\beta_1\beta_2-1+\sqrt{(1+\beta_1^2)(1+\beta_2^2)}]$ | $\dfrac{\beta_3^2}{(\beta_1+\beta_2)^2(1+\beta_1^2)(1+\beta_2^2)}[(1+\beta_2^2)^2\,\widehat{\text{Var}}\,(\beta_1)$ $+\,2(1+\beta_1^2)(1+\beta_2^2)\,\widehat{\text{Cov}}\,(\beta_1,\beta_2)+(1+\beta_1^2)^2\,\widehat{\text{Var}}\,(\beta_2)]$ |
| Orthogonal regression ....... | $\beta_4 = \tfrac{1}{2}[(\beta_2-\beta_1^{-1})+\text{Sign}\,(S_{xy})\sqrt{4+(\beta_2-\beta_1^{-1})^2}]$ | $\dfrac{\beta_4^2}{4\beta_1^2+(\beta_1\beta_2-1)^2}[\beta_1^{-2}\,\widehat{\text{Var}}\,(\beta_1)+2\,\widehat{\text{Cov}}\,(\beta_1,\beta_2)+\beta_1^2\,\widehat{\text{Var}}\,(\beta_2)]$ |
| Reduced major-axis .......... | $\beta_5 = \text{Sign}\,(S_{xy})(\beta_1\beta_2)^{1/2}$ | $\dfrac{1}{4}\left[\dfrac{\beta_2}{\beta_1}\,\widehat{\text{Var}}\,(\beta_1)+2\,\widehat{\text{Cov}}\,(\beta_1,\beta_2)+\dfrac{\beta_1}{\beta_2}\,\widehat{\text{Var}}\,(\beta_2)\right]$ |

NOTE.—An estimate of covariance term is given by:

$$\widehat{\text{Cov}}\,(\beta_1,\beta_2) = (\beta_1 S_{xx}^2)^{-1}\left\{\sum_{i=1}^{n}(x_i-\bar{x})(y_i-\bar{y})[y_i-\bar{y}-\beta_1(x_i-\bar{x})][y_i-\bar{y}-\beta_2(x_i-\bar{x})]\right\}.$$

**Isobe, Feigelson, Akritas & Babu, ApJ 364, 105 1990**

# Comments

- Standard estimates of variances of slopes are valid strictly under a very restrictive assumption: errors are independent of X values

- The estimates are valid even when this condition is violated. These are derived using the so called `delta method'
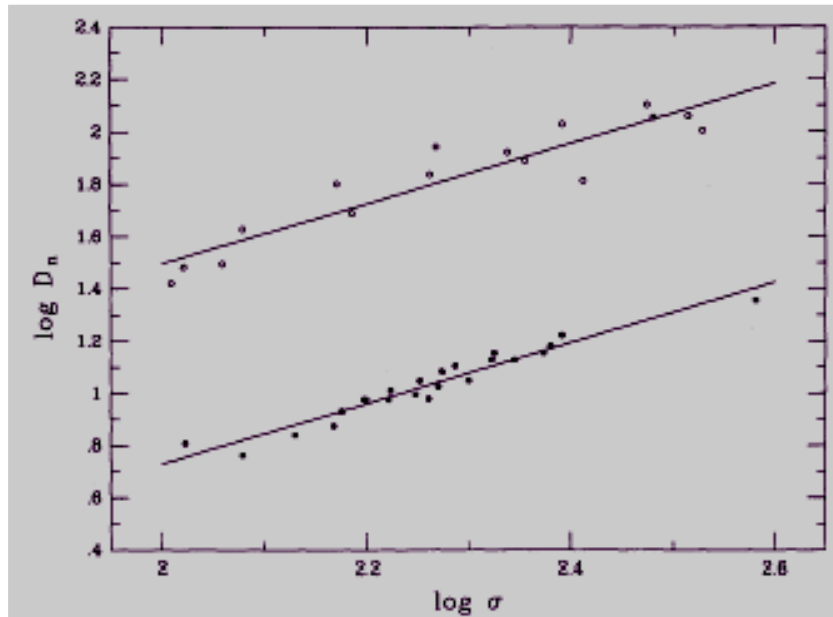
# Relations among the slopes

Suppose $S_{XY} > 0$

- If $\beta_5 < 1$, then
$$\beta_3 \leq 1 \text{ and } \beta_1 \leq \beta_4 \leq \beta_5 \leq \beta_3 \leq \beta_2 \,.$$

- If $\beta_5 > 1$, then
$$\beta_3 \geq 1 \text{ and } \beta_1 \leq \beta_3 \leq \beta_5 \leq \beta_4 \leq \beta_2 \,.$$

- If $\beta_5 = 1$, then
$$\beta_3 = \beta_4 = \beta_5 \,.$$

$\beta_5$ is the slope of the reduced major axis

- Feigelson & Babu, ApJ 397, p.55, 1992

**Example: Faber-Jackson relation between diameter and stellar velocity dispersion of elliptical galaxies**

### TABLE 4

REGRESSIONS FOR COMA AND VIRGO $\log D_n'$ VERSUS $\log \sigma^\bullet$

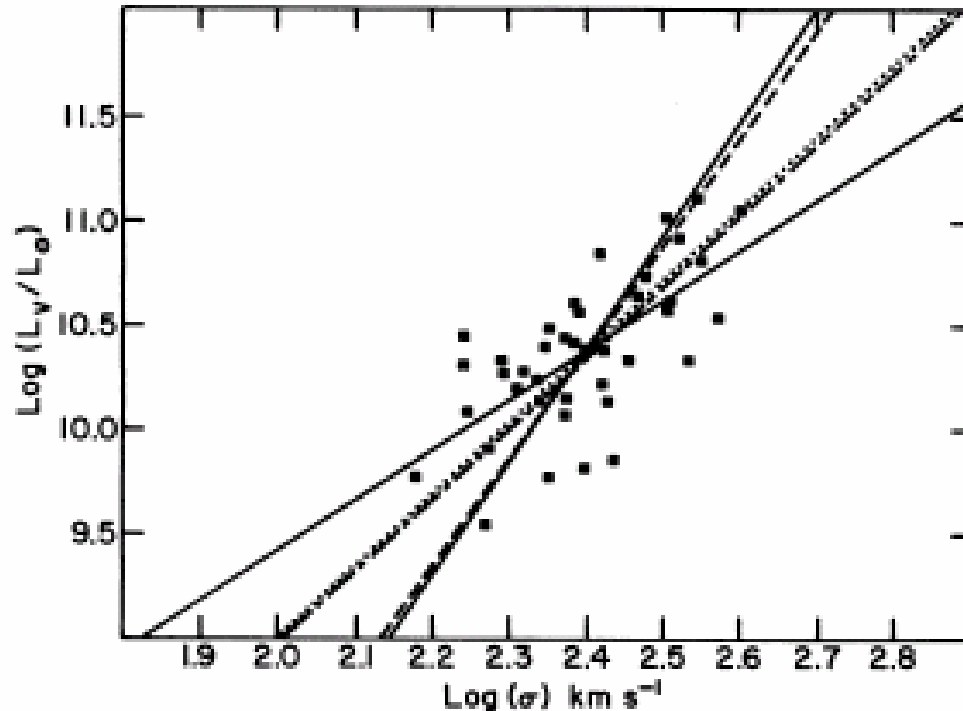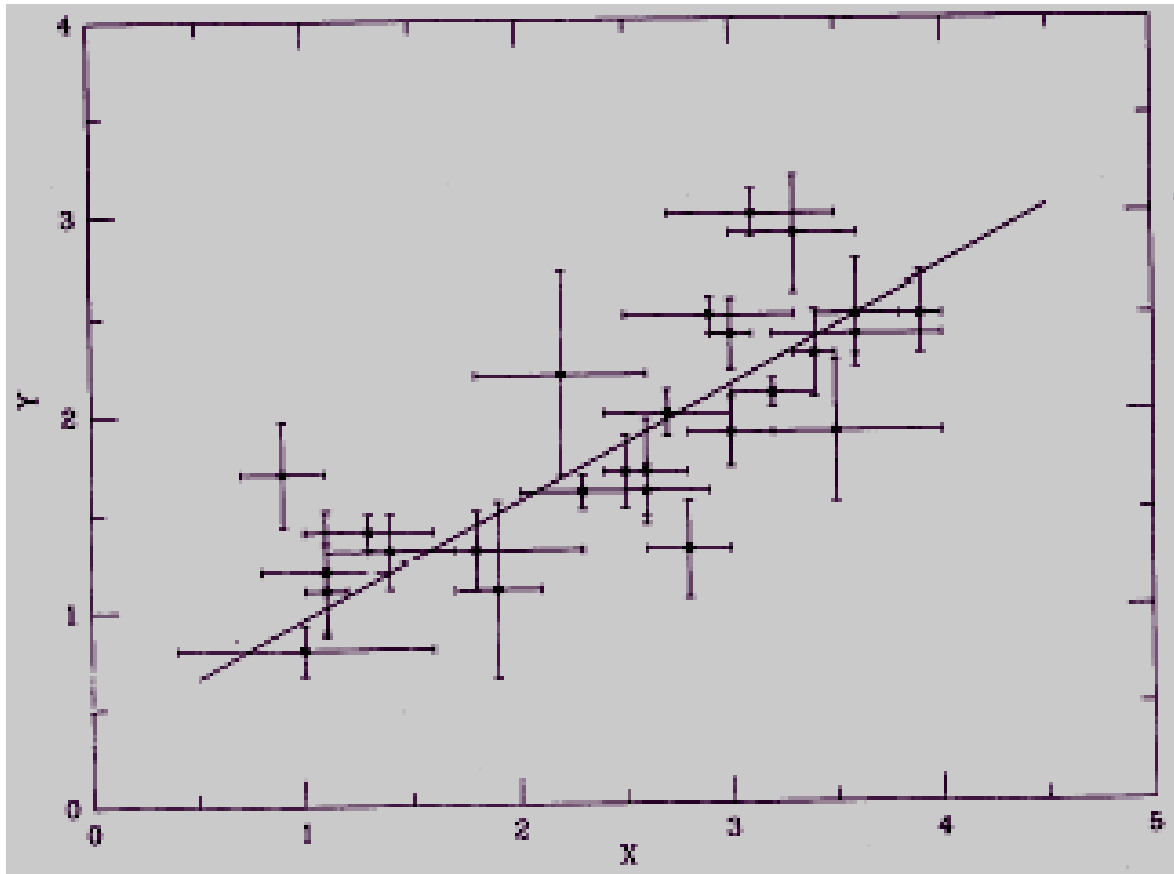| METHOD (1) | ASYMPTOTIC FORMULAE Intercept (2) | ASYMPTOTIC FORMULAE Slope (3) | BOOTSTRAP SLOPE (4) | JACKKNIFE SLOPE (5) |
|---|---|---|---|---|
| 23 Coma Ellipticals | | | | |
| OLS(Y\|X) | $-1.595 \pm 0.186$ | $1.162 \pm 0.082$ | $1.186 \pm 0.094$ | $1.164 \pm 0.111$ |
| OLS(X\|Y) | $-1.765 \pm 0.216$ | $1.238 \pm 0.096$ | $1.261 \pm 0.104$ | $1.239 \pm 0.128$ |
| OLS bisector | $-1.678 \pm 0.200$ | $1.199 \pm 0.088$ | $1.223 \pm 0.099$ | $1.201 \pm 0.119$ |
| Orthogonal | $-1.694 \pm 0.209$ | $1.206 \pm 0.092$ | $1.231 \pm 0.102$ | $1.208 \pm 0.124$ |
| Reduced major axis | $-1.679 \pm 0.200$ | $1.199 \pm 0.088$ | $1.223 \pm 0.099$ | $1.201 \pm 0.119$ |
| OLS mean | $-1.680 \pm 0.200$ | $1.200 \pm 0.088$ | $1.224 \pm 0.099$ | $1.201 \pm 0.119$ |
| 16 Virgo Ellipticals | | | | |
| OLS(Y\|X) | $-0.790 \pm 0.230$ | $1.144 \pm 0.101$ | $1.143 \pm 0.127$ | $1.114 \pm 0.118$ |
| OLS(X\|Y) | $-1.183 \pm 0.180$ | $1.316 \pm 0.082$ | $1.322 \pm 0.132$ | $1.316 \pm 0.093$ |
| OLS bisector | $-0.978 \pm 0.190$ | $1.227 \pm 0.085$ | $1.227 \pm 0.107$ | $1.226 \pm 0.099$ |
| Orthogonal | $-1.021 \pm 0.198$ | $1.245 \pm 0.089$ | $1.246 \pm 0.121$ | $1.245 \pm 0.104$ |
| Reduced major axis | $-0.979 \pm 0.190$ | $1.227 \pm 0.085$ | $1.228 \pm 0.108$ | $1.227 \pm 0.099$ |
| OLS mean | $-0.986 \pm 0.188$ | $1.230 \pm 0.084$ | $1.233 \pm 0.110$ | $1.230 \pm 0.098$ |

FIG. 2.—Example of a data set with large scatter obtained from Schechter's (1980) measurements of the Faber-Jackson relation in elliptical galaxies. The luminosity is in solar luminosity units. The two solid lines present OLS($Y \mid X$) (*shallowest line*) and OLS($X \mid Y$) (*steepest line*). The dot-dashed line, dashed line, and dotted line represent the OLS bisector, OR, and RMA, respectively.

The calculated slopes are 2.4 $\pm$ 0.4 and 5.4 $\pm$ 0.8 for the extrema OLS(L/$\sigma$) and OLS ($\sigma$/L), respectively, and 3.4 $\pm$ 0.4, 3.6 $\pm$ 0.4 and 5.2 $\pm$ 0.8 for the OLS bisector, reduced major axis, and orthogonal regression respectively. The scientific conclusions regarding distances and galaxy formation models obviously depend greatly on the regression method adopted. The dispersion of the five estimates is larger than the variance of any one estimate. The astronomer should calculate all the regression lines and be cautious about the confidence intervals and conclusions.

# Heteroscedastic measurement errors in both variables



**Homoscedastic functional**
Deeming (Vistas Astr 1968)
Fuller "Measurement Error Models" (1987)

**Heteroscedastic functional**
York (Can J Phys 1966)
ODRPACK  Boggs et al. (ACM Trans Math Soft 1990)

**Heteroscedastic structural**
BCES (Akritas & Bershady ApJ 1996)

# Functional Regression

$$Y_i = y_i + \varepsilon_i$$

$$X_i = x_i + \tau_i$$

$\varepsilon_i$ and $\tau_i$ are measurement errors

- We are interested in the real (regression) relation

$$y_i = bx_i + a$$

- $x_i$ are fixed.

# Fitting Power Law

- $f(z) = c\,z^{-\alpha}$ for $z > h > 0$ and for some $\alpha > 1$.
- $Y = \log(f(x))$, $\quad X = \log z$
- $Y = a + b\,X$
- Fitting the curve is equivalent to estimating a and b by linear regression
- Clearly we use OLS(Y|X)
- X is independent variable and Y is dependent variable

# Structural Regression

$$Y_i = y_i + \varepsilon_i$$

$$X_i = x_i + \tau_i$$

$\varepsilon_i$ and $\tau_i$ are measurement errors

- We are interested in the real (regression) relation

$$y_i = bx_i + a$$

- For any i, $x_i$ is a random variable, it has its own intrinsic variability

# Regression with measurement errors and intrinsic scatter

**Y = observed data**
**V = measurement errors**

$$(Y_{1i}, Y_{2i}, V_i), \quad i = 1, \ldots n$$

**X = intrinsic variables**
**e = intrinsic scatter**

$$Y_{1i} = X_{1i} + \epsilon_{1i} \quad \text{and} \quad Y_{2i} = X_{2i} + \epsilon_{2i}$$

**Regression model**
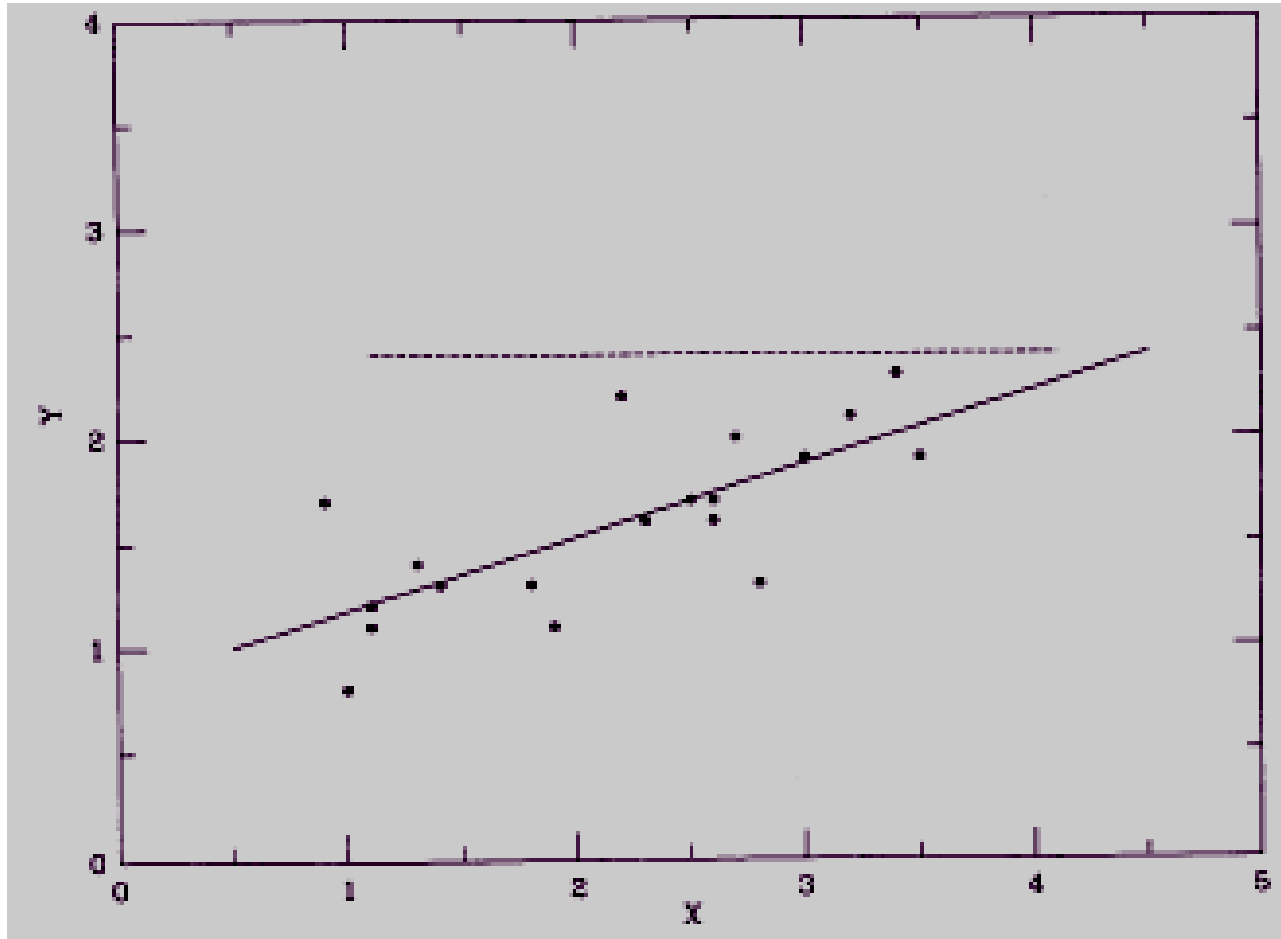
$$X_{2i} = \alpha_1 + \beta_1 X_{1i} + \epsilon_i$$

**Slope estimator**

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n}(Y_{1i} - \bar{Y}_1)(Y_{2i} - \bar{Y}_2) - \sum_{i=1}^{n} V_{12,i}}{\sum_{i=1}^{n}(Y_{1i} - \bar{Y}_1)^2 - \sum_{i=1}^{n} V_{11,i}}$$

$$\hat{\alpha}_1 = \bar{Y}_2 - \beta_1 \bar{Y}_1 .$$

**Slope variance**

$$\hat{\sigma}_{\beta_1}^2 = n^{-1}\sum_{i=1}^{n}(\xi_{1i} - \bar{\xi}_1)^2 \qquad \xi_{1i} = \frac{[Y_{1i} - E(Y_{1i})](Y_{2i} - \beta_1 Y_{1i} - \alpha_1) + \beta_1 V_{11,i} - V_{12,i}}{V(Y_{1i}) - E(V_{11,i})}$$
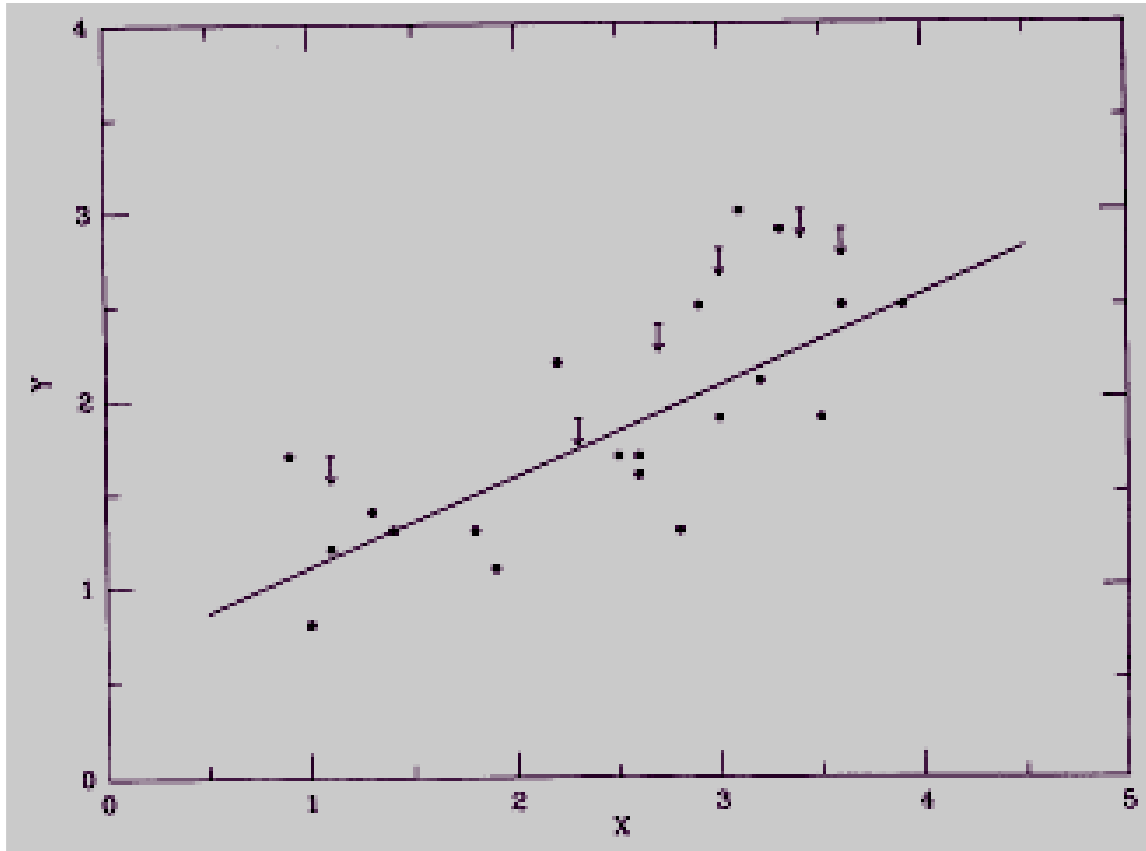
**Akritas & Bershady, ApJ 470, 706 1996**

# Truncation due to flux limits



**Econometrics:** Tobit & LIMDEP models (Amemiya, Advanced econometrics 1985; Maddala, Limited-dependent & Quantitative Variables in Econometrics 1983)
**Astronomy:** Malmquist bias in Hubble diagram (Deeming, Vistas Astr 1968, Segal, PNAS 1975)

# Censoring due to non-detections



**Correlation coefficients:**
Generalized Kendall's $\tau$ (Brown, Hollander & Korwar 1974)

**Linear regression with normal residuals:**
EM Algorithm (Wolynetz Appl Stat 1979)

**Linear regression with Kaplan-Meier residuals:**
Buckley & James (Biometrika 1979)    Schmitt (ApJ 1985)

**Isobe, Feigelson & Nelson (ApJ 1986)**
**Implemented in Astronomy Survival Analysis (ASURV) package**

# Conclusions

**Bivariate linear regression in astronomy can be surprisingly complex. Pay attention to precise question being asked, and details of situation. Several codes are available through http://astrostatistics.psu.edu/statcodes.**

- **Functional vs. structural regression**
- **Symmetrical vs. dependent regression**
- **Weighting by measurement error**
- **Truncation & censoring due to flux limits**