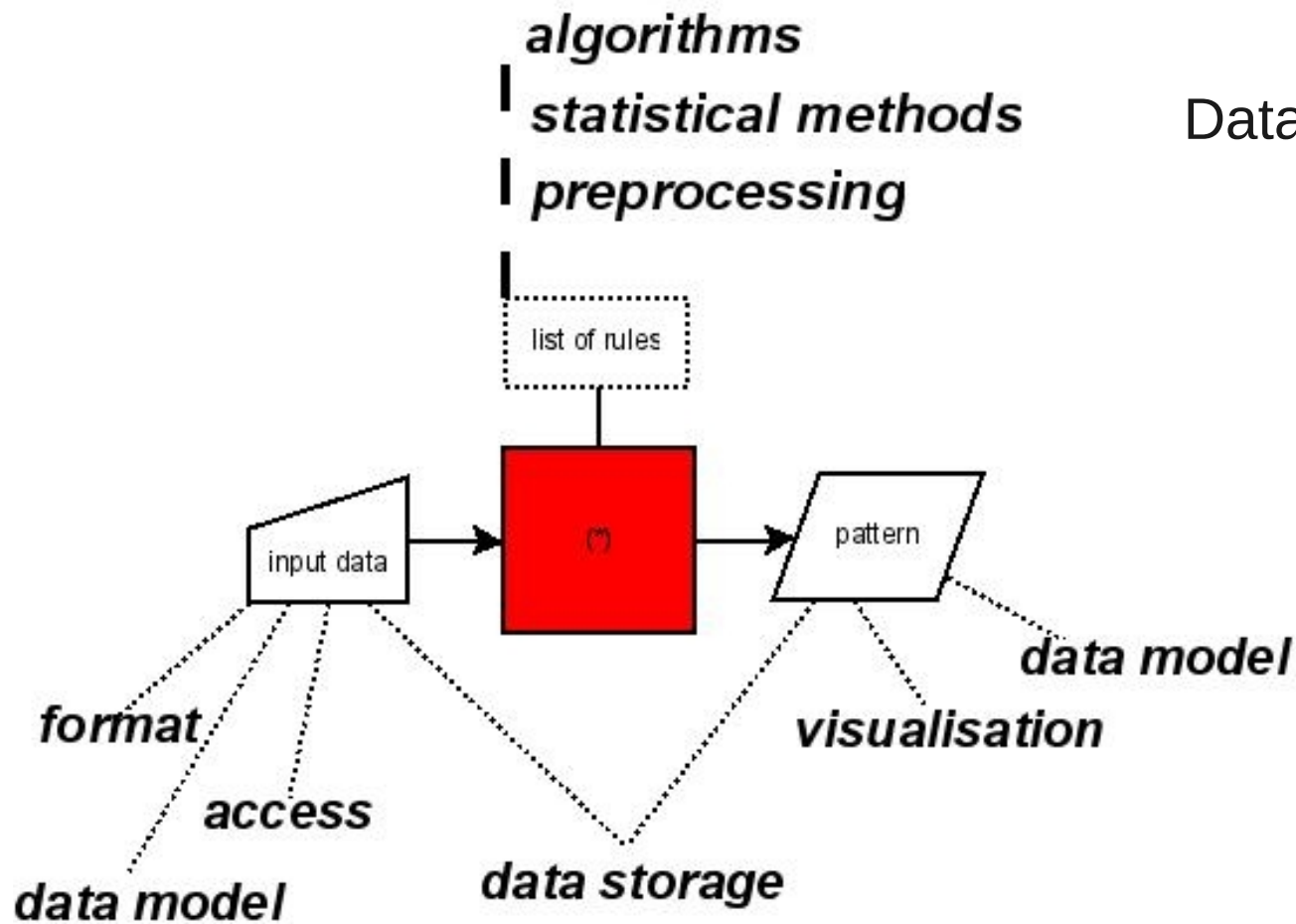


Virtual Observations 2012

Statistics in Data Mining



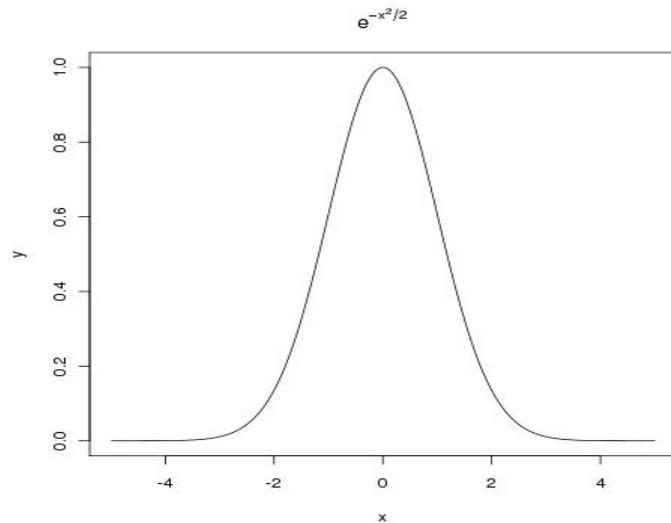
Database: Data format
Access
Data model
Preprocessing
Processing

- List of rules: verification of model or parametric search
- Verification of model: parametric and non-parametric
- PDF
- Bayesian statistic
- Tests

$$dP = f(x)dx$$

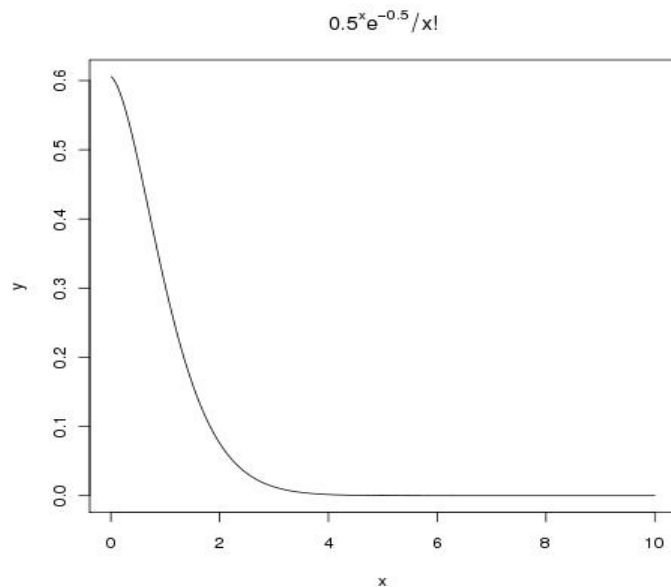
$$P(x_1 < x < x_2) = \int_{x_1}^{x_2} f(x) dx$$

$$P(x < x_1) = \int_{-\infty}^{x_1} f(x) dx$$



$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-x_0)^2}{2\sigma^2}}$$

$$f(x) = \lambda^x e^{-x} / x!$$

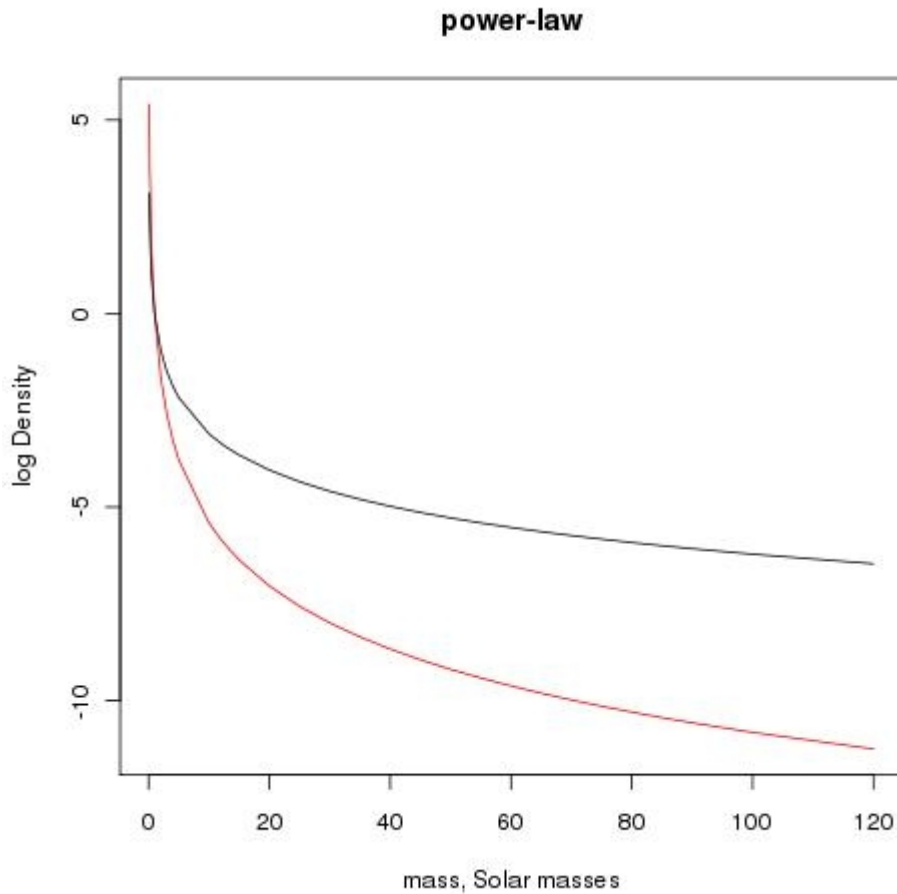


- Gauss
- Poisson

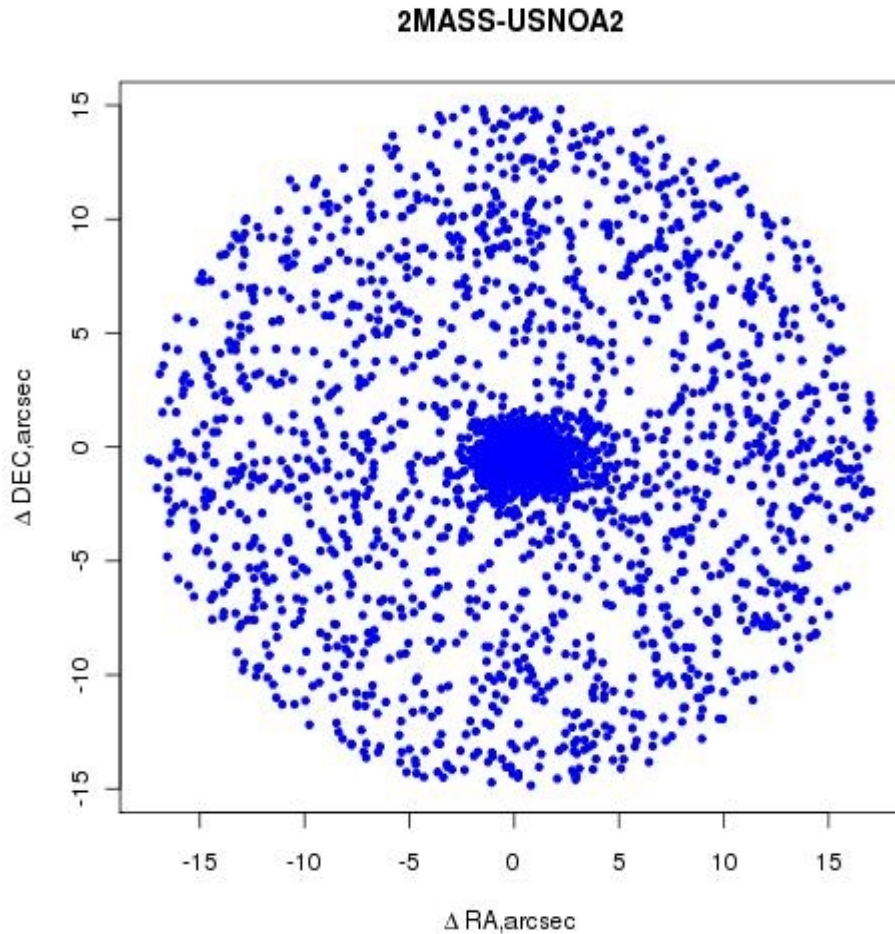
$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = N(\bar{x}, \sigma^2/n)$$

$$\lim_{n \rightarrow \infty} \frac{1/n \sum_{i=1}^n X_i - \bar{x}}{\sigma/\sqrt{n}} = N(0, 1)$$

Power-law distribution



$$dN = m^{-\alpha} dm$$



- Cross-identify 2MASS and USNO-A2
- Select coordinate differences
- Find mean, median, variance, MAD

$$\bar{x} = E(x) = \int_{-\infty}^{\infty} f(x) x dx$$

$$\bar{x} = \sum_i f(x_i) x_i$$

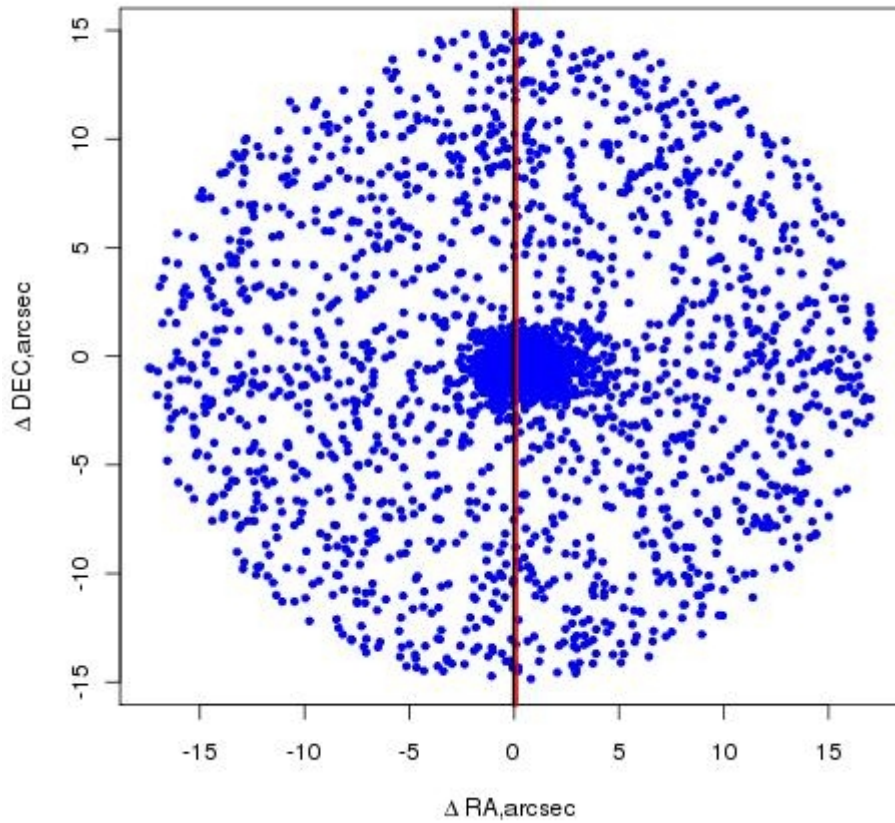
$$\text{median}(x) = \tilde{x} : P(x < \tilde{x}) = \int_{-\infty}^{\tilde{x}} f(x) dx = 1/2$$

$$P(x > \tilde{x}) = \int_{\tilde{x}}^{\infty} f(x) dx = 1/2$$

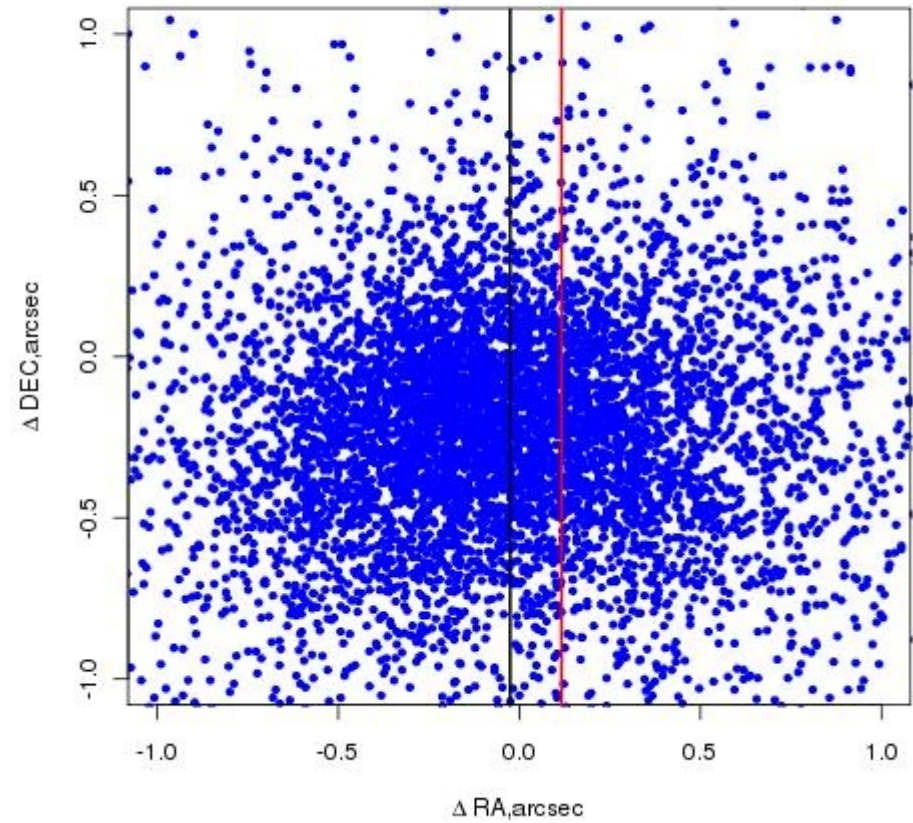
Mean, median



2MASS-USNOA2



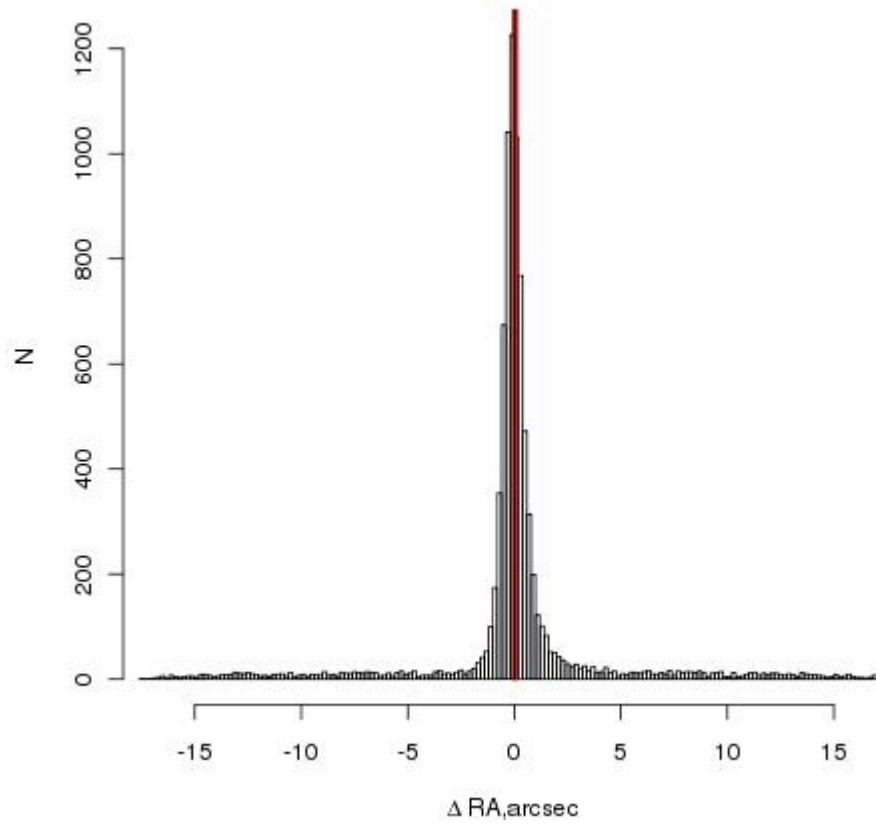
2MASS-USNOA2



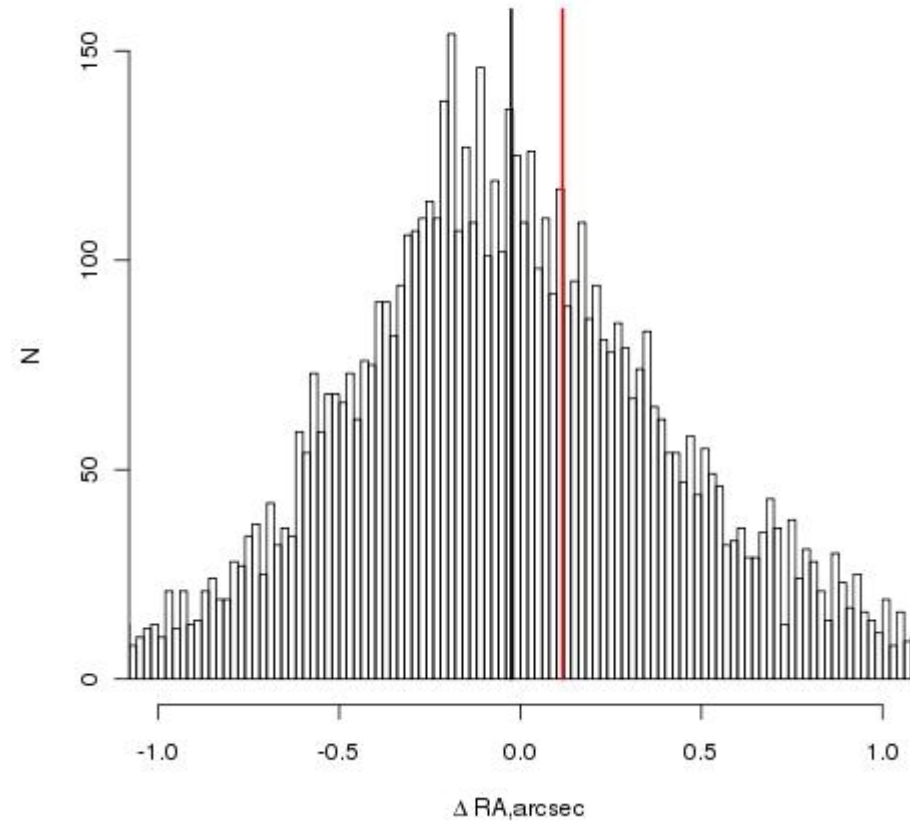
Mean, median



2MASS-USNOA2



2MASS-USNOA2



$$\sigma^2 = \int (x - \bar{x})^2 f(x) dx$$

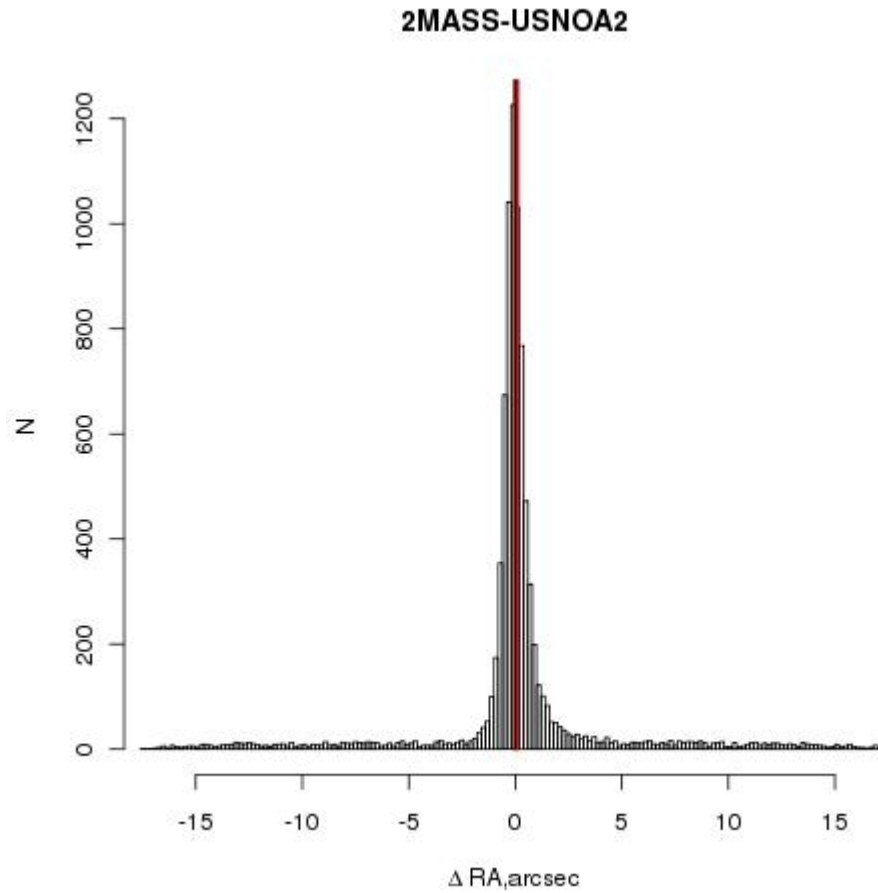
$$\sigma^2 = \bar{x^2} - \bar{x}^2$$

$$\text{Var}(x) = E(x - \bar{x})^2 = \int (x - \bar{x})^2 f(x) dx = \sigma^2$$

$$x_{RMS}^2 = \int x^2 f(x) dx$$

$$x_{RMS}^2 = \bar{x}^2 = \sigma^2 + \bar{x}^2$$

K-sigma clipping

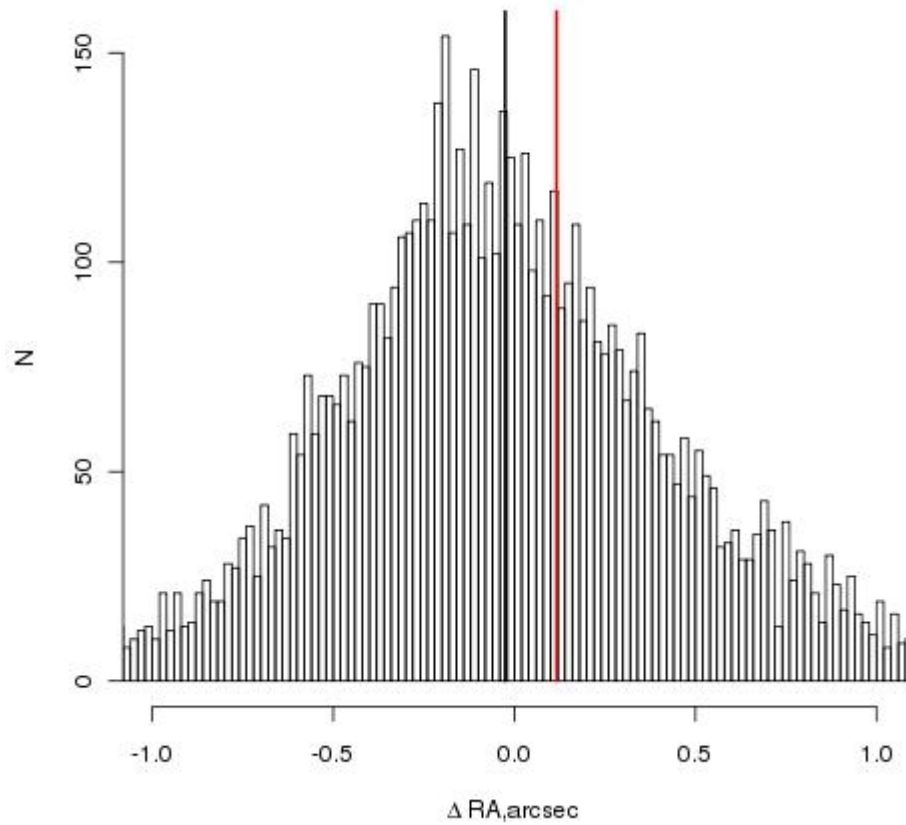


$$|x_i - \bar{x}| > k \sigma$$

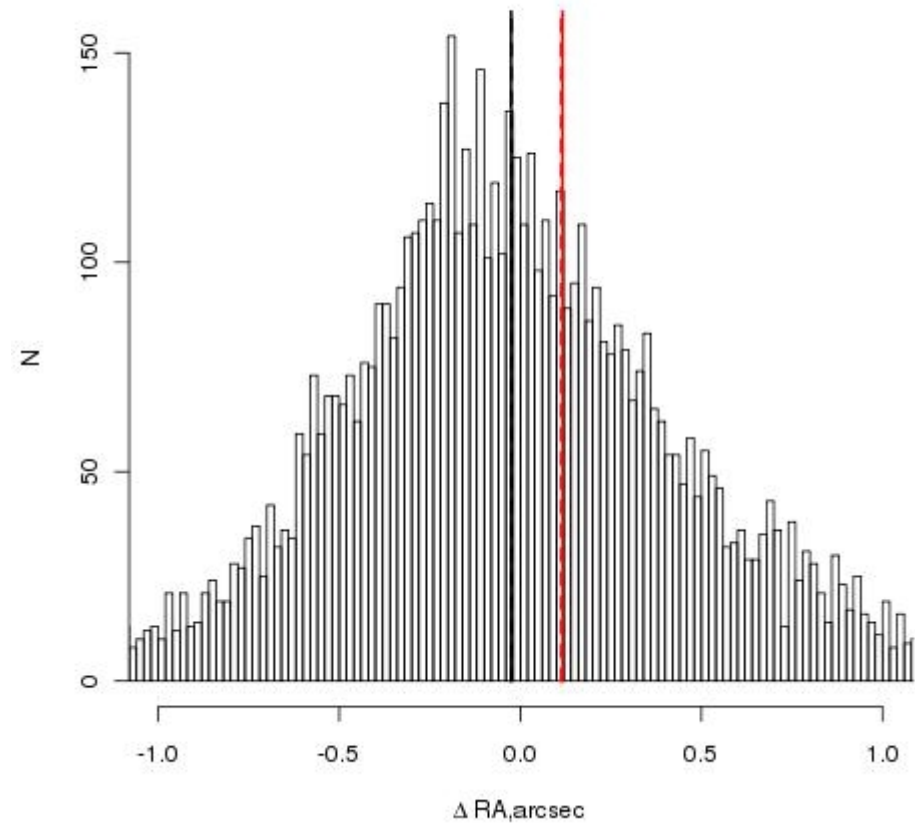
K-sigma clipping



2MASS-USNOA2



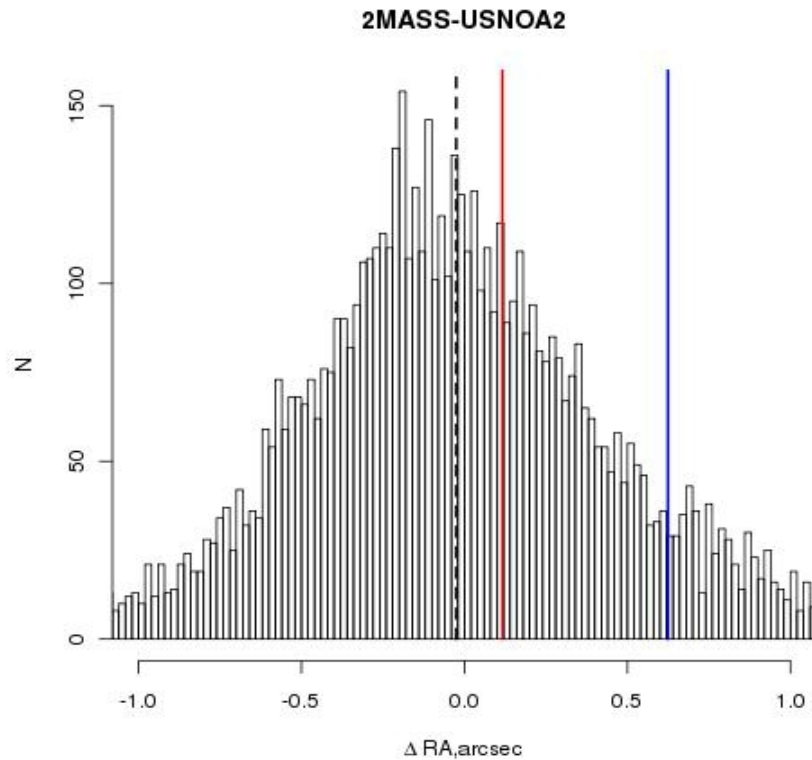
2MASS-USNOA2



MAD



$$\text{MAD}(\mathbf{x}) = \int_{-\infty}^{\infty} (x - \tilde{x}) f(x) dx$$



- MySQL:
 - AVG() - mean
 - COUNT() - N
 - STD(), STDDEV(), STDDEV_POP()
 - STDDEV_SAMP()
 - VARIANCE(), VAR_POP()
 - VAR_SAMP()

Distributions to test



0.025

0.52

-0.025

0.52

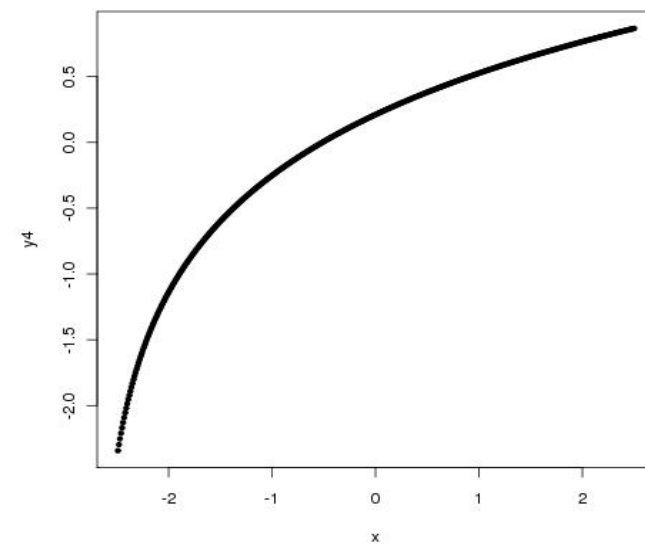
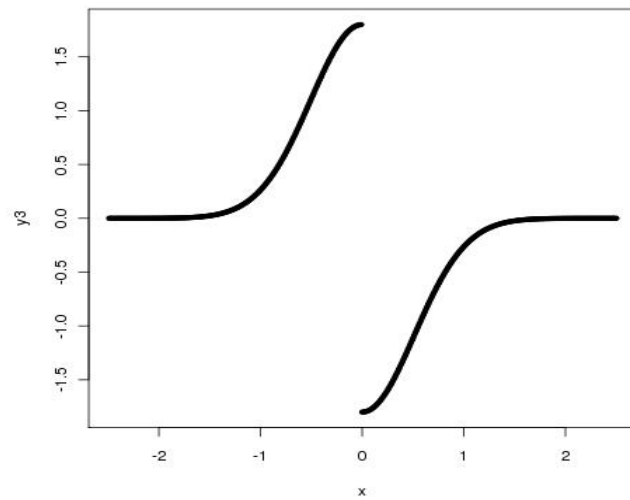
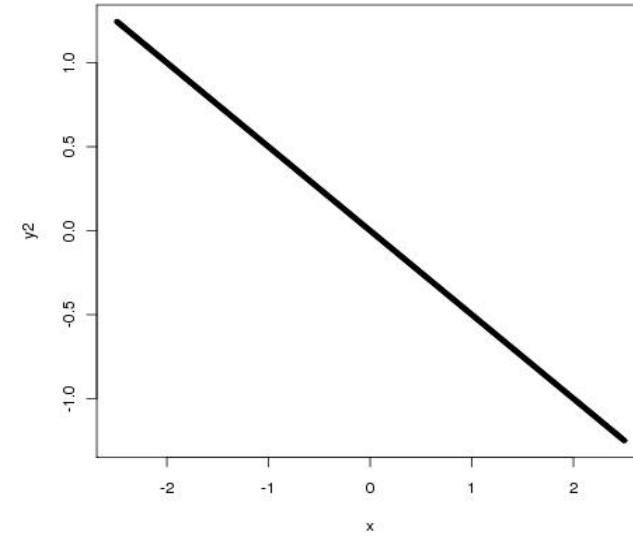
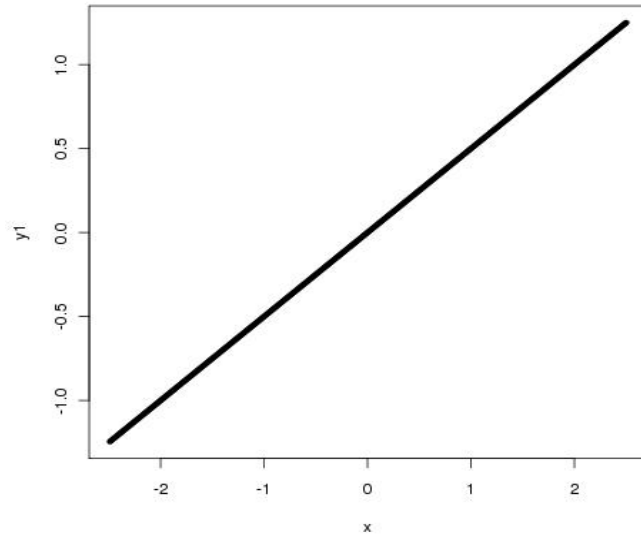
-0.003

0.52

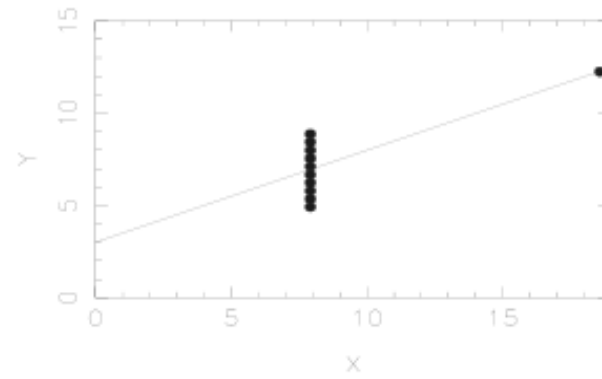
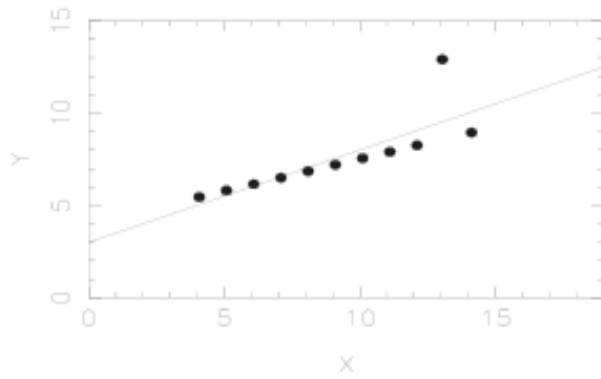
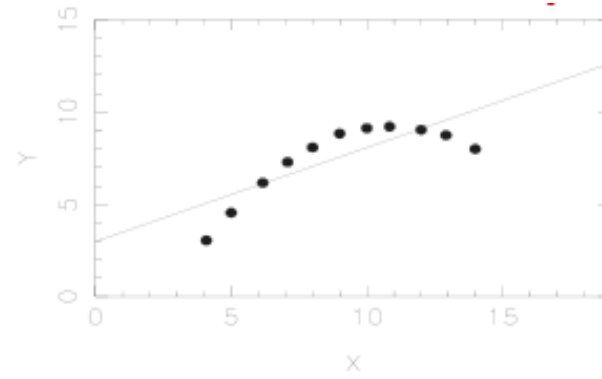
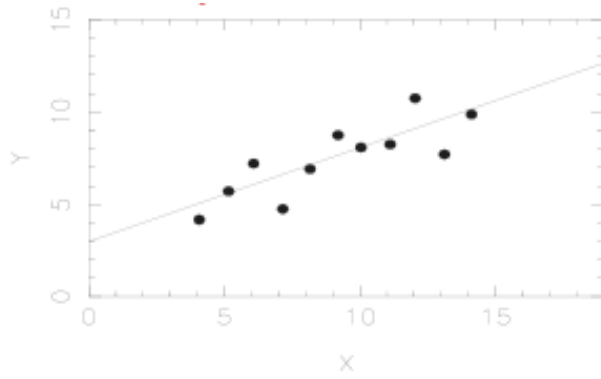
0.0003

0.52

Real situation



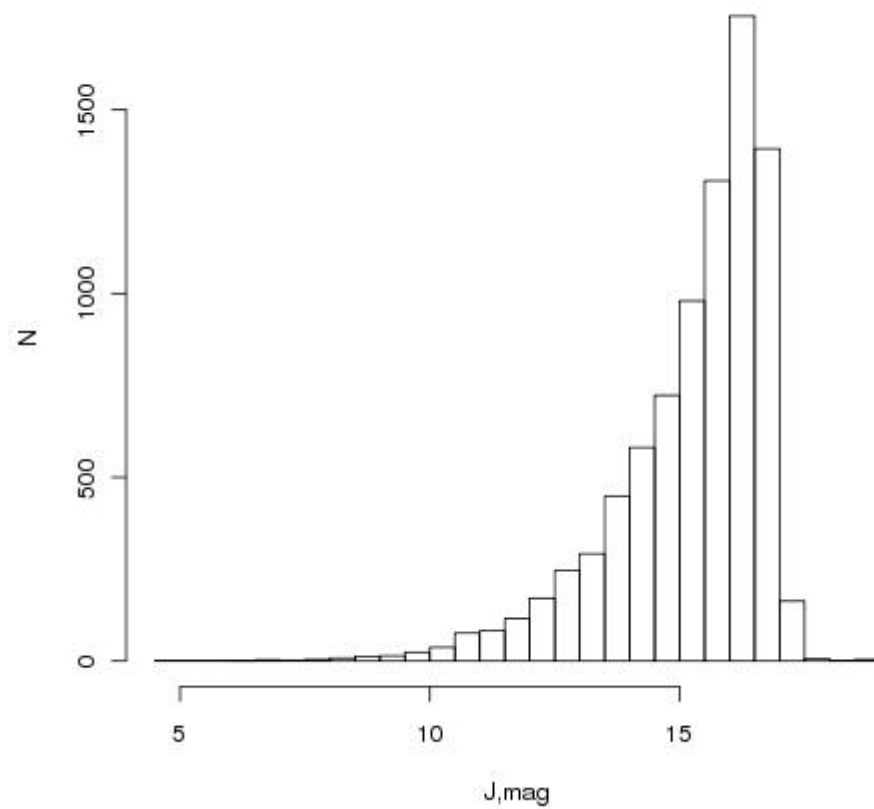
Anscombe quartet



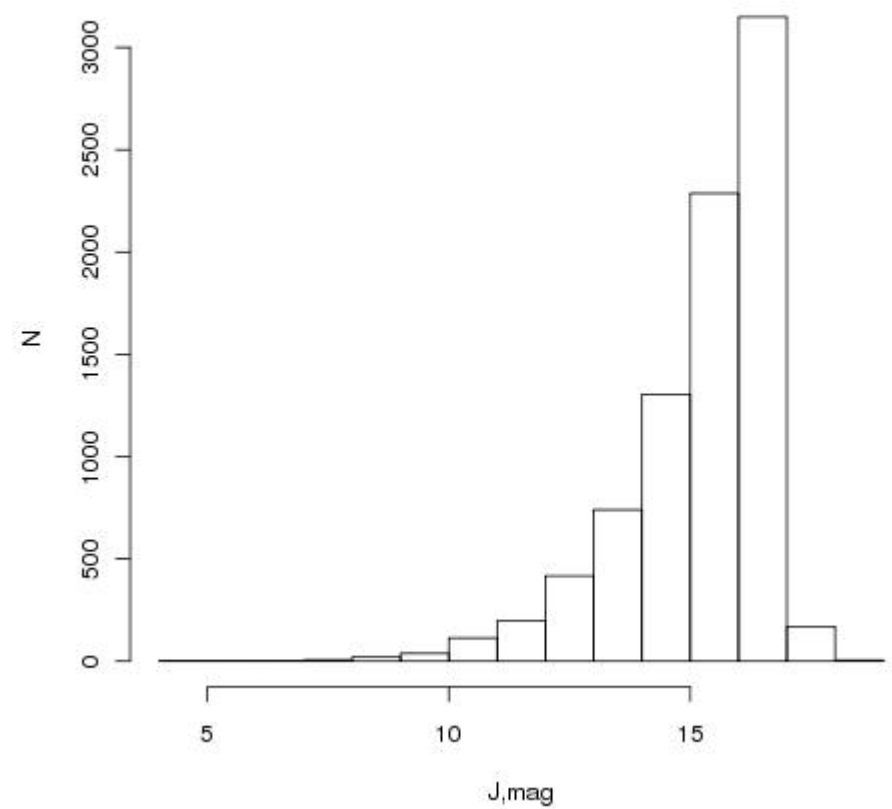
Histogram



2MASS



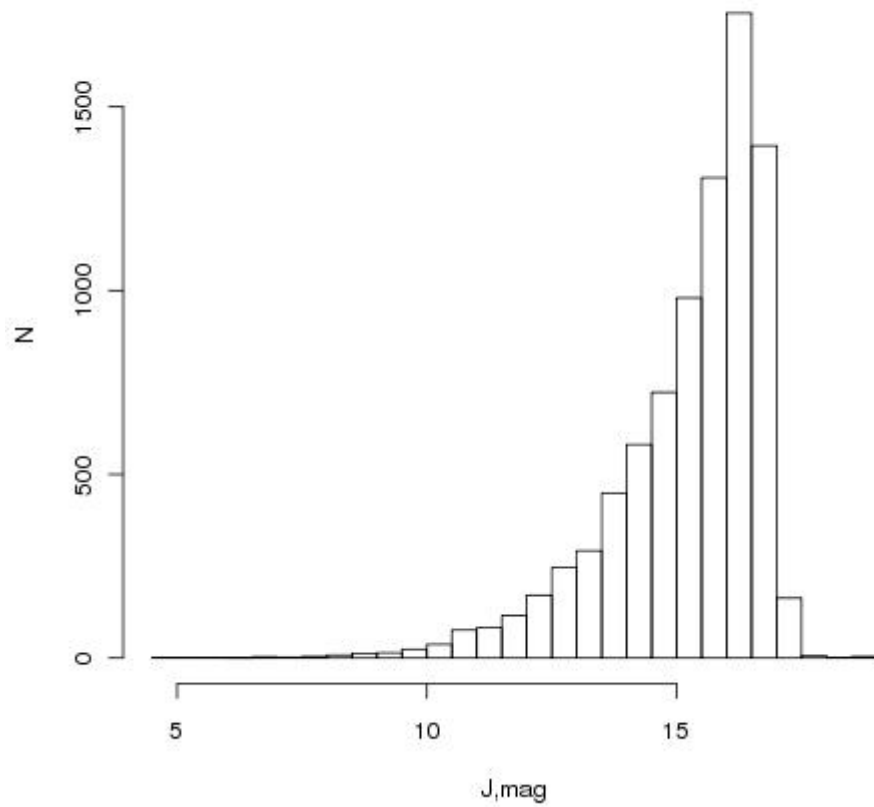
2MASS



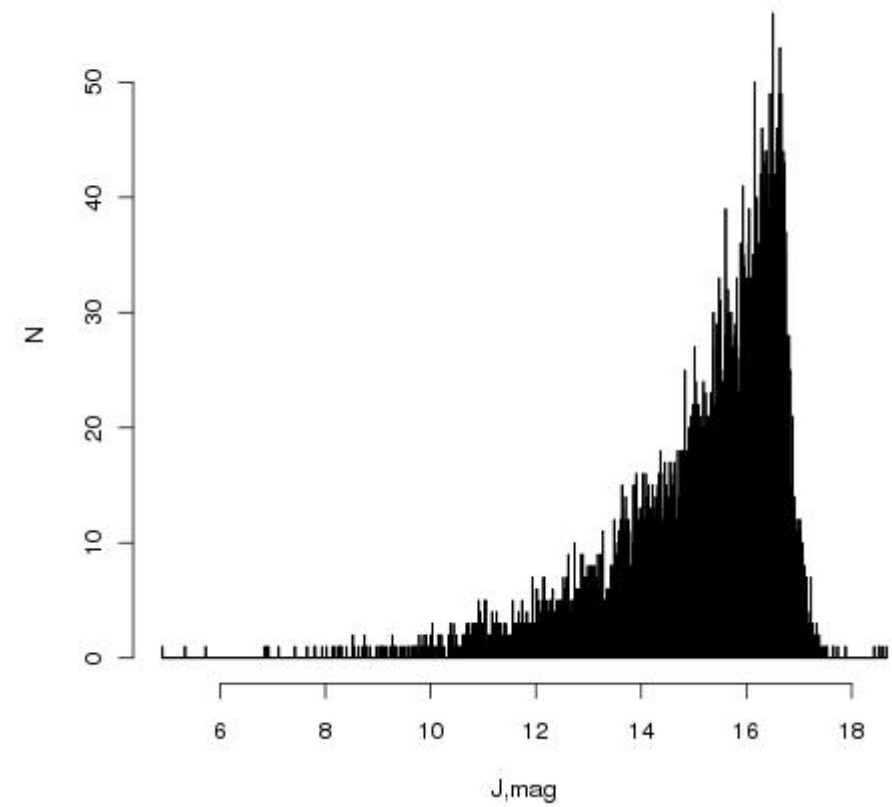
Histogram



2MASS



2MASS



- Biased
- Oversmoothed/Undersmoothed

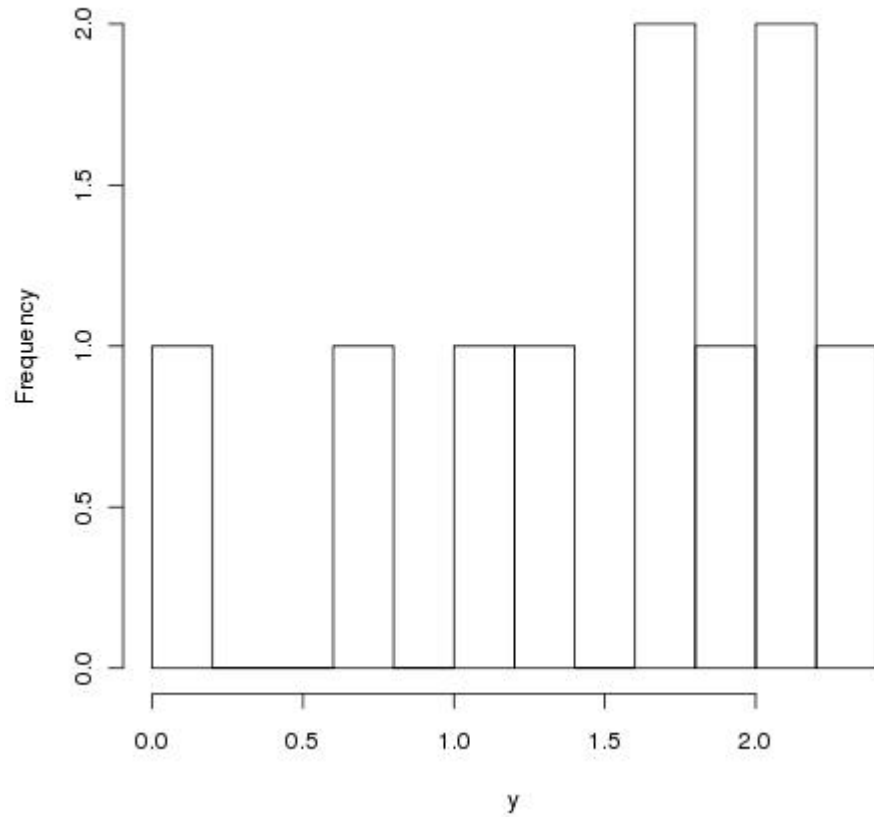
$$p_i = \frac{n_i}{\Delta_i \sum_i n_i}$$

$$h_{opt} = \frac{3.5 \sigma}{n^{1/3}}$$

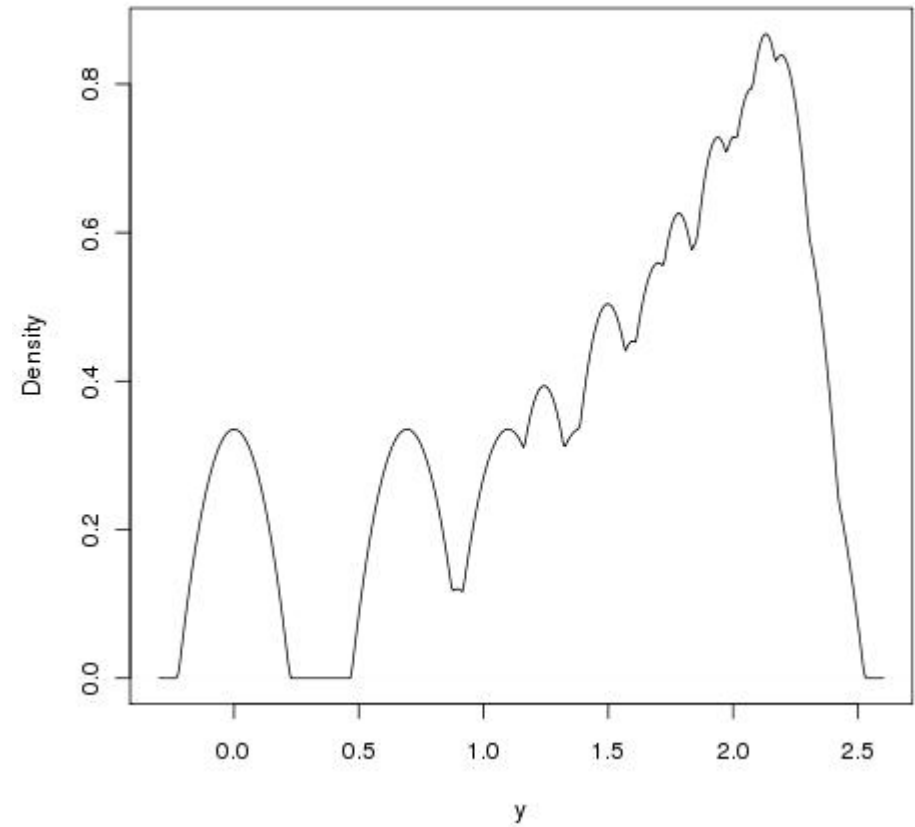
Kernel smoothing



Histogram of y

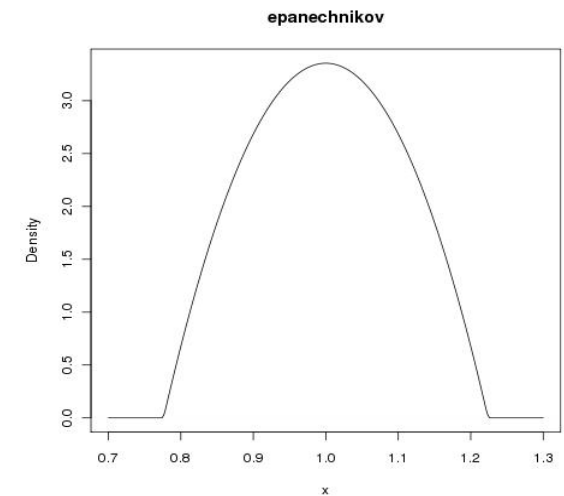
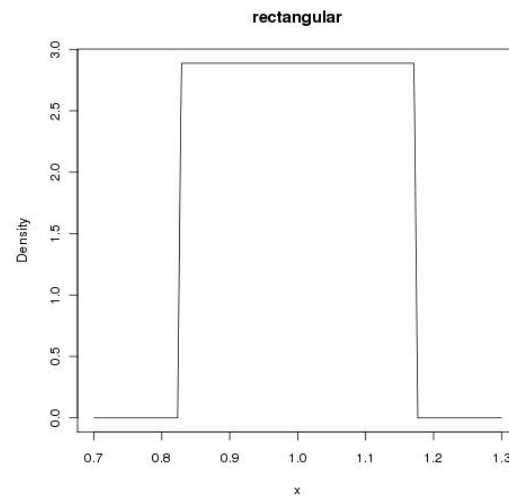
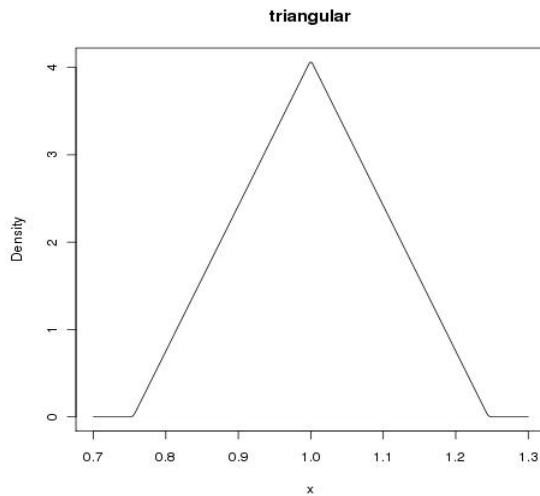


epanechnikov



$$\hat{f}(x; h) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)$$

Kernel smoothing



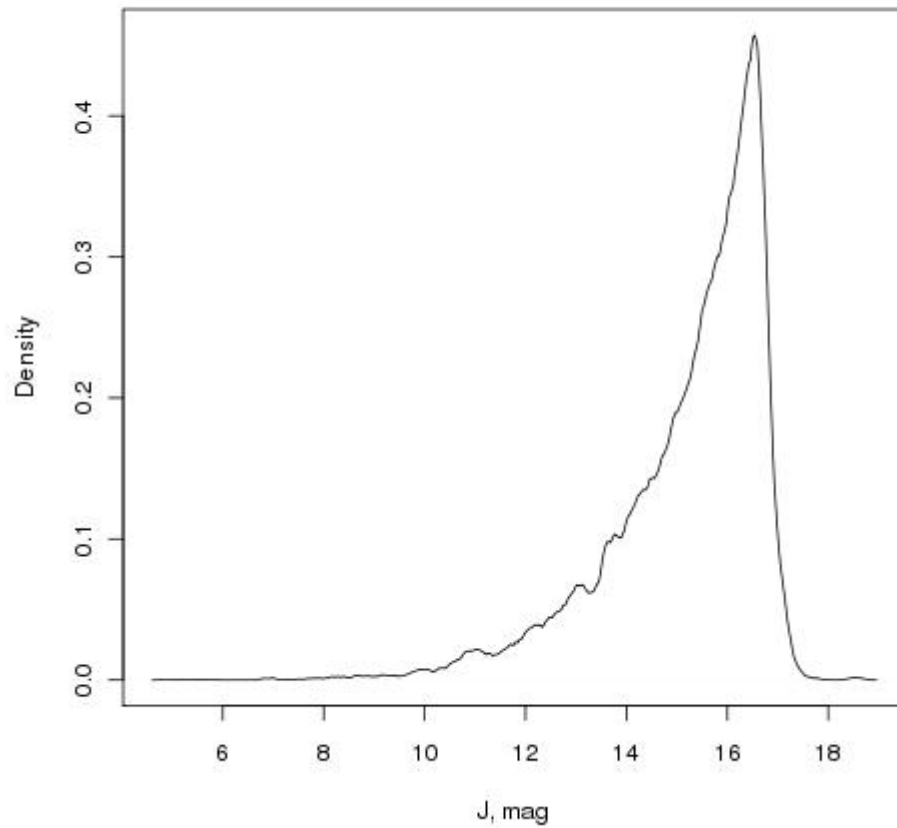
Coordinates are shifted on 0.5 on figures!

$$K(x) = 1/2, |x| < 1 \quad K(x) = (1 - |x|), |x| < 1 \quad K(x) = \frac{3}{4}(1 - x^2), |x| < 1$$

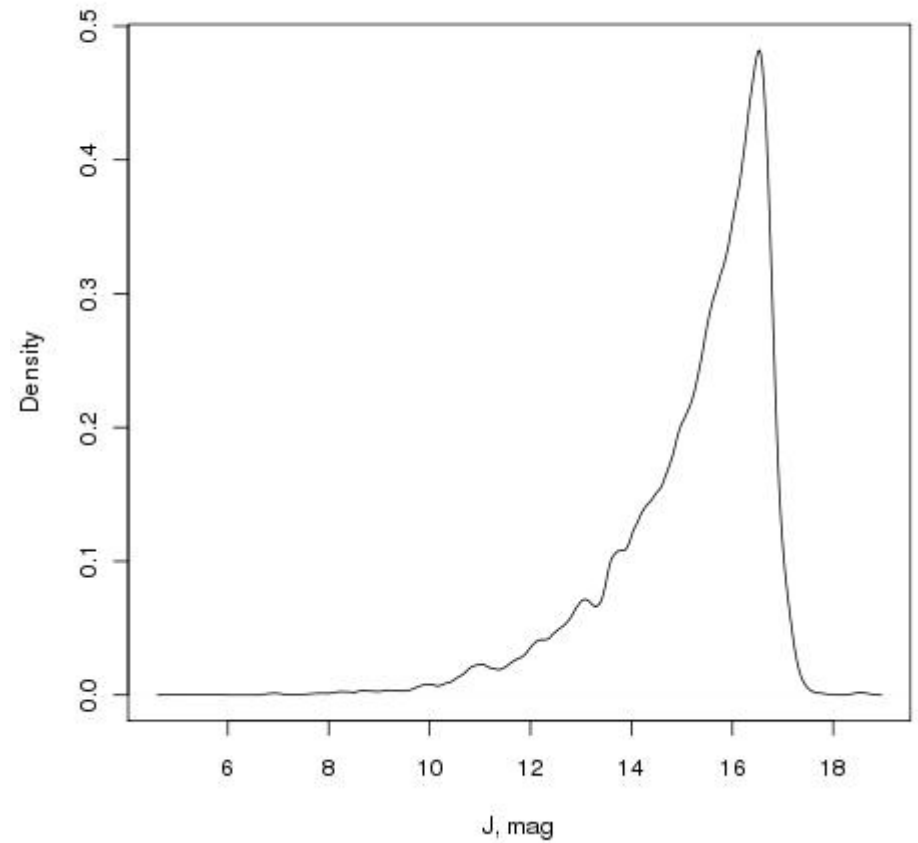
Kernel smoothing



rectangular



epanechnikov



$$h_{opt} = \left[\frac{R(K)}{\mu_2(K)^2 R(\ddot{f})n} \right]^{1/5}$$

$$\mu_2(K) = \int x^2 K(x) dx$$

$$R(f) = \int f^2(x) dx$$

- $P(B|A) = P(A|B)P(B)/P(A)$
 - $P(A)$ – normalization
 - $P(B)$ – prior probability
 - $P(A|B)$ – likelihood
 - $P(B|A)$ – posterior probability

Maximum likelihood



$$f(p_k | p_1, \dots, p_{k-1}, p_{k+1}, \dots, p_n) = \int_{p_1, \dots, p_{k-1}, p_{k+1}, \dots, p_n} f(\vec{P}) dp_1, \dots, dp_{k-1}, dp_{k+1}, \dots, dp_n$$

$$L = \prod_{i=1}^n f(p_i)$$

$$\chi_k^2 = -2 \ln \frac{L_{-k}}{L}$$

$$\chi^2 = \sum_i \frac{(O_i - T_i)^2}{T_i}$$

$$\chi^2 < \chi^2(\alpha, n - k - 1)$$

- Data sample
- Mean, median, variance, rms
- Visualisation
- Model
- Test

Covariance

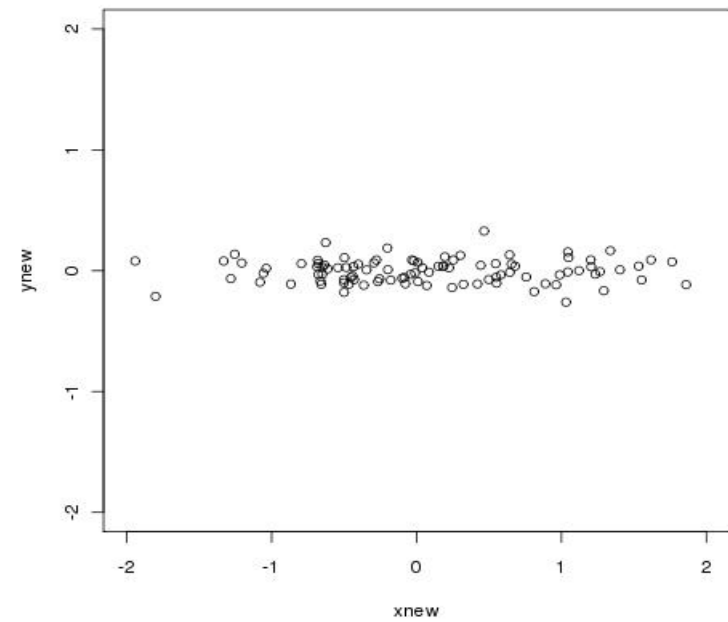
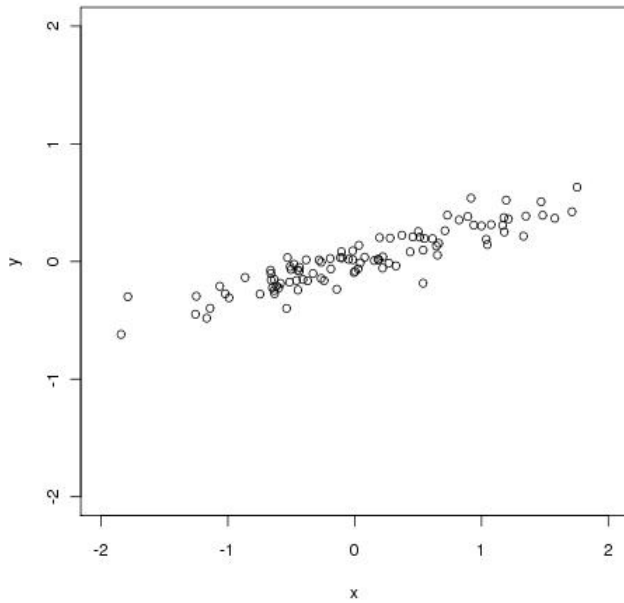


$$X = \{X_1, X_2, X_3, \dots, X_n\}$$

$$\bar{X} = \{\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n\}$$

$$\text{cov}(X_i, X_j) = E((X_i - \bar{X}_i)(X_j - \bar{X}_j))$$

PCA



$$\text{cov}(X, Y) = C$$

$$TC T^{-1} = \Lambda$$

Lutz-Kelker effect



$$g(\pi_0 | \pi) = \frac{1}{\sqrt{(2\pi)\sigma}} \exp\left(-\frac{(\pi_0 - \pi)^2}{2\sigma^2}\right)$$

$$g(\pi | \pi_0) \propto \frac{1}{\pi^4} \exp\left(-\frac{(\pi - \pi_0)^2}{2\sigma^2}\right)$$

$$Z \equiv \frac{\pi}{\pi_0}$$

$$g(\pi | \pi_0) \propto G(Z) \equiv \frac{1}{Z^4} \exp\left(-\frac{(Z - 1)^2}{2(\sigma/\pi_0)^2}\right)$$

$$\Delta M = M_{\text{true}} - M_{\text{observed}} = 5 \log \frac{\pi}{\pi_0} = 5 \log Z$$

$$\langle \Delta M \rangle = \frac{5 \int_0^\infty \log Z G(Z) dZ}{\int_0^\infty G(Z) dZ}$$

