# MULTIPLE EXPOSURES IN LARGE SURVEYS
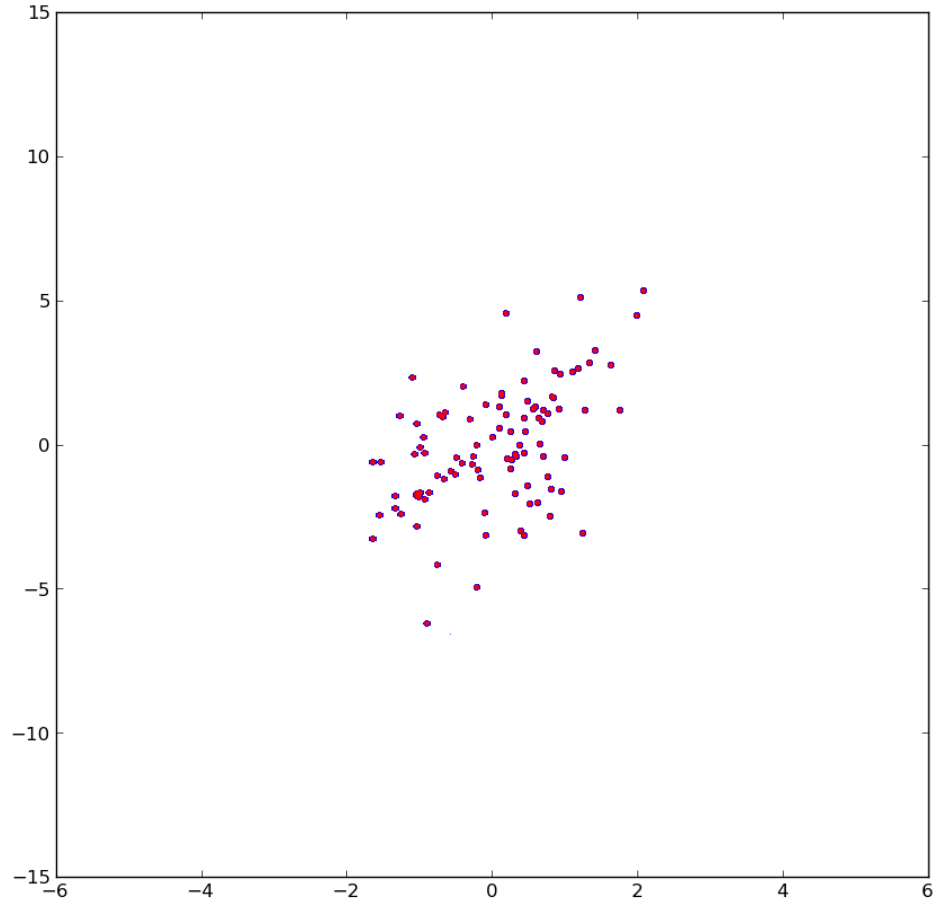
Tamás Budavári / Johns Hopkins University
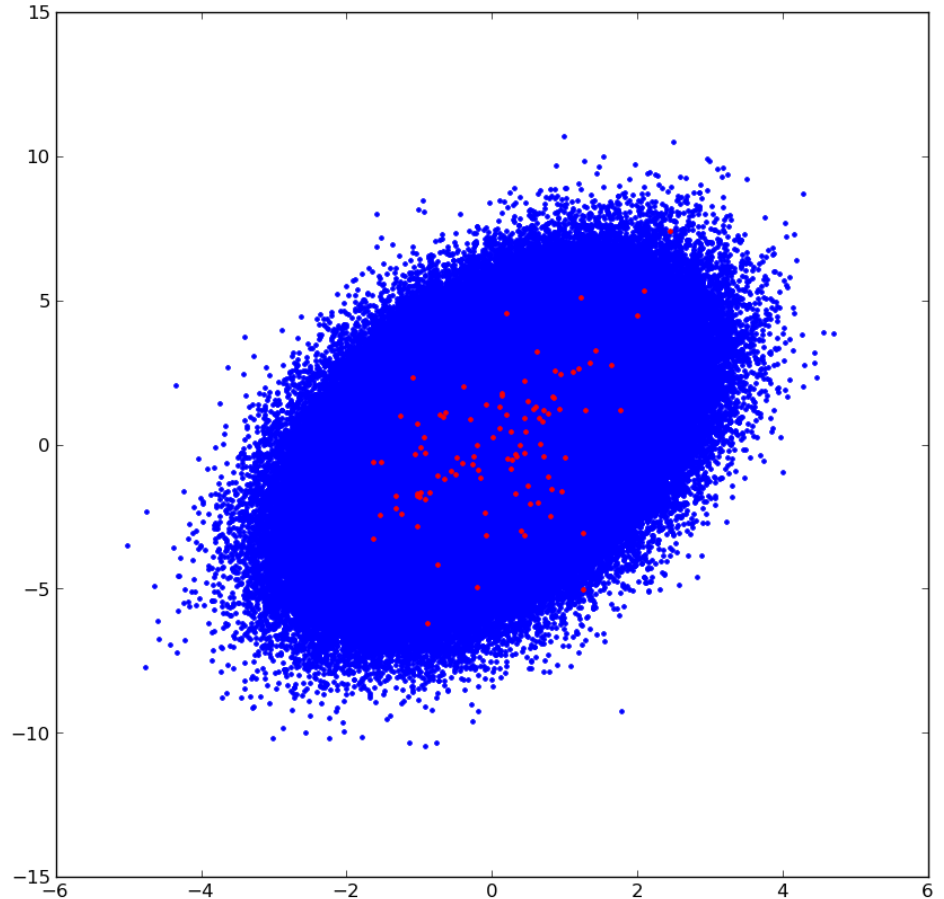
# Big Data?

- Noisy
- Skewed
- Artifacts

# Big Data?

- Noisy
- Skewed
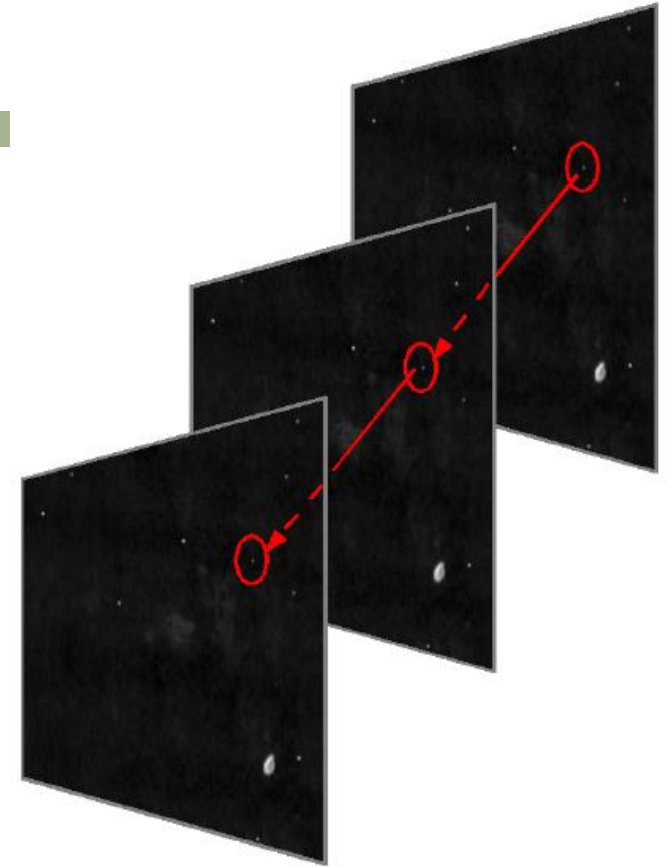- Artifacts
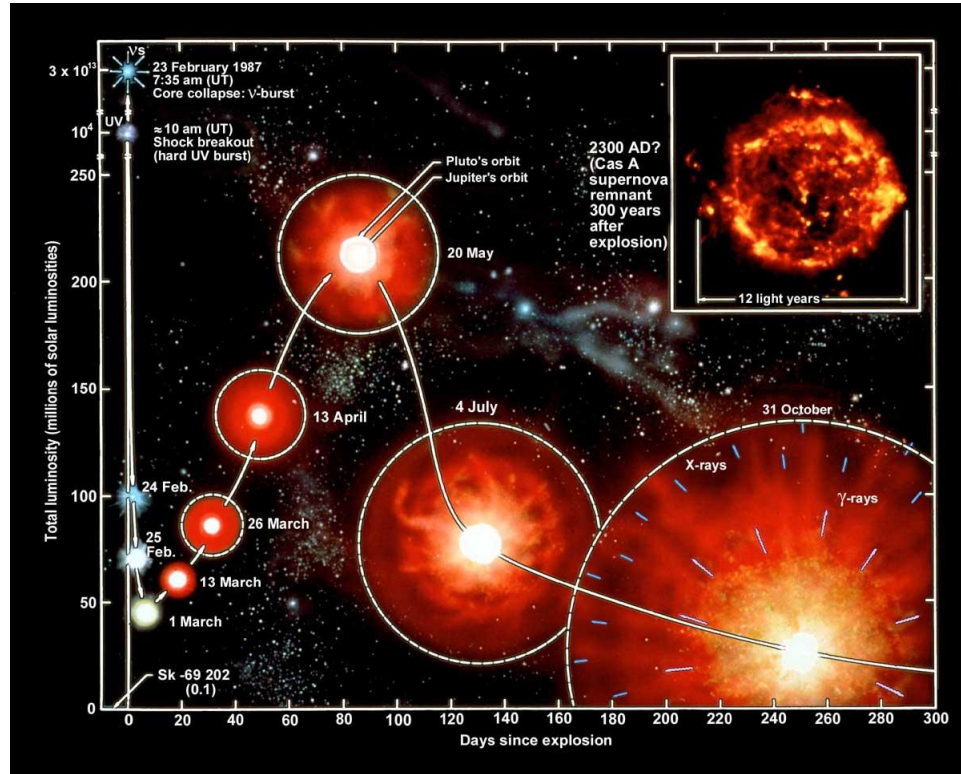
# Serious Issues

- Significant fraction of catalogs is junk
  - GALEX         ~50%
  - PS1 3PI         50-80%
  - PS1 MDS        >95%

- Textbook methods often fail due to artifacts
  - What are the good techniques?

12/12/2013

# Time Domain

# Time Series of Faint Sources?

**Tamás Budavári**

- ☐ Co-add images and do forced photometry
  - ☐ Ideal if we have all observations but we never do
- ☐ Independent catalogs as we go
  - ☐ Need to dig in the noise to build good timeseries
- ☐ Goal is an incremental strategy to weed out noise
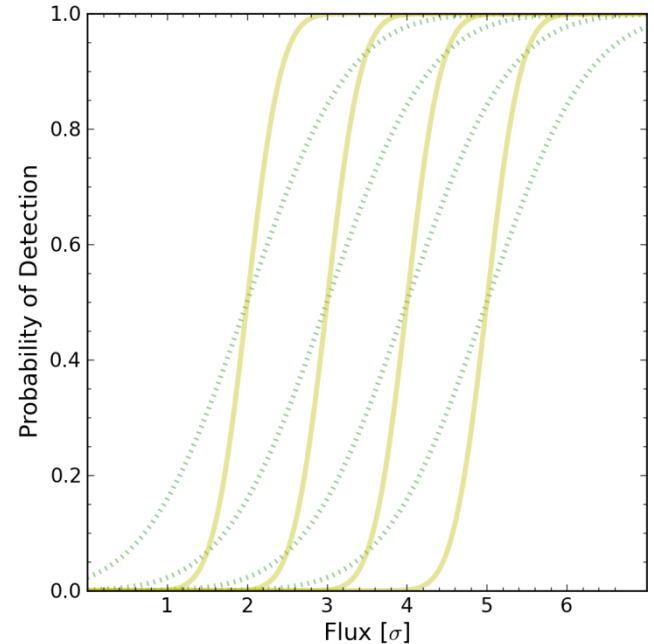  - ☐ Otherwise catalogs are overwhelmed by junk

12/12/2013

# Detection Probability

**Tamás Budavári**

☐ Measured flux is true + normal error $f_i = f + \epsilon_i$

☐ Probability of detection

$$P_f \equiv P(f_i > f_D | f) = \frac{1}{2} \text{erfc} \left( \frac{f_D - f}{\sigma \sqrt{2}} \right)$$
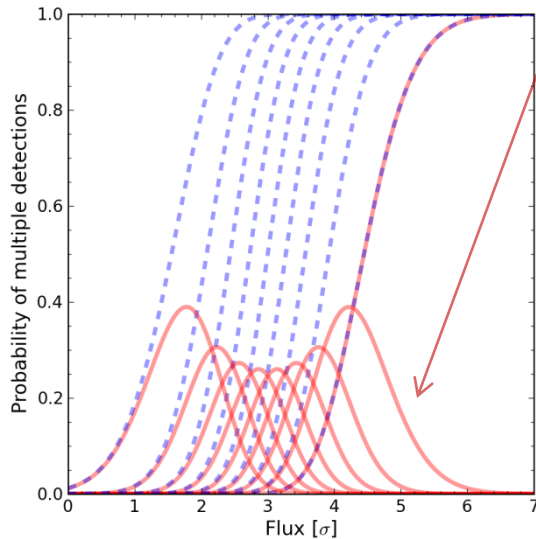
12/12/2013

# Detection Probability

☐ Measured flux is true + normal error $f_i = f + \epsilon_i$

☐ Probability of detection

    ☐ As a function of the true flux

        ■ Thresholds at *2-, 3-, 4- & 5σ*

    ☐ Sharper for 9-way stacks

# Detection Probability

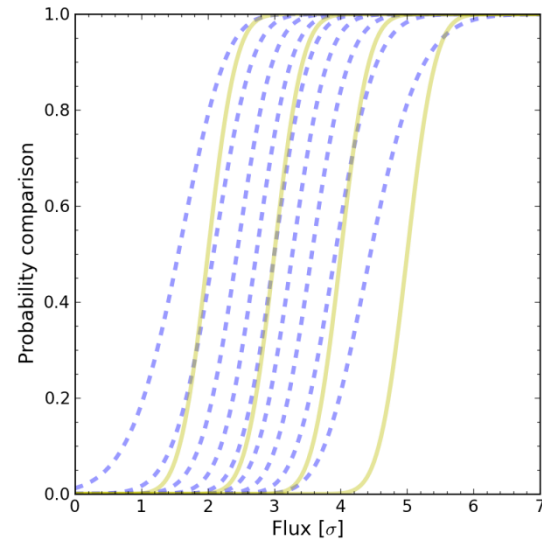□ Multiple exposures

■ Binomial $$P(n|k,f) = \binom{k}{n} P_f^n \left(1 - P_f\right)^{k-n}$$

# Detection Probability

## ☐ Multiple exposures

### ◻ Binomial

$$P(n|k,f) = \binom{k}{n} P_f^n \left(1 - P_f\right)^{k-n}$$

# What is a Real Source?

□ Is it "real" or just "noise" ?

    □ Bayesian hypothesis testing

$$B = \frac{L_{\text{real}}}{L_{\text{noise}}}$$

# What is a Real Source?

- Is it "real" or just "noise" ?
  - Bayesian hypothesis testing

$$B = \frac{L_{\mathrm{real}}}{L_{\mathrm{noise}}}$$

$$L_{\mathrm{real}} = \int df\, \pi(f)\, L(f)$$

$$L(f) = (1 - P_f)^{k-n} \prod_i^n G(f_i; f, \sigma^2)$$

- Out of $k$ observations $n$ detections of $f_i$ fluxes

12/12/2013
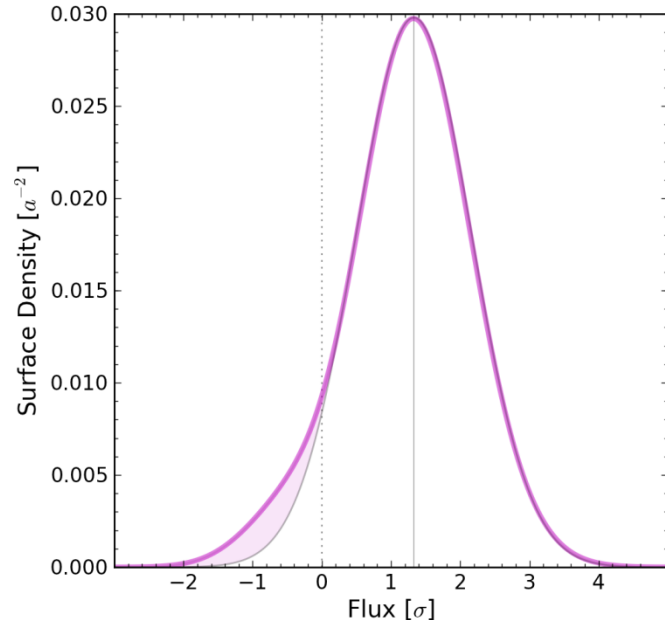
# Apparent Flux Distribution

Tamás Budavári

□ Galaxy number-counts as fn of magnitude

    ■ Empirical relation approximately shows

$$dN \propto 10^{0.4m}\, dm \propto \frac{df}{f^2}$$

□ More and more fainter and fainter sources!

    ■ But there is a limit, cf. Olbers' paradox
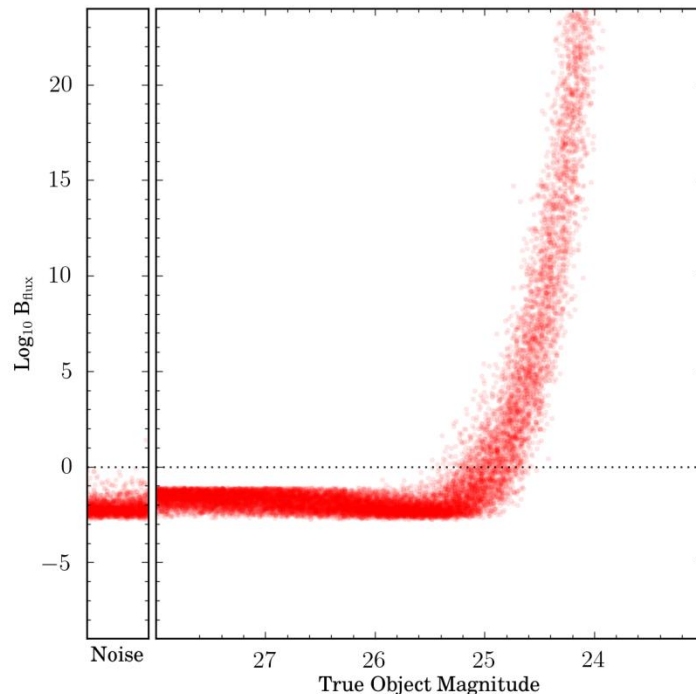
12/12/2013

# Distribution of Noise Peaks

☐ Local maxima of continuous Gaussian random field

  ☐ Cf., $P(\mathbf{k})$ by Barden, Bond, Kaiser, Szalay (BBKS; 1986)

  ☐ Now in 2D:

# Something Like LSST

**Tamás Budavári**

☐ Simulation

  ☐ Sky at $5\sigma$ is 24 mag

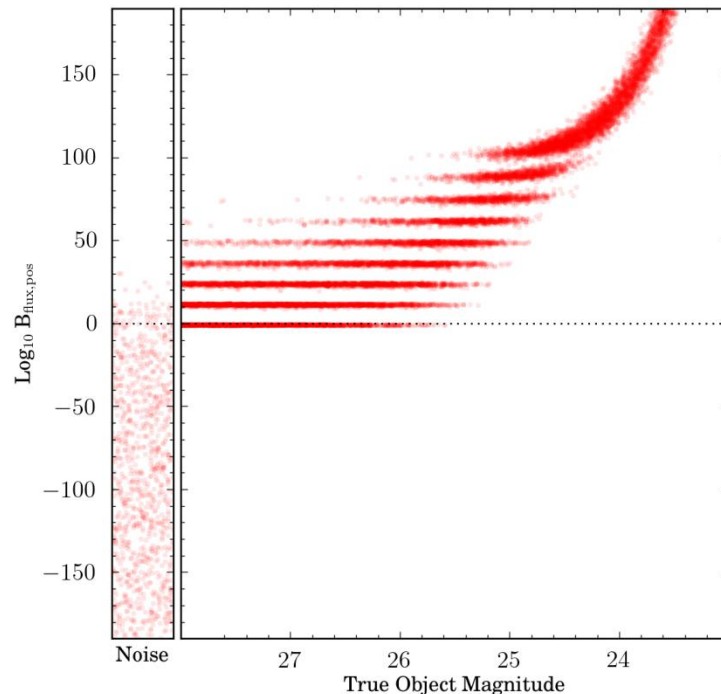  ☐ Object limit is at 28

☐ Bayes factor

  ☐ Considering only fluxes

# Adding Directions

☐ Bayes factor from cross-id

  ☐ As TB & Szalay (2008)

  ☐ Faint sources can be distinguished based on their celestial coordinates

*Always at "same" place!*

# Cross-Identification

- Hard problem
  - Computationally, Scientifically & Statistically
  - Need symmetric *n*-way solution
  - Need reliable quality measure

- Same or not?
  - Distance threshold? Maximum likelihood?

12/12/2013

# Same or Not?

**OR** ☐ The Bayes factor

$$B(H, K|D) \quad = \quad \frac{p(D|H)}{p(D|K)}$$

**SAME** ☐ *H:* all observations of the same object

**NOT** ☐ *K:* might be from separate objects

12/12/2013

# Same or Not?

**OR** ☐ The Bayes factor

$$B(H, K|D) \quad = \quad \frac{p(D|H)}{p(D|K)}$$

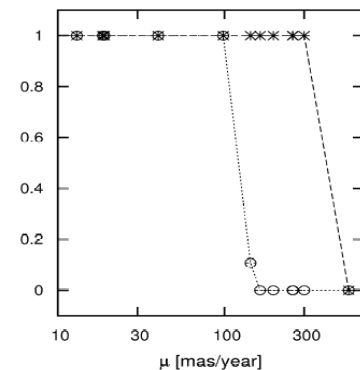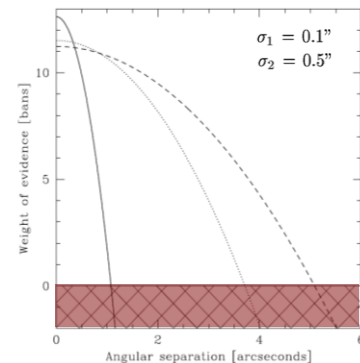**SAME** ☐ *H:* all observations of the same object

    ▫ Same properties, e.g., coordinates, brightness

**NOT** ☐ *K:* might be from separate objects

    ▫ Properties could be different

12/12/2013

# Works in General

□ Analytic results for Gaussian errors

  ▪ Incremental $n$-way strategy

□ We can find moving stars

  ▪ With unknown velocities

□ Matching events in time

  ▪ E.g., supernovae





2013

# SkyQuery – the new generation!

**Tamás Budavári**

- Dynamic federation of astronomy databases
  - Query the collection as if they were one

- The 3$^{rd}$ generation tool coming in December
  - Cluster of machines running partitioned jobs
  - Proper probabilistic exec with variable errors

# SkyQuery – the new generation!

□ Dynamic                                                 ses

■ Query

□ The 3rd g                                              ber

■ Cluster                                                s

■ Proper                                                rs

**SkyQuery**
🔒 username: budavari | account | sign out

| home | schema | query | jobs | my db | docs |

| syntax check | quick execute | execute | Output table: xmatch1 | Comments: |

```
1  SELECT s.ObjID, g.ObjID, t.ObjID, ...,
2      x.RA, x.Dec, x.LogBF
3  FROM SDSS:PhotoObjAll AS s
4      CROSS JOIN GALEX:PhotoObjAll AS g
5      CROSS JOIN TwoMASS:PhotoXSC AS t
6  XMATCH BAYESFACTOR AS x
7      EXIST s ON POINT(s.Cx, s.Cy, s.Cz), 0.1
8      EXIST g ON POINT(g.RA, g.Dec), 0.2
9      MAY   t ON POINT(t.RA, t.Dec), 0.5
10     HAVING LIMIT 1e6
11 WHERE s.Galaxy = 1
```
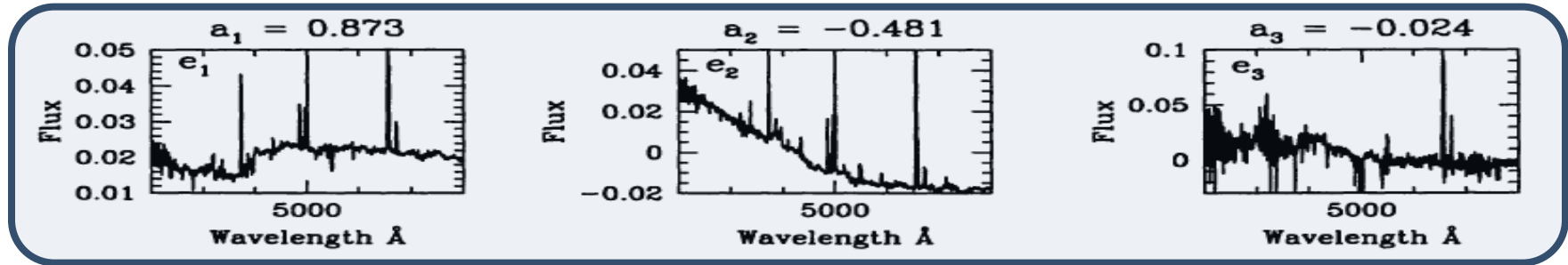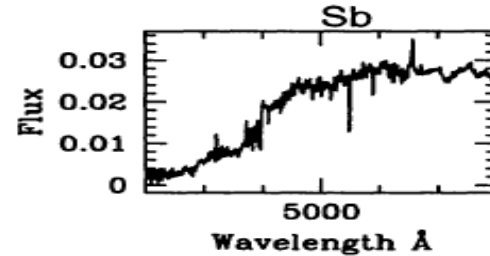
12/12/2013

# Only the first steps…

☐ Resolved shapes: radio morphology  *(Fan, TB+ 2014)*

☐ Colors to augment matches *(Marquez, TB, Sarro 2014)*
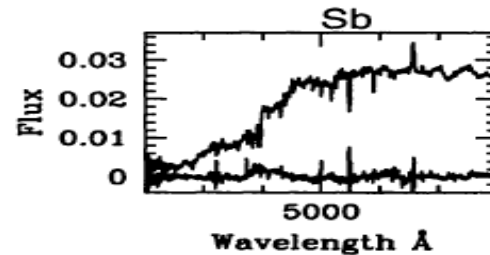
# Galaxy Light ~ Linear Combination


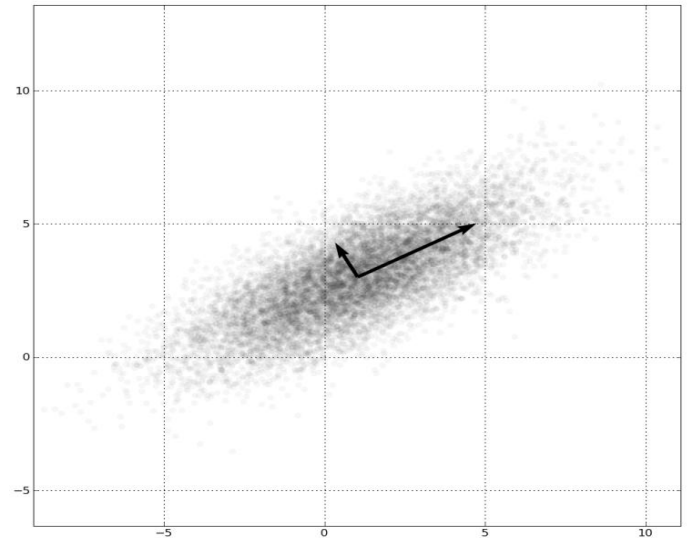
☐ Principal Components Analysis (PCA)

☐ SDSS galaxy type

☐ On a big memory machine

12/12/2013

# Principal Component Analysis

- Principal directions
  - Directions of largest variations
  - Eigenproblem of covariances
  - Singular Value Decomposition

- Problems
  - Needs lots of memory
  - Only need largest ones
  - Very sensitive to outliers



12/12/2013

# Science is Interactive

*"Too much to be accurate"*

By the time you do the calculations,
the answer may have changed...



12/12/2013

# Streams of Data

□ Mean

$$\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$$

$$\mu_n = \frac{n-1}{n} \mu_{n-1} + \frac{1}{n} x_n$$

$$\mu = \gamma \mu_{\text{prev}} + (1 - \gamma) x$$

12/12/2013

# Streams of Data

☐ Mean

☐ Covariance

$$\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$$

$$C = \gamma C_{\text{prev}} + (1 - \gamma) y y^{\text{T}}$$

$$\mu_n = \frac{n-1}{n} \mu_{n-1} + \frac{1}{n} x_n$$

$$y = x - \mu_{\text{prev}}$$

$$\mu = \gamma \mu_{\text{prev}} + (1 - \gamma) x$$

*Iterative evaluation!*

# Streaming PCA

**Tamás Budavári**

□ Initialization

    ◘ Eigensystem of a small, random subset

    ◘ Truncate at *p* largest eigenvalues
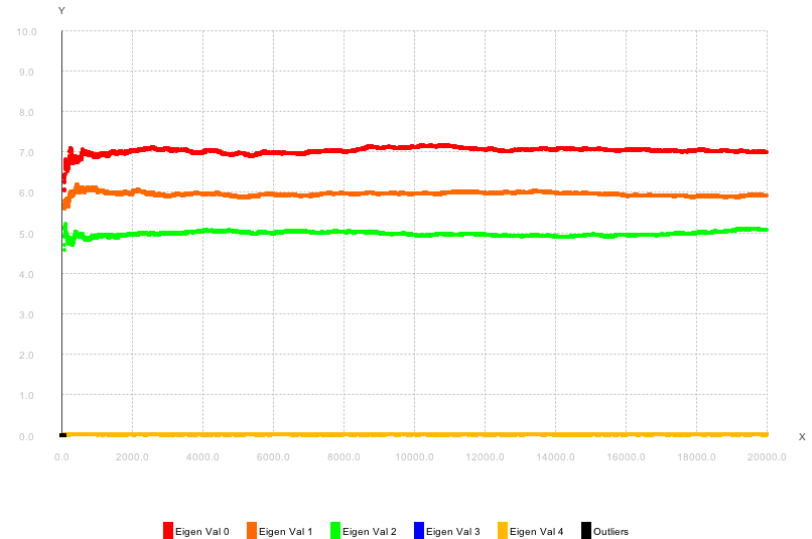
$$C \approx E_p \Lambda_p E_p^{\mathrm{T}}$$

□ Incremental updates

    ◘ Mean and the low-rank *A* matrix

    ◘ SVD of *A* yields new eigensystem

$$C \approx \gamma E_p \Lambda_p E_p^{\mathrm{T}} + (1 - \gamma) y y^{\mathrm{T}}$$

$$\approx A A^{\mathrm{T}}$$
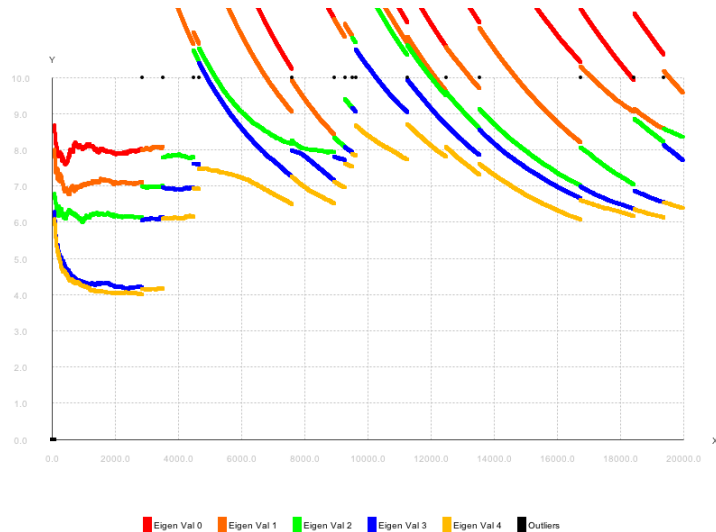
□ Randomized algorithm!

12/12/2013

# Streaming PCA

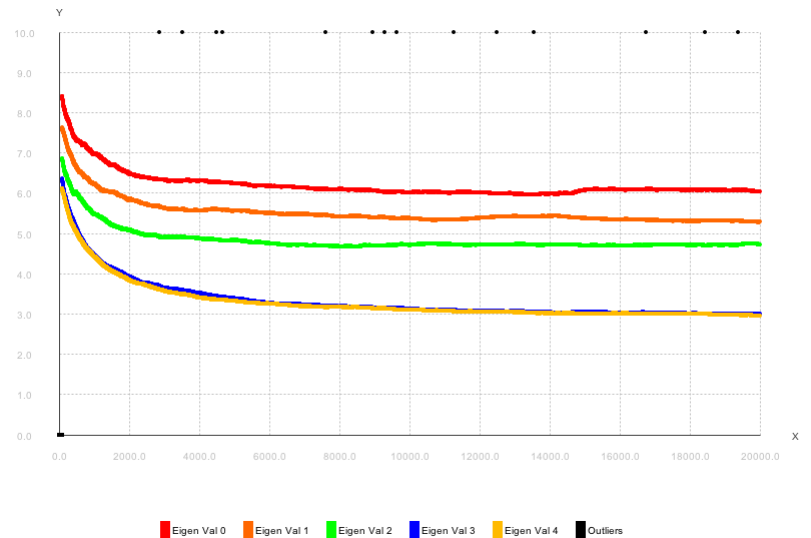☐ 3D Gaussian rotated into 50D

  ☐ Stretches: 7, 6, 5

  ☐ Total Var = 110

# With Outliers

- Adding 0.1% outliers
    - $\sigma = 100$ in each bin
- Outliers take over the PCs
    - Instability, no convergence



Eigen Val 0  Eigen Val 1  Eigen Val 2  Eigen Val 3  Eigen Val 4  Outliers

12/12/2013

# Robust Algorithm

☐ Outliers under control

  ◻ Marked on top

☐ Initialized with SVD

  ◻ On a set of 100 vectors

# Summary

- Plan for the junk
  - Proper statistics save money and gain speed
- Incremental randomized strategies scale
  - Crossmatching, embeddings, ML, etc.
- Not there, yet
  - Need new methods and tools

12/12/2013