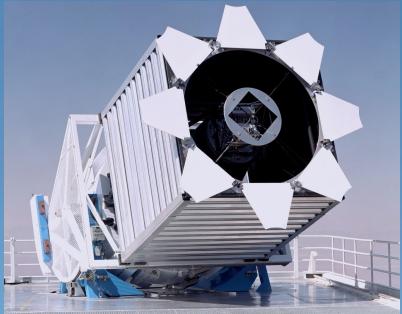


Big Data Sets in Astronomy - day 3

Željko Ivezić, University of Washington

LSST



SDSS



Gaia



Main Topics:

Day 1: Introduction

- who I think you are?
- who I am?
- why do astronomers need Big Data?
- Large Synoptic Survey Telescope: Big Data!
- astroML

Day 2: Density Estimation, Clustering and Classification in Astronomy

Day 3a: Dimensionality reduction, Regression and Time Series Analysis in Astronomy

Day 3b: Schedule reserve and free-form discussions

Outline

- Dimensionality Reduction
 - Principal Component Analysis
 - Non-negative Matrix Factorization
 - Independent Component Analysis
 - Manifold learning (Locally Linear Embedding)
- Regression
 - (Gaussian) errors in both variables
 - regression with non-Gaussian errors and/or outliers

Principal Component Analysis (PCA)

Motivation:

2D case: what is the direction of largest variance in the data?
Equivalently, what is the rotation of the coordinate system in which results in no covariance between the variables?

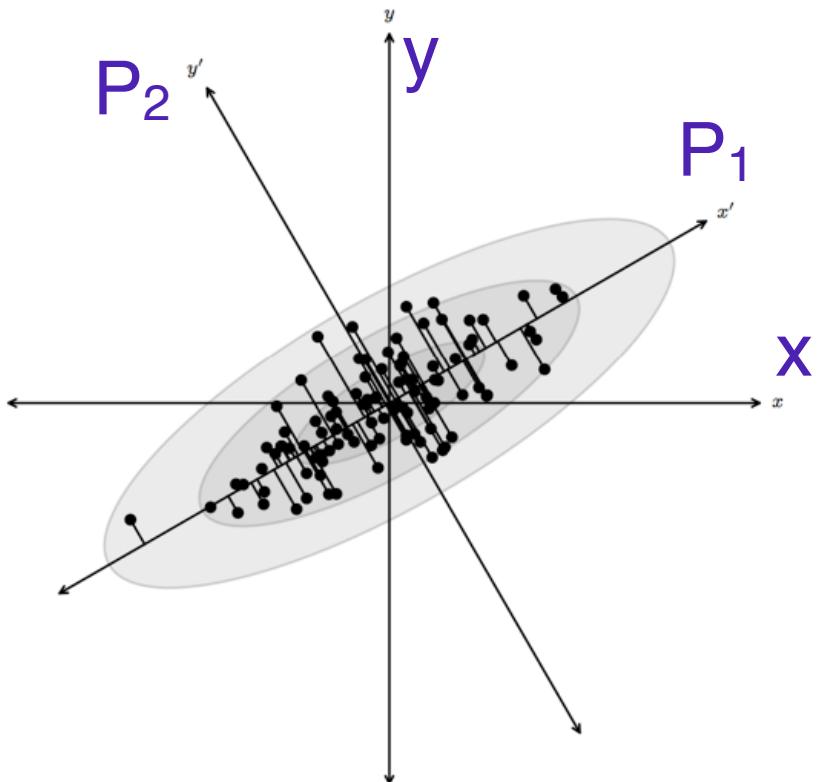


Figure 7.2.: A distribution of points drawn from a bivariate Gaussian and centered on the origin of x and y . PCA defines a rotation such that the new axes (x' and y') are aligned along the directions of maximal variance (the principal components) with zero covariance. This is equivalent to minimizing the square of the perpendicular distances between the points and the principal components.

The two variables x and y have a non-vanishing covariance; the covariance in the P_1 vs. P_2 coordinate system is 0 by construction (recall discussion of 2-D Gaussian in Lecture 1)

Principal Component Analysis (PCA)

Motivation:

High-D case (driven by the “curse of dimensionality”): there are ~4000 data points in SDSS spectra. Can we represent these spectra as linear combinations of **much smaller** number of eigen-components?

The curse of dimensionality:
in high-D space, the ratio of the volume of the unit hyper-sphere and the volume of the hyper-cube that encloses it goes to 0 as D increases

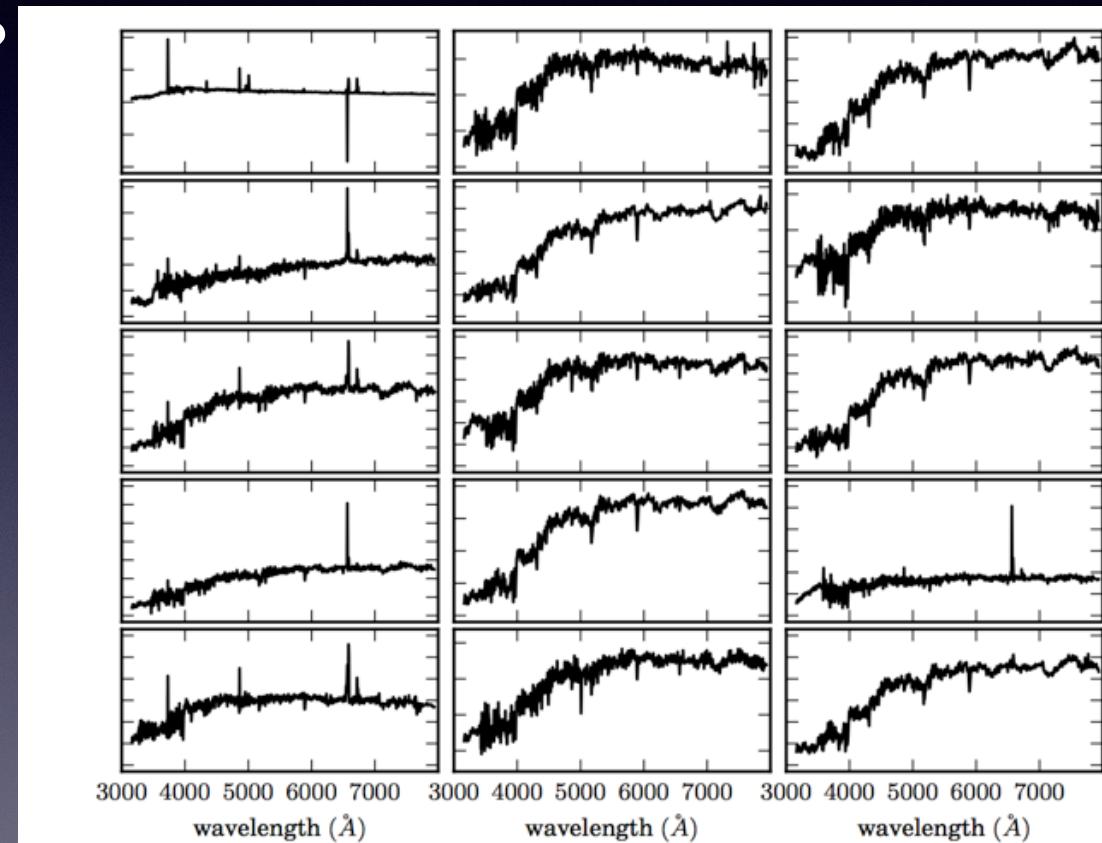


Figure 7.1.: A sample of 15 galaxy spectra selected from the SDSS spectroscopic data set (see §1.5.5). These spectra span a range of galaxy types, from star-forming to passive galaxies. Each spectrum has been shifted to its rest frame and covers the wavelength interval 3000–8000 Å. The specific fluxes, $F_\lambda(\lambda)$, on the ordinate axes have an arbitrary scaling.

Principal Component Analysis (PCA)

Each spectrum $\mathbf{x}_i(k)$ can be described by

$$P_1 = (x - \mu_x) \cos \alpha + (y - \mu_y) \sin \alpha,$$

$$P_2 = -(x - \mu_x) \sin \alpha + (y - \mu_y) \cos \alpha.$$

$$\mathbf{x}_i(k) = \boldsymbol{\mu}(k) + \sum_j^R \theta_{ij} \mathbf{e}_j(k), \quad (7.17)$$

where i represents the number of the input spectrum, j represents the number of the eigenspectrum, and, for the case of a spectrum, k represents the wavelength. Here, $\boldsymbol{\mu}(k)$ is the mean spectrum and θ_{ij} are the linear expansion coefficients derived from

$$\theta_{ij} = \sum_k \mathbf{e}_j(k) (\mathbf{x}_i(k) - \boldsymbol{\mu}(k)). \quad (7.18)$$

R is the total number of eigenvectors (given by the rank of X , $\min(N, K)$). If the summation is over all eigenvectors, the input spectrum is fully described with no loss of information. Truncating this expansion (i.e., $r < R$),

$$\mathbf{x}_i(k) = \sum_i^{r < R} \theta_i \mathbf{e}_i(k), \quad (7.19)$$

will exclude those eigencomponents with smaller eigenvalues. These components will, predominantly, reflect the noise within the data set. This is reflected in figure 7.5: truncating the recon-

Principal Component Analysis (PCA)

sometimes called Karhunen-Loeve and Hotelling transform

How is it done:

Data matrix is formed by N sets (objects) of K measured features (attributes); e.g. $K \sim 4000$ for SDSS spectra of, say, $N=100,000$ quasars; we subtract the mean of each feature (and sometimes normalized by their st. deviation, called whitening; in case of spectra, flux is normalized to 1 to avoid uninteresting correlations with source brightness)

Using matrix algebra, the first eigen-component is found by maximizing variance and other eigen-components by requiring covariance to vanish.

There are different approaches, whose performance depends on N and K .

Principal Component Analysis (PCA)

How is it done:

PCA can be performed easily using Scikit-learn:

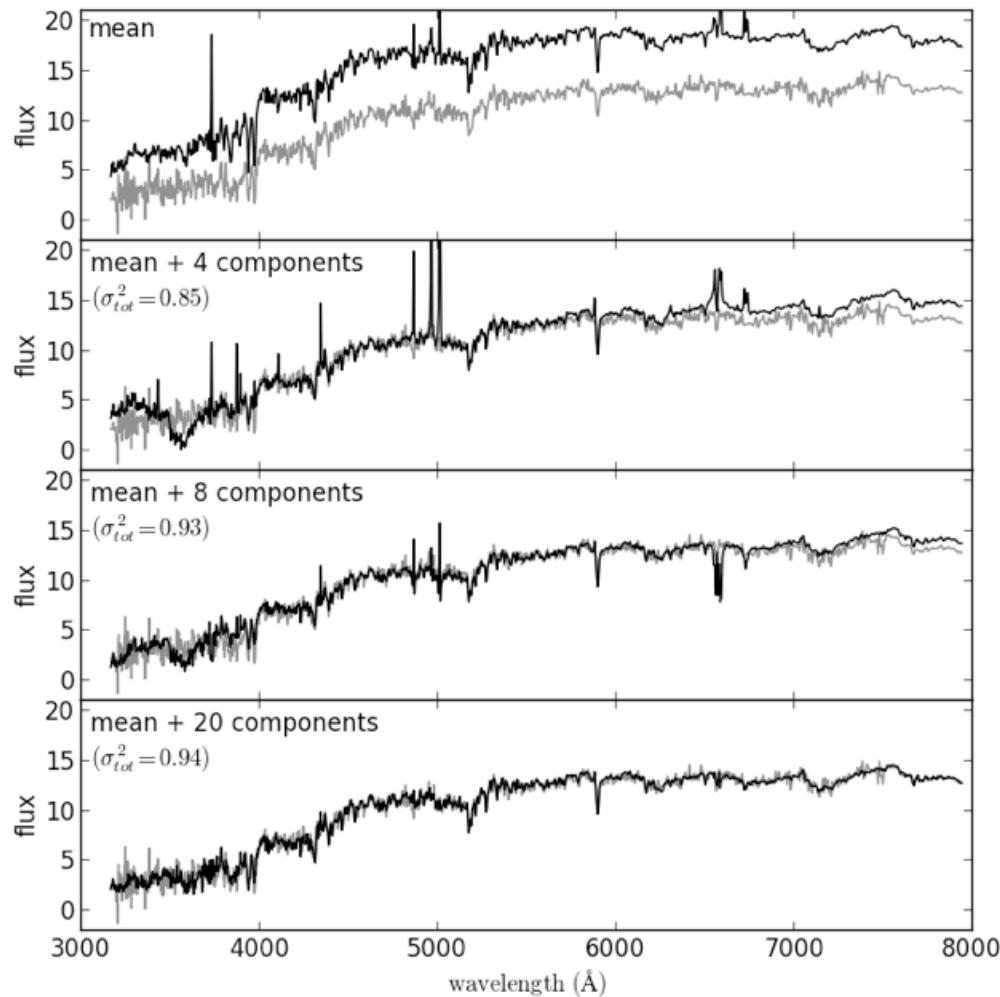
```
import numpy as np
from sklearn.decomposition import PCA

X = np.random.normal(size=(100, 3)) # 100 points in 3 dimensions
R = np.random.random((3, 10)) # projection matrix
X = np.dot(X, R) # X is now 10-dim, with 5 intrinsic dims
pca = PCA(n_components=4) # n_components can be optionally set
pca.fit(X)
comp = pca.transform(X) # compute the subspace projection of X

mean = pca.mean_ # length 10 mean of the data
components = pca.components_ # 4 x 10 matrix of components
var = pca.explained_variance_ # the length 4 array of eigenvalues
```

Principal Component Analysis (PCA)

- made this plot by running
%run fig_spec_reconstruction.py



How to choose the number of eigencomponents?

Mean spectrum

4 eigencomponents

8 eigencomponents

20 eigencomponents

Principal Component Analysis (PCA)

How to choose the number of eigencomponents?

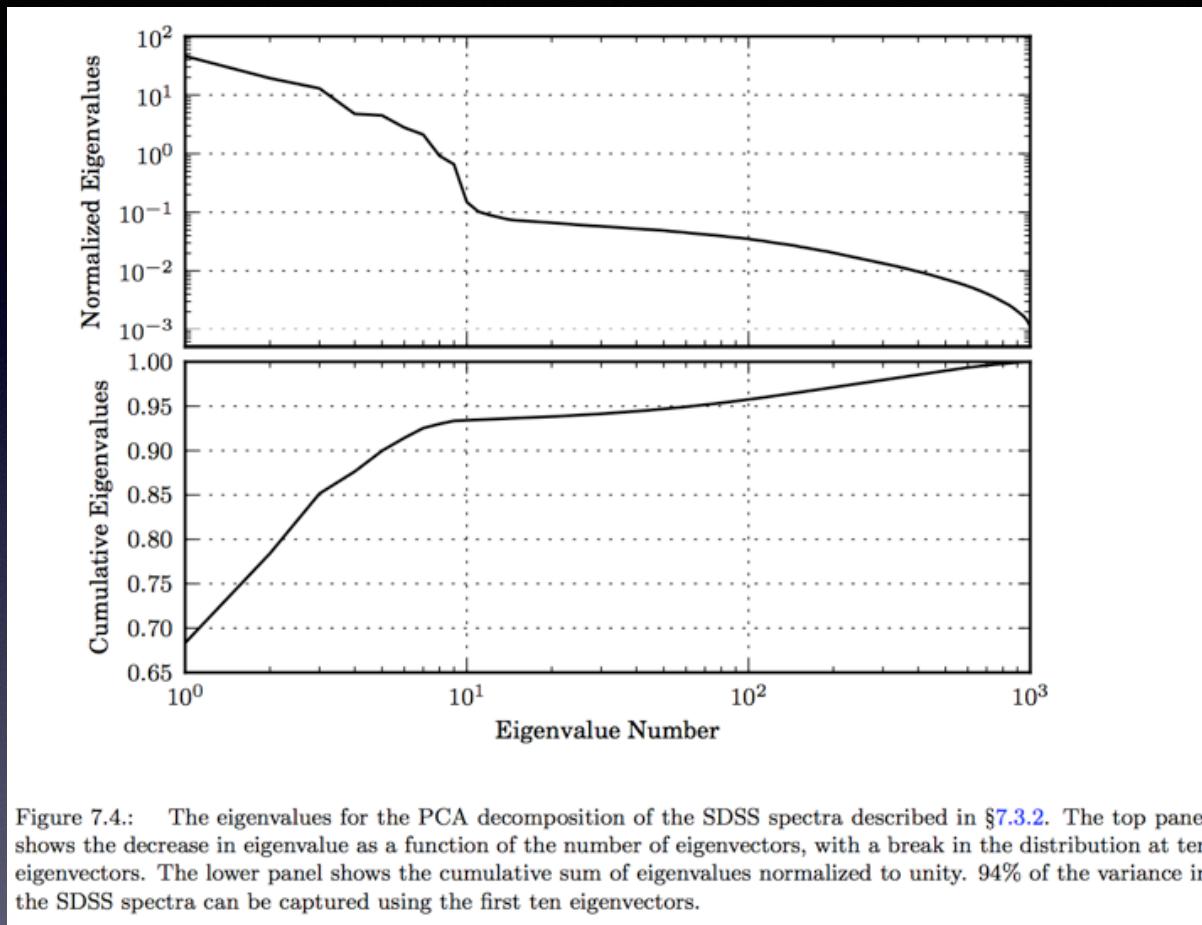


Figure 7.4.: The eigenvalues for the PCA decomposition of the SDSS spectra described in §7.3.2. The top panel shows the decrease in eigenvalue as a function of the number of eigenvectors, with a break in the distribution at ten eigenvectors. The lower panel shows the cumulative sum of eigenvalues normalized to unity. 94% of the variance in the SDSS spectra can be captured using the first ten eigenvectors.

There is no universal method, but typically we require that some fraction of data variance is captured by truncated series.

This plot (called the scree plot) shows that the first ten eigencomponents already capture about 94% of the variance in the dataset. Higher terms attempt to describe high-frequency measurement noise.

Principal Component Analysis (PCA)

A few more points about PCA

It can treat missing data by introducing weights for each point
(Connolly & Szalay, 1999, AJ 117, 2052).

If the dataset cannot fit in memory, then PCA computation can be challenging. One way to address this is the use of iterative online algorithms.

In some problems, linear decomposition might not be appropriate (e.g. a set of Planck functions with varying temperature)

Eigencomponents can be negative, which in some applications provides only limited insight into underlying physics (e.g. galaxy spectra).

Non-negative Matrix Factorization (NMF)

Attempts to address the fact that PCA eigencomponents can be negative even when we have physical reasons (a priori knowledge) to expect them to be non-negative

Assumes that the data matrix X is a product of two positive matrices, $X=Y^*W$

Solved iteratively for Y and W by minimizing the reconstruction error $\|X-WY\|^2$; sometimes problems with local minima are encountered (Lee & Seung, 2001)

We will see an example after introducing ICA

Independent Component Analysis (ICA)

Also known as the “cocktail party problem”, it assumes that the signal is a linear combination of components:

Each spectrum, $x_i(k)$, can now be described by

$$x_1(k) = a_{11}s_1(k) + a_{12}s_2(k) + a_{13}s_3(k) + \dots, \quad (7.34)$$

$$x_2(k) = a_{21}s_1(k) + a_{22}s_2(k) + a_{23}s_3(k) + \dots, \quad (7.35)$$

$$x_3(k) = a_{31}s_1(k) + a_{32}s_2(k) + a_{33}s_3(k) + \dots, \quad (7.36)$$

where $s_i(k)$ are the individual stellar spectra and a_{ij} the appropriate mixing amplitudes. In matrix format we can write this as

$$X = AS, \quad (7.37)$$

where X and S are matrices for the set of input spectra and stellar spectra, respectively. Extracting these signal spectra is equivalent to estimating the appropriate weight matrix, W , such that

$$S = WX. \quad (7.38)$$

This appears very similar to PCA, but...

Independent Component Analysis (ICA)

The principle that underlies ICA comes from the observation that the input signals, $s_i(k)$, should be statistically independent. Two random variables are considered statistically independent if their joint probability distribution, $f(x, y)$, can be fully described by a combination of their marginalized probabilities, that is,

$$f(x^p, y^q) = f(x^p)f(y^q), \quad (7.39)$$

where p and q represent arbitrary higher-order moments of the probability distributions. For the case of PCA, $p = q = 1$ and the statement of independence simplifies to the weaker condition of uncorrelated data (see §7.3.1 on the derivation of PCA).

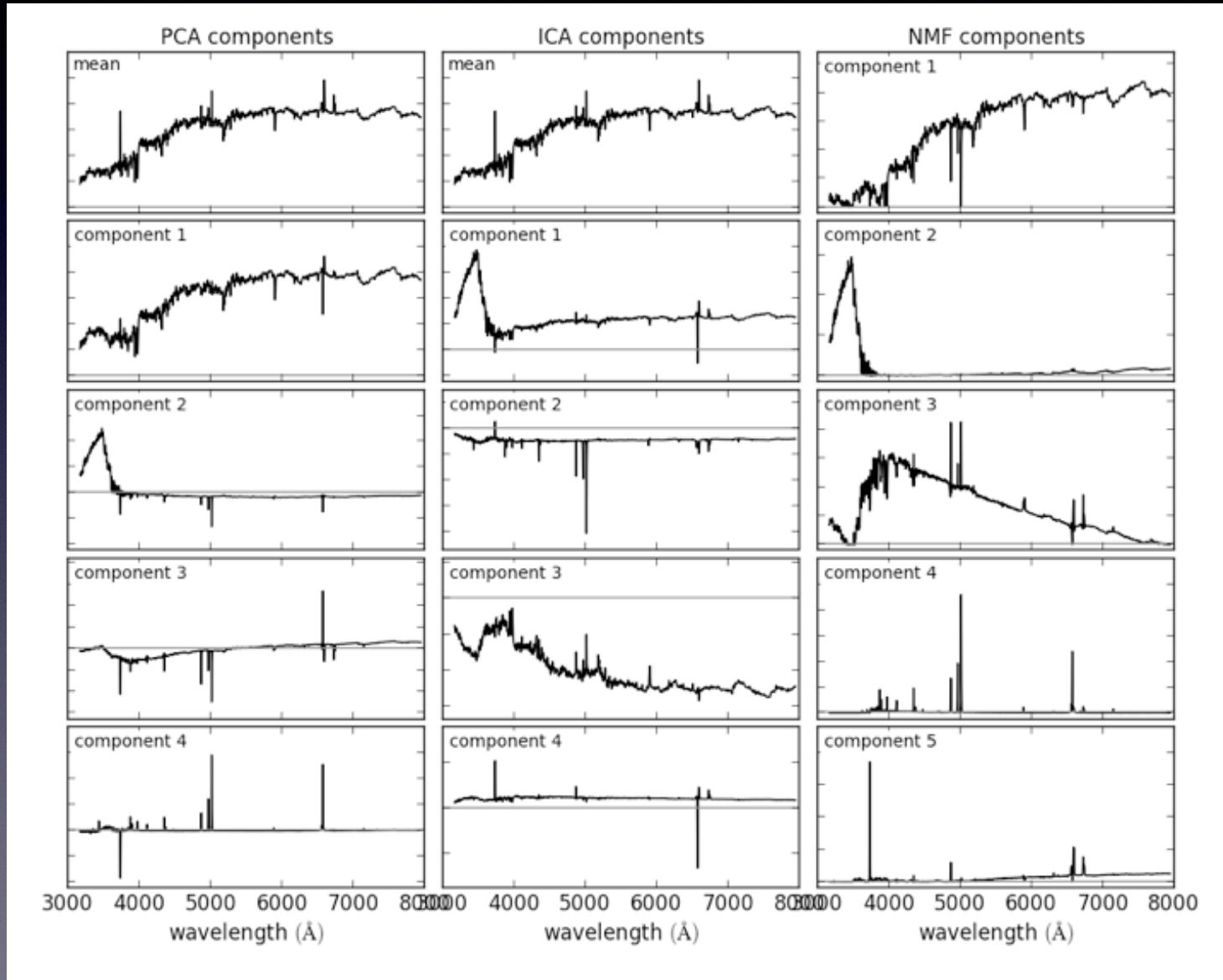
In most implementations of ICA algorithms the requirement for statistical independence is expressed in terms of the non-Gaussianity of the probability distributions. The rationale for this is that the sum of any two independent random variables will always be more Gaussian than either

Let's now compare PCA, NMF and ICA using our sample of SDSS galaxy spectra.

Comparison of PCA, ICA and NMF

Which one is
the most
astrophysical?

- made this plot by running (it takes a minute)
`%run fig_spec_decompositions.py`



Manifold Learning

What do we do if the assumption of linearity does not hold?

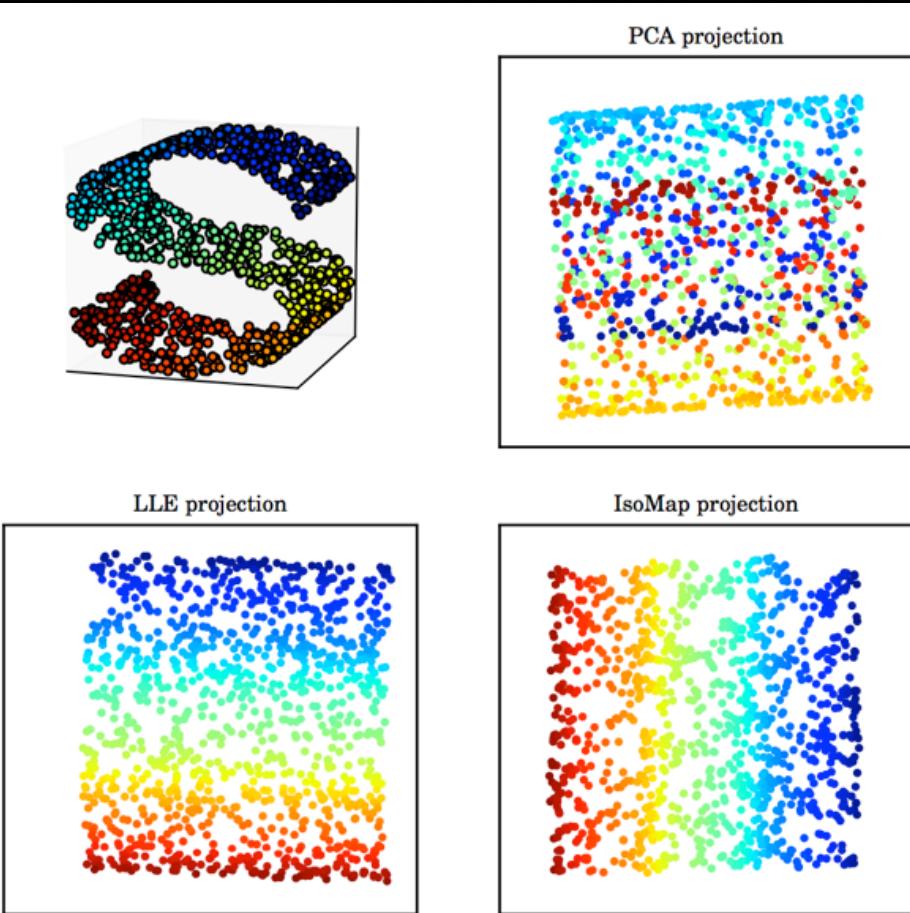


Figure 7.8.: A comparison of PCA and manifold learning. The top-left panel shows an example S-shaped data set (a two-dimensional manifold in a three-dimensional space). PCA identifies three principal components within the data. Projection onto the first two PCA components results in a mixing of the colors along the manifold. Manifold learning (LLE and IsoMap) preserves the local structure when projecting the data, preventing the mixing of the colors.

Here we have a 2D manifold in 3D space. PCA and other linear methods cannot uncover this structure.

Manifold learning techniques “unfold” or “unwrap” this manifold so that its structure becomes clear

Manifold Learning

Locally Linear Embedding

One of more popular techniques among a number of recent methods for non-linear dimensionality reduction.

In non-linear problems, it can have significantly better performance than PCA; e.g. it takes only two components to describe the same fraction of variance in SDSS galaxy spectra that requires dozens of PCA components (Vanderplas & Connolly 2009, AJ, 138, 1365)

The local manifold is determined by analyzing the nearest neighbor distribution around a given point; in some sense, one can think of a tangential hyper-plane approximation in high-D geometry

SDSS galaxy spectra

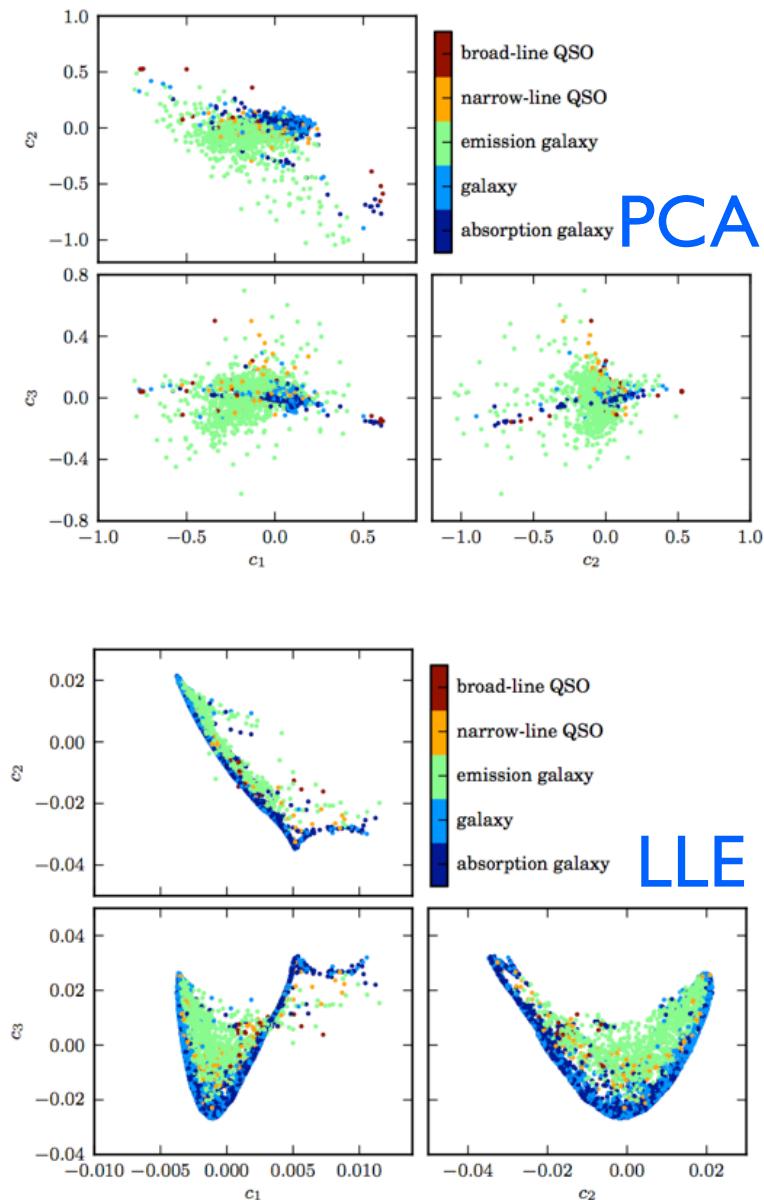


Figure 7.9.: A comparison of the classification of quiescent galaxies and sources with strong line emission using LLE and PCA. The top panel shows the segregation of galaxy types as a function of the first three PCA components. The lower panel shows the segregation using the first three LLE dimensions. The preservation of locality in LLE enables nonlinear features within a spectrum (e.g., variation in the width of an emission line) to be captured with fewer components. This results in better segregation of spectral types with fewer dimensions.

Similarly to the “S curve” example:

PCA: the structure in eigencoefficients is scrambled

LLE: the structure in eigencoefficients preserves information (from independent classification)

To reproduce: astroML, book figures, chapter 9

Which dimensionality reduction technique to use in practice?

Simple summary. Table 7.1 is a simple summary of the trade-offs along our axes of accuracy, interpretability, simplicity, and speed in dimension reduction methods, expressed in terms of high (H), medium (M), and low (L) categories.

Table 7.1.: Summary of the practical properties of the main dimensionality reduction techniques.

Method	Accuracy	Interpretability	Simplicity	Speed
Principal component analysis	H	H	H	H
Locally linear embedding	H	M	H	M
Nonnegative matrix factorization	H	H	M	M
Independent component analysis	M	M	L	L

Break here, go to notebook...

Regression

Motivation:

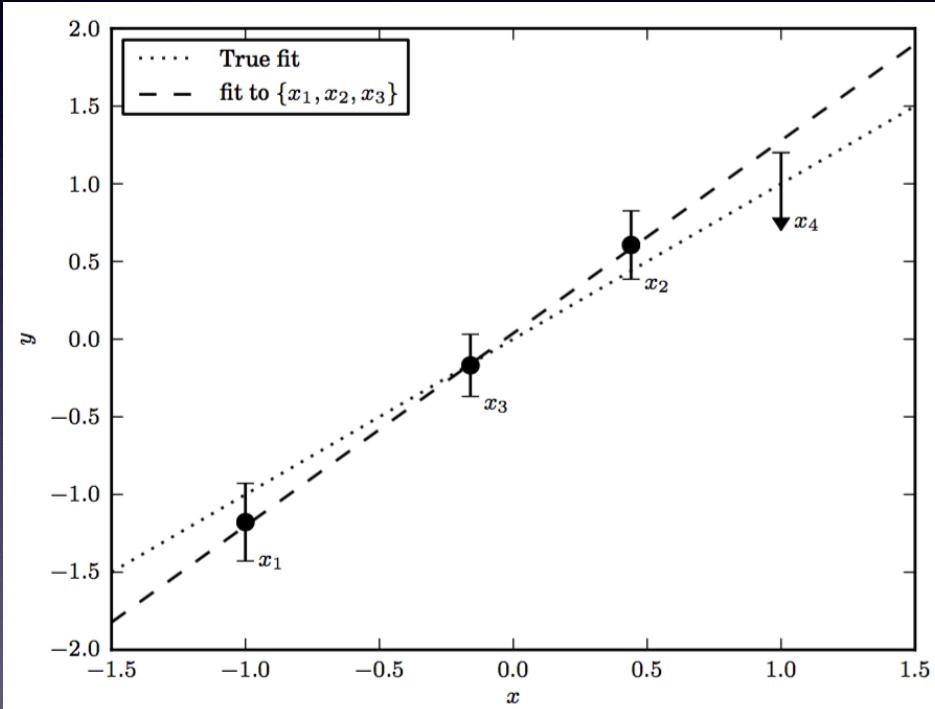
While e.g. least-square regression (LSQ) is often discussed, there are a few additional exceedingly useful and related tools in astroML for regression problems that often appear in practice:

- (Gaussian) errors in both variables
- regression with non-Gaussian errors and/or outliers

Least Squares Method is an application of the Maximum Likelihood method:

$$p(\{y_i\} | \{x_i\}, \theta, I) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(\frac{-(y_i - (\theta_0 + \theta_1 x_i))^2}{2\sigma_i^2}\right)$$

The origin of “Least Squares”: 2 comes from Gaussian errors



$$\begin{aligned}\theta_1 &= \frac{\sum_i^N x_i y_i - \bar{x}\bar{y}}{\sum_i^N (x_i - \bar{x})^2}, \\ \theta_0 &= \bar{y} - \theta_1 \bar{x},\end{aligned}$$

$$\begin{aligned}\sigma^2 &= \sum_{i=1}^N (y_i - \theta_0 + \theta_1 x_i)^2, \\ \sigma_{\theta_1}^2 &= \sigma^2 \frac{1}{\sum_i^N (x_i - \bar{x})^2}, \\ \sigma_{\theta_0}^2 &= \sigma^2 \left(\frac{1}{N} + \frac{\bar{x}^2}{\sum_i^N (x_i - \bar{x})^2} \right)\end{aligned}$$

Main assumptions:

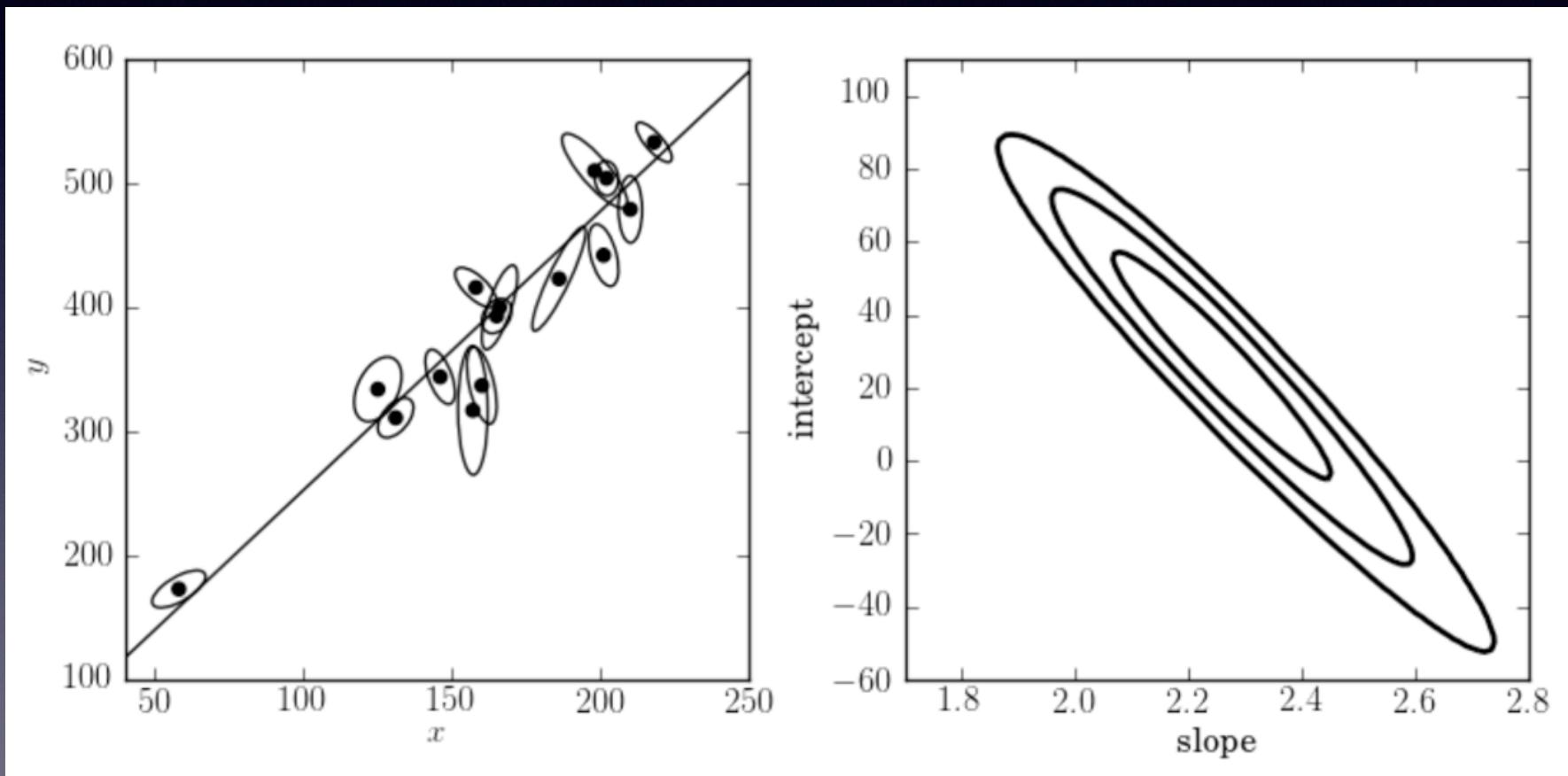
- 1) Gaussian errors for y
- 2) No errors in x

Gaussian errors in both variables

(see Hogg, Bovy & Lang, 2010, arXiv:1008.4686)

Example code:

[http://www.astroml.org/book_figures/chapter8/
fig_total_least_squares.html](http://www.astroml.org/book_figures/chapter8/fig_total_least_squares.html)



Regression with non-Gaussian errors and outliers

M estimators (M stands for “maximum-likelihood-type”) approach the problem of outliers by modifying the underlying likelihood estimator to be less sensitive than the classic L_2 norm. M estimators are a class of estimators that include many maximum-likelihood approaches (including least squares). They replace the standard least squares, which minimizes the sum of the squares of the residuals between a data value and the model, with a different function. Ideally the M estimator has the property that it increases less than the square of the residual and has a unique minimum at zero.

Huber loss function

An example of an M estimator that is common in robust regression is that of the Huber loss (or cost) function [9]. The Huber estimator minimizes

$$\sum_{i=1}^N e(y_i|y), \quad (8.65)$$

where

$$e(t) = \begin{cases} \frac{1}{2}t^2 & \text{if } |t| \leq c, \\ c|t| - \frac{1}{2}c^2 & \text{if } |t| \geq c, \end{cases} \quad (8.66)$$

Regression with non-Gaussian errors and outliers

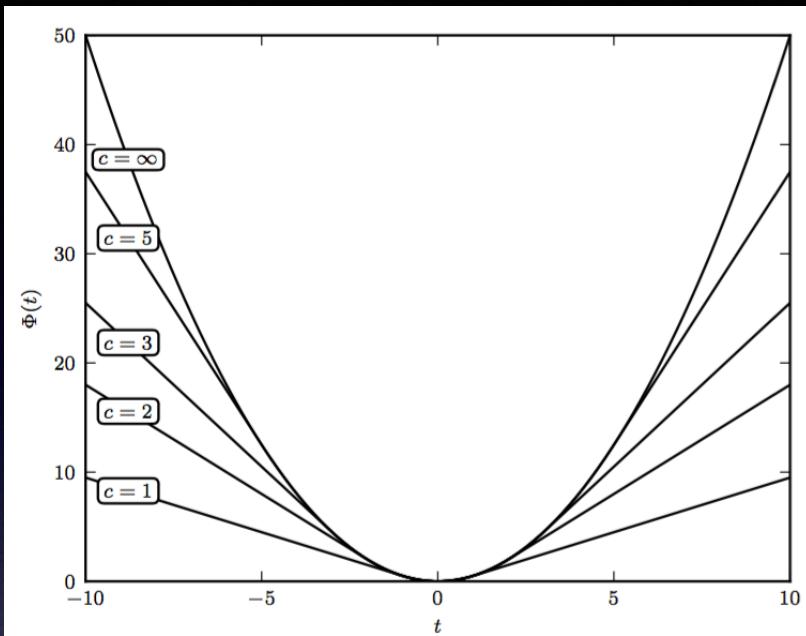


Figure 8.7.: The Huber loss function for various values of c .

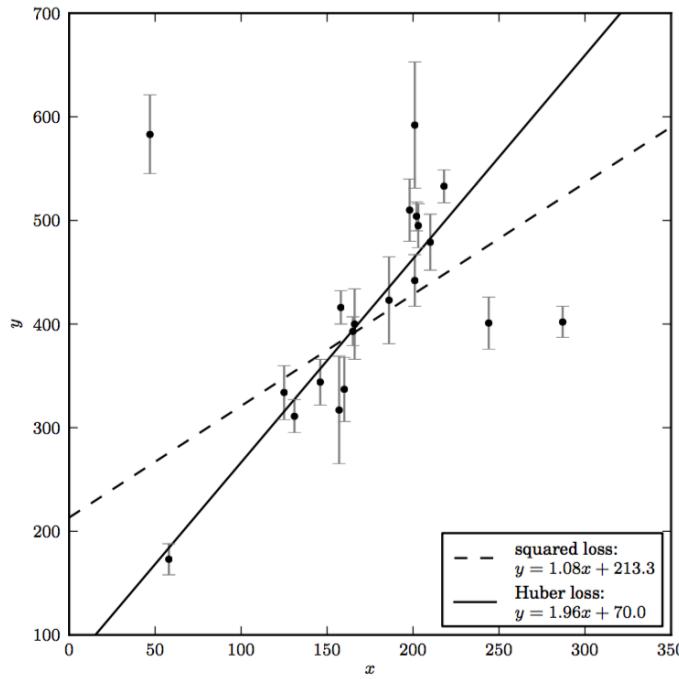


Figure 8.8.: An example of fitting a simple linear model to data which includes outliers (data is from table 1 of [8]). A comparison of linear regression using the squared-loss function (equivalent to ordinary least-squares regression) and the Huber loss function, with $c = 1$ (i.e., beyond 1 standard deviation, the loss becomes linear).

Regression with non-Gaussian errors and outliers

Bayesian mixture model:

distribution. The mixture model includes three additional parameters: μ_b and V_b , the mean and standard deviation of the background, and p_b , the probability that any point is an outlier. With this model, the likelihood becomes (cf. eq. 5.83; see also [8])

$$p(\{y_i\}|\{x_i\}, \{\sigma_i\}, \theta_0, \theta_1, \mu_b, V_b, p_b) \propto \prod_{i=1}^N \left[\frac{1-p_b}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(y_i - \theta_1 x_i - \theta_0)^2}{2\sigma_i^2}\right) + \frac{p_b}{\sqrt{2\pi(V_b + \sigma_i^2)}} \exp\left(-\frac{(y_i - \mu_b)^2}{2(V_b + \sigma_i^2)}\right) \right]. \quad (8.67)$$

Using MCMC sampling and marginalizing over the background parameters yields the dashed-line fit in figure 8.9. The marginalized posterior for this model is shown in the lower-left panel. This fit is much less affected by the outliers than is the simple regression model used above.

Finally, we can go further and perform an analysis analogous to that of §5.6.7, in which we attempt to identify bad points individually. In analogy with eq. 5.94 we can fit for nuisance parameters g_i , such that if $g_i = 1$, the point is a “good” point, and if $g_i = 0$ the point is a “bad” point. With this addition our model becomes

$$p(\{y_i\}|\{x_i\}, \{\sigma_i\}, \{g_i\}, \theta_0, \theta_1, \mu_b, V_b) \propto \prod_{i=1}^N \left[\frac{g_i}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(y_i - \theta_1 x_i - \theta_0)^2}{2\sigma_i^2}\right) + \frac{1-g_i}{\sqrt{2\pi(V_b + \sigma_i^2)}} \exp\left(-\frac{(y_i - \mu_b)^2}{2(V_b + \sigma_i^2)}\right) \right]. \quad (8.68)$$

This model is very powerful: by marginalizing over all parameters but a particular g_i , we obtain a posterior estimate of whether point i is an outlier. Using this procedure, the “bad” points have been marked with a circle in the upper-left panel of figure 8.9. If instead we marginalize over the

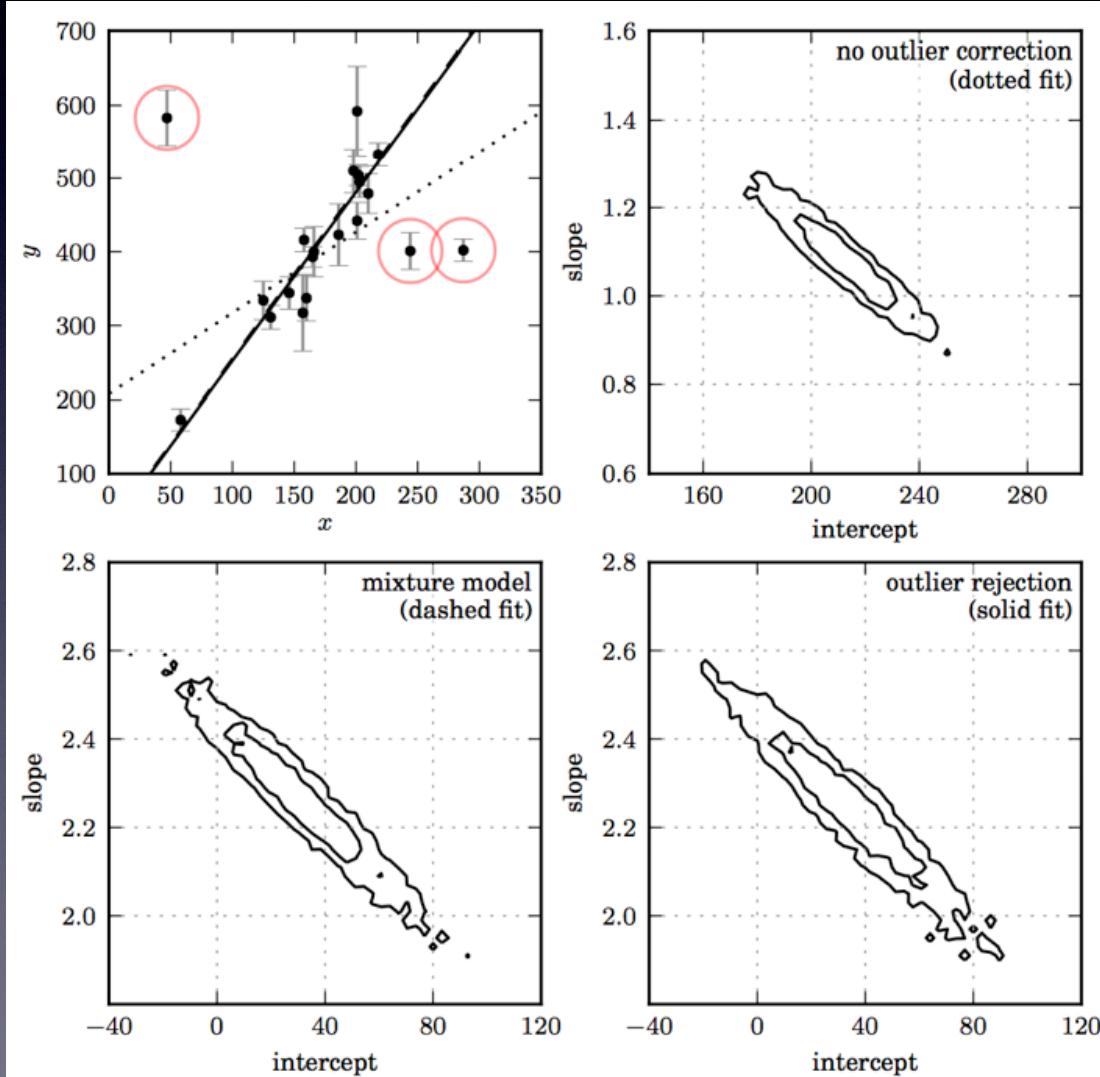
Regression with non-Gaussian errors and outliers (see Hogg, Bovy & Lang, 2010, arXiv:1008.4686)

Example code:

http://www.astroml.org/book_figures/chapter8/fig_outlier_rejection.html

The code uses MCMC

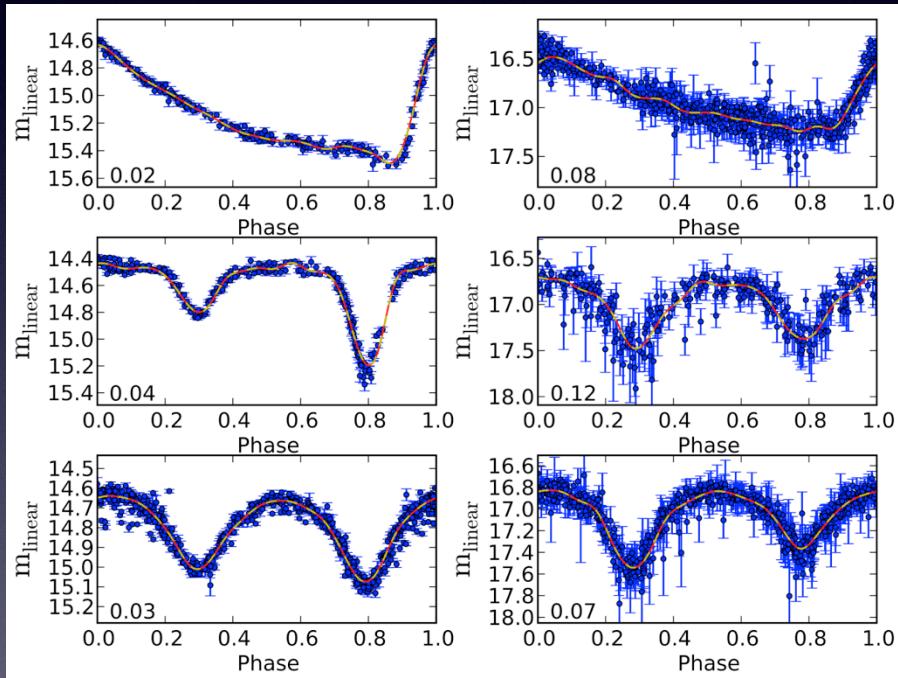
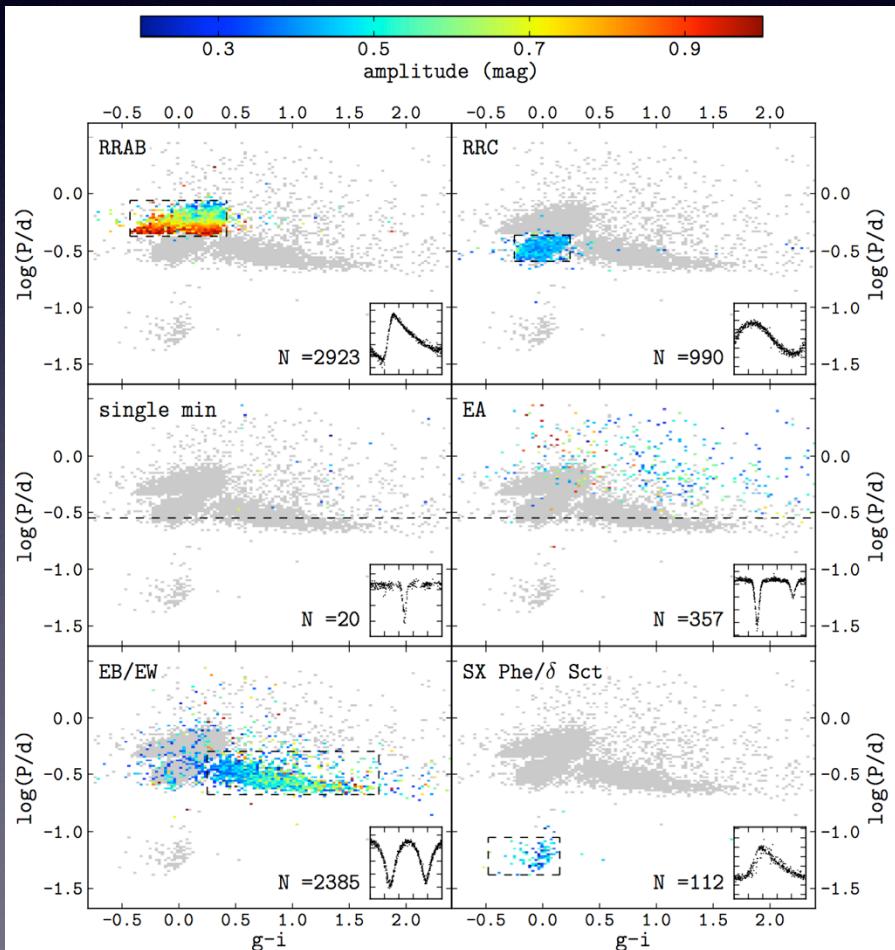
It's easy to change the outlier model (the mixture likelihood)...



Time series analysis

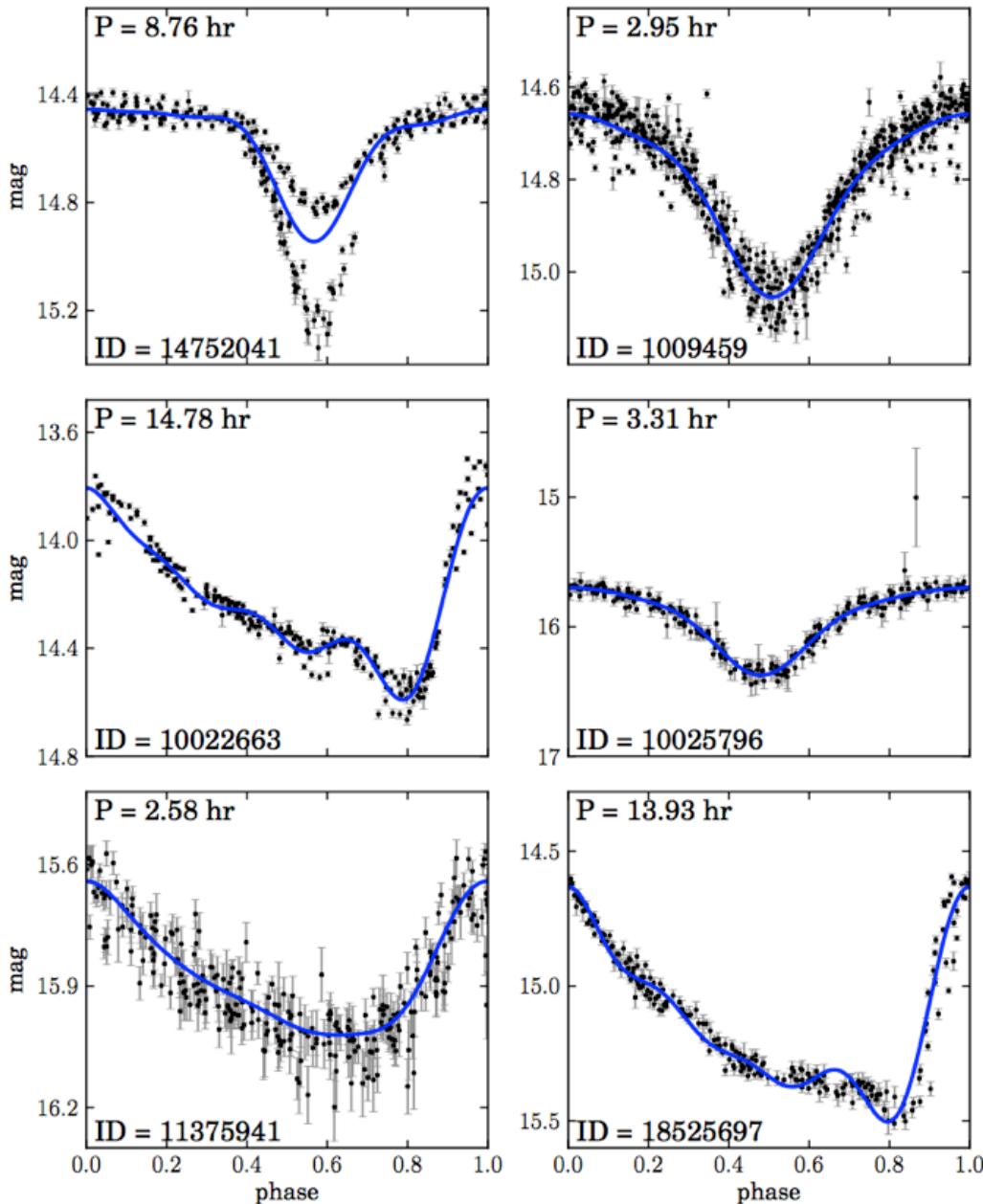
Motivation:

Use temporal changes to find interesting objects, and to learn more about the underlying physics of astrophysical objects (variable stars, exploding stars, quasars, comets...)



Above: light curves of variable stars; **Left:** their period-color diagrams

Time series analysis

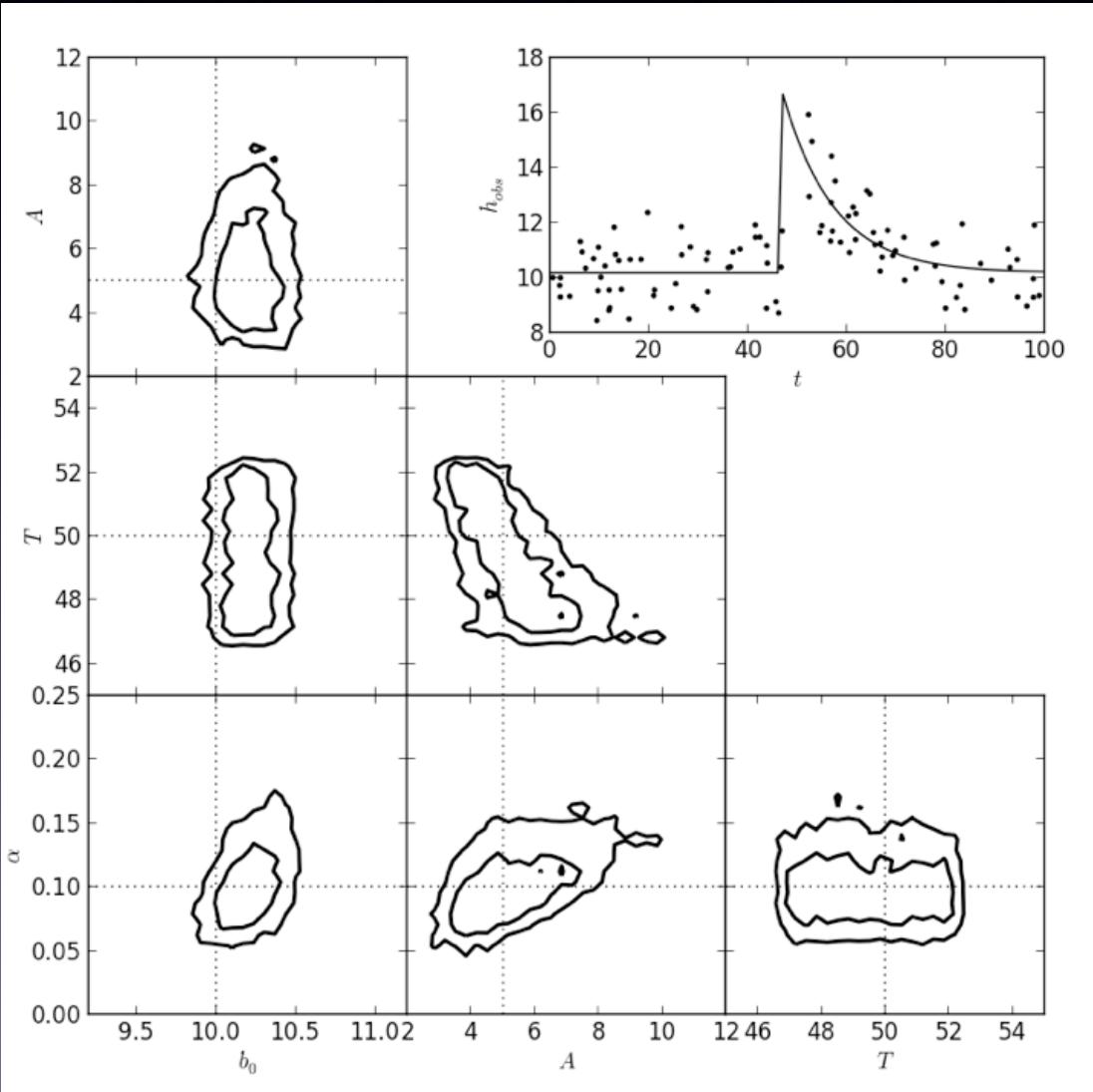


Light-curve analysis:

- 1) find the best period using periodogram
- 2) for this period, find the best-fit k-term Fourier expansion
- 3) compute the chi2, amplitude, mean flux, mean color, etc.

- **Matched filter analysis of a low-SNR signal**

- made this plot by running (note “pickling” of results):
%run fig_matchedfilt_burst.py



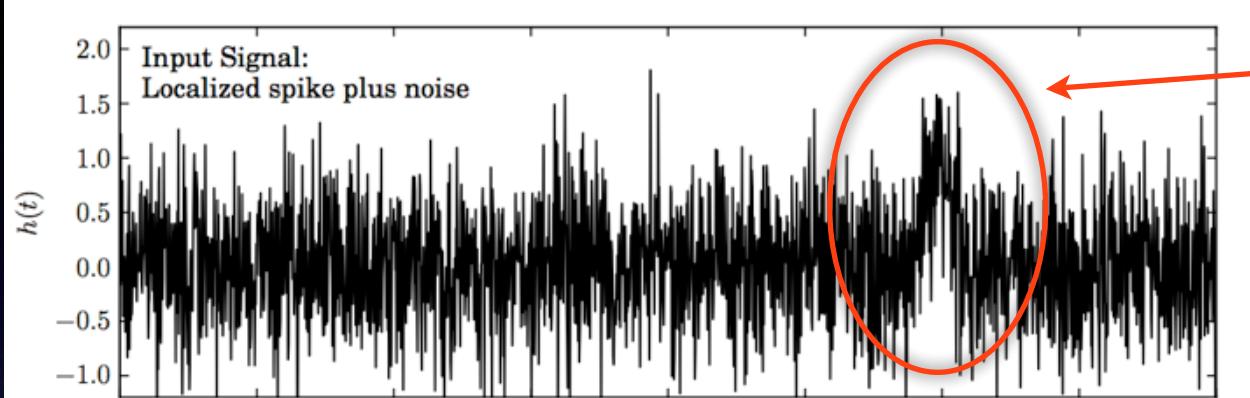
It calls pyMC module
Easy to change model
It tells you about
parameter covariances

- **Find a low-SNR burst (wavelet analysis)**
- What if we don't know the signal shape?
- And we know that the signal is not periodic.
- The discrete wavelet transform (DWT) can be used to analyze **the power spectrum of a time series as a function of time**. While similar analysis could be performed using the Fourier transform evaluated in short sliding windows, the DWT is superior. If a time series contains a localized event in time and frequency, DWT may be used to discover the event and characterize its power spectrum. **PyWavelets** is a toolkit with wavelet analysis implemented in Python.

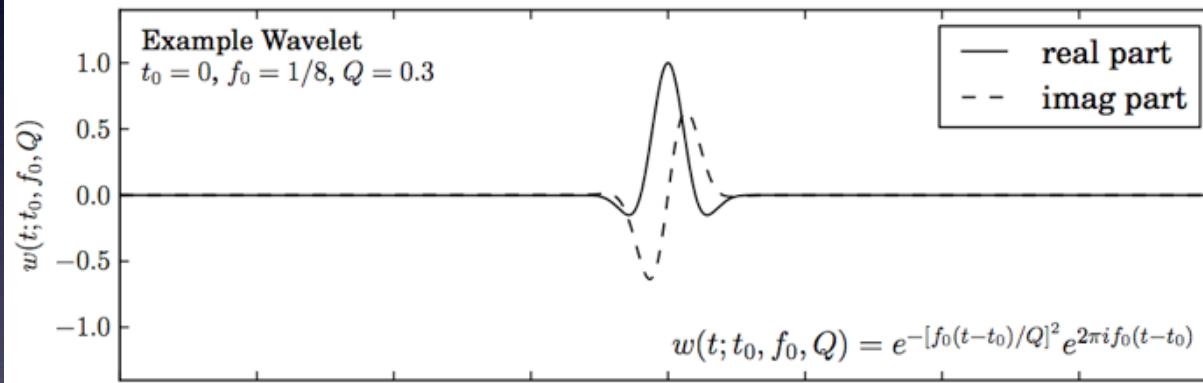
$$w(t|t_0, f_0, Q) = A \exp[i2\pi f_0(t - t_0)] \exp[-f_0^2(t - t_0)^2/Q^2]$$

- made this plot by running:

```
%run fig_line_wavelet_PSD.py
```



Signal



Found it!

