

Principal Component Analysis (PCA)

BIO337 Systems Biology / Bioinformatics – Spring 2014

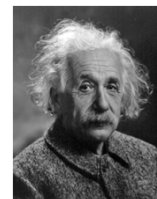
Edward Marcotte, Univ of Texas at Austin

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

What is Principal Component Analysis? What does it do?

So, first let's build some intuition.

“You do not really understand something unless you can explain it to your grandmother”, Albert Einstein



wikipedia

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

With thanks for many of these explanations to <http://stats.stackexchange.com/questions/2691/making-sense-of-principal-component-analysis-eigenvectors-eigenvalues>

What is Principal Component Analysis? What does it do?

A general (and imprecise) political example:

Suppose you conduct a political poll with 30 questions, each answered by 1 (*strongly disagree*) through 5 (*strongly agree*). Your data is the answers to these questions from many people, so it's 30-dimensional, and you want to understand what the major trends are.

You run PCA and discover 90% of your variance comes from one direction, corresponding not to a single question, but to a specific weighted combination of questions. This new hybrid axis corresponds to the political left-right spectrum, *i.e.* democrat/republican spectrum.

Now, you can study that, or factor it out & look at the remaining more subtle aspects of the data.

So, PCA is a method for discovering the major trends in data, simplifying the data to focus only on those trends, or removing those trends to focus on the remaining information.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

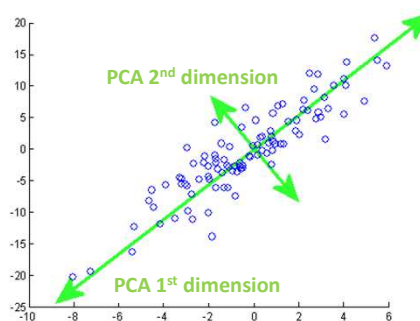
Example: Christian Bueno, <http://stats.stackexchange.com/questions/2691/making-sense-of-principal-component-analysis-eigenvectors-eigenvalues>

What is Principal Component Analysis? What does it do?

A more precise graphical example:

In a general sense, PCA rotates your axes to “line up” better with your data.

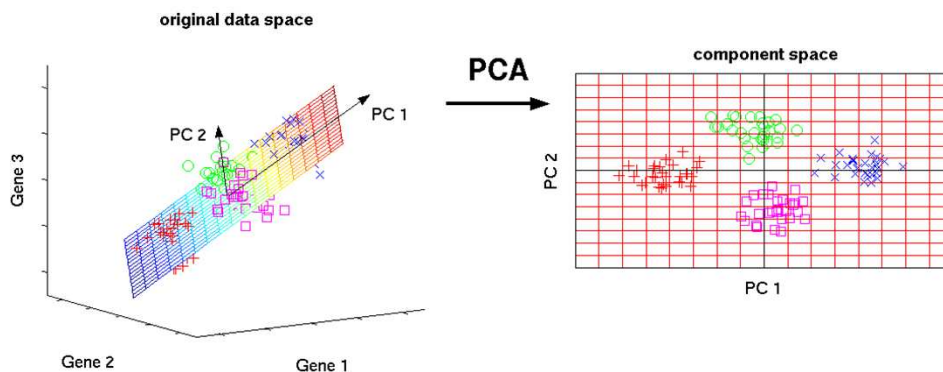
Because rotation is a kind of linear transformation, your new dimensions will be weighted sums of the old ones, like $\langle 1 \rangle = 23\% \cdot [1] + 46\% \cdot [2] + 39\% \cdot [3]$



Quotes & image adapted from isomorphisms, <http://stats.stackexchange.com/questions/2691/making-sense-of-principal-component-analysis-eigenvectors-eigenvalues>

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

PCA finds new variables which are linear combinations of the original variables such that in the new space, the data has fewer dimensions.



Think of a data set consisting of points in 3D on the surface of a flat plate held up at an angle. In the original x, y, z axes you need 3 dimensions to represent the data, but with the correct linear transformation, you only need 2.

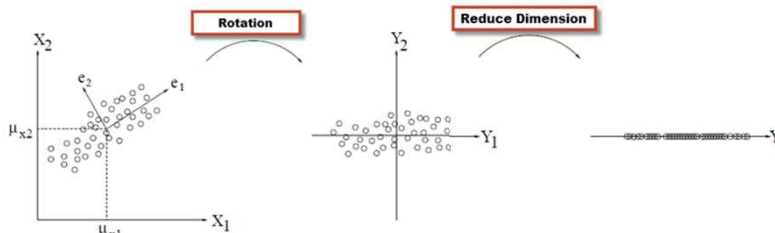
Quotes: Shlomo Argamon, from <http://stats.stackexchange.com/questions/2691/making-sense-of-principal-component-analysis-eigenvectors-eigenvalues>
Image: <http://phdthesis-bioinformatics-maxplanckinstitute-molecularplantphys.matthias-scholz.de/>

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

To summarize so far:

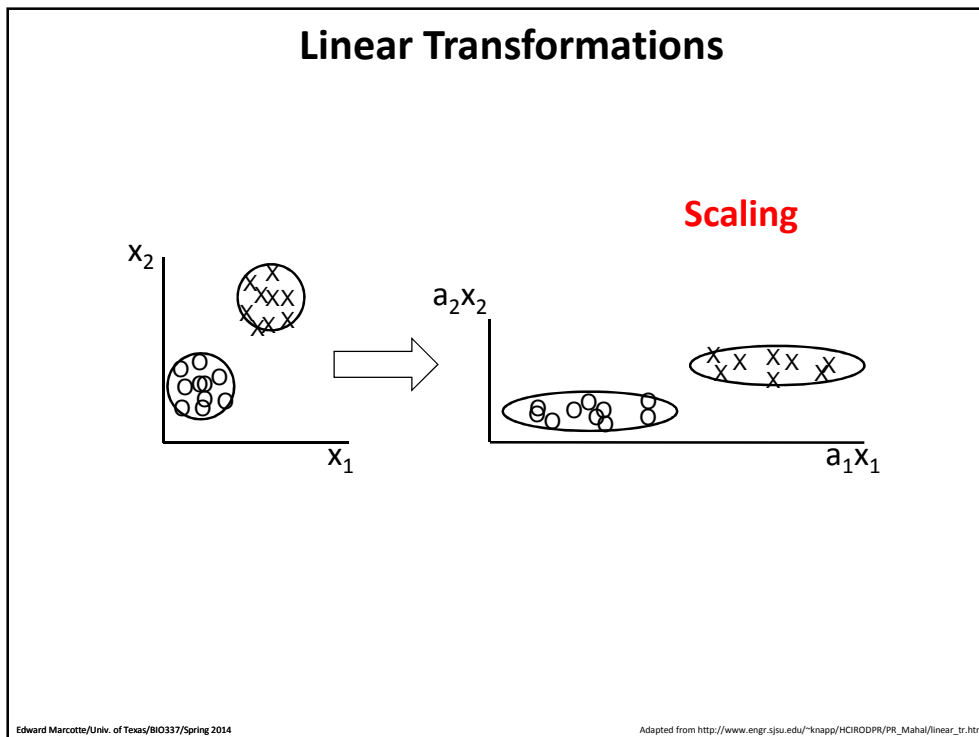
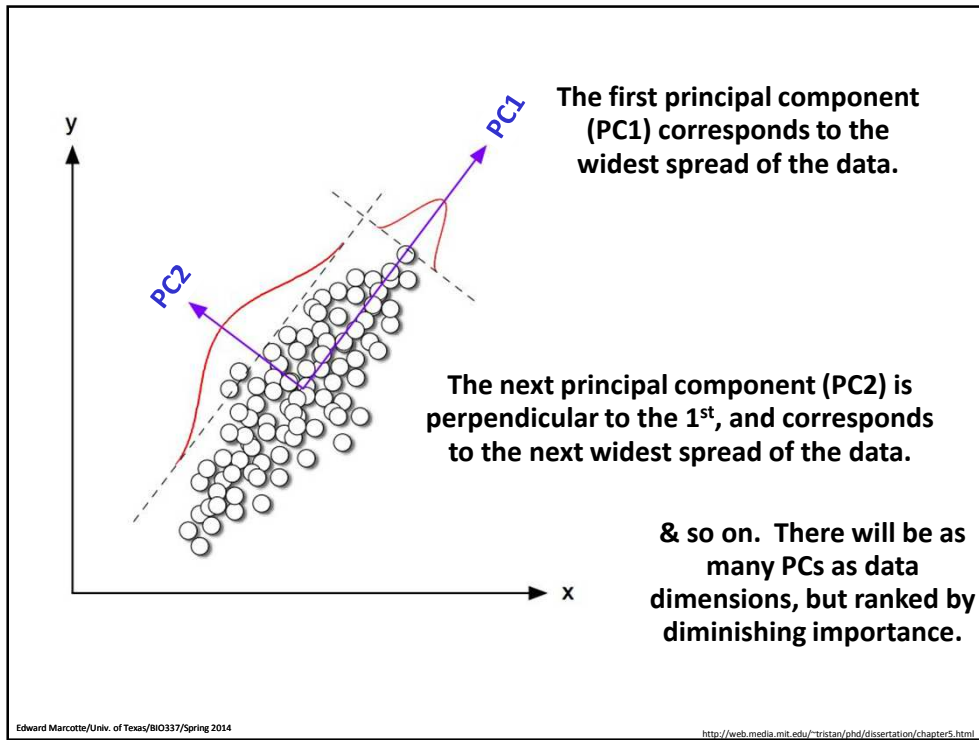
PCA is a technique to reduce dimension by:

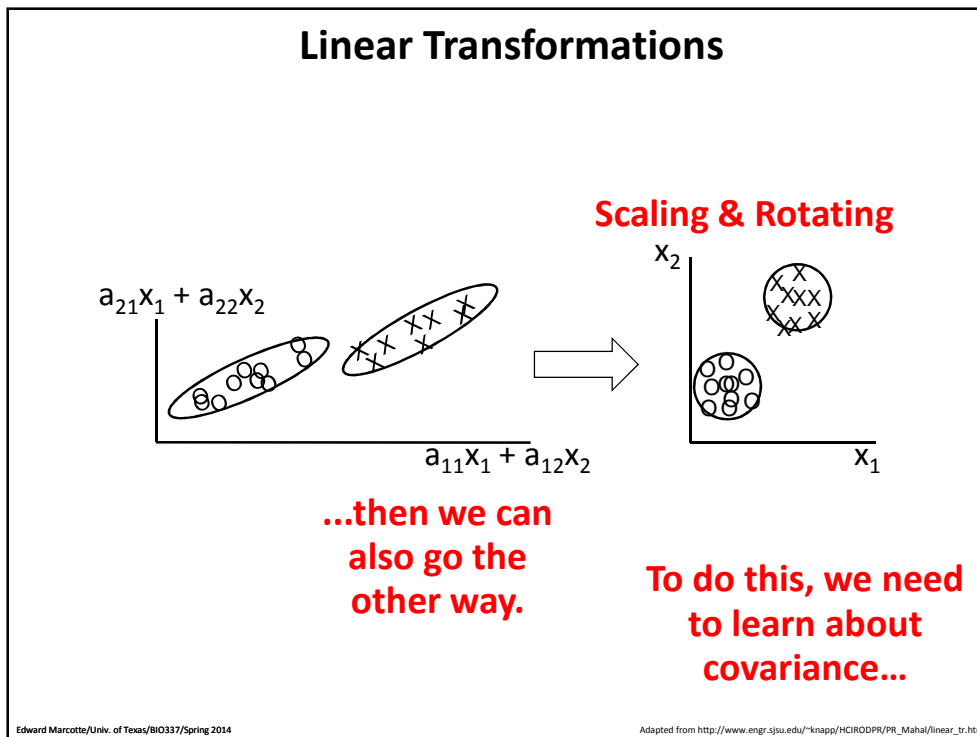
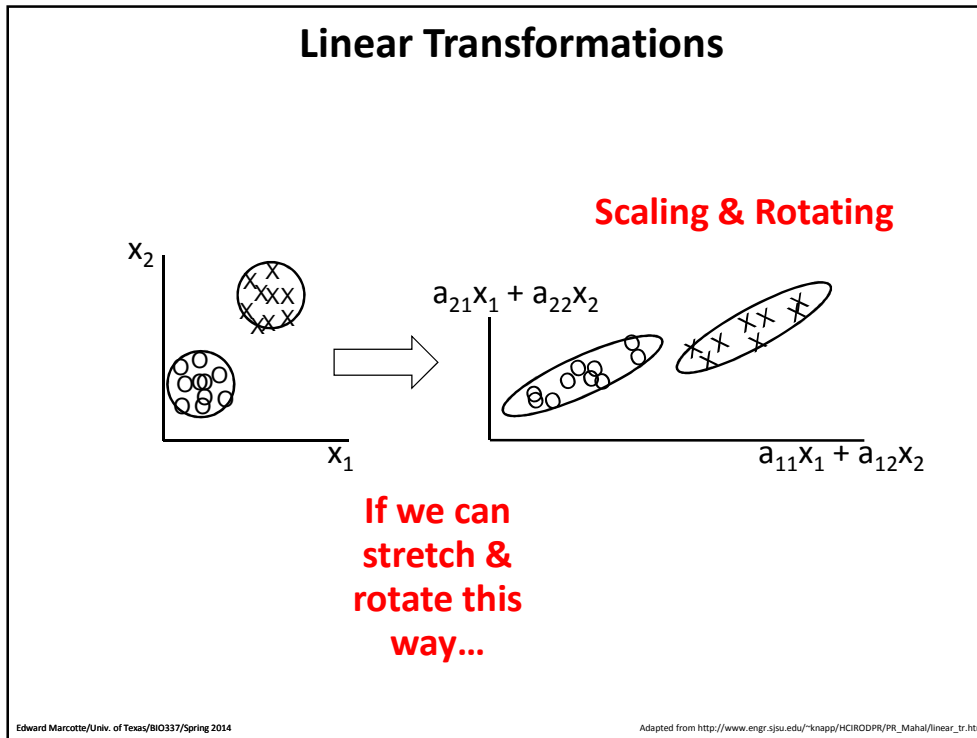
1. Taking **linear combinations** of the original variables.
2. Each linear combination explains the **most variance** in the data it can.
3. Each linear combination is uncorrelated (**orthogonal**) with the others
4. Plot the data in terms of only the most important (**principal**) dimensions



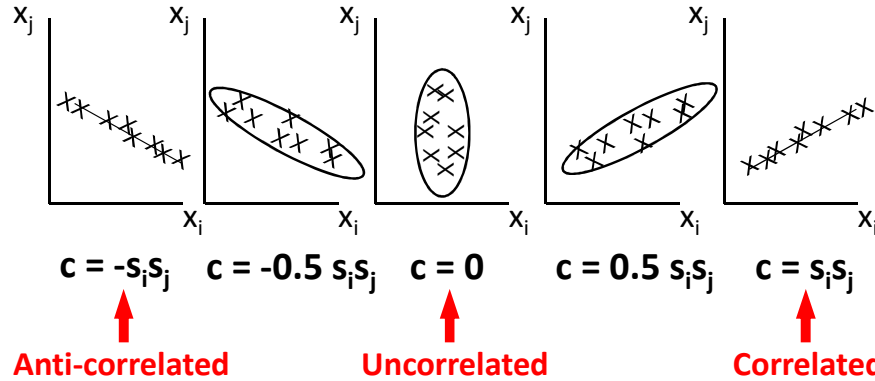
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Image: <http://blog.equametrics.com/2013/02/an-introduction-to-principal-component-analysis>
Quote: <http://stats.stackexchange.com/questions/2691/making-sense-of-principal-component-analysis-eigenvectors-eigenvalues>





Covariance = measures tendency to vary together (to co-vary). Similar to correlation.



Variance = average of the squared deviations of a feature from its mean
Covariance = average of the products of the deviations of feature values from their means

$$c(i,j) = \frac{[x(1,i) - m(i)] [x(1,j) - m(j)] + \dots + [x(n,i) - m(i)] [x(n,j) - m(j)]}{(n - 1)}$$

Adapted from http://www.engr.sjsu.edu/~knapp/HCIRODPR/PR_Mahal/cov.htm

All of the **covariances** $c(i,j)$ between features can be collected together into a **covariance matrix C**.

This summarizes all of the correlation structure among all pairs of features.

$$C = \begin{bmatrix} c(1,1) & c(1,2) & \dots & c(1,n) \\ c(2,1) & c(2,2) & \dots & c(2,n) \\ \vdots & \vdots & & \vdots \\ c(n,1) & c(n,2) & \dots & c(n,n) \end{bmatrix}$$

covariance between feature 1 and 2
 covariance between feature 2 and n
 etc.

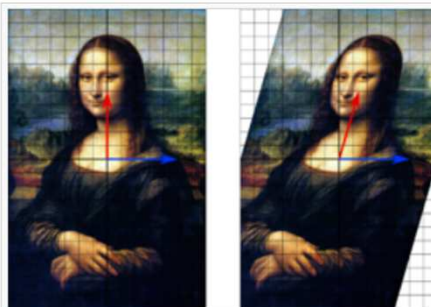
We need one last concept: Eigenvectors and eigenvalues

An **eigenvector** of a square matrix A is a non-zero vector v that, when the matrix is multiplied by v , yields a constant multiple of v , the multiplier being commonly denoted by λ . That is:

$$Av = \lambda v$$

(Because this equation uses **post-multiplication** by v , it describes a **right eigenvector**.)

The number λ is called the **eigenvalue** of A corresponding to v .

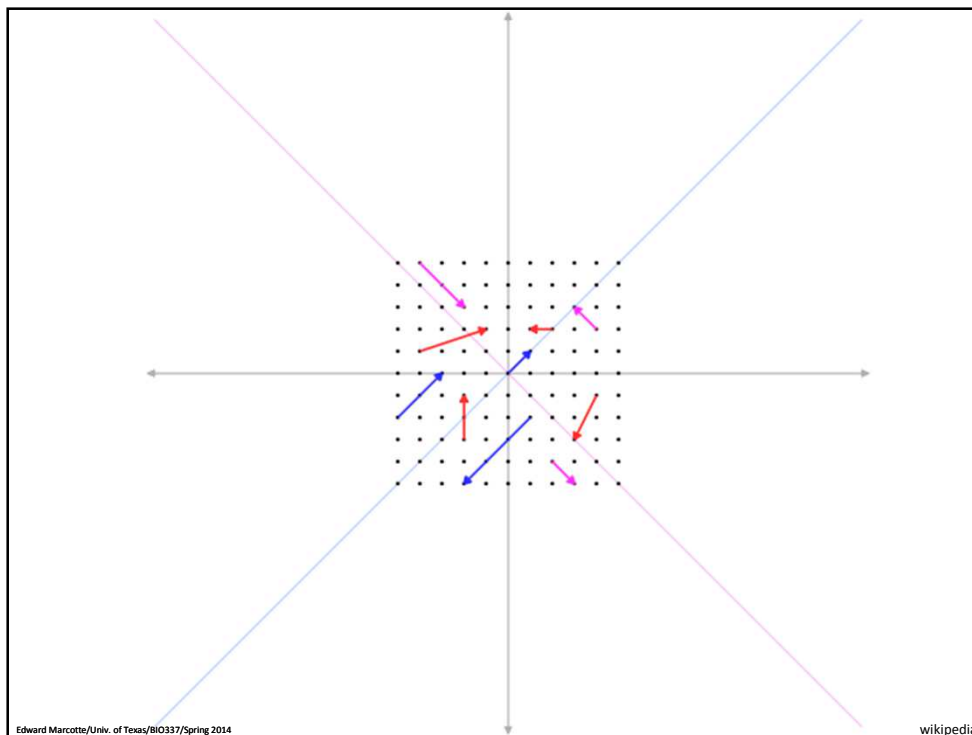


The **blue arrow** is an eigenvector of this linear transformation matrix, since it doesn't change direction.

Its scale is also unchanged, so its eigenvalue is 1.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Wikipedia



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

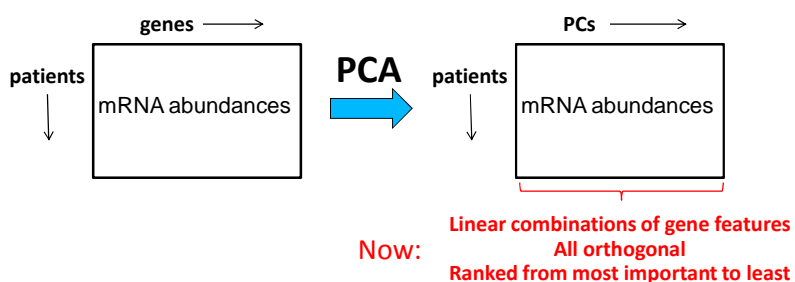
wikipedia

Calculating the PCA

1. Calculate the covariation matrix C between features of the data
2. Calculate the eigenvectors and eigenvalues of C
3. Order the eigenvectors according to the eigenvalues

PC1 is the eigenvector corresponding to the largest eigenvalue, PC2 is the eigenvector corresponding to the next largest, etc.

The data can be plotted as projections along the PCs of interest.



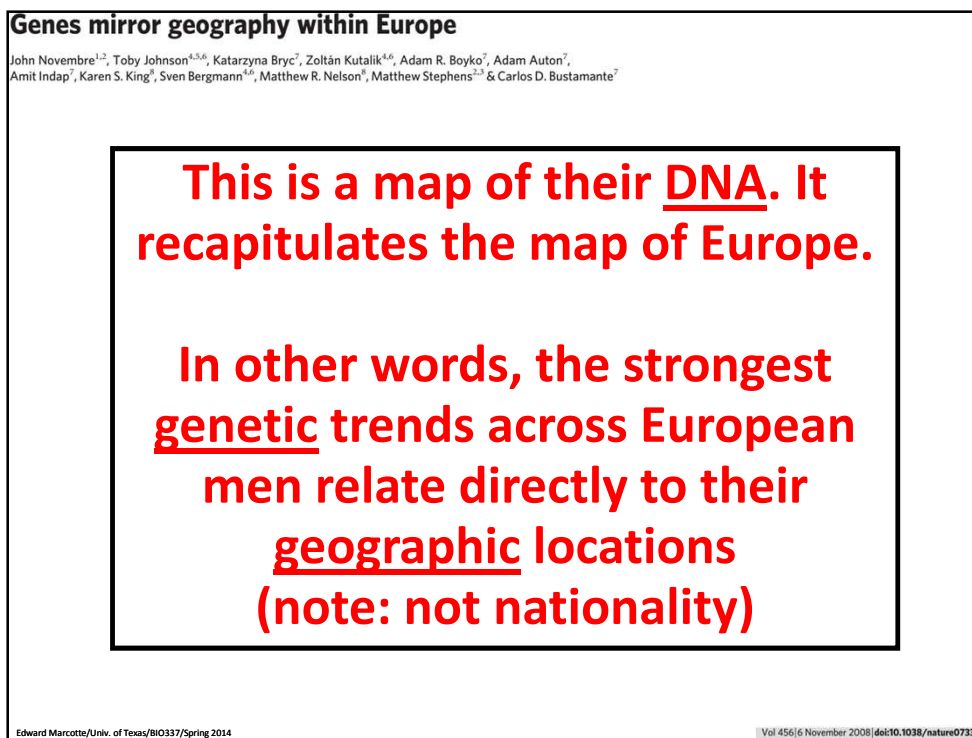
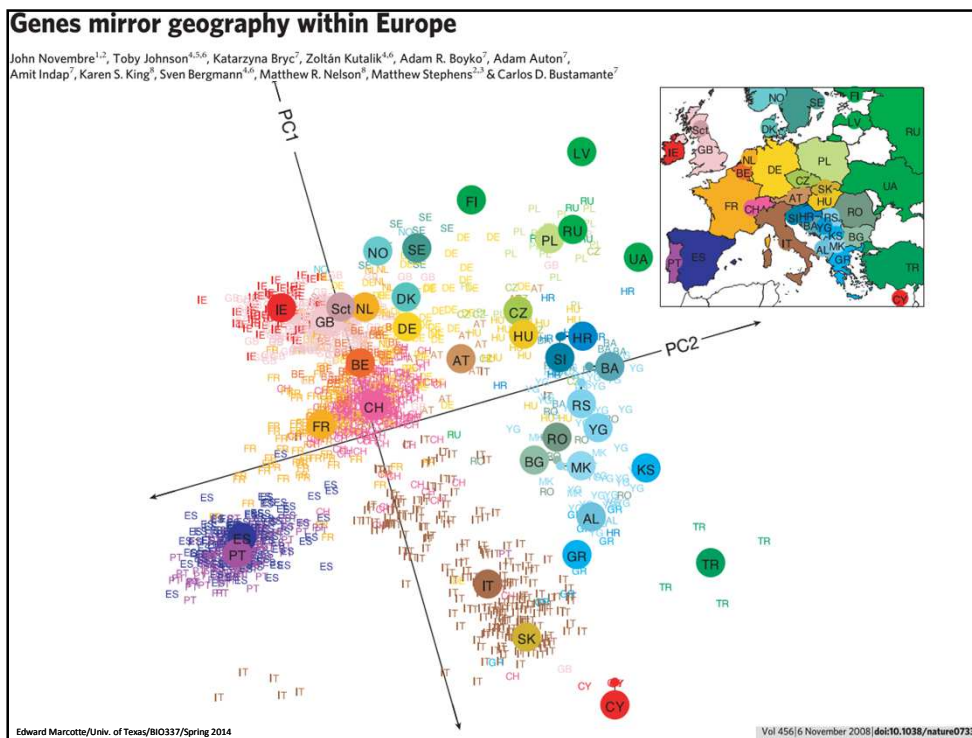
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Here's a beautiful application of PCA that illustrates how it can simplify complicated datasets:

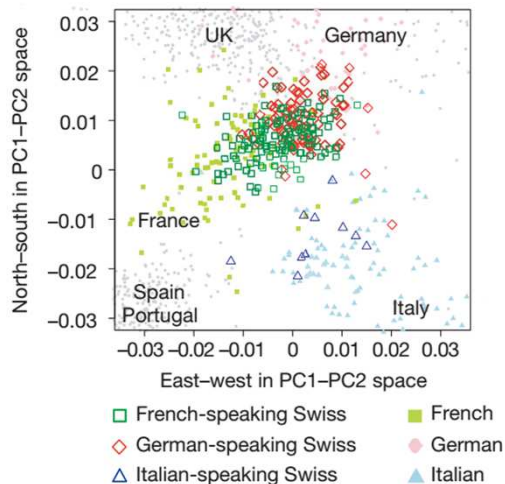


Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Vol 456 | 6 November 2008 | doi:10.1038/nature07331



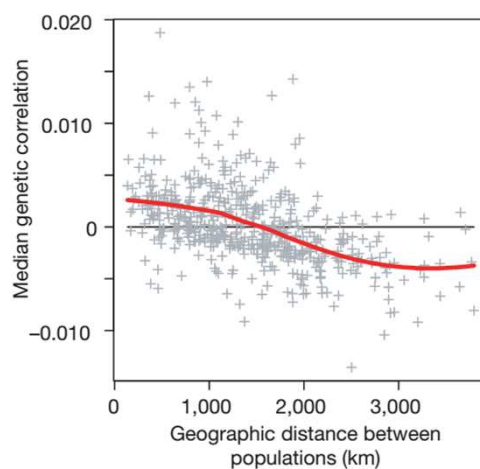
The trend even holds up within a country, e.g. Switzerland:



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Vol 456|6 November 2008|doi:10.1038/nature07331

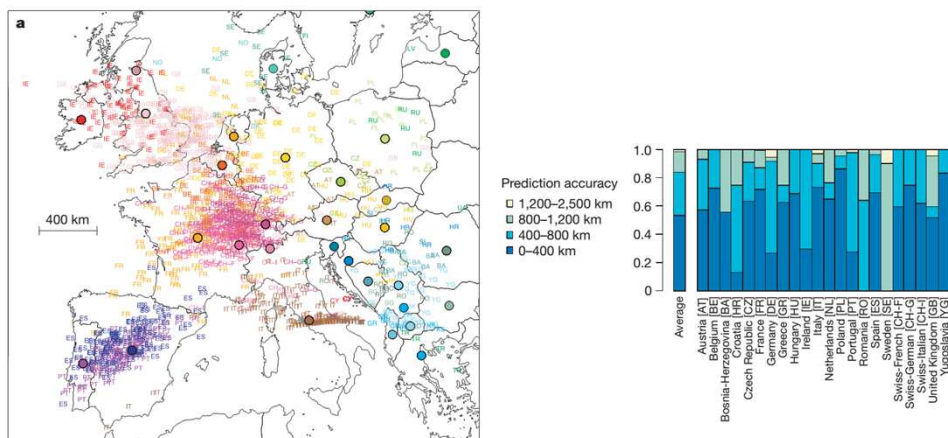
Why? Genetic similarity is strongly distance dependent in Europe, presumably because people tend to live near where they were born and tend to marry locally.



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Vol 456|6 November 2008|doi:10.1038/nature07331

The trend is so strong that it can be used to *predict* where the men are from, e.g. shown here by leave-one-out cross-validation:



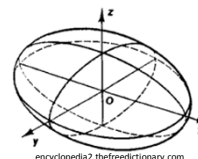
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Vol 456|6 November 2008|doi:10.1038/nature07331

SUMMARY

In a sense, PCA fits a (multidimensional) ellipsoid to the data

- Described by directions and lengths of principal (semi-)axes, e.g. the axis of a cigar or egg or the plane of a pancake



- No matter how an ellipsoid is turned, the eigenvectors point in those principal directions. The eigenvalues give the lengths.
- The biggest eigenvalues correspond to the fattest directions (having the most data variance). The smallest eigenvalues correspond to the thinnest directions (least data variance).
- Ignoring the smallest directions (*i.e.*, collapsing them) loses relatively little information.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Adapted from whuber, <http://stats.stackexchange.com/questions/2691/making-sense-of-principal-component-analysis-eigenvectors-eigenvalues>