

Principal Components Analysis

36-490

Spring 2010

Abstract

EXERCISE: Step through the `pca.R` file on the class website. Then replicate the analysis of the cars data given below.

Contents

| | | |
|----------|---|-----------|
| 1 | Mathematics of Principal Components | 2 |
| 1.1 | Minimizing Projection Residuals | 2 |
| 1.2 | Maximizing Variance | 3 |
| 1.3 | More Geometry; Back to the Residuals | 5 |
| 1.4 | Statistical Inference, or Not | 6 |
| 2 | Example: Cars | 7 |
| 3 | Latent Semantic Analysis | 10 |
| 3.1 | Principal Components of the New York <i>Times</i> | 11 |
| 4 | PCA for Visualization | 14 |
| 5 | PCA Cautions | 15 |

Principal components analysis (PCA) is one of a family of techniques for taking high-dimensional data, and using the dependencies between the variables to represent it in a more tractable, lower-dimensional form, without losing too much information. PCA is one of the simplest and most robust ways of doing such **dimensionality reduction**. It is also one of the oldest, and has been rediscovered many times in many fields, so it is also known as the Karhunen-Loève transformation, the Hotelling transformation, the method of empirical orthogonal functions, and singular value decomposition¹. We will call it PCA.

¹Strictly speaking, singular value decomposition is a matrix algebra trick which is used in the most common algorithm for PCA.

1 Mathematics of Principal Components

We start with p -dimensional feature vectors, and want to summarize them by projecting down into a q -dimensional subspace. Our summary will be the projection of the original vectors on to q directions, the **principal components**, which span the sub-space.

There are several equivalent ways of deriving the principal components mathematically. The simplest one is by finding the projections which maximize the variance. The first principal component is the direction in feature space along which projections have the largest variance. The second principal component is the direction which maximizes variance among all directions orthogonal to the first. The k^{th} component is the variance-maximizing direction orthogonal to the previous $k - 1$ components. There are p principal components in all.

Rather than maximizing variance, it might sound more plausible to look for the projection with the smallest average (mean-squared) distance between the original vectors and their projections on to the principal components; this turns out to be equivalent to maximizing the variance.

Throughout, assume that the data have been “centered”, so that every feature has mean 0. If we write the centered data in a matrix \mathbf{X} , where rows are objects and columns are features, then $\mathbf{X}^T \mathbf{X} = n \mathbf{V}$, where \mathbf{V} is the covariance matrix of the data. (You should check that last statement!)

1.1 Minimizing Projection Residuals

We’ll start by looking for a one-dimensional projection. That is, we have p -dimensional feature vectors, and we want to project them on to a line through the origin. We can specify the line by a unit vector along it, \vec{w} , and then the projection of a data vector \vec{x}_i on to the line is $\vec{x}_i \cdot \vec{w}$, which is a scalar. (Sanity check: this gives us the right answer when we project on to one of the coordinate axes.) This is the distance of the projection from the origin; the actual coordinate in p -dimensional space is $(\vec{x}_i \cdot \vec{w})\vec{w}$. The mean of the projections will be zero, because the mean of the vectors \vec{x}_i is zero:

$$\frac{1}{n} \sum_{i=1}^n (\vec{x}_i \cdot \vec{w}) \vec{w} = \left(\left(\frac{1}{n} \sum_{i=1}^n \vec{x}_i \right) \cdot \vec{w} \right) \vec{w} \quad (1)$$

If we try to use our projected or **image** vectors instead of our original vectors, there will be some error, because (in general) the images do not coincide with the original vectors. (When do they coincide?) The difference is the error or **residual** of the projection. How big is it? For any one vector, say \vec{x}_i , it’s

$$\|\vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w}\|^2 = \|\vec{x}_i\|^2 - 2(\vec{w} \cdot \vec{x}_i)(\vec{w} \cdot \vec{x}_i) + \|\vec{w}\|^2 \quad (2)$$

$$= \|\vec{x}_i\|^2 - 2(\vec{w} \cdot \vec{x}_i)^2 + 1 \quad (3)$$

(This is the same trick used to compute distance matrices in the solution to the first homework; it’s really just the Pythagorean theorem.) Add those residuals

up across all the vectors:

$$RSS(\vec{w}) = \sum_{i=1}^n \|\vec{x}_i\|^2 - 2(\vec{w} \cdot \vec{x}_i)^2 + 1 \quad (4)$$

$$= \left(n + \sum_{i=1}^n \|\vec{x}_i\|^2 \right) - 2 \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2 \quad (5)$$

The term in the big parenthesis doesn't depend on \vec{w} , so it doesn't matter for trying to minimize the residual sum-of-squares. To make RSS small, what we must do is make the second sum big, i.e., we want to maximize

$$\sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2 \quad (6)$$

Equivalently, since n doesn't depend on \vec{w} , we want to maximize

$$\frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2 \quad (7)$$

which we can see is the sample mean of $(\vec{w} \cdot \vec{x}_i)^2$. The mean of a square is always equal to the square of the mean plus the variance:

$$\frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2 = \left(\frac{1}{n} \sum_{i=1}^n \vec{x}_i \cdot \vec{w} \right)^2 + \text{Var}[\vec{w} \cdot \vec{x}_i] \quad (8)$$

Since we've just seen that the mean of the projections is zero, minimizing the residual sum of squares turns out to be equivalent to maximizing the variance of the projections.

(Of course in general we don't want to project on to just one vector, but on to multiple principal components. If those components are orthogonal and have the unit vectors $\vec{w}_1, \vec{w}_2, \dots, \vec{w}_k$, then the image of \vec{x}_i is its projection into the space spanned by these vectors,

$$\sum_{j=1}^k (\vec{x}_i \cdot \vec{w}_j) \vec{w}_j \quad (9)$$

The mean of the projection on to each component is still zero. If we go through the same algebra for the residual sum of squares, it turns out that the cross-terms between different components all cancel out, and we are left with trying to maximize the sum of the variances of the projections on to the components. EXERCISE: Do this algebra.)

1.2 Maximizing Variance

Accordingly, let's maximize the variance! Writing out all the summations grows tedious, so let's do our algebra in matrix form. If we stack our n data vectors

into an $n \times p$ matrix, \mathbf{X} , then the projections are given by $\mathbf{X}\mathbf{w}$, which is an $n \times 1$ matrix. The variance is

$$\sigma_{\vec{w}}^2 = \frac{1}{n} \sum_i (\vec{x}_i \cdot \vec{w})^2 \quad (10)$$

$$= \frac{1}{n} (\mathbf{X}\mathbf{w})^T (\mathbf{X}\mathbf{w}) \quad (11)$$

$$= \frac{1}{n} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \quad (12)$$

$$= \mathbf{w}^T \frac{\mathbf{X}^T \mathbf{X}}{n} \mathbf{w} \quad (13)$$

$$= \mathbf{w}^T \mathbf{V} \mathbf{w} \quad (14)$$

We want to choose a unit vector \vec{w} so as to maximize $\sigma_{\vec{w}}^2$. To do this, we need to make sure that we only look at unit vectors — we need to **constrain** the maximization. The constraint is that $\vec{w} \cdot \vec{w} = 1$, or $\mathbf{w}^T \mathbf{w} = 1$. This needs a brief excursion into constrained optimization.

We start with a function $f(w)$ that we want to maximize. (Here, that function is $\mathbf{w}^T \mathbf{V} \mathbf{w}$.) We also have an equality constraint, $g(w) = c$. (Here, $g(w) = \mathbf{w}^T \mathbf{w}$ and $c = 1$.) We re-arrange the constraint equation so its right-hand side is zero, $g(w) - c = 0$. We now add an extra variable to the problem, the **Lagrange multiplier** λ , and consider $\mathcal{L}(w, \lambda) = f(w) - \lambda(g(w) - c)$. This is our new objective function, so we differentiate with respect to both arguments and set the derivatives equal to zero:

$$\frac{\partial \mathcal{L}}{\partial w} = 0 = \frac{\partial f}{\partial w} - \lambda \frac{\partial g}{\partial w} \quad (15)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = 0 = -g(w) + c \quad (16)$$

That is, maximizing with respect to λ gives us back our constraint equation, $g(w) = c$. However, we have modified the equation we use to find the optimal w , because it now includes the $\lambda \frac{\partial g}{\partial w}$ term. If λ were zero, we'd be back to the unconstrained optimization problem; the bigger λ is, the more the constraint “bites”, and forces modifies the optimal w . At the same time, when we have the constraint satisfied, the *value* our new objective function \mathcal{L} is the same as the value of the old objective function f . (If we had more than one constraint, we would just need more Lagrange multipliers.)²

For our projection problem,

$$u = \mathbf{w}^T \mathbf{V} \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w} - 1) \quad (17)$$

$$\frac{\partial u}{\partial \mathbf{w}} = 2\mathbf{V}\mathbf{w} - 2\lambda\mathbf{w} = 0 \quad (18)$$

$$\mathbf{V}\mathbf{w} = \lambda\mathbf{w} \quad (19)$$

²To learn more about Lagrange multipliers, read Boas (1983) or (more compactly) Klein (2001).

Thus, desired vector \mathbf{w} is an **eigenvector** of the covariance matrix \mathbf{V} , and the maximizing vector will be the one associated with the largest **eigenvalue** λ . This is good news, because finding eigenvectors is something which can be done comparatively rapidly, and because eigenvectors have many nice mathematical properties, which we can use as follows.

We know that \mathbf{V} is a $p \times p$ matrix, so it will have p different eigenvectors.³ We know that \mathbf{V} is a covariance matrix, so it is symmetric, and then linear algebra tells us that the eigenvectors must be orthogonal to one another. Again because \mathbf{V} is a covariance matrix, it is a **positive matrix**, in the sense that $\vec{x} \cdot \mathbf{V} \vec{x} \geq 0$ for any \vec{x} . This tells us that the eigenvalues of \mathbf{V} must all be ≥ 0 .

The eigenvectors of \mathbf{V} are the **principal components** of the data. We know that they are all orthogonal to each other from the previous paragraph, so together they span the whole p -dimensional feature space. The first principal component, i.e. the eigenvector which goes the largest value of λ , is the direction along which the data have the most variance. The second principal component, i.e. the second eigenvector, is the direction orthogonal to the first component with the most variance. Because it is orthogonal to the first eigenvector, their projections will be uncorrelated. In fact, projections on to all the principal components are uncorrelated with each other. If we use q principal components, our weight matrix \mathbf{w} will be a $p \times q$ matrix, where each column will be a different eigenvector of the covariance matrix \mathbf{V} . The eigenvalues will give the total variance described by each component. The variance of the projections on to the first q principal components is then $\sum_{i=1}^q \lambda_i$.

1.3 More Geometry; Back to the Residuals

Suppose that the data really are q -dimensional. Then \mathbf{V} will have only q positive eigenvalues, and $p - q$ zero eigenvalues. If the data fall near a q -dimensional subspace, then $p - q$ of the eigenvalues will be nearly zero.

If we pick the top q components, we can define a projection operator \mathbf{P}_q . The images of the data are then $\mathbf{X}\mathbf{P}_q$. The **projection residuals** are $\mathbf{X} - \mathbf{X}\mathbf{P}_q$ or $\mathbf{X}(\mathbf{I} - \mathbf{P}_q)$. (Notice that the residuals here are vectors, not just magnitudes.) If the data really are q -dimensional, then the residuals will be zero. If the data are *approximately* q -dimensional, then the residuals will be small. In any case, we can define the R^2 of the projection as the fraction of the original variance kept by the image vectors,

$$R^2 \equiv \frac{\sum_{i=1}^q \lambda_i}{\sum_{j=1}^p \lambda_j} \quad (20)$$

just as the R^2 of a linear regression is the fraction of the original variance of the dependent variable retained by the fitted values.

The $q = 1$ case is especially instructive. We know, from the discussion of projections in the last lecture, that the residual vectors are all orthogonal to the projections. Suppose we ask for the first principal component of the residuals. This will be the direction of largest variance which is perpendicular

³Exception: if $n < p$, there are only n distinct eigenvectors and eigenvalues.

to the first principal component. In other words, it will be the second principal component of the data. This suggests a recursive algorithm for finding all the principal components: the k^{th} principal component is the leading component of the residuals after subtracting off the first $k - 1$ components. In practice, it is faster to use eigenvector-solvers to get all the components at once from \mathbf{V} , but we will see versions of this idea later.

This is a good place to remark that if the data really fall in a q -dimensional subspace, then \mathbf{V} will have only q positive eigenvalues, because after subtracting off those components there will be no residuals. The other $p - q$ eigenvectors will all have eigenvalue 0. If the data cluster around a q -dimensional subspace, then $p - q$ of the eigenvalues will be very small, though how small they need to be before we can neglect them is a tricky question.⁴

Projections on to the first two or three principal components can be visualized; however they may not be enough to really give a good summary of the data. Usually, to get an R^2 of 1, you need to use all p principal components.⁵ How many principal components you should use depends on your data, and how big an R^2 you need. In some fields, you can get better than 80% of the variance described with just two or three components. A sometimes-useful device is to plot $1 - R^2$ versus the number of components, and keep extending the curve it until it flattens out.

1.4 Statistical Inference, or Not

You may have noticed, and even been troubled by, the fact that I have said nothing at all yet like “assume the data are drawn at random from some distribution”, or “assume the different rows of the data frame are statistically independent”. This is because no such assumption is required for principal components. *All* it does is say “these data can be summarized using projections along these directions”. It says nothing about the larger population or stochastic process the data came from; it doesn’t even suppose the latter exist.

However, we could *add* a statistical assumption and see how PCA behaves under those conditions. The simplest one is to suppose that the data are IID draws from a distribution with covariance matrix \mathbf{V}_0 . Then the sample covariance matrix $\mathbf{V} \equiv \frac{1}{n} \mathbf{X}^T \mathbf{X}$ will converge on \mathbf{V}_0 as $n \rightarrow \infty$. Since the principal components are smooth functions of \mathbf{V} (namely its eigenvectors), they will tend to converge as n grows⁶. So, along with that additional assumption about the

⁴One tricky case where this can occur is if $n < p$. Any two points define a line, and three points define a plane, etc., so if there are fewer data points than features, it is *necessarily* true that the fall on a low-dimensional subspace. If we look at the bags-of-words for the *Times* stories, for instance, we have $p \approx 4400$ but $n \approx 102$. Finding that only 102 principal components account for all the variance is not an empirical discovery but a mathematical artifact.

⁵The exceptions are when some of your features are linear combinations of the others, so that you don’t really have p *different* features, or when $n < p$.

⁶There is a wrinkle if \mathbf{V}_0 has “degenerate” eigenvalues, i.e., two or more eigenvectors with the same eigenvalue. Then any linear combination of those vectors is also an eigenvector, with the same eigenvalue. (EXERCISE: show this.) For instance, if \mathbf{V}_0 is the identity matrix, then

| Variable | Meaning |
|------------|---|
| Sports | Binary indicator for being a sports car |
| SUV | Indicator for sports utility vehicle |
| Wagon | Indicator |
| Minivan | Indicator |
| Pickup | Indicator |
| AWD | Indicator for all-wheel drive |
| RWD | Indicator for rear-wheel drive |
| Retail | Suggested retail price (US\$) |
| Dealer | Price to dealer (US\$) |
| Engine | Engine size (liters) |
| Cylinders | Number of engine cylinders |
| Horsepower | Engine horsepower |
| CityMPG | City gas mileage |
| HighwayMPG | Highway gas mileage |
| Weight | Weight (pounds) |
| Wheelbase | Wheelbase (inches) |
| Length | Length (inches) |
| Width | Width (inches) |

Table 1: Features for the 2004 cars data.

data-generating process, PCA does make a prediction: in the future, the principal components will look like they do now.

2 Example: Cars

Let's work an example. The data⁷ consists of 388 cars from the 2004 model year, with 18 features. Eight features are binary indicators; the other 11 features are numerical (Table 1). All of the features except **Type** are numerical. Table 2 shows the first few lines from the data set. PCA only works with numerical features, so we have ten of them to play with.

There are *two* R functions for doing PCA, `princomp` and `prcomp`, which differ in how they do the actual calculation.⁸ The latter is generally more robust, so we'll just use it.

```
cars04 = read.csv("cars-fixed04.dat")
cars04.pca = prcomp(cars04[,8:18], scale.=TRUE)
```

every vector is an eigenvector, and PCA routines will return an essentially arbitrary collection of mutually perpendicular vectors. Generically, however, any arbitrarily small tweak to \mathbf{V}_0 will break the degeneracy.

⁷On the course website; from <http://www.amstat.org/publications/jse/datasets/04cars.txt>, with incomplete records removed.

⁸`princomp` actually calculates the covariance matrix and takes its eigenvalues. `prcomp` uses a different technique called “singular value decomposition”.

```

Sports, SUV, Wagon, Minivan, Pickup, AWD, RWD, Retail, Dealer, Engine,
  Cylinders, Horsepower, CityMPG, HighwayMPG, Weight, Wheelbase, Length, Width
Acura 3.5 RL, 0, 0, 0, 0, 0, 0, 0, 43755, 39014, 3.5, 6, 225, 18, 24, 3880, 115, 197, 72
Acura MDX, 0, 1, 0, 0, 0, 1, 0, 36945, 33337, 3.5, 6, 265, 17, 23, 4451, 106, 189, 77
Acura NSX S, 1, 0, 0, 0, 0, 0, 1, 89765, 79978, 3.2, 6, 290, 17, 24, 3153, 100, 174, 71

```

Table 2: The first few lines of the 2004 cars data set.

The second argument to `prcomp` tells it to first scale all the variables to have variance 1, i.e., to standardize them. You should experiment with what happens with this data when we don't standardize.

We can now extract the loadings or weight matrix from the `cars04.pca` object. For comprehensibility I'll just show the first two components.

```

> round(cars04.pca$rotation[,1:2], 2)
      PC1    PC2
Retail  -0.26 -0.47
Dealer  -0.26 -0.47
Engine  -0.35  0.02
Cylinders -0.33 -0.08
Horsepower -0.32 -0.29
CityMPG   0.31  0.00
HighwayMPG 0.31  0.01
Weight   -0.34  0.17
Wheelbase -0.27  0.42
Length   -0.26  0.41
Width    -0.30  0.31

```

This says that all the variables *except* the gas-mileages have a negative projection on to the first component. This means that there is a negative correlation between mileage and everything else. The first principal component tells us about whether we are getting a big, expensive gas-guzzling car with a powerful engine, or whether we are getting a small, cheap, fuel-efficient car with a wimpy engine.

The second component is a little more interesting. Engine size and gas mileage hardly project on to it at all. Instead we have a contrast between the physical size of the car (positive projection) and the price and horsepower. Basically, this axis separates mini-vans, trucks and SUVs (big, not so expensive, not so much horse-power) from sports-cars (small, expensive, lots of horse-power).

To check this interpretation, we can use a useful tool called a **biplot**, which plots the data, along with the projections of the original features, on to the first two components (Figure 1). Notice that the car with the lowest value of the second component is a Porsche 911, with pick-up trucks and mini-vans at the other end of the scale. Similarly, the highest values of the first component all belong to hybrids.

3 Latent Semantic Analysis

Information retrieval systems (like search engines) and people doing computational text analysis often represent documents as what are called **bags of words**: documents are represented as vectors, where each component counts how many times each word in the dictionary appears in the text. This throws away information about word order, but gives us something we can work with mathematically. Part of the representation of one document might look like:

| | | | | | | | | | | |
|----|-----------|-----|---------|------|-------|-------|--------|----------|-----------|----------|
| a | abandoned | abc | ability | able | about | above | abroad | absorbed | absorbing | abstract |
| 43 | 0 | 0 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 1 |

and so on through to “zebra”, “zoology”, etc. at the end of the dictionary. These vectors are very, very large! At least in English and similar languages, these bag-of-word vectors have three outstanding properties:

1. Most words do not appear in most documents; the bag-of-words vectors are very **sparse** (most entries are zero).
2. A small number of words appear many times in almost all documents; these words tell us almost nothing about what the document is about. (Examples: “the”, “is”, “of”, “for”, “at”, “a”, “and”, “here”, “was”, etc.)
3. Apart from those hyper-common words, most words’ counts are correlated with some but not all other words; words tend to come in bunches which appear together.

Taken together, this suggests that we do not really get a lot of information from keeping around *all* the words. We would be better off if we could project down a smaller number of new features, which we can think of as combinations of words that tend to appear together in the documents, or not at all. But this tendency needn’t be absolute — it can be partial because the words mean slightly different things, or because of stylistic differences, etc. This is *exactly* what principal components analysis does.

To see how this can be useful, imagine we have a collection of documents (a **corpus**), which we want to search for documents about agriculture. It’s entirely possible that many documents on this topic don’t actually contain the *word* “agriculture”, just closely related words like “farming”. A simple feature-vector search on “agriculture” will miss them. But it’s very likely that the occurrence of these related words is well-correlated with the occurrence of “agriculture”. This means that all these words will have similar projections on to the principal components, and will be easy to find documents whose principal components projection is like that for a query about agriculture. This is called **latent semantic indexing**.

To see why this is *indexing*, think about what goes into coming up with an index for a book by hand. Someone draws up a list of topics and then goes through the book noting all the passages which refer to the topic, and maybe a little bit of what they say there. For example, here’s the start of the entry for “Agriculture” in the index to Adam Smith’s *The Wealth of Nations*:

AGRICULTURE, the labour of, does not admit of such subdivisions as manufactures, 6; this impossibility of separation, prevents agriculture from improving equally with manufactures, 6; natural state of, in a new colony, 92; requires more knowledge and experience than most mechanical professions, and yet is carried on without any restrictions, 127; the terms of rent, how adjusted between landlord and tenant, 144; is extended by good roads and navigable canals, 147; under what circumstances pasture land is more valuable than arable, 149; gardening not a very gainful employment, 152–3; vines the most profitable article of culture, 154; estimates of profit from projects, very fallacious, *ib.*; cattle and tillage mutually improve each other, 220; . . .

and so on. (Agriculture is an important topic in *The Wealth of Nations*.) It’s asking a lot to hope for a computer to be able to do something like this, but we could at least hope for a list of pages like “6, 92, 126, 144, 147, 152 – –3, 154, 220, . . .”. One could imagine doing this by treating each page as its own document, forming its bag-of-words vector, and then returning the list of pages with a non-zero entry for the feature “agriculture”. This will fail: only two of those nine pages actually contains that word, and this is pretty typical. On the other hand, they are full of words strongly correlated with “agriculture”, so asking for the pages which are most similar in their principal components projection to that word will work great.⁹

At first glance, and maybe even second, this seems like a wonderful trick for extracting meaning, or **semantics**, from pure correlations. Of course there are also all sorts of ways it can fail, not least from spurious correlations. If our training corpus happens to contain lots of documents which mention “farming” and “Kansas”, as well as “farming” and “agriculture”, latent semantic indexing will not make a big distinction between the relationship between “agriculture” and “farming” (which is genuinely semantic) and that between “Kansas” and “farming” (which is accidental, and probably wouldn’t show up in, say, a corpus collected from Europe).

Despite this susceptibility to spurious correlations, latent semantic indexing is an *extremely* useful technique in practice, and the foundational papers (??) are worth reading.

3.1 Principal Components of the New York *Times*

To get a more concrete sense of how latent semantic analysis works, and how it reveals semantic information, let’s apply it to some data. The accompanying R file and R workspace contains some news stories taken from the New York *Times* Annotated Corpus (?), which consists of about 1.8 million stories from the *Times*, from 1987 to 2007, which have been hand-annotated by actual human beings with standardized machine-readable information about their contents. From this corpus, I have randomly selected 57 stories about art and 45 stories

⁹Or it should anyway; I haven’t actually done the experiment with this book.

about music, and turned them into a bag-of-words data frame, one row per story, one column per word; plus an indicator in the first column of whether the story is one about art or one about music.¹⁰ The original data frame thus has 102 rows, and 4432 columns: the categorical label, and 4431 columns with counts for every distinct word that appears in at least one of the stories.¹¹

The PCA is done as it would be for any other data:

```
nyt.pca <- prcomp(nyt.frame[,-1])
nyt.latent.sem <- nyt.pca$rotation
```

We need to omit the first column in the first command because it contains categorical variables, and PCA doesn't apply to them. The second command just picks out the matrix of projections of the features on to the components — this is called `rotation` because it can be thought of as rotating the coordinate axes in feature-vector space.

Now that we've done this, let's look at what the leading components are.

```
> signif(sort(nyt.latent.sem[,1],decreasing=TRUE)[1:30],2)
      music      trio      theater orchestra composers      opera
0.110      0.084      0.083      0.067      0.059      0.058
theaters      m      festival      east      program      y
0.055      0.054      0.051      0.049      0.048      0.048
jersey      players      committee      sunday      june      concert
0.047      0.047      0.046      0.045      0.045      0.045
symphony      organ      matinee      misstated instruments      p
0.044      0.044      0.043      0.042      0.041      0.041
X.d      april      samuel      jazz      pianist      society
0.041      0.040      0.040      0.039      0.038      0.038

> signif(sort(nyt.latent.sem[,1],decreasing=FALSE)[1:30],2)
      she      her      ms      i      said      mother      cooper
-0.260     -0.240     -0.200     -0.150     -0.130     -0.110     -0.100
my      painting      process      paintings      im      he      mrs
-0.094     -0.088     -0.071     -0.070     -0.068     -0.065     -0.065
me      gagosian      was      picasso      image      sculpture      baby
-0.063     -0.062     -0.058     -0.057     -0.056     -0.056     -0.055
artists      work      photos      you      nature      studio      out
-0.055     -0.054     -0.051     -0.051     -0.050     -0.050     -0.050
says      like
-0.050     -0.049
```

These are the thirty words with the largest positive and negative projections

¹⁰Actually, following standard practice in language processing, I've normalized the bag-of-word vectors so that documents of different lengths are comparable, and used "inverse document-frequency weighting" to de-emphasize hyper-common words like "the" and emphasize more informative words. See the lecture notes for 36-350 if you're interested.

¹¹If we were trying to work with the complete corpus, we should expect at least 50000 words, and perhaps more.

on to the first component.¹² The words with positive projections are mostly associated with music, those with negative components with the visual arts. The letters “m” and “p” show up with music because of the combination “p.m”, which our parsing breaks into two single-letter words, and because stories about music give show-times more often than do stories about art. Personal pronouns appear with art stories because more of those quote people, such as artists or collectors.¹³

What about the second component?

```
> signif(sort(nyt.latent.sem[,2],decreasing=TRUE)[1:30],2)
      art      museum      images      artists      donations      museums
0.150      0.120      0.095      0.092      0.075      0.073
painting      tax      paintings      sculpture      gallery      sculptures
0.073      0.070      0.065      0.060      0.055      0.051
painted      white      patterns      artist      nature      service
0.050      0.050      0.047      0.047      0.046      0.046
decorative      feet      digital      statue      color      computer
0.043      0.043      0.043      0.042      0.042      0.041
paris      war      collections      diamond      stone      dealers
0.041      0.041      0.041      0.041      0.041      0.040

> signif(sort(nyt.latent.sem[,2],decreasing=FALSE)[1:30],2)
      her      she      theater      opera      ms
-0.220      -0.220      -0.160      -0.130      -0.130
      i      hour      production      sang      festival
-0.083      -0.081      -0.075      -0.075      -0.074
music      musical      songs      vocal      orchestra
-0.070      -0.070      -0.068      -0.067      -0.067
      la      singing      matinee      performance      band
-0.065      -0.065      -0.061      -0.061      -0.060
awards      composers      says      my      im
-0.058      -0.058      -0.058      -0.056      -0.056
      play      broadway      singer      cooper      performances
-0.056      -0.055      -0.052      -0.051      -0.051
```

Here the positive words are about art, but more focused on acquiring and trading (“collections”, “dealers”, “donations”, “dealers”) than on talking with artists or about them. The negative words are musical, specifically about musical theater and vocal performances.

I could go on, but by this point you get the idea.

¹²Which direction is positive and which negative is of course arbitrary; basically it depends on internal choices in the algorithm.

¹³You should check out these explanations for yourself. The raw stories are part of the R workspace.

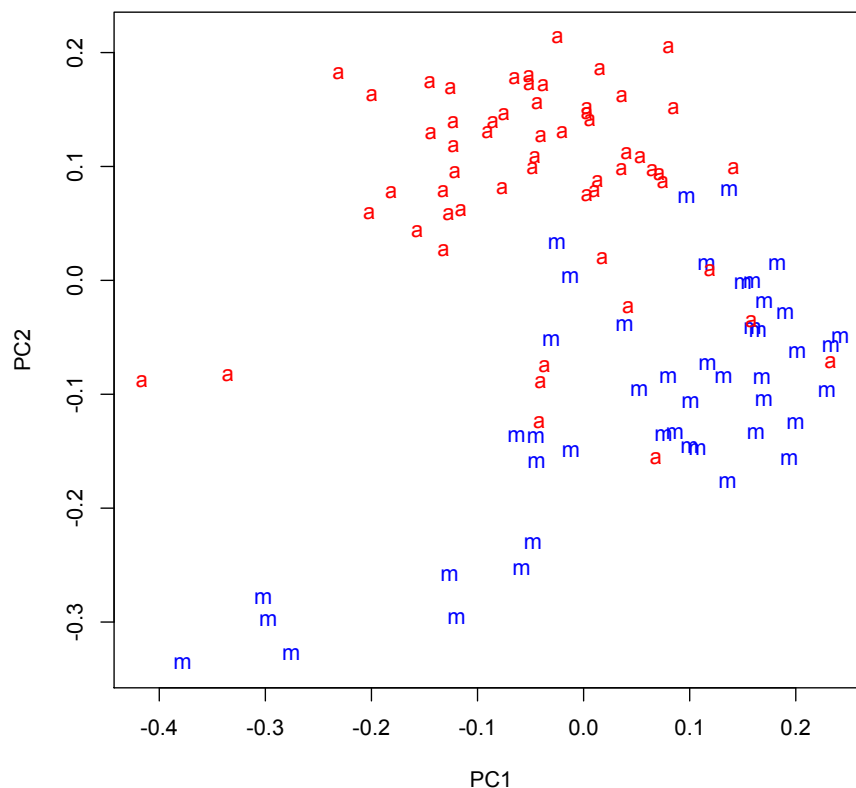


Figure 2: Projection of the *Times* stories on to the first two principal components. Labels: “a” for art stories, “m” for music.

4 PCA for Visualization

Let’s try displaying the *Times* stories using the principal components. (Assume that the objects from just before are still in memory.)

```
plot(nyt.pca$x[,1:2],type="n")
points(nyt.pca$x[nyt.frame[,"class.labels"]=="music",1:2],pch="m",col="blue")
points(nyt.pca$x[nyt.frame[,"class.labels"]=="art",1:2],pch="a",col="red")
```

The first command makes an empty plot — I do this just to set up the axes nicely for the data which will actually be displayed. The second and third commands plot a blue “m” at the location of each music story, and a red “a” at the location of each art story. The result is Figure 2.

Notice that even though we have gone from 4431 dimensions to 2, and so thrown away a lot of information, we could draw a line across this plot and have most of the art stories on one side of it and all the music stories on the other. If we let ourselves use the first four or five principal components, we'd still have a thousand-fold savings in dimensions, but we'd be able to get almost-perfect separation between the two classes. This is a sign that PCA is really doing a good job at summarizing the information in the word-count vectors, and in turn that the bags of words give us a lot of information about the meaning of the stories.

The figure also illustrates the idea of **multidimensional scaling**, which means finding low-dimensional points to represent high-dimensional data by preserving the distances between the points. If we write the original vectors as $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$, and their images as $\vec{y}_1, \vec{y}_2, \dots, \vec{y}_n$, then the MDS problem is to pick the images to minimize the difference in distances:

$$\sum_i \sum_{j \neq i} (\|\vec{y}_i - \vec{y}_j\| - \|\vec{x}_i - \vec{x}_j\|)^2$$

This will be small if distances between the images are all close to the distances between the original points. PCA accomplishes this precisely because \vec{y}_i is itself close to \vec{x}_i (on average).

5 PCA Cautions

Trying to guess at what the components might mean is a good idea, but like many good ideas it's easy to go overboard. Specifically, once you attach an idea in your mind to a component, and especially once you attach a *name* to it, it's very easy to forget that those are names and ideas you made up; to **reify** them, as you might reify clusters. Sometimes the components actually do measure real variables, but sometimes they just reflect patterns of covariance which have many different causes. If I did a PCA of the same features but for, say, European cars, I might well get a similar first component, but the second component would probably be rather different, since SUVs are much less common there than here.

A more important example comes from population genetics. Starting in the late 1960s, L. L. Cavalli-Sforza and collaborators began a huge project of mapping human genetic variation — of determining the frequencies of different genes in different populations throughout the world. (Cavalli-Sforza *et al.* (1994) is the main summary; Cavalli-Sforza has also written several excellent popularizations.) For each point in space, there are a very large number of features, which are the frequencies of the various genes among the people living there. Plotted over space, this gives a map of that gene's frequency. What they noticed (unsurprisingly) is that many genes had similar, but not identical, maps. This led them to use PCA, reducing the huge number of features (genes) to a few components. Results look like Figure 3. They interpreted these components, very reasonably, as signs of large population movements. The first principal component for Europe and the Near East, for example, was supposed to show

the expansion of agriculture out of the Fertile Crescent. The third, centered in steppes just north of the Caucasus, was supposed to reflect the expansion of Indo-European speakers towards the end of the Bronze Age. Similar stories were told of other components elsewhere.

Unfortunately, as Novembre and Stephens (2008) showed, spatial patterns like this are what one should expect to get when doing PCA of any kind of spatial data with local correlations, because that essentially amounts to taking a Fourier transform, and picking out the low-frequency components.¹⁴ They simulated genetic diffusion processes, without any migration or population expansion, and got results that looked very like the real maps (Figure 4). This doesn't mean that the stories of the maps *must be* wrong, but it does undercut the principal components as evidence for those stories.

References

- Boas, Mary L. (1983). *Mathematical Methods in the Physical Sciences*. New York: Wiley, 2nd edn.
- Cavalli-Sforza, L. L., P. Menozzi and A. Piazza (1994). *The History and Geography of Human Genes*. Princeton: Princeton University Press.
- Klein, Dan (2001). “Lagrange Multipliers without Permanent Scarring.” Online tutorial. URL <http://dbpubs.stanford.edu:8091/~klein/lagrange-multipliers.pdf>.
- Novembre, John and Matthew Stephens (2008). “Interpreting principal component analyses of spatial population genetic variation.” *Nature Genetics*, **40**: 646–649. doi:10.1038/ng.139.

¹⁴Remember that PCA re-writes the original vectors as a weighted sum of new, orthogonal vectors, just as Fourier transforms do. When there is a lot of spatial correlation, values at nearby points are similar, so the low-frequency modes will have a lot of amplitude, i.e., carry a lot of the variance. So first principal components will tend to be similar to the low-frequency Fourier modes.

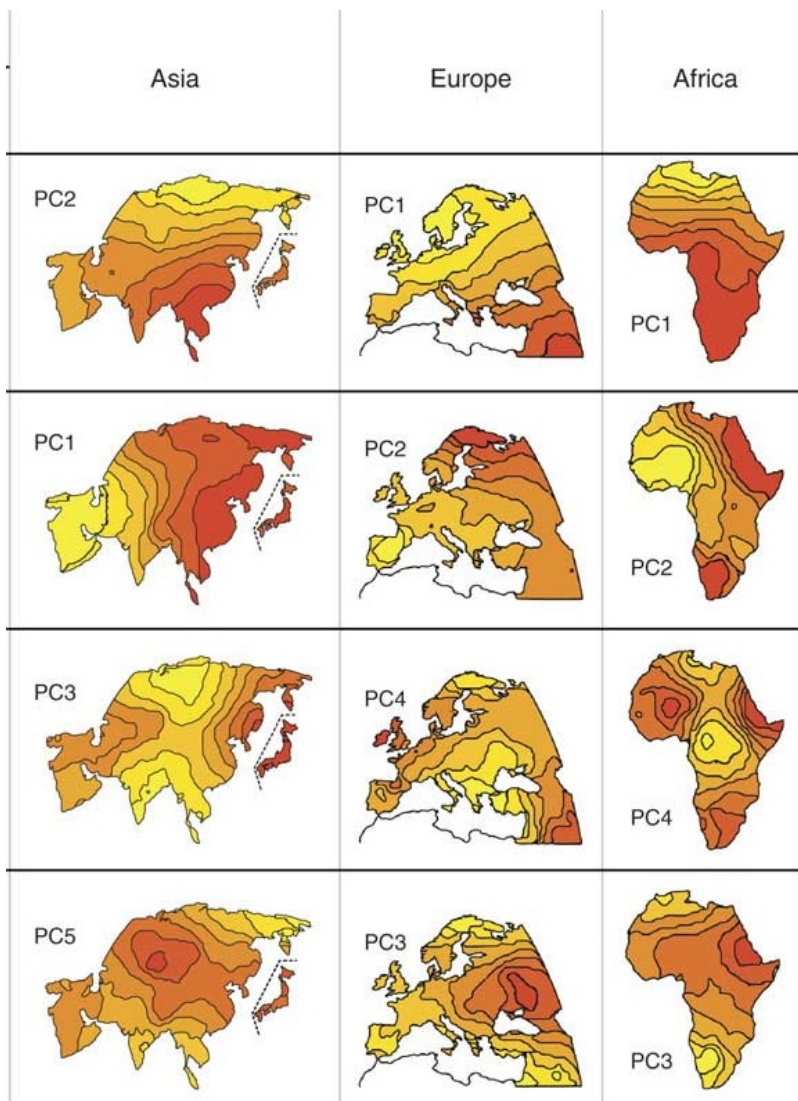


Figure 3: Principal components of genetic variation in the old world, according to Cavalli-Sforza *et al.* (1994), as re-drawn by Novembre and Stephens (2008).

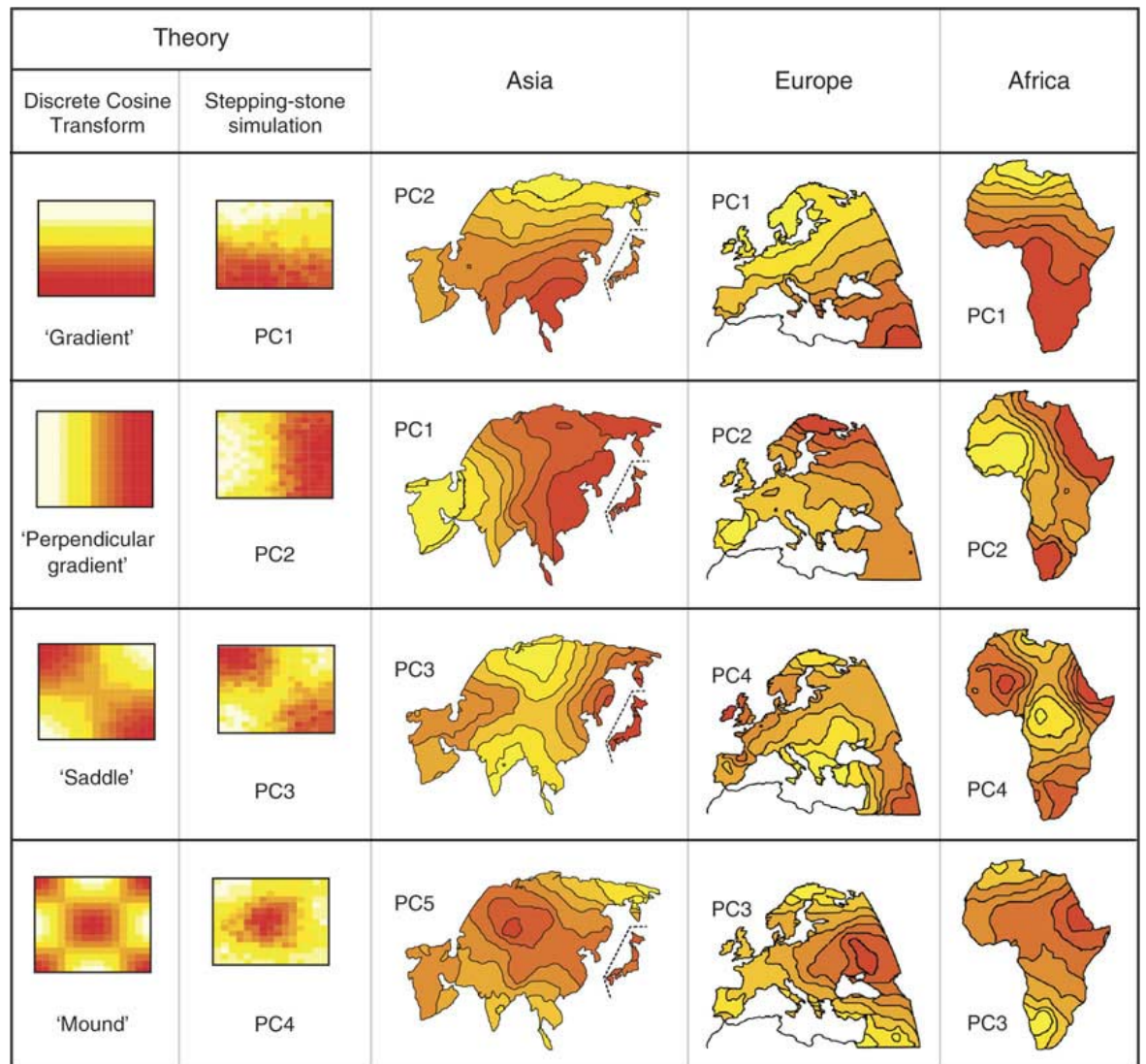


Figure 4: How the PCA patterns can arise as numerical artifacts (far left column) or through simple genetic diffusion (next column). From Novembre and Stephens (2008).