Displays for Statistics 5401/8401


Lecture 24


November 2, 2005


Christopher Bingham, Instructor


612-625-1024, `kb@umn.edu`
372 Ford Hall


Class Web Page

---

# Principal components

*Principal components* are specific
<u>linear combinations</u> $z_j \equiv \mathbf{v}_j{}'\mathbf{x} = \sum_{1 \le \ell \le p} v_{\ell j} x_\ell$
of variables $x_1, \ldots, x_p$. The vectors
$\mathbf{v}_j = [v_{1j}, \ldots, v_{pj}]'$, $1 \le j \le p$, of coefficients
are chosen to have certain properties.

There are at least two ways to motivate
principal components.

- Principal components are linear com-
  binations $\mathbf{v}_j{}'\mathbf{x}$ of variables which have
  the <u>largest variances subject to cons-</u>
  <u>traints</u> on the coefficients $v_{\ell j}$:

$$\|\mathbf{v}_j\|^2 = \mathbf{v}_j{}'\mathbf{v}_j = \sum_\ell v_{\ell j}{}^2 = 1 \ (\underline{\text{normalized}})$$
$$\mathbf{v}_j{}'\mathbf{v}_k = \sum_\ell v_{\ell j} v_{\ell k} = 0, \ j \ne k \ (\underline{\text{orthogonal}})$$

- Principal components are linear com-
  binations from which you can <u>approx-</u>
  <u>imately reconstruct</u> a data matrix $\mathbf{X}$
  using a <u>least squares</u> criterion.

2

---

The two approaches agree in an important
way.

Both view the first few principal com-
ponents as a set, preferably small, of
*new* variables $z_1, z_2, \ldots, z_m$, linearly
related to $x_1, x_2, \ldots, x_p$, and which lose as
little as possible "important information"
in the complete data.

- The variance maximization approach
  equates "important information" with
  high variability.
- The data matrix approximation ap-
  proach equates "important information"
  with having an approximation with
  small errors.

I focus on their use in *approximating* a N
by p data matrix $\mathbf{X}$, because I think that
is closer to the way principal components
are usually used.

---

In many cases, the single most important
description or summary of a data matrix
$\mathbf{X}$ is its sample mean vector $\overline{\mathbf{x}} = \sum_i \mathbf{x}_i / N$.

I think this is because the N by p matrix

$$\mathbf{1}_N \overline{\mathbf{x}}' = \begin{bmatrix} \overline{\mathbf{x}}' \\ \overline{\mathbf{x}}' \\ \ldots \\ \overline{\mathbf{x}}' \end{bmatrix}$$

often "explains" or predicts $\mathbf{X}$ well in the
following sense:

　Elements of $\tilde{\mathbf{X}} \equiv \mathbf{X} - \mathbf{1}_N \overline{\mathbf{x}}'$ are often
　much smaller than the elements of $\mathbf{X}$.

Principal components are actually derived
from a process which attempts to ap-
proximate $\tilde{\mathbf{X}}$ rather than $\mathbf{X}$. By adding
$\mathbf{1}_n \overline{\mathbf{x}}'$ to an approximation $\hat{\tilde{\mathbf{X}}}$ to $\tilde{\mathbf{X}}$, you can
then get an approximation $\hat{\mathbf{X}}$ for $\mathbf{X}$:

$$\hat{\mathbf{X}} = \hat{\tilde{\mathbf{X}}} + \mathbf{1}_N \overline{\mathbf{x}}'.$$

The matrix of <u>residuals</u> of $X$ from $\overline{x}$ is,

$$\tilde{X} = X - 1_N\overline{x}' = \begin{bmatrix} (x_1 - \overline{x})' \\ (x_2 - \overline{x})' \\ (x_3 - \overline{x})' \\ \cdots \cdots \\ (x_N - \overline{x})' \end{bmatrix}, \text{ N by p}$$

When N > p, usually rank$(\tilde{X})$ = p. That is, you can always find p pairs of vectors,

$$U_k \text{ (N × 1) and } v_k \text{ (p × 1)}$$

so that $\tilde{X}$ is the sum of p *outer products* $U_k v_k'$ of $U_k$ and $v_k$:

$$\tilde{X} = \sum_{1 \leq k \leq p} U_k v_k'$$

The element in row i (case i) and column $\ell$ (variable $\ell$) of $\tilde{X}$ is

$$\tilde{x}_{i\ell} \equiv x_{i\ell} - \overline{x}_\ell = \sum_{1 \leq k \leq p} u_{ik}v_{\ell k} = \sum_{1 \leq k \leq p} v_{\ell k}u_{ik}$$

There are an infinite number of such sets of vectors $\{U_k\}$ and $\{v_k\}$.

In this representation, $\tilde{X} = \sum_{1 \leq k \leq p} U_k v_k'$, think of $U_k$ and $v_k$ in the following way:

• Each N by 1 $U_k$ is a new *variable* with a value for each of the N cases

• Element $v_{\ell k}$ of the p by 1 vector $v_k$ is a *coefficient* of $U_k$, in a representation for column $\tilde{X}_\ell$ of $\tilde{X}$.

  Specifically, if $\tilde{X}_\ell$ is column $\ell$ of $\tilde{X}$,
$$\tilde{X}_\ell = \sum_{1 \leq k \leq p} v_{\ell k}U_k.$$

Except for there being no constant term (intercept), this looks a little like a multiple regression of $\tilde{X}_\ell$ on $U_1, ..., U_p$.

When rank$(\tilde{X})$ = p, when m < p you can't *exactly* <u>reproduce</u> $\tilde{X}$ by $\sum_{1 \leq k \leq m} U_k v_k'$, but it may be possible to get a good fit that is, have $\tilde{X} - \sum_{1 \leq k \leq m} U_k v_k'$.

That is, you *may* be able closely to *approximate* $\tilde{X}$ by a sum of m < p outer products:

$$\tilde{X} \cong \sum_{1 \leq k \leq m} U_k v_k' \text{ or } \tilde{X}_\ell = \sum_{1 \leq k \leq m} v_{\ell k}U_k, \ell = 1,...,p$$

in the sense that the elements of $\tilde{X} - \sum_{1 \leq k \leq m} U_k v_k'$ are small.

This would be a *reduced rank approx-imation* (specifically a <u>rank m approx-imation</u>) to $\tilde{X}$ because

$$\text{rank}(\sum_{1 \leq i \leq m} U_i v_i') = m < p$$

---
| How do you find such $U_i$'s and $v_i$'s? |
---

That is, for a specific m < p, how do you find

• Variables $U_k$, k = 1,...,m

• Coefficient vectors $v_k$, k = 1,...,m

such that $\tilde{X} - \sum_{1 \leq k \leq m} U_k v_k'$ is small?

Let's start with m = 1, that is, find a N×1 $U$ and p×1 $v$ such that

$$\tilde{X}^{(1)} = Uv' = [v_1U \ v_2U \ ... \ v_pU] \cong \tilde{X} \text{ (rank 1)}$$

that is, find numbers

• $\{u_i\}$, i = 1,...,N

• $\{v_\ell\}$, $\ell$ = 1,...,p

so that $\tilde{x}_{i\ell}^{(1)} - u_iv_\ell$ is small.

Let's drop the $\tilde{}$ and just use $X$, since what we are doing does not depend on working with a matrix with 0 means, although that is the usual case.

If you view $U$ as a new variable, this is like finding one predictor variable $U$ so that the regressions (without intercept) of each column of $X$ on $U$ is a good fit.

Contrast this with the usual regression situation where you are *given* a predictor variable $Z$ and seek to find coefficients. Here you need to find both predictor and coefficients.

The first thought many statisticians would have would be to find $U$ and $v$ so that $Uv'$ is close to $X$ by the *least squares* criterion.

That is, find $U$ and $v$ so as to minimize

$$\sum_{1 \le i \le N}\sum_{1 \le \ell \le p}(x_{i\ell} - u_i v_\ell)^2 = \| X - Uv' \|^2$$

That's what we're going to do.

**Notation**: When $A = [a_{ij}]$ is a matrix

$$\| A \|^2 = \sum_i \sum_j a_{ij}^2 = \text{trace}(A'A)$$

The *Singular Value Decomposition* (SVD) of $X$ is the key to finding $U$ and $v$ to minimize $\| X - Uv' \|^2$.

The SVD of $X$ is a mathematical representation of $X$ which is useful in many contexts.

9

## The Singular Value Decomposition

For *any* N×p matrix $X$ with $N \ge p$, there are always <u>three matrices</u> $L$, $R$ and $T$ such that $X = LTR'$ where

- $L = [L_1, L_2, ..., L_p]$ is $N \times p$ with $L'L = I_p$. That is, the columns $L_1, ... , L_p$ of $L$ are *orthonormal*:

$$L_j'L_j = 1, \quad L_j'L_k = 0, \; j \ne k$$

  **Note**: When $N > p$, this does not mean that $L^{-1} = L'$, since $L$ is not square.

- $R = [r_1, r_2, ... , r_p]$ is $p \times p$ *square* with $R'R = I_p$, that is, the columns $r_1, ... , r_p$ of $R$ are *orthonormal*:

$$r_j'r_j = 1, \quad r_j'r_k = 0, \; j \ne k$$

  Since $R$ is p×p, this means that $R^{-1} = R'$ and $RR' = I_p$. $R$ is an <u>orthogonal matrix</u>.

- $T = \text{diag}[t_1, ..., t_p]$, $p \times p$, *diagonal* with $t_1 \ge t_2 \ge ... \ge t_p \ge 0$

10

## Vocabulary

$X = LTR'$ is the *singular value decomposition* (SVD) of $X$

## Facts

- $X = LTR' = \sum_{1 \le k \le p} t_k L_k r_k' = \sum_{1 \le k \le p}(t_k L_k)r_k'$, a sum of *p outer products* of $t_k L_k$ and $r_k$

- The $t_i$'s are unique

- When $t_i \ne t_j$, all $j \ne i$, $L_i$ and $r_i$ are unique (except for multiplication of both by $-1$ : $(-L_i)(-r_i)' = L_i r_i'$)

Thus the *SVD* $X = LTR'$ of $X$ is <u>essentially unique</u>.

- The N by 1 vectors $L_j$, $j = 1,...,p$ are the **left singular vectors** of $X$.

- The p by 1 vectors $r_j$, $j = 1,...,p$ are the **right singular vectors** of $X$.

- The p scalars $t_1 \ge t_2 \ge ... \ge t_p \ge 0$ are the **singular values** of $X$.

11

When there are only $s < p$ singular values $t_i \ne 0$, so that $t_{s+1} = ... = t_p = 0$,

$$X = \sum_{1 \le k \le s}(t_k L_k)r_k'$$

a sum of only s outer products.

## Fact:

- $\text{Rank}(X) = s =$ number of non-zero singular values. When $s < p$,

$$t_s > t_{s+1} = t_{s+2} = ... = t_p = 0.$$

The SVD is often the best way numerically to determine the Rank($X$):

- Compute $T$ from $X$

- Count how many diagonal elements are non-zero except for rounding error. This should be Rank($X$).

12

## Computing the SVD in MacAnova: `svd()`

Suppose `x` is a REAL matrix.  Then
$svd(x)$ computes the vector $[t_1,...,t_p]'$ of
<u>singular values</u> of x (diag($T$), *not* $T$)

$svd(x,left:T)$ computes a structure with
two components:
- `values`,  vector of <u>singular values</u>
- `leftvectors` matrix $L$ whose columns
  are $L_1, ..., L_p$, the *left* <u>singular vectors</u>

$svd(x,right:T)$ computes a structure
with two components:
- `values`, vector of singular values
- `rightvectors` matrix $R$ whose columns
  are $r_1, ..., r_p$, *right* <u>singular vectors</u>

$svd(x,right:T,left:T)$ or $svd(x,all:T)$
computes a 3 component structure:
- `values`: <u>singular values</u>
- `leftvectors`: *left* <u>singular vectors</u>
- `rightvectors`: *right* <u>singular vectors</u>

## Example

```
Cmd> x <- run(10)^run(0,3)' # powers of run(10)

Cmd> setlabels(x,\
       structure("@",vector("i^0","i^1","i^2","i^3")))

Cmd> x # 10 ny 4 matrix
          i^0        i^1        i^2        i^3
(1)         1          1          1          1
(2)         1          2          4          8
(3)         1          3          9         27
(4)         1          4         16         64
(5)         1          5         25        125
(6)         1          6         36        216
(7)         1          7         49        343
(8)         1          8         64        512
(9)         1          9         81        729
(10)        1         10        100       1000

Cmd> vals <- svd(x); vals #just Sing values
(1)      1415.4      27.14     2.2961     0.41587
```

$X$ has rank 4, but since the two smallest
singular value are so small, it is close to
having rank 3, or possibly even rank 2.

```
Cmd> results <- svd(x,left:T,right:T)# or all:T

Cmd> compnames(results)
(1) "values"              Structure
(2) "leftvectors"         component
(3) "rightvectors"        names
```
- `values`　　　　　$p$-vector, $(t_1, t_2, ..., t_p)$
- `leftvectors`　　$N$ by $p$ $L = [L_1 ... L_p]$
- `rightvectors`　$p$ by $p$ square $R = [r_1...r_p]$

You can construct the diagonal matrix $T$
in the SVD by

```
Cmd> tmatrix <- dmat(results$values); tmatrix#Diagonal matrix T
(1,1)     1415.4          0          0          0
(2,1)          0      27.14          0          0
(3,1)          0          0     2.2961          0
(4,1)          0          0          0     0.41587
```

Here are <u>numerical checks</u> that the right
and left singular vectors are orthonormal
($R'R = I_p$ and $L'L = I_p$):

```
Cmd> R <- results$rightvectors # right singular vectors

Cmd> list(R) # size is p by p
R                    REAL    4    4      (labels)

Cmd> R' %*% R # = I_4
           (1)        (2)         (3)         (4)
(1) _____1  4.0533e-17 -6.9218e-17  4.9043e-17
(2)  4.0533e-17 _____1  6.2095e-17  4.8526e-17
(3) -6.9218e-17  6.2095e-17 _____1 -1.4827e-16
(4)  4.9043e-17  4.8526e-17 -1.4827e-16 _____1
```

This is $I_4$.

```
Cmd> L <- results$leftvectors # left singular vectors

Cmd> list(L) # size is N by p
L                    REAL   10    4      (labels)

Cmd> L' %*% L # = I_4
           (1)        (2)         (3)         (4)
(1) _____1 -6.9389e-18 -1.1102e-16 -2.3592e-16
(2) -6.9389e-18 _____1  1.6653e-16  3.8858e-16
(3) -1.1102e-16  1.6653e-16 _____1 -4.4409e-16
(4) -2.3592e-16  3.8858e-16 -4.4409e-16 _____1
```

This is also $I_4$.

## Relationship with Eigenvalues and Eigenvectors

- Each <u>right</u> singular vector $r_j$ is an
  eigenvector of $X'X$ with eigenvalue $t_j^2$.

  That is, $r_j$ satisfies $X'Xr_j = t_j^2 r_j$.

  <u>Check</u>:  $X'Xr_j = RTL'LTR'r_j = t_j^2 r_j$,
  because $L'L = I_p$ and $R'R = I_p$.

```
Cmd> sqrt(eigenvals(x' %*% x)) # numerical check
(1)      1415.4      27.14     2.2961     0.41587
```

$X'X$ is the $p{\times}p$ matrix of sums of
squares (SS) and sums of products (SP)
of the *columns* of $X$.

And, for the case we apply this to, $\tilde{X}'\tilde{X}$
consists of sums of squares $\sum_i(x_{i\ell}-\overline{x}_\ell)^2$
and products $\sum_i(x_{i\ell} - \overline{x}_\ell)(x_{ik} - \overline{x}_k)$.  An,
of course, $S_x = (1/(N-1))\tilde{X}'\tilde{X}$.

- Each <u>left</u> singular vector $L_j$ is an eigenvector of the N by N matrix $XX'$ with eigenvalue $t_j^2$.

  That is, $L_j$ satisfies $XX'L_j = t_j^2 L_j$.

  <u>Check</u>: $XX'L_j = LTR'RTL'L_j = t_j^2 L_j$ since $R'R = I_p$ and $L'L = I_p$.

```
Cmd> sqrt(round(eigenvals(x %*% x'),9)) # numerical check
(1)      1415.4        27.14       2.2961      0.41587           0
(6)           0            0            0            0           0
```

  There are N – p = 6 zero eigenvalues.

  $XX'$ is the N×N matrix of SS and SP of the *rows* of $X$.

---

## Summary

If $v_1$, $v_2$, ..., $v_p$ are eigenvectors of $X'X$ with eigenvalues $\lambda_1 \geq \lambda_2 \geq \lambda_p \geq 0$ then

- the $j^{th}$ *singular value* is
  $$t_j = \sqrt{\lambda_j}$$
- the $j^{th}$ *right singular vector* is
  $$r_j = v_j \text{ (could be } -v_j\text{), } j = 1, ..., p$$

If $\ell_1$, $\ell_2$, ..., $\ell_N$ (all N by 1) are eigenvectors of $XX'$ with eigenvalues $\lambda_1 \geq \lambda_2 \geq \lambda_p \geq ... \geq \lambda_N$ then

- the $j^{th}$ *singular value* is
  $$t_j = \sqrt{\lambda_j}, \; j = 1, ..., p$$
- the $j^{th}$ *left singular vector* is
  $$L_j = \ell_j \text{ (could be } -\ell_j\text{), } j = 1, ..., p$$
- The remaining eigenvalues $\lambda_{p+1}$, ..., $\lambda_N$ are 0
- The remaining eigenvectors $\ell_{p+1}$, ..., $\ell_N$ are irrelevant

---

Now define p linear combinations of the columns of $X$ with coefficients from the <u>right</u> singular vectors $r_j$, j = 1,... p:

$$U_j \equiv Z_j = Xr_j = LTR'r_j = t_j L_j, \; j = 1,...,p$$

$Z_j = Xr_j = \sum_{1 \leq k \leq p} r_{kj} X_k$ is a linear combination of the columns of $X$.

The coefficients (weights) are the elements of

  $r_j = j^{th}$ right singular vector of $X$

  $= j^{th}$ eigenvector of $X'X$.

Because $R'R = I_p$, $R'r_j = [0 \; 0 \; ... \; 1 \; ... \; 0]'$ and so

  $Z_j = LTR'r_j = t_j L_j$

is proportional to a left singular vector of $X$.

---

## Back to low rank approximations

Q. With m = 1, what N×1 $U$ and p×1 $v$ minimize (make smallest) the "residual SS"

$$\|X - Uv'\|^2 = \sum_{1 \leq i \leq N} \sum_{1 \leq \ell \leq p} (x_{i\ell} - u_i v_\ell)^2 ?$$

A.    $U = Z_1 = t_1 L_1$ and $v = r_1$

  That is

  $$\hat{X}^{(1)} \equiv Z_1 r_1' = t_1 L_1 r_1' = X r_1 r_1'$$

  is the best rank 1 approximation to $X$ in the least squares sense.

This generalizes to rank m > 1:

$$\hat{X}^{(m)} = \sum_{1 \leq j \leq m} Z_j r_j' = \sum_{1 \leq j \leq m} t_j L_j r_j'$$
$$= X\left(\sum_{1 \leq j \leq m} r_j r_j'\right)$$

is the <u>best rank m approximation</u> to $X$ in the least squares sense.

How good is the approximation?

- The "residual sum of squares" is

$$\| \mathbf{X} - \hat{\mathbf{X}}^{(m)} \|^2 = \sum_{1 \le i \le N} \sum_{1 \le \ell \le p} (x_{i\ell} - \hat{x}_{i\ell}^{(m)})^2$$
$$= \sum_{m+1 \le k \le p} t_k^2 = \sum_{m+1 \le k \le p} \lambda_k$$

  = the sum of the squared smallest
  p – m singular values of **X**

  = sum of the smallest p – m eigen-
  values of **X′X**.

- The "total sum of squares" is

$$\| \mathbf{X} \|^2 = \sum_\ell \sum_i x_{i\ell}^2 = \sum_{1 \le k \le p} t_k^2 = \sum_{1 \le k \le p} \lambda_k = \mathrm{tr}\ \mathbf{X'X}$$

Therefore, when the ratio

$$\frac{\| \mathbf{X} - \hat{\mathbf{X}}^{(m)} \|^2}{\| \mathbf{X} \|^2} = \frac{\sum_{1 \le i \le N} \sum_{1 \le \ell \le p} (x_{i\ell} - \hat{x}_{i\ell}^{(m)})^2}{\sum_{1 \le i \le N} \sum_{1 \le \ell \le p} x_{i\ell}^2} = \frac{\sum_{m+1 \le k \le p} t_k^2}{\sum_{1 \le k \le p} t_k^2}$$

is small, the approximation is pretty
good.  This ratio is analogous to

$$SS_{residual} / SS_{total} = 1 - R^2$$

in regression.

In the rank m approximation,

$$\hat{\mathbf{X}}^{(m)} = \sum_{1 \le k \le m} \mathbf{Z}_k \mathbf{r}_k' = \sum_{1 \le k \le m} t_k \mathbf{L}_k \mathbf{r}_k' ,$$

column $\mathbf{X}_\ell$ of **X** is approximated by

$$\hat{\mathbf{X}}_\ell^{(m)} = \sum_{1 \le k \le m} \mathbf{Z}_k r_{\ell k} = \sum_{1 \le k \le m} r_{\ell k} \mathbf{Z}_k ,$$

a linear combination of $\mathbf{Z}_1$, $\mathbf{Z}_2$, ..., $\mathbf{Z}_m$.

Since the $\mathbf{Z}_k$'s themselves are linear
combinations of the columns of **X** ($\mathbf{Z}_k$ =
$\mathbf{X}\mathbf{r}_k$), so are the columns of $\hat{\mathbf{X}}^{(m)}$:

$$\hat{\mathbf{X}}_\ell^{(m)} = \sum_{1 \le k \le m} r_{\ell k} \mathbf{X} \mathbf{r}_k$$
$$= \sum_{1 \le \ell \le m} (\sum_{1 \le k \le m} r_{\ell k} r_{jk}) \mathbf{X}_j$$

$\sum_{1 \le k \le m} r_{\ell k} r_{jk}$ is a partial sum of squares (j
= ℓ) or sum of products (j ≠ ℓ) of <u>rows</u>
of **R**.  Since **RR′** = $\mathbf{I}_p$,

$$\sum_{1 \le k \le m} r_{\ell k}^2 = 1 - \sum_{m+1 \le k \le p} r_{\ell k}^2$$
$$\sum_{1 \le k \le m} r_{\ell k} r_{jk} = -\sum_{m+1 \le k \le p} r_{\ell k} r_{jk}$$