

# STAD29: Statistics for the Life and Social Sciences

Lecture notes

# Section 1

## Principal components

# Principal Components

- Have measurements on (possibly large) number of variables on some individuals.
- Question: can we describe data using fewer variables (because original variables correlated in some way)?
- Look for direction (linear combination of original variables) in which values *most spread out*. This is *first principal component*.
- Second principal component then direction uncorrelated with this in which values then most spread out. And so on.

# Principal components

- See whether small number of principal components captures most of variation in data.
- Might try to interpret principal components.
- If 2 components good, can make plot of data.
- (Like discriminant analysis, but no groups.)
- “What are important ways that these data vary?”

# xxx Packages

You might not have installed the first of these. See over for instructions.

```
library(ggbiplot) # see over  
library(tidyverse)  
library(ggrepel)
```

## xxx Installing ggbiplot

- ggbiplot not on CRAN, so usual `install.packages` will not work. This is same procedure you used for `smmr` in C32:
- Install package `devtools` first (once):

```
install.packages("devtools")
```

- Then install `ggbiplot` (once):

```
library(devtools)  
install_github("vqv/ggbiplot")
```

## xxx Small example: 2 test scores for 8 people

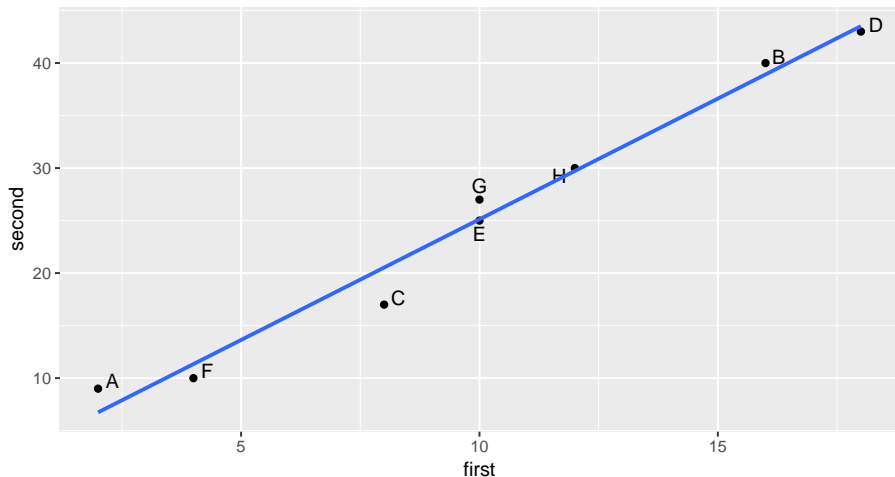
```
my_url <- "http://www.uts.utoronto.ca/~butler/d29/test12.txt"
test12 <- read_table2(my_url)
test12
```

```
## # A tibble: 8 x 3
##   first second id
##   <dbl>  <dbl> <chr>
## 1      2      9 A
## 2     16     40 B
## 3      8     17 C
## 4     18     43 D
## 5     10     25 E
## 6      4     10 F
## 7     10     27 G
## 8     12     30 H
```

```
g <- ggplot(test12, aes(x = first, y = second, label = id)) +
  geom_point() + geom_text_repel()
```

## xxx The plot

```
g + geom_smooth(method = "lm", se = F)
```





# xxx Principal component analysis

- Grab just the numeric columns:

```
test12 %>% select_if(is.numeric) -> test12_numbers
```

- Strongly correlated, so data nearly 1-dimensional:

```
cor(test12_numbers)
```

```
##           first    second
## first  1.000000  0.989078
## second 0.989078  1.000000
```

# Finding principal components xxx

- Make a score summarizing this one dimension. Like this:

```
test12.pc <- princomp(test12_numbers, cor = T)
summary(test12.pc)
```

```
## Importance of components:
```

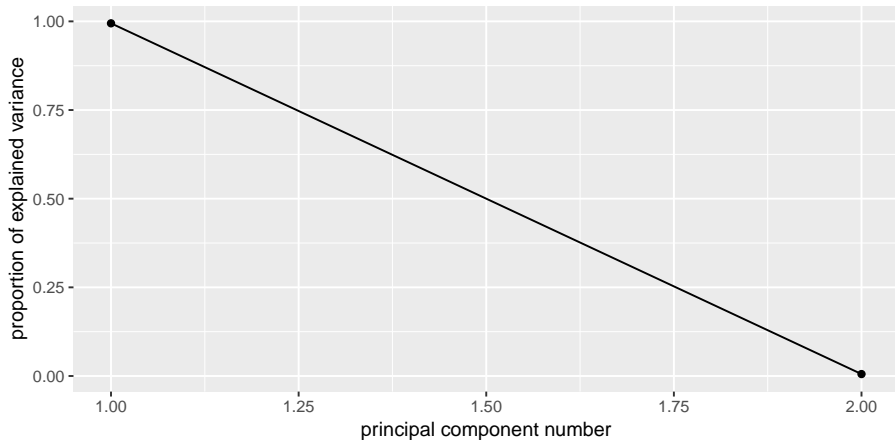
##	Comp.1	Comp.2
## Standard deviation	1.410347	0.104508582
## Proportion of Variance	0.994539	0.005461022
## Cumulative Proportion	0.994539	1.000000000

# Comments

- “Standard deviation” shows relative importance of components (as for LDs in discriminant analysis)
- Here, first one explains almost all (99.4%) of variability.
- That is, look only at first component and ignore second.
- $\text{cor=T}$  standardizes all variables first. Usually wanted, because variables measured on different scales. (Only omit if variables measured on same scale and expect similar variability.)

# Scree plot xxx

```
ggscreeplot(test12.pc)
```



## xxx Component loadings

explain how each principal component depends on (standardized) original variables (test scores):

```
test12.pc$loadings
```

```
##
## Loadings:
##          Comp.1 Comp.2
## first    0.707  0.707
## second   0.707 -0.707
##
##          Comp.1 Comp.2
## SS loadings      1.0    1.0
## Proportion Var   0.5    0.5
## Cumulative Var   0.5    1.0
```

First component basically negative sum of (standardized) test scores. That is, person tends to score similarly on two tests, and a composite score would summarize performance.

## xxx Component scores

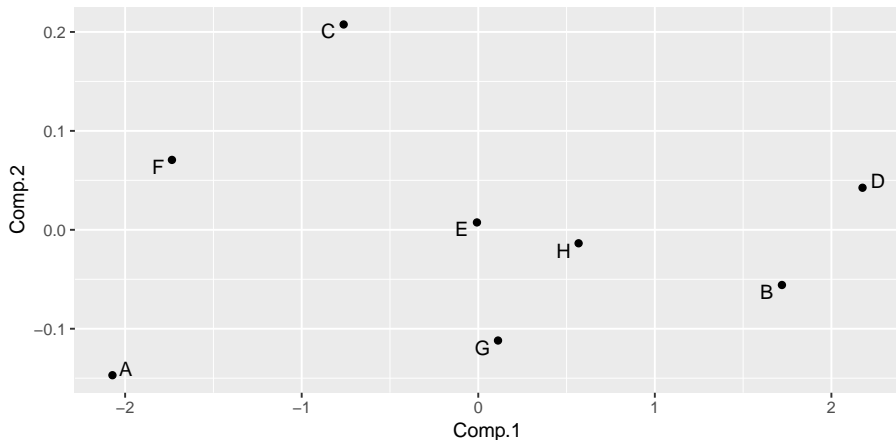
```
d <- data.frame(test12, test12.pc$scores)
d
```

##	first	second	id	Comp.1	Comp.2
## 1	2	9	A	-2.071819003	-0.146981782
## 2	16	40	B	1.719862811	-0.055762223
## 3	8	17	C	-0.762289708	0.207589512
## 4	18	43	D	2.176267535	0.042533250
## 5	10	25	E	-0.007460609	0.007460609
## 6	4	10	F	-1.734784030	0.070683441
## 7	10	27	G	0.111909141	-0.111909141
## 8	12	30	H	0.568313864	-0.013613668

- Person A is a low scorer, high positive comp.1 score.
- Person D is high scorer, high negative comp.1 score.
- Person E average scorer, near-zero comp.1 score.

## xxx Plot of scores

```
ggplot(d, aes(x = Comp.1, y = Comp.2, label = id)) +  
  geom_point() + geom_text_repel()
```



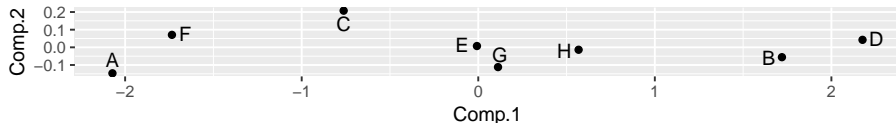
# xxx Comments

- Vertical scale exaggerates importance of comp.2.
- Fix up to get axes on same scale:

```
g <- ggplot(d, aes(x = Comp.1, y = Comp.2, label = id)) +  
  geom_point() + geom_text_repel() +  
  coord_fixed()
```

- Shows how exam scores really spread out along one dimension:

g



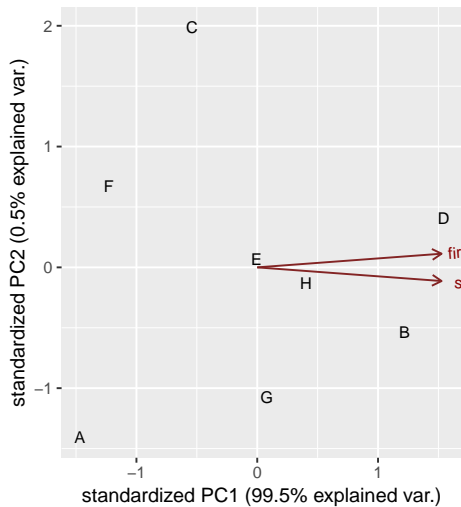


## xxx The biplot

- Plotting variables and individuals on one plot.
- Shows how components and original variables related.
- Shows how individuals score on each component, and therefore suggests how they score on each variable.
- Add `labels` option to identify individuals:

```
g <- ggbiplot(test12.pc, labels = test12$id)
```

## xxx The biplot



## xxx Comments

- Variables point almost same direction (left). Thus very negative value on comp.1 goes with high scores on both tests, and test scores highly correlated.
- Position of individuals on plot according to scores on principal components, implies values on original variables. Eg.:
- D very negative on comp.1, high scorer on both tests.
- A and F very positive on comp.1, poor scorers on both tests.
- C positive on comp.2, high score on first test relative to second.
- A negative on comp.2, high score on second test relative to first.

# xxx Track running data

- (1984) track running records for distances 100m to marathon, arranged by country. Countries labelled by (mostly) Internet domain names (ISO 2-letter codes):

```
my_url <- "http://www.utsc.utoronto.ca/~butler/d29/men_track_1
track <- read_table(my_url)
track %>% sample_n(12)
```

```
## # A tibble: 12 x 9
##      m100  m200  m400  m800 m1500 m5000 m10000 marathon coun
##      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>      <dbl> <chr>
##  1  12.2  23.2  52.9  2.02  4.24  16.7   35.4      165. ck
##  2  10.1   20    44.6  1.75  3.59  13.2   27.5      131. ru
##  3  10.6  20.5  45.9  1.78  3.61  13.5   28.1      131. dk
##  4  10.3  20.9  46.9  1.79  3.77  14.0   29.2      136. kr
##  5  10.6  21.5  47.8  1.84  3.92  14.7   30.8      149. id
##  6  10.6  21.3  46.8  1.79  3.77  14.1   30.1      139. tw
```

## xxx Country names

Also read in a table to look country names up in later:

```
my_url <- "http://www.utsc.utoronto.ca/~butler/d29/isocodes.csv"
iso <- read_csv(my_url)
iso
```

```
## # A tibble: 251 x 4
```

##	Country	ISO2	ISO3	M49
##	<chr>	<chr>	<chr>	<dbl>
## 1	<NA>	<NA>	<NA>	NA
## 2	Afghanistan	af	afg	4
## 3	Aland Islands	ax	ala	248
## 4	Albania	al	alb	8
## 5	Algeria	dz	dza	12
## 6	American Samoa	as	asm	16
## 7	Andorra	ad	and	20
## 8	Angola	ao	ago	24

## xxx Data and aims

- Times in seconds 100m–400m, in minutes for rest (800m up).
- This taken care of by standardization.
- 8 variables; can we summarize by fewer and gain some insight?
- In particular, if 2 components tell most of story, what do we see in a plot?

## xxx Fit and examine principal components

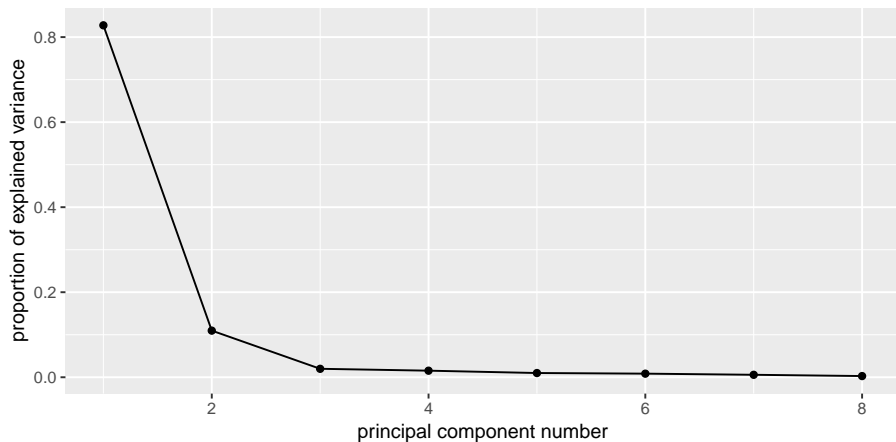
```
track_num <- track %>% select_if(is.numeric)
track.pc <- princomp(track_num, cor = T)
summary(track.pc)
```

```
## Importance of components:
```

```
##              Comp.1      Comp.2
## Standard deviation      2.5733531 0.9368128
## Proportion of Variance 0.8277683 0.1097023
## Cumulative Proportion 0.8277683 0.9374706
##              Comp.3      Comp.4
## Standard deviation      0.39915052 0.35220645
## Proportion of Variance 0.01991514 0.01550617
## Cumulative Proportion 0.95738570 0.97289187
##              Comp.5      Comp.6
## Standard deviation      0.282630981 0.260701267
## Proportion of Variance 0.009985034 0.008495644
## Cumulative Proportion 0.999976000 0.999976000
```

# xxx Scree plot

```
ggscreeplot(track.pc)
```





## xxx How many components?

- As for discriminant analysis, look for “elbow” in scree plot.
- See one here at 3 components; everything 3 and beyond is “scree”.
- So take 2 components.
- Note difference from discriminant analysis: want “large” rather than “small”, so go 1 step left of elbow.
- Another criterion: any component with eigenvalue bigger than about 1 is worth including. 2nd one here has eigenvalue just less than 1.
- Refer back to summary: cumulative proportion of variance explained for 2 components is 93.7%, pleasantly high. 2 components tell almost whole story.

## xxx How do components depend on original variables?

Loadings:

```
track.pc$loadings
```

```
##
```

```
## Loadings:
```

```
##          Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7
## m100      0.318  0.567  0.332  0.128  0.263  0.594  0.136
## m200      0.337  0.462  0.361 -0.259 -0.154 -0.656  -
0.113
## m400      0.356  0.248 -0.560  0.652 -0.218 -0.157
## m800      0.369          -0.532 -0.480  0.540          -
0.238
## m1500     0.373 -0.140 -0.153 -0.405 -0.488  0.158  0.610
## m5000     0.364 -0.312  0.190          -0.254  0.141  -
0.591
## m10000    0.367 -0.307  0.182          -0.133  0.219  -
```

## xxx Comments

- comp.1 loads about equally (has equal weight) on times over all distances.
- comp.2 has large positive loading for long distances, large negative for short ones.
- comp.3: large negative for middle distance, large positive especially for short distances.
- Country overall good at running will have lower than average record times at all distances, so comp.1 *large*. Conversely, for countries bad at running, comp.1 very negative.
- Countries relatively better at sprinting (low times) will be *positive* on comp.2; countries relatively better at distance running *negative* on comp.2.

## xxx Commands for plots

- Principal component scores (first two). Also need country names.

```
d <- data.frame(track.pc$scores,
  country = track$country
)
names(d)
```

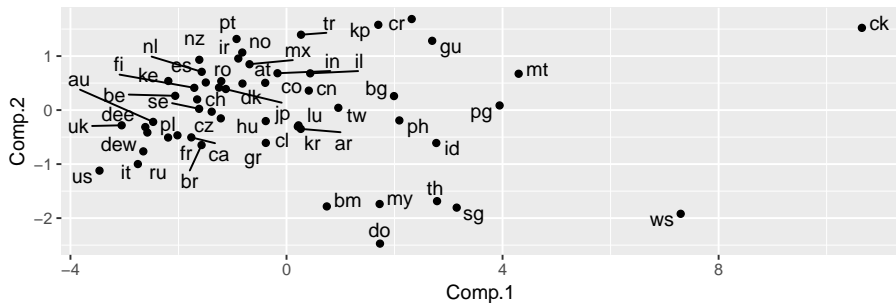
```
## [1] "Comp.1" "Comp.2" "Comp.3" "Comp.4" "Comp.5"
## [6] "Comp.6" "Comp.7" "Comp.8" "country"
```

```
g1 <- ggplot(d, aes(
  x = Comp.1, y = Comp.2,
  label = country
)) +
  geom_point() + geom_text_repel() +
  coord_fixed()
```

- Biplot:

## xxx Principal components plot

g1

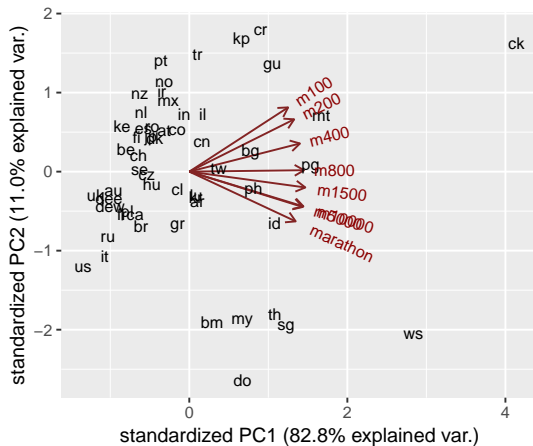


## xxx Comments on principal components plot

- Good running countries at right of plot: US, UK, Italy, Russia, East and West Germany.
- Bad running countries at left: Western Samoa, Cook Islands.
- Better sprinting countries at bottom: US, Italy, Russia, Brazil, Greece. do is Dominican Republic, where sprinting records relatively good, distance records very bad.
- Better distance-running countries at top: Portugal, Norway, Turkey, Ireland, New Zealand, Mexico. ke is Kenya.

## xxx Biplot

g2



## xxx Comments on biplot

- Had to do some pre-work to interpret PC plot. Biplot more self-contained.
- All variable arrows point left; countries on left have large (bad) record times overall, countries on right good overall.
- Variable arrows extend negatively as well. Top left = bad at distance running, bottom right = good at distance running.
- Bottom left = bad at sprinting, top right = good at sprinting.
- Doesn't require so much pre-interpretation of components.



# xxx How do I know which country is which?

Need to look up two-letter abbreviations in ISO table, eg. for best 8 running countries:

```
d %>%
  arrange(desc(Comp.1)) %>%
  left_join(iso, by = c("country" = "ISO2")) %>%
  select(Comp.1, country, Country) %>%
  slice(1:8)
```

##	Comp.1	country	Country
## 1	10.652914	ck	Cook Islands
## 2	7.297865	ws	Samoa
## 3	4.297909	mt	Malta
## 4	3.945224	pg	Papua New Guinea
## 5	3.150886	sg	Singapore
## 6	2.787273	th	Thailand
## 7	2.773125	id	Indonesia

## xxx Best 8 running countries

```
d %>%
  arrange(Comp.1) %>%
  left_join(iso, by = c("country" = "ISO2")) %>%
  select(Comp.1, country, Country) %>%
  slice(1:8)
```

	Comp.1	country	Country
## 1	-3.462175	us	United States of America
## 2	-3.052104	uk	United Kingdom
## 3	-2.752084	it	Italy
## 4	-2.651062	ru	Russian Federation
## 5	-2.613964	dee	East Germany
## 6	-2.576272	dew	West Germany
## 7	-2.468919	au	Australia
## 8	-2.191917	fr	France

## xxx Worst 8 running countries

```
d %>%
  arrange(desc(Comp.1)) %>%
  left_join(iso, by = c("country" = "ISO2")) %>%
  select(Comp.1, country, Country) %>%
  slice(1:8)
```

	##	Comp.1	country	Country
	## 1	10.652914	ck	Cook Islands
	## 2	7.297865	ws	Samoa
	## 3	4.297909	mt	Malta
	## 4	3.945224	pg	Papua New Guinea
	## 5	3.150886	sg	Singapore
	## 6	2.787273	th	Thailand
	## 7	2.773125	id	Indonesia
	## 8	2.697066	gu	Guam

## xxx Better at distance running

```
d %>%
  arrange(desc(Comp.2)) %>%
  left_join(iso, by = c("country" = "ISO2")) %>%
  select(Comp.2, country, Country) %>%
  slice(1:8)
```

	Comp.2	country	Country
## 1	1.6860391	cr	Costa Rica
## 2	1.5791490	kp	Korea (North)
## 3	1.5226742	ck	Cook Islands
## 4	1.3957839	tr	Turkey
## 5	1.3167578	pt	Portugal
## 6	1.2829272	gu	Guam
## 7	1.0663756	no	Norway
## 8	0.9547437	ir	Iran, Islamic Republic of

## xxx Better at sprinting

```
d %>%
  arrange(Comp.2) %>%
  left_join(iso, by = c("country" = "ISO2")) %>%
  select(Comp.2, country, Country) %>%
  slice(1:10)
```

	Comp.2	country	Country
## 1	-2.4715736	do	Dominican Republic
## 2	-1.9196130	ws	Samoa
## 3	-1.8055052	sg	Singapore
## 4	-1.7832229	bm	Bermuda
## 5	-1.7386063	my	Malaysia
## 6	-1.6851772	th	Thailand
## 7	-1.1204235	us	United States of America
## 8	-0.9989821	it	Italy
## 9	-0.7639385	ru	Russian Federation
## 10	-0.6170624	br	Brazil

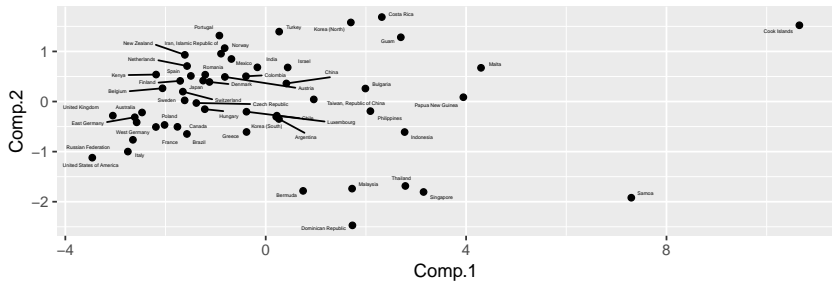
## xxx Plot with country names

```
g <- d %>%
  left_join(iso, by = c("country" = "ISO2")) %>%
  select(Comp.1, Comp.2, Country) %>%
  ggplot(aes(x = Comp.1, y = Comp.2, label = Country)) +
  geom_point() + geom_text_repel(size = 1) +
  coord_fixed()
```

```
## Warning: Column `country`/`ISO2` joining factor and
## character vector, coercing into character vector
```

## xxx The plot

g



# xxx Principal components from correlation matrix

Create data file like this: cov.txt and read in like this:

```
my_url <- "http://www.utsc.utoronto.ca/~butler/d29/cov.txt"
mat <- read_table(my_url, col_names = F)
mat
```

```
## # A tibble: 3 x 3
##       X1      X2      X3
##   <dbl> <dbl> <dbl>
## 1  1      0.970 -0.96
## 2  0.970  1      -0.998
## 3 -0.96  -0.998  1
```



## xxx Pre-processing

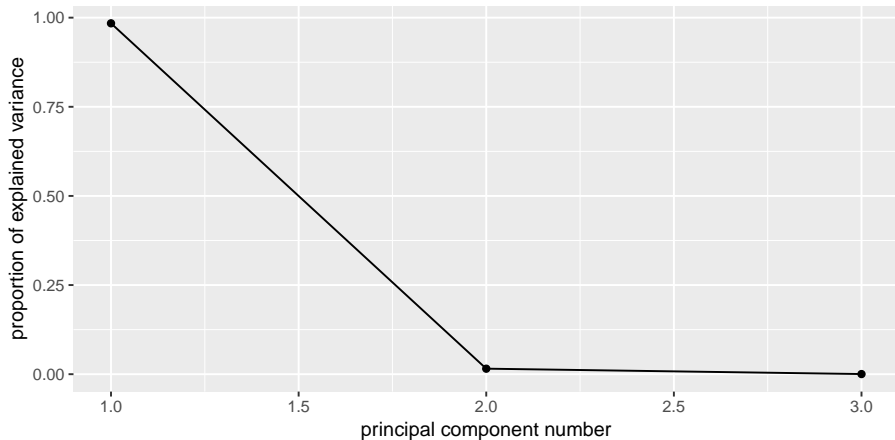
A little pre-processing required:

- Turn into matrix (from data frame)
- Feed into princomp as covmat=

```
mat.pc <- mat %>%  
  as.matrix() %>%  
  princomp(covmat = .)
```

# xxx Scree plot: one component fine

```
ggscreeplot(mat.pc)
```



# xxx Component loadings

Compare correlation matrix:

```
""r mat ""
```

```
"" A tibble: 3 x 3 X1 X2 X3 <dbl> <dbl>
<dbl> 1 1 0.970 -0.96 2 0.970 1 -0.998 3
-0.96 -0.998 1 ""
```

with component loadings

```
""r mat.pcloadings""
```

```
"" Loadings: Comp.1 Comp.2 Comp.3 X1
0.573 0.812 0.112 X2 0.581 -0.306 -0.755 X3
-0.578 0.498 -0.646 Comp.1 Comp.2 Comp.3
SS loadings 1.000 1.000 1.000 Proportion
Var 0.333 0.333 0.333 Cumulative Var 0.333
0.667 1.000 ""
```

\* When X1 large, X2 also large, X3 small.

\* Then 'comp.1' \*negative\*.

\* When X1 small, X2 small, X3 large.

\* Then 'comp.1' \*positive\*.

## xxx No scores

- With correlation matrix rather than data, no component scores
- So no principal component plot
- and no biplot.

```
## Error in FUN(X[[i]], ...): invalid 'name' argument
```