

```
In [ ]: # Initialize Otter
import otter
grader = otter.Notebook("lab5-pca.ipynb")

In [ ]: import numpy as np
import pandas as pd
import altair as alt
from scipy import linalg
from statsmodels.multivariate.pca import PCA
# disable row limit for plotting
alt.data_transformers.disable_max_rows()
# uncomment to ensure graphics display with pdf
# alt.renderers.enable('mimetype')
```

Lab 5: Principal components

Principal components analysis (PCA) is a widely-used multivariate analysis technique. Depending on the application, PCA is variously described as:

- a dimension reduction method
- an approximation method
- a latent factor model
- a filtering or compression method

The core technique of PCA is *finding linear data transformations that preserve variance*.

What does it mean to say that 'principal components are linear data transformations'?
Suppose you have a dataset with n observations

and p variables. We can represent the values as a data matrix \mathbf{X} with n rows and p columns:

$$\mathbf{X} = \underbrace{\begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_p \end{bmatrix}}_{\text{column vectors}} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$

To say that the principal components are linear data transformations means that each principal component is of the form:

$$\text{PC} = \mathbf{X}\mathbf{v} = v_1\mathbf{x}_1 + v_2\mathbf{x}_2 + \cdots + v_p\mathbf{x}_p$$

for some vector \mathbf{v} . In PCA, the following terminology is used:

- linear combination coefficients v_j are known as *loadings*
- values of the linear combinations are known as *scores*
- the vector of loadings \mathbf{v} is known as a *principal axis*

As discussed in lecture, the values of the loadings are found by decomposing the correlation structure.

Objectives

In this lab, you'll focus on computing and interpreting principal components:

- finding the loadings (linear combination coefficients) for each PC;
- quantifying the variation captured by each PC;
- visualization-based techniques for selecting a number of PC's to A(nalyze);
- plotting and interpreting loadings.

You'll work with a selection of county summaries from the 2010 U.S. census. The first few rows of the dataset are shown below:

```
In [ ]: # import tidy county-level 2010 census data
census = pd.read_csv('data/census2010.csv', enc
census.head()
```

The observational units are U.S. counties, and each row is an observation on one county. The values are, for the most part, percentages of the county population. You can find variable descriptions in the metadata file `census2010metadata.csv` in the data directory.

Correlations

PCA identifies variable combinations that capture covariation by decomposing the correlation matrix. So, to start with, let's examine the correlation matrix for the 2010 county-level census data to get a sense of which variables tend to vary together.

The correlation matrix is a matrix of all pairwise correlations between variables. If x_{ij} denotes the value for the i th observation of variable j , then the entry at row j and column k of the correlation matrix \mathbf{R} is:

$$r_{jk} = \frac{\sum_i (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)}{S_j S_k}$$

In the census data, the `State` and `County` columns indicate the geographic region for each observation; essentially, they are a row index. So we'll drop them before computing the matrix \mathbf{R} :

```
In [ ]: # store quantitative variables separately
x_mx = census.drop(columns = ['State', 'County'])
```

From here, the matrix is simple to compute in pandas using `.corr()` :

```
In [ ]: # correlation matrix
corr_mx = x_mx.corr()
```

The matrix can be inspected directly to determine which variables vary together. For example, we could look at the correlations between employment rate and every other variable in the dataset by extracting the `Employed` column from the correlation matrix and sorting the correlations:

```
In [ ]: # correlation between employment rate and other
corr_mx.loc[:, 'Employed'].sort_values()
```

Recall that correlation is a number in the interval $[-1, 1]$ whose magnitude indicates the strength of the linear relationship between variables:

- correlations near -1 are *strongly negative*, and mean that the variables *tend to vary in opposition*
- correlations near 1 are *strongly positive*, and mean that the variables *tend to vary together*

From examining the output above, it can be seen that the percentage of the county population that is employed is:

- strongly *negatively* correlated with child poverty, poverty, and unemployment, meaning it *tends to vary in opposition* with these variables
- strongly *positively* correlated with income per capita, meaning it *tends to vary together* with this variable

If instead we wanted to look up the correlation between just two variables, we could retrieve the relevant entry directly using `.loc[...]` with the variable names:

```
In [ ]: # correlation between employment and income per  
corr_mx.loc['Employed', 'IncomePerCap']
```

So across U.S. counties employment is, perhaps unsurprisingly, strongly and positively correlated with income per capita, meaning that higher

employment rates tend to coincide with higher incomes per capita.

Question 1

Find the correlation between the poverty rate and demographic minority rate and store the value as `pov_dem_rate` . Interpret the value in context.

Type your answer here, replacing this text.

```
In [ ]: # correlation between poverty and percent minor
pov_dem_rate = ...

# print
pov_dem_rate
```

```
In [ ]: grader.check("q1")
```

While direct inspection is useful, it can be cumbersome to check correlations for a large number of variables this way. A heatmap -- a colored image of the matrix -- provides a (sometimes) convenient way to see what's going on without having to examine the numerical values directly. The cell below shows one way of constructing this plot. Notice the diverging color scale; this should always be used.

```
In [ ]: # melt corr_mx
corr_mx_long = corr_mx.reset_index().rename(
    columns = {'index': 'row'}
).melt(
    id_vars = 'row',
```

```

    var_name = 'col',
    value_name = 'Correlation'
)

# construct plot
alt.Chart(corr_mx_long).mark_rect().encode(
  x = alt.X('col', title = '', sort = {'field
  y = alt.Y('row', title = '', sort = {'field
  color = alt.Color('Correlation',
                    scale = alt.Scale(scheme
                                     domain
                                     type =
                    legend = alt.Legend(tickCo
).properties(width = 300, height = 300)

```

Question 2

Which variable is self employment rate most *positively* correlated with? Refer to the heatmap.

Type your answer here, replacing this text.

Computing principal components

Each principal component is of the form:

$$PC_i = \sum_j v_j x_{ij} \quad (\text{PC score for observation } i)$$

The loading v_j for each component indicate which variables are most influential (heavily weighted) on that principal axis, and thus offer an indirect

picture of which variables are driving variation and covariation in the original data.

Loadings and scores

In `statsmodels`, the module `multivariate.pca` contains an easy-to-use implementation.

```
In [ ]: # compute principal components
pca = PCA(data = x_mx, standardize = True)
```

Most quantities you might want to use in PCA can be retrieved as attributes of `pca`. In particular:

- `.loadings` contains the loadings
- `.scores` contains the scores
- `.eigenvals` contains the variances along each principal axis (see lecture notes)

Examine the loadings below. Each column gives the loadings for one principal component; components are ordered from largest to smallest variance.

```
In [ ]: # inspect loadings
pca.loadings
```

Similarly, inspect the scores below and check your understanding; each row is an observation and the columns give the scores on each principal axis.


```
In [ ]: # inspect scores
pca.scores
```

Importantly, `statsmodels` rescales the scores so that they have unit inner product; in other words, so that the variances are all $\frac{1}{n-1}$.

```
In [ ]: # variance of scores
pca.scores.var()
```

```
In [ ]: # for comparison
1/(x_mx.shape[0] - 1)
```

To change this behavior, set `normalize = False` when computing the principal components.

Question 3

Check your understanding. Which variable contributes most to the sixth principal component? Store the variable name exactly as it appears among the original column names as `pc6_most_influential_variable`, and store the corresponding loading as `pc6_most_influential_variable_loading`. Print the variable name.

```
In [ ]: # find most influential variable
pc6_most_influential_variable = ...

# find loading
pc6_most_influential_variable_loading = ...
```

```
# print  
...
```

```
In [ ]: grader.check("q3")
```

Variance ratios

The *variance ratios* indicate the proportions of total variance in the data captured by each principal axis. You may recall from lecture that the variance ratios are computed from the eigenvalues of the correlation (or covariance, if data are not standardized) matrix.

When using `statsmodels`, these need to be computed manually.

```
In [ ]: # compute variance ratios  
var_ratios = pca.eigenvals/pca.eigenvals.sum()  
  
# print  
var_ratios
```

Note again that the principal components have been computed in order of *decreasing* variance.

Question 4

Check your understanding. What proportion of variance is captured *jointly* by the first three components taken together? Provide a calculation to justify your answer.

Type your answer here, replacing this text.

In []: ...

Selecting a subset of PCs

PCA generally consists of choosing a small subset of components. The basic strategy for selecting this subset is to determine how many are needed to capture some analyst-chosen minimum portion of total variance in the original data.

Most often this assessment is made graphically by inspecting the variance ratios and their cumulative sum, *i.e.*, the amount of total variation captured jointly by subsets of successive components. We'll store these quantities in a data frame.

```
In [ ]: # store proportion of variance explained as a d
pca_var_explained = pd.DataFrame({
    'Component': np.arange(1, 23),
    'Proportion of variance explained': var_rat

# add cumulative sum
pca_var_explained['Cumulative variance explaine

# print
pca_var_explained.head()
```

Now we'll make a dual-axis plot showing, on one side, the proportion of variance explained (y) as a function of component (x), and on the other side, the cumulative variance explained (y) also as a

function of component (x). Make sure that you've completed Q1(a) before running the next cell.

```
In [ ]: # encode component axis only as base layer
base = alt.Chart(pca_var_explained).encode(
    x = 'Component')

# make a base layer for the proportion of variance explained
prop_var_base = base.encode(
    y = alt.Y('Proportion of variance explained',
              axis = alt.Axis(titleColor = '#57
    )

# make a base layer for the cumulative variance explained
cum_var_base = base.encode(
    y = alt.Y('Cumulative variance explained',
    )

# add points and lines to each base layer
prop_var = prop_var_base.mark_line(stroke = '#57
cum_var = cum_var_base.mark_line() + cum_var_base

# layer the layers
var_explained_plot = alt.layer(prop_var, cum_var)

# display
var_explained_plot
```

The purpose of making this plot is to quickly determine the fewest number of principal components that capture a considerable portion of variation and covariation. 'Considerable' here is a bit subjective.

Question 5

How many principal components explain more than 6% of total variation individually? Store this number as `num_pc`, and store the proportion of variation that they capture jointly as `var_explained`.

```
In [ ]: # number of selected components
num_pc = ...

# variance explained
var_explained = ...

#print
print('number selected: ', num_pc)
print('proportion of variance captured: ', var_
```

```
In [ ]: grader.check("q5")
```

Interpreting loadings

Now that you've chosen the number of components to work with, the next step is to examine loadings to understand just *which* variables the components combine with significant weight.

We'll store the scores for the components you selected as a dataframe.

```
In [ ]: # subset loadings
loading_df = pca.loadings.iloc[:, 0:num_pc]

# rename columns
loading_df = loading_df.rename(columns = dict(z
```

```
# print
loading_df.head()
```

Again, the loadings are the *weights* with which the variables are combined to form the principal components. For example, the **PC1** column tells us that this component is equal to:

$$(-0.020055 \times \text{women}) + (0.289614 \times \text{white}) + (0.$$

Since the components together capture over half the total variation, the heavily weighted variables in the selected components are the ones that drive variation in the original data.

By visualizing the loadings, we can see which variables are most influential for each component, and thereby also which variables seem to drive total variation in the data.



```
In [ ]: # melt from wide to long
loading_plot_df = loading_df.reset_index().melt(
    id_vars = 'index',
    var_name = 'Principal Component',
    value_name = 'Loading'
).rename(columns = {'index': 'Variable'})

# add a column of zeros to encode for x = 0 line
loading_plot_df['zero'] = np.repeat(0, len(load

# create base layer
base = alt.Chart(loading_plot_df)

# create lines + points for loadings
```

```
loadings = base.mark_line(point = True).encode(  
  y = alt.X('Variable', title = ''),  
  x = 'Loading',  
  color = 'Principal Component'  
)  
  
# create line at zero  
rule = base.mark_rule().encode(x = alt.X('zero'  
  
# layer  
loading_plot = (loadings + rule).properties(wid  
  
# show  
loading_plot.facet(column = alt.Column('Princip
```

Look first at PC1: the variables with the largest loadings (points farthest in either direction from the zero line) are Child Poverty (positive), Employed (negative), Income per capita (negative), Poverty (positive), and Unemployment (positive). We know from exploring the correlation matrix that employment rate, unemployment rate, and income per capita are all related, and similarly child poverty rate and poverty rate are related. Therefore, the positively-loaded variables are all measuring more or less the same thing, and likewise for the negatively-loaded variables.

Essentially, then, PC1 is predominantly (but not entirely) a representation of income and poverty. In particular, counties have a higher value for PC1 if they have lower-than-average income per capita and higher-than-average poverty rates, and a smaller value for PC1 if they have higher-than-

average income per capita and lower-than-average poverty rates.

A system for loading interpretation

Often interpreting principal components can be difficult, and sometimes there's no clear interpretation available! That said, it helps to have a system instead of staring at the plot and scratching our heads. Here is a semi-systematic approach to interpreting loadings:

1. Divert your attention away from the zero line.
2. Find the largest positive loading, and list all variables with similar loadings.
3. Find the largest negative loading, and list all variables with similar loadings.
4. The principal component represents the difference between the average of the first set and the average of the second set.
5. Try to come up with a description of less than 4 words.

This system is based on the following ideas:

- a high loading value (negative or positive) indicates that a variable strongly influences the principal component;
- a negative loading value indicates that
 - increases in the value of a variable *decrease* the value of the principal

component

- and decreases in the value of a variable *increase* the value of the principal component;
- a positive loading value indicates that
 - increases in the value of a variable *increase* the value of the principal component
 - and decreases in the value of a variable *decrease* the value of the principal component;
- similar loadings between two or more variables indicate that the principal component reflects their *average*;
- divergent loadings between two sets of variables indicates that the principal component reflects their *difference*.

Question 6

Work with your neighbor to interpret PC2. Come up with a name for the component and explain which variables are most influential.

Type your answer here, replacing this text.

Standardization

Data are typically standardized because otherwise the variables on the largest scales tend to dominate the principal components, and most of

the time PC1 will capture the majority of the variation. However, that is artificial. In the census data, income per capita has the largest magnitudes, and thus, the highest variance.

```
In [ ]: # three largest variances
x_mx.var().sort_values(ascending = False).head()
```

When PCs are computed without normalization, the total variation is mostly just the variance of income per capita because it is orders of magnitude larger than the variance of any other variable. But that's just because of the *scale* of the variable -- incomes per capita are large numbers -- not a reflection that it varies more or less than the other variables.

Run the cell below to see what happens to the variance ratios if the data are not normalized.

```
In [ ]: # recompute pcs without normalization
pca_unscaled = PCA(data = x_mx, standardize = F

# show variance ratios for first three pcs
pca_unscaled.eigenvals[0:3]/pca_unscaled.eigenv
```

Further, let's look at the loadings when data are not standardized:

```
In [ ]: # subset loadings
unscaled_loading_df = pca_unscaled.loadings.ilo

# rename columns
unscaled_loading_df = unscaled_loading_df.renam
columns = dict(zip(unscaled_loading_df.colu
```

```

)

# melt from wide to long
unscaled_loading_plot_df = unscaled_loading_df.
    id_vars = 'index',
    var_name = 'Principal Component',
    value_name = 'Loading'
).rename(
    columns = {'index': 'Variable'}
)

# add a column of zeros to encode for x = 0 line
unscaled_loading_plot_df['zero'] = np.repeat(0,

# create base layer
base = alt.Chart(unscaled_loading_plot_df)

# create lines + points for loadings
loadings = base.mark_line(point = True).encode(
    y = alt.X('Variable', title = ''),
    x = 'Loading',
    color = 'Principal Component'
)

# create line at zero
rule = base.mark_rule().encode(x = alt.X('zero'

# layer
loading_plot = (loadings + rule).properties(wid

# show
loading_plot.facet(column = alt.Column('Princip

```

Notice that the variables with nonzero loadings in unscaled PCA are simply the three variables with the largest variances.

```

In [ ]: # three largest variances
x_mx.var().sort_values(ascending = False).head(

```

Exploratory analysis based on PCA

Now that we have the principal components, we can use them for exploratory data visualizations. To this end, let's retrieve the scores from the components you selected:

```
In [ ]: # subset scores
score_df = pca.scores.iloc[:, 0:num_pc]

# rename columns
score_df = score_df.rename(
    columns = dict(zip(score_df.columns, ['PC'
    ]

# add state and county
score_df[['State', 'County']] = census[['State'

# print
score_df.head()
```

The PC's can be used to construct scatterplots of the data and search for patterns. We'll illustrate that by identifying some outliers. The cell below plots PC2 (employment type) against PC4 (carpooling?):

```
In [ ]: # base chart
base = alt.Chart(score_df)

# data scatter
scatter = base.mark_point(opacity = 0.2).encode
    x = alt.X('PC2:Q', title = 'Self-employment
```

```

    y = alt.Y('PC4:Q', title = 'Carpooling PC')
)

# show
scatter

```

Notice that there are a handful of outlying points in the upper right region away from the dense scatter. What are those?

In order to inspect the outlying counties, we first need to figure out how to identify them. The outlying values have a large *sum* of PC2 and PC4. We can distinguish them by finding a cutoff value for the sum; a simple quantile will do.

```

In [ ]: # find cutoff value
pc2_pc4_sum = (score_df.PC2 + score_df.PC4)
cutoff = pc2_pc4_sum.quantile(0.99999)

# store outlying rows using cutoff
outliers = score_df[(-score_df.PC2 + score_df.PC4) > cutoff]

# plot outliers in red
pts = alt.Chart(outliers).mark_circle(
    color = 'red',
    opacity = 0.3
).encode(
    x = 'PC2',
    y = 'PC4'
)

# layer
scatter + pts

```

Notice that almost all the outlying counties are remote regions of Alaska:

```
In [ ]: outliers
```

What sets them apart? The cell below retrieves the normalized data and county name for the outlying rows, and then plots the Standardized values of each variable for all 9 counties as vertical ticks, along with a point indicating the mean for the outlying counties. This plot can be used to determine which variables are over- or under-average for the outlying counties relative to the nation by simply locating means that are far from zero in either direction.

```
In [ ]: x_ctr = (x_mx - x_mx.mean())/x_mx.std()

# retrieve normalized data for outlying rows
outlier_data = x_ctr.loc[outliers.index.values]
              census.loc[outliers.index.values, ['County'
          )

# melt to long format for plotting
outlier_plot_df = outlier_data.melt(
    id_vars = 'County',
    var_name = 'Variable',
    value_name = 'Standardized value'
)

# plot ticks for values (x) for each variable (
ticks = alt.Chart(outlier_plot_df).mark_tick().
    x = 'Standardized value',
    y = 'Variable'
)

# shade out region within 3SD of mean
grey = alt.Chart(
    pd.DataFrame(
```

```

        {'Variable': x_ctr.columns,
         'upr': np.repeat(3, 22),
         'lwr': np.repeat(-3, 22)}
    )
).mark_area(opacity = 0.2, color = 'gray').encode(
    y = 'Variable',
    x = alt.X('upr', title = 'Standardized value'),
    x2 = 'lwr'
)

# compute means of each variable across counties
means = alt.Chart(outlier_plot_df).transform_aggregate(
    group_mean = 'mean(Standardized value)',
    groupby = ['Variable']
).transform_calculate(
    large = 'abs(datum.group_mean) > 3'
).mark_circle(size = 80).encode(
    x = 'group_mean:Q',
    y = 'Variable',
    color = alt.Color('large:N', legend = None)
)

# Layer
ticks + grey + means

```

Question 7

The two variables that clearly set the outlying counties apart from the nation are the percentage of the population using alternative transportation (extremely above average) and the percentage that drive to work (extremely below average). What about those counties explains this?

(Hint: take a peek at the [Wikipedia page on transportation in Alaska](#).)

Type your answer here, replacing this text.

Submission

1. Save the notebook.
2. Restart the kernel and run all cells. (**CAUTION:** if your notebook is not saved, you will lose your work.)
3. Carefully look through your notebook and verify that all computations execute correctly and all graphics are displayed clearly. You should see **no errors**; if there are any errors, make sure to correct them before you submit the notebook.
4. Download the notebook as an `.ipynb` file. This is your backup copy.
5. Export the notebook as PDF and upload to Gradescope.

To double-check your work, the cell below will rerun all of the autograder tests.

```
In [ ]: grader.check_all()
```