

Data Science Concepts and Analysis

Week 5: Principal components

(EXPLORATORY DATA ANALYSIS II)

- Covariation between variables
- Eigendecomposition
- Principal components analysis

This week: principal components

Objective: introduce an exploratory technique for exploring covariation among several variables.

- **Covariation between variables**

- What is covariation and how is it measured quantitatively?
- Covariance and correlation
- Correlation matrix
- Heatmap visualization

- **Eigendecomposition**

- The eigenvalue problem
- Geometric interpretation
- Computations

- **Principal components analysis**

- What is PCA exactly?
- PCA in the low-dimensional setting
- Variation capture and loss
- Interpreting principal components
- Example application

Covariation between variables

- What is covariation and how is it measured quantitatively?
- Covariance and correlation matrices
- Visualizations

City sustainability data

Throughout this week's lecture we'll use the following dataset on the sustainability of U.S. cities:

In [2]:

```
city_sust.head(4)
```

Out[2]:

	GEOID_MSA	Name	Econ_Domain	Social_Domain	Env_Domain	Sustain_Index
0	310M300US10100	Aberdeen, SD Micro Area	0.565264	0.591259	0.444472	1.600995
1	310M300US10140	Aberdeen, WA Micro Area	0.427671	0.520744	0.429274	1.377689
2	310M300US10180	Abilene, TX Metro Area	0.481092	0.496874	0.454192	1.432157
3	310M300US10220	Ada, OK Micro Area	0.466566	0.526206	0.425964	1.418737

For each **Metropolitan Statistical Area** (MSA), a sustainability index is calculated based on economic, social, and environmental indicators (also indices).

$$\text{sustainability index} = \text{economic} + \text{social} + \text{environmental}$$

The domain indices are computed from a large number of development indicator variables. To keep the application simple, we'll work with this derived data instead of the underlying variables.

If you're interested, you can dig deeper on the [Sustainable Development Report website](#), which provides detailed data reports related to the U.N.'s 2030 sustainable development goals.

What is covariation?

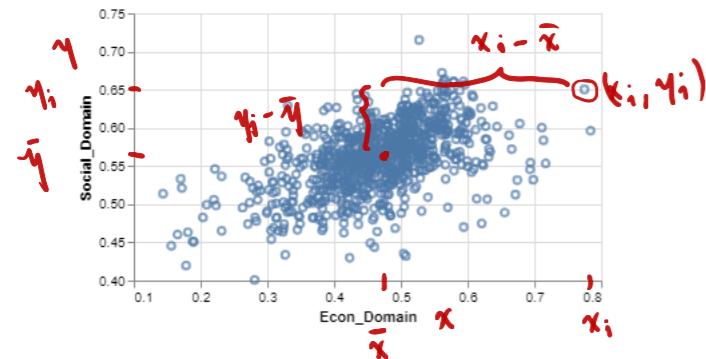
Last week we talked about exploring variation among individual variables: the tendency of values to change between observations.

Covariation refers to ***the tendency of two variables to change together between observations***. Covariation is about relationships.

In [4]:

```
econ_social
```

Out[4]:



The social and economic indices do seem to *vary together*: higher values of the economic index coincide with higher values of the social index.

That's all there is to it. ***You've already thought about covariation, even if you didn't know it!***

How is covariation measured?

Let $(x_1, y_1), \dots, (x_n, y_n)$ denote n values of two variables, X and Y .

If X and Y tend to vary together, then whenever X is far from its mean, so is Y .

- In other words, *their deviations from typical values coincide*.

This coincidence (or lack thereof) can be captured quantitatively by the (sample) **covariance**:

$$\text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

deviation of
 y from its
 mean
 ↓
 deviation of
 x from its
 mean
 ↓
 sum over
 observations

In [5]:

```
xy = city_sust.iloc[:, 2:4] # x = econ, y = social
xy_ctr = xy - xy.mean() # center by subtracting column means
xy_ctr.head(3)
```

Out[5]:

	$x_i - \bar{x}$	$y_i - \bar{y}$
Econ_Domain	Social_Domain	
0	0.098892	0.029696
1	-0.038701	-0.040819
2	0.014720	-0.064690

$$\begin{aligned}
 & \text{original} && \text{column means} \\
 & \begin{bmatrix} \text{econ} & \text{social} \\ x_1 & y_1 \\ x_2 & y_2 \\ x_3 & y_3 \\ \vdots & \vdots \end{bmatrix} & - & \begin{bmatrix} \text{econ} & \text{social} \\ \bar{x} & \bar{y} \\ \bar{x} & \bar{y} \\ \bar{x} & \bar{y} \\ \vdots & \vdots \end{bmatrix} \\
 & & & \text{sum of product of deviations}
 \end{aligned}$$

In [6]:

```
xy_cov = xy_ctr.prod(axis = 1).sum() / (len(xy) - 1) # cov(x, y)
xy_cov
```

Out[6]:

Econ x Social
 0.0988 x 0.0297
 -0.0387 x -0.0408
 ...

0.0019850231634313846

Correlation: standardized covariance

Covariance is a little tricky to interpret. Is 0.00199 large or small?

It is a little more useful to compute the (sample) **correlation**:

$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{S_x S_y}$$

This is simply **covariance standardized to [-1, 1]**.

Properties:

- $\text{corr}(x, x) = 1$, since any variable's deviations coincide *perfectly* with themselves
- the sign indicates whether x and y vary together or in opposition
- $\text{corr}(x, y) = 1, -1$ are the *strongest possible* correlations
- $\text{corr}(x, y) = 0$ is the weakest possible correlation
- moderate to large correlations indicate that x and y covary
- small correlations *do not indicate* that x and y don't covary

Correlation of social and economic indices

Standardizing the covariance makes it more interpretable:

In [7]:

```
xy_cov/xy.std().prod() # cov(x, y)/(sx*sy)
```

Out[7]:

0.5314906371658936

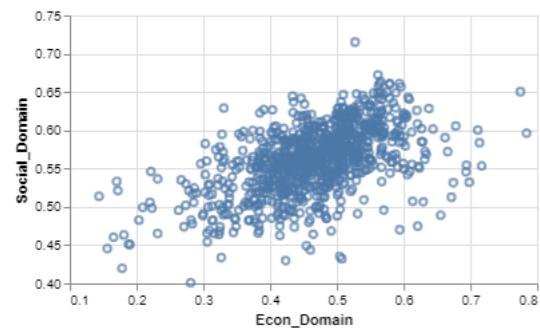
The correlation indicates that the social and economic indices vary together (positive) weakly (moderate magnitude on [0, 1] scale).

This is just a number that quantifies what you already knew from the graphic: there is a positive relationship.

In [8]:

```
econ_social
```

Out[8]:



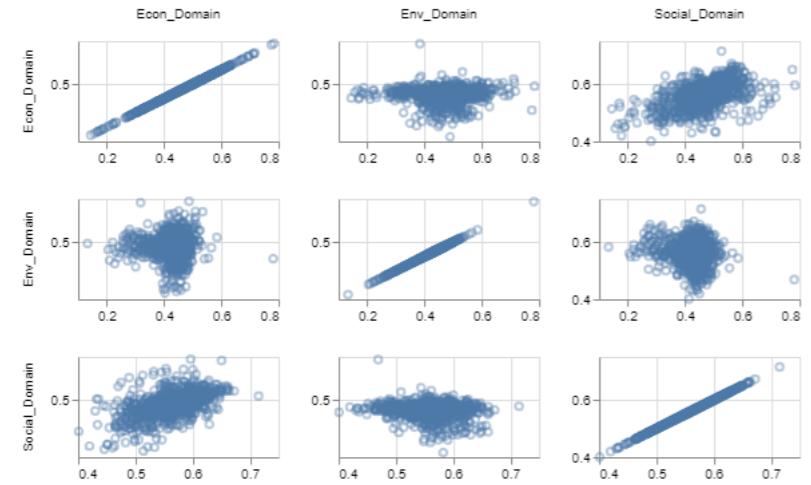
What is the use of a quantitative measure of covariation?

It helps to have a number so that we can talk precisely about the *strength* of a relationship, and compare relationships quantitatively.
Consider looking at all pairwise relationships for patterns of covariation:

In [10]:

```
scatter_panel
```

Out[10]:



Is the economic index *more* related to the environmental index or the social index?

Correlation matrix (example)

The pairwise correlations among all three variables can be represented in a simple square matrix:

In [11]:

```
x_mx.corr() # correlation matrix in pandas
```

Out[11]:

	Econ_Domain	Social_Domain	Env_Domain
Econ_Domain	1.000000	0.531491	0.011206
Social_Domain	0.531491	1.000000	-0.138674
Env_Domain	0.011206	-0.138674	1.000000

The strongest relationship is between social and economic indices; the weakest is between environmental and economic indices.

Notice the correspondence between this matrix and the scatterplot panel -- there is one number per plot that describes the strength of the relationship shown visually.

$$X = \{x_{ij}\} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ x_{21} & \dots & x_{2p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}$$

Correlation matrix (definition)

The (sample) correlation matrix can be expressed as a simple matrix product; let $X \in \mathbb{R}^{n \times p}$ denote the data values for n observations of p variables.

Consider the centered and scaled version of the data: $Z = \left\{ \frac{x_{ij} - \bar{x}_j}{s_{x_j}} \right\}$.

The (sample) **correlation matrix** is defined as:

$$Z = \begin{bmatrix} \frac{x_{11} - \bar{x}_1}{s_1} & \dots & \frac{x_{1p} - \bar{x}_p}{s_p} \\ \vdots & \ddots & \vdots \\ \frac{x_{n1} - \bar{x}_1}{s_1} & \dots & \frac{x_{np} - \bar{x}_p}{s_p} \end{bmatrix}$$

think of as data frame:

$$\begin{array}{c|cccc} & x_1 & x_2 & \dots & x_p \\ \hline 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 2 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ n & x_{n1} & x_{n2} & \dots & x_{np} \end{array}$$

X is an $n \times p$ real-valued matrix

$$\text{corr}(X) = \frac{1}{n-1} Z' Z$$

call this R

$$\begin{bmatrix} z_{11} & z_{12} & \dots & z_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ z_{n1} & z_{n2} & \dots & z_{np} \end{bmatrix} \begin{bmatrix} z_{11} & z_{12} & \dots & z_{1p} \\ z_{21} & z_{22} & \dots & z_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ z_{n1} & z_{n2} & \dots & z_{np} \end{bmatrix} = \begin{bmatrix} r_{11} & \dots & r_{1p} \\ \vdots & \ddots & \vdots \\ r_{p1} & \dots & r_{pp} \end{bmatrix}$$

In [12]:

```
# correlation matrix 'by hand'
n = len(x_mx) # sample size
z_mx = (x_mx - x_mx.mean()) / x_mx.std() # (xi - xbar)/sx
z_mx.transpose().dot(z_mx)/(n - 1) # Z'Z/(n - 1)
```

Out[12]:

	Econ_Domain	Social_Domain	Env_Domain
Econ_Domain	1.000000	0.531491	0.011206
Social_Domain	0.531491	1.000000	-0.138674
Env_Domain	0.011206	-0.138674	1.000000

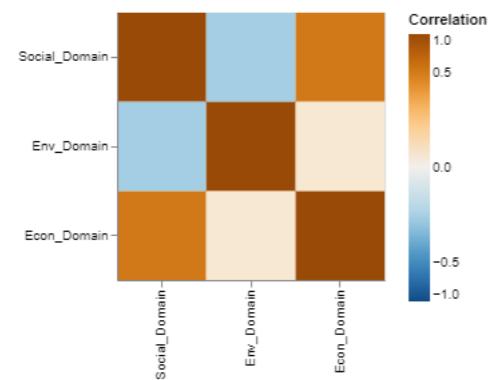
Heatmap visualization

Often it's useful to visualize the correlation matrix as a heatmap.

In [14]:

```
heatmap
```

Out[14]:



- Each cell corresponds to a pair of variables
- Cells are colored according to the magnitude of correlation between the pair
- Rows and columns are sorted in order of correlation strength

Notice again the correspondence with the scatterplot panel.

Eigendecomposition

- The eigenvalue problem
- Geometric interpretation
- Computations

Why are we talking about eigendecomposition?

We are going to factor the correlation matrix into components. This is, in essence, PCA.

This is useful for two reasons:

- identifying which variables drive variation and covariation
- reducing correlation structure among many variables to simpler components;
- visualizing multivariate data.

The eigendecomposition provides a means of factoring the matrix.

The eigenvalue problem

Let A be a square ($n \times n$) matrix.

The **eigenvalue problem** refers to finding nonzero λ and x that satisfy the equation:

$$Ax = \lambda x$$

For any such solutions:

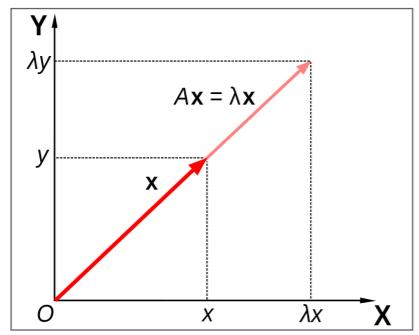
- λ is an **eigenvalue** of A ;
- x is an **eigenvector** of A .

Geometry

For a simple example, suppose $n = 2$ and x is an eigenvalue of A . Then:

$$Ax = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 \\ a_{21}x_1 + a_{22}x_2 \end{bmatrix} = \begin{bmatrix} \lambda x_1 \\ \lambda x_2 \end{bmatrix} = \lambda x$$

So the eigenvalue problem equation says that the *linear transformation of x by A* is simply a rescaling of x by a factor of λ .



The decomposition

If λ, x are an eigenvalue and eigenvector of A , then so are $c\lambda, cx$ for any constant c ; so usually attention is constrained to solutions such that $\|x\| = 1$; these are unique to within the sign of x .

The **eigendecomposition** of a square matrix consists in ***finding its unique nonzero eigenvalues and (unit-length) eigenvectors***.

This task is considered a 'decomposition' because the eigenvectors v and eigenvalues Λ satisfy

$$\begin{array}{c} \text{A} \\ \text{original matrix} \end{array} = \begin{array}{c} \text{V} \Lambda \text{V}' \\ \text{eigendecomposition} \end{array}$$

So the original matrix can be *reconstructed* from the eigenvalues and eigenvectors.

$$\begin{array}{c} \text{A} \\ \text{original} \end{array} = \underbrace{\text{V} \Lambda \text{V}'}_{\text{decomposition}} = \begin{array}{c} \text{V} \\ \left[v_1 \ v_2 \ \dots \ v_n \right] \end{array} \begin{array}{c} \Lambda \\ \left[\begin{matrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_n & \end{matrix} \right] \end{array} \begin{array}{c} \text{V}' \\ \left[v'_1 \ v'_2 \ \dots \ v'_n \right] \end{array}$$

eigenvectors
of unit length

Special case

We will be applying eigendecomposition to correlation matrices. These are of the form $A = Z'Z$.

Such matrices have some special properties, with the consequences:

- $V'V = I$, in other words, the eigenvectors are an orthonormal basis;
- $\lambda_i \geq 0$, in other words, all eigenvalues are nonnegative.

Computations

The eigendecomposition is computed numerically using iterative methods. Luckily, these are very easy to implement:

In [15]:

```
# compute decomposition
decomp = linalg.eig(corr_mx)
decomp
```

Out[15]:

```
(array([0.44788559+0.j, 1.54664121+0.j, 1.0054732 +0.j]),
 array([[ 0.68276586, -0.68571102,  0.2522522 ],
       [-0.70523279, -0.7087523 , -0.01780118],
       [-0.19099079,  0.16574249,  0.96749778]]))
```

In [16]:

```
# eigenvalues
decomp[0]
```

Out[16]:

```
array([0.44788559+0.j, 1.54664121+0.j, 1.0054732 +0.j])
```

In [17]:

```
# eigenvectors
decomp[1]
```

Out[17]:

```
array([[ 0.68276586, -0.68571102,  0.2522522 ],
       [-0.70523279, -0.7087523 , -0.01780118],
       [-0.19099079,  0.16574249,  0.96749778]])
```

Does the decomposition really work?

Let's check that in fact $A = V\Lambda V'$:

In [19]:

```
eigenvecs.dot(eigenvals).dot(eigenvecs.transpose()) # V Lambda V'
```

Out[19]:

	Econ_Domain	Social_Domain	Env_Domain
Econ_Domain	1.000000	0.531491	0.011206
Social_Domain	0.531491	1.000000	-0.138674
Env_Domain	0.011206	-0.138674	1.000000

In [20]:

```
corr_mx # correlation matrix (our 'A')
```

Out[20]:

	Econ_Domain	Social_Domain	Env_Domain
Econ_Domain	1.000000	0.531491	0.011206
Social_Domain	0.531491	1.000000	-0.138674
Env_Domain	0.011206	-0.138674	1.000000

Orthogonality

Let's check also that in fact $\mathbf{v}'\mathbf{v} = \mathbf{I}$:

In [21]:

```
eigenvecs.transpose().dot(eigenvecs) # v'v
```

Out[21]:

	v1	v2	v3
v1	1.000000e+00	-1.510283e-16	1.241330e-17
v2	-1.510283e-16	1.000000e+00	3.849861e-17
v3	1.241330e-17	3.849861e-17	1.000000e+00

The significance of this property is that $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$ form an *orthonormal basis*.

What's a basis?

A **basis** in linear algebra is a **set of vectors that span a linear space**. Think of a basis as a set of axes.

The usual basis for \mathbb{R}^3 are the unit vectors:

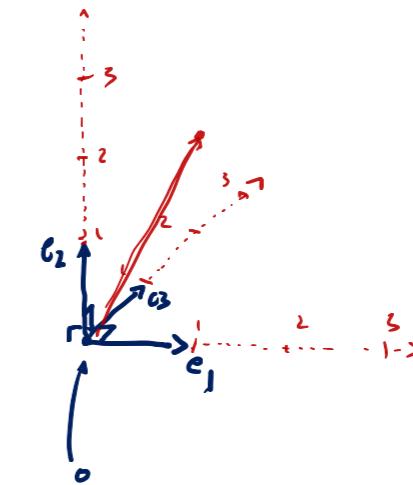
$$\mathbf{e}_1 = (1, 0, 0) \quad \mathbf{e}_2 = (0, 1, 0) \quad \mathbf{e}_3 = (0, 0, 1)$$

Coordinates in \mathbb{R}^3 are usually represented as multiples of these basis vectors:

$$(1, 2, 1) = 1\mathbf{e}_1 + 2\mathbf{e}_2 + 1\mathbf{e}_3$$

Terminology:

- A basis is *orthogonal* if all basis vectors are at right angles to one another.
- A basis is *orthonormal* if it is orthogonal and basis vectors are of unit length.



$$c_i' e_j = 0 \quad \forall j \neq i$$

$$\|\mathbf{e}_1\| = \sqrt{1^2 + 0^2 + 0^2} = \sqrt{1^2} = \sqrt{1} = 1$$

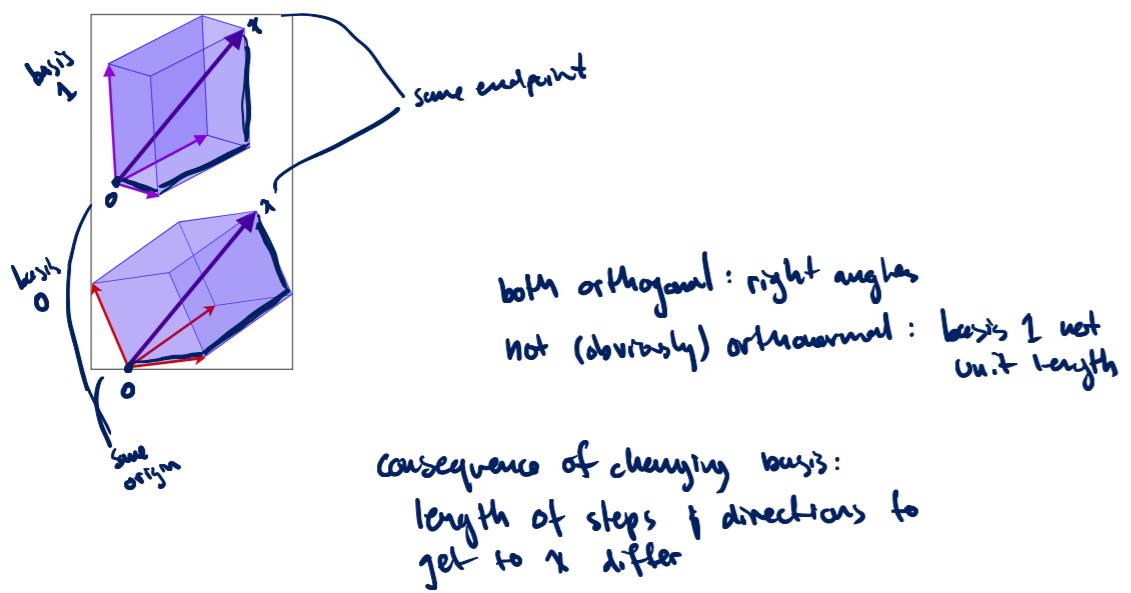
$$\|\mathbf{e}_2\| = 1$$

$$\|\mathbf{e}_3\| = 1$$

Nonstandard bases

To an extent the usual choice of e_1, e_2, e_3 is arbitrary.

It's possible to choose *any* system based on three directions to represent the same point relative to the origin.



Interpreting a change of basis

Different bases are simply different coordinate systems.

It's akin to the possibility of locating the distance to the corner store by different units and relative to different directions. For example:

- half a block to the right from the front door; or
- sixty steps diagonally in a straight line; or
- fifteen steps north, fifty steps east.

Eigendecomposition of the correlation matrix

So, decomposing the correlation matrix yields a basis.

It turns out that when data are represented on that basis, patterns of covariation vanish.

Let $Y = XV$; this is the representation of X on the basis given by V . Then

$$Y'Y = V'X'XV = V'V\Lambda V'V = \Lambda$$

which implies that $\text{cov}(y_j, y_k) = 0$ since Λ is diagonal.

So the change of basis to V in a sense 'captures' the covariation in the data.

Notice that the change of basis is a linear transformation XV .

This is the essence of PCA: using eigendecomposition to obtain a linear transformation that captures covariation.

Principal components analysis

- What is PCA exactly?
- PCA in the low-dimensional setting
- Variation capture and loss
- Interpreting principal components
- Example application

Principal components analysis

Principal components analysis (PCA) consists in ***finding a linear transformation of one's data that captures covariation.***

Applications are varied, and so are descriptions of just what the method is supposed to do, but the core technique is simple:

1. Compute the eigendecomposition of the correlation matrix.
2. Choose a subset of eigenvectors.
3. Then do other stuff (the 'analysis' part).

PCA terminology:

- Eigenvectors are called 'principal component directions'
- The values of eigenvectors are called 'loadings'
- The projected data consist of the 'principal components'

PCA in the low-dimensional setting

Let's consider what happens if we follow the steps of PCA outlined above with just two variables, say the social index and the economic index, which were the most correlated pair in the example data.

Denote the observations on these two variables by

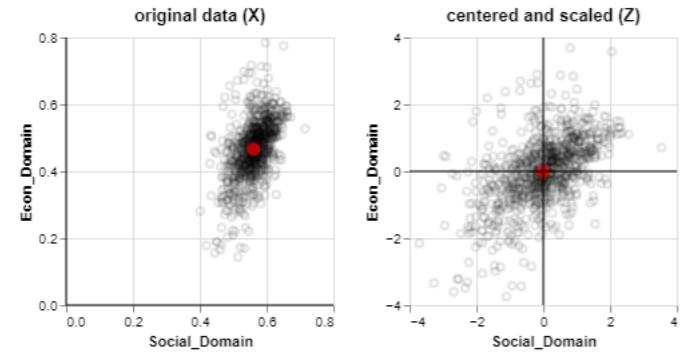
$$X = [x_1 \ x_2] \quad \text{where} \quad x_1 = \text{social index and } x_2 = \text{economic index}$$

To get the correlation matrix, first compute $Z = \left\{ \frac{x_i - \bar{x}}{s_x} \right\}$.

In [23]:

```
plot.configure_axis(domain = False)
```

Out[23]:



The red dots indicate the means.

PCA in the low-dimensional setting

Now we'll compute the eigendecomposition. This is pretty much all there is to it.

In [25]:

```
# compute correlation mx
r_mx = z_mx.transpose().dot(z_mx)/(len(z_mx) - 1)
r_mx
```

Out[25]:

	Econ_Domain	Social_Domain
Econ_Domain	1.000000	0.531491
Social_Domain	0.531491	1.000000

In [26]:

```
# decomposition
r_eig = linalg.eig(r_mx.values)

# show PC directions
directions = pd.DataFrame(r_eig[1], columns = ['PC1_Direction', 'PC2_Direction'], index = r_mx.columns)
directions
```

Out[26]:

	PC1_Direction	PC2_Direction
Econ_Domain	0.707107	-0.707107
Social_Domain	0.707107	0.707107

Remember that the principal component directions are just the eigenvectors.

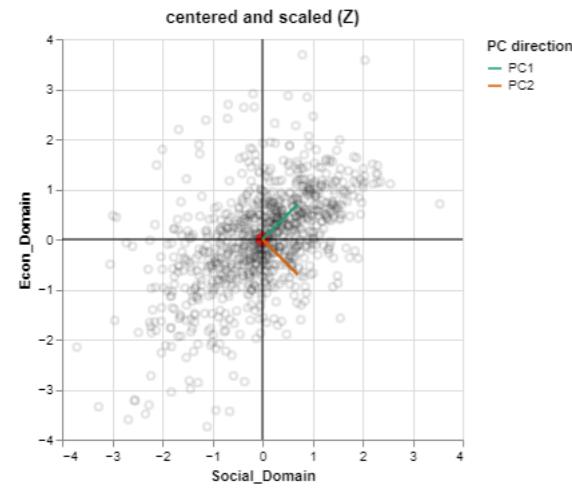
Geometry of PCA

What did we just do? Let's plot the principal component directions on the centered and scaled data.

In [28]:

```
centered_plot + eigenbasis
```

Out[28]:



Do you notice any relationship between the directions and the scatter?

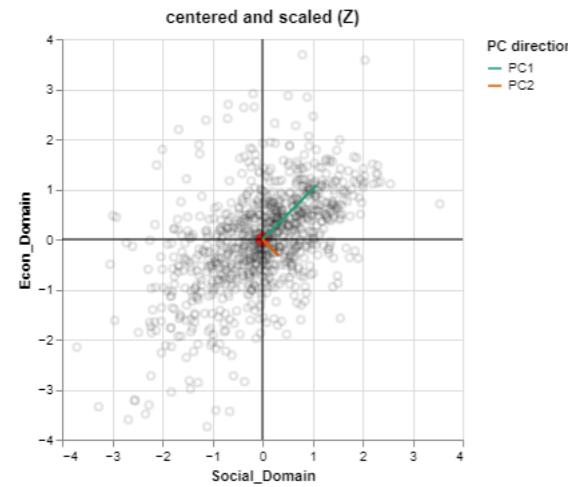
Geometry of PCA

How about now?

In [30]:

```
centered_plot + eigenbasis
```

Out[30]:



The difference here comes from scaling the directions by the corresponding eigenvalues (plotting $\lambda_1 v_1$ and $\lambda_2 v_2$).

The principal component directions are the axes along which the data vary most, and the eigenvalues give the magnitude of that variation.

Geometry of PCA: projection

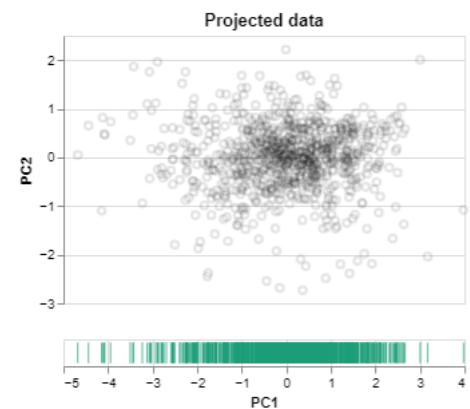
So if we wanted to look at just *one* quantity that retains variation and covariation in the original data, we could:

- project the data onto the new basis
- and look at the values on the axis of the first PC direction (the one with the most variation, *i.e.*, largest λ_i).

In [32]:

```
projection & pc1
```

Out[32]:



If we select just the first principal component to 'do other stuff' with, we lose the variation in PC2. Is this a problem?

Quantifying variation capture/loss

To figure out how much variation/covariation is captured and lost, we need to know how much is available in the first place.

The **total variation in x** is given in terms of the correlation matrix by:

$$\text{total variation} = \det(R) = \sum_{i=1}^p \lambda_i$$

Now let $y_j = Zv_j$ be the j th principal component. Its variance is:

$$\frac{1}{n-1}y_j'y_j = \frac{1}{n-1}v_j'Z'Zv_j = v_j'V'\Lambda V v_j = e_j'\Lambda e_j = \lambda_j$$

Quantifying variation capture/loss

So the total variation is the sum of the eigenvalues, and the variance of each PC is the corresponding eigenvalue.

We can therefore define the **proportion of total variation explained** by the j th principal component as:

$$\frac{\lambda_j}{\sum_{j=1}^p \lambda_j}$$

This is sometimes also called the *variance ratio* for the j th principal component.

So in our example:

In [33]:

```
eigenvalues = r_eig[0].real # store eigenvalues as real array
eigenvalues[0]/eigenvalues.sum() # variance ratio
```

Out[33]:

0.765745318582947

The first principal component captures 77% of total variation.

Interpreting principal components

So we have obtained a single derived variable that captures 3/4 of all variation and covariation in both variables. But what is this?

The values of the first principal component (green ticks) are given by:

$$PC1_i = \mathbf{x}'_i \mathbf{v}_1 = 0.7071 \times \text{economic}_i + 0.7071 \times \text{social}_i$$

So the principal component is a linear combination of the underlying variables.

Those coefficients are known as **loadings**: they determine the **relative weights** by which the variables contribute to the principal component.

In this case, the PC1 loadings are equal; so this principal component is simply the average of the social and economic indices.

Example application

Now that we've reviewed the basic technique, we can consider cases with a larger number of variables.

This is really where PCA is most useful.

- It can help answer the question: which variables are driving variation in the data?
- It can help reduce data dimensions by finding combinations of variables that preserve variation.
- It can provide a means of visualizing high-dimensional data.

So, let's look at a different example: world development indicators.

In [34]:

```
wdi = pd.read_csv('data/wdi-data.csv').iloc[:, 2: ].set_index('country')
wdi.head(3)
```

Out[34]:

country	inequality_adjusted_life	inequality_adjusted_education	maternal_mortality	parliament_pct_women	labor_participation_women	labor_participation_men	total_pop	urban_pct_pop	pop_under5	pop_15to64	...	total_unemployment	youth_unemployment	prison_5yr	trade	foreign_investment	net_migration	imm
Norway	0.931	0.908	2	40.8	60.4	67.2	5.4	82.6	0.055556	0.648148	...	3.3	9.3	80	72.1	0.4	5.3	
Ireland	0.926	0.892	5	24.3	56.0	68.4	4.9	63.4	0.061224	0.653061	...	4.9	13.1	82	239.2	-20.4	4.9	
Switzerland	0.947	0.883	5	38.6	62.9	73.8	8.6	73.8	0.046512	0.662791	...	4.6	7.4	77	119.4	-2.6	6.1	

3 rows × 31 columns

Computation via sklearn.decomposition.PCA

Luckily, `sklearn` contains an easy-to-use implementation that saves the trouble of going through all the calculations we did before 'by hand'.

In [35]:

```
from sklearn.decomposition import PCA
```

The implementation simply requires centering and scaling the data first.

In [36]:

```
# center and scale
wdi_ctr = (wdi - wdi.mean()) / wdi.std()

# compute principal components
pca = PCA(31)
pca.fit(wdi_ctr)
```

Out[36]:

```
PCA(n_components=31)
```

Variance ratios

Selecting the number of principal components to use is somewhat subjective, but always based on the variance ratios.

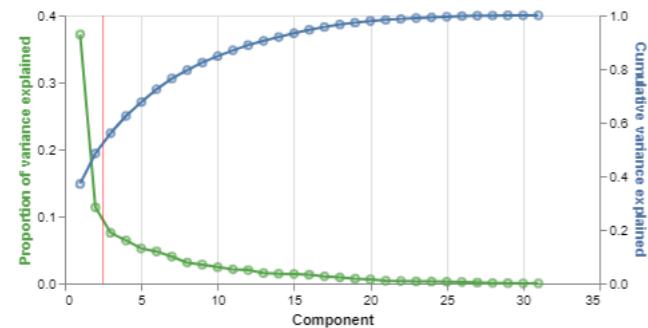
- how many PCs capture substantial variation individually before the variance ratio 'tails off'?
- how many PCs are needed to capture a certain proportion of the total variation?

This can be assessed by plotting the variance ratios and their cumulative sum:

In [38]:

```
variance_plot.properties(height = 200, width = 400)
```

Out[38]:



In this case, the variance ratios drop off after 2 components. These first two capture about half of total variation in the data.

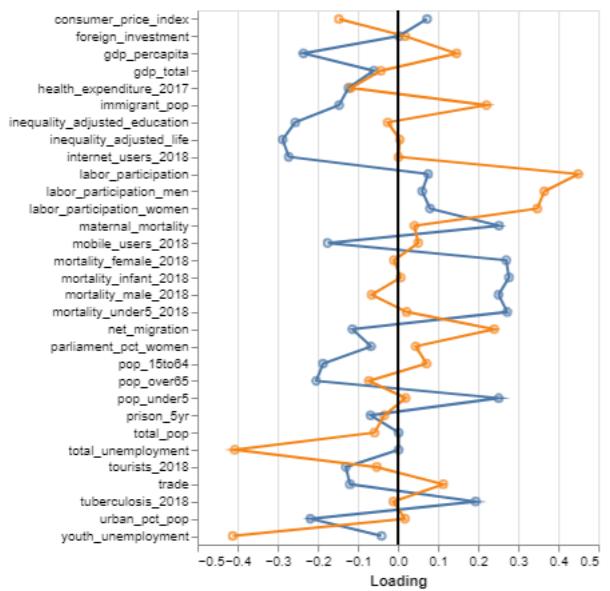
Loading plots

Examining the loadings can help address the question: *which variables drive variation in the data?*

In [40]:

```
loading_plot.properties(width = 300, height = 400)
```

Out[40]:



Visualization

Finally, we might want to use the first two principal components to visualize the data. Often it's helpful to merge the principal components with the original data and apply visualization techniques you already know to search for interesting patterns.

In [42]:

```
scatter + text
```

Out[42]:

