

Principal Component Analysis

EDP 619 Week 10

Dr. Abhik Roy

Welcome!



There are a lot of things going on behind the scenes when using PCAs and this is just a very brief introduction without any audio. I have tried to minimize the jargon and complexity, though some items may not be as clear as others. If you have questions, please feel free to reach out.

Additionally you may notice the following icons in the footnotes. These contain links to external sites that provide extra materials that may be of interest to you.



HTML
EXTRA
CONTENT



PDF
PAPER
MATERIAL



GUIDED
EXAMPLE



R
SCRIPT



ONLINE
VIDEO



Prerequisites

This slideshow assumes that you have a basic understanding of variance and correlations. For a refresher, please take a look at both reviews below



Prerequisites

This slideshow assumes that you have a basic understanding of variance and correlations. For a refresher, please take a look at both reviews below

Variance is essentially a measure of the spread between points in a data set. Specifically it tells us how far each data point in a set is from the mean and by proxy from every other data point in that set.

The infographic features a large blue hand-drawn style letter 'V' followed by the word 'VARIANCE' in a matching blue font. Below this, the mathematical formula for variance is shown: $\text{Var}[\hat{f}(x)] = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$. At the bottom, a green text box contains the following explanatory text: "Variance is the amount our predicted values would change if we had a different training dataset. It is the 'flexibility' of our model, balanced against bias." The author's name, "BY CHRIS ALBON", is at the very bottom.

VARIANCE

$$\text{Var}[\hat{f}(x)] = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$$

Variance is the amount our predicted values would change if we had a different training dataset. It is the "flexibility" of our model, balanced against bias.

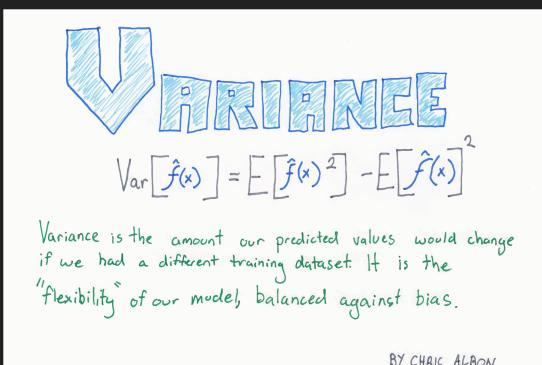
BY CHRIS ALBON



Prerequisites

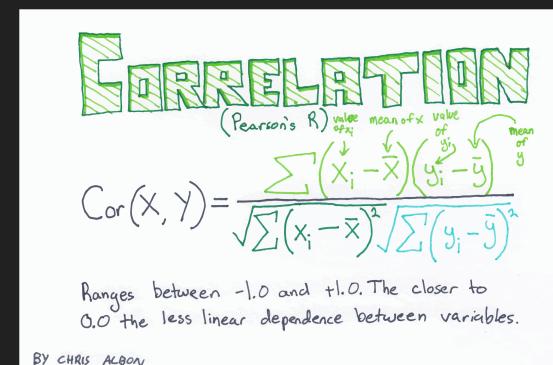
This slideshow assumes that you have a basic understanding of variance and correlations. For a refresher, please take a look at both reviews below

Variance is essentially a measure of the spread between points in a data set. Specifically it tells us how far each data point in a set is from the mean and by proxy from every other data point in that set.



A slide titled "VARIANCE" with a blue hand-drawn style title. Below it is the mathematical formula for variance: $\text{Var}[\hat{f}(x)] = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$. A green note at the bottom states: "Variance is the amount our predicted values would change if we had a different training dataset. It is the 'flexibility' of our model, balanced against bias." The slide is attributed to "BY CHRIS ALBON".

Correlation gives you an idea of the strength or weakness of the relationship between two variables. In a survey where each item is set to measure a single construct, these are essentially the applicable questions.



A slide titled "CORRELATION" with a green hand-drawn style title. Above the title, it says "(Pearson's R) value". Below the title is the formula for Pearson's Correlation coefficient: $\text{Cor}(x, y) = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}}$. A green note at the bottom states: "Ranges between -1.0 and +1.0. The closer to 0.0 the less linear dependence between variables." The slide is attributed to "BY CHRIS ALBON".



More Review

If you would like a deeper dive on either area, take a look at the videos below



More Review

If you would like a deeper dive on either area, take a look at the videos below

Variance

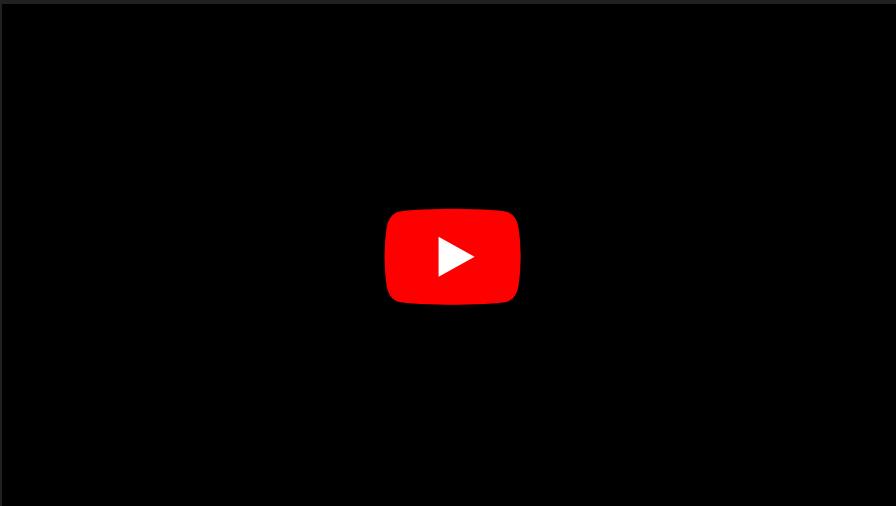




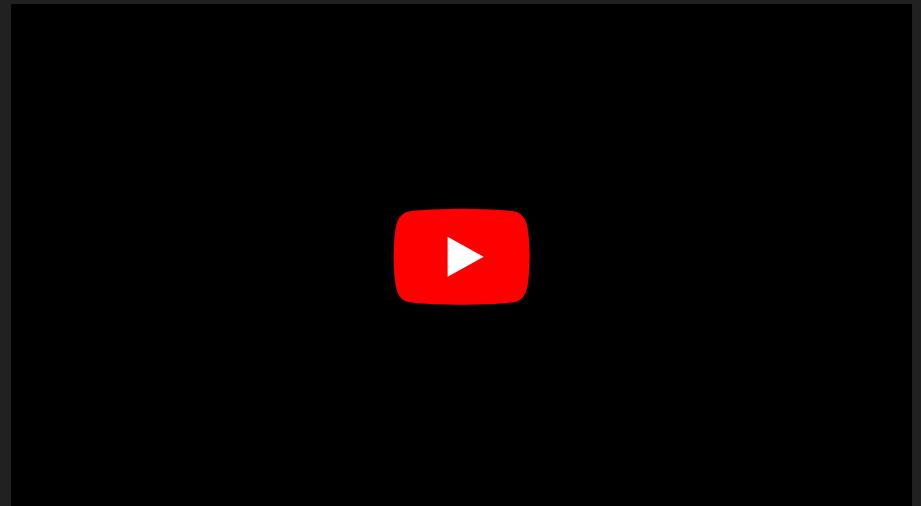
More Review

If you would like a deeper dive on either area, take a look at the videos below

Variance



Correlation



From Basic to Better



Survey Research
Methods

From Basic to Better



correlations are great but they don't...

From Basic to Better



correlations are great but they don't...

1. tell you how every question is related to every other question
2. differentiate between relevant data and noise

From Basic to Better



correlations are great but they don't...

1. tell you how every question is related to every other question
2. differentiate between relevant data and noise

...enter a method called ***Principle Component Analysis***

Principle Component Analysis (PCA)



Steps in a Nutshell





Steps in a Nutshell

The basic idea of a PCA can be broken into two steps



Steps in a Nutshell

The basic idea of a PCA can be broken into two steps

Locate the directions, or *components*, in a data set with high variance

PRINCIPAL COMPONENTS

Principal components are the linear combination of features that have the maximum variance out of all linear combinations.

Alternative interpretation: Principal components are low dimensional linear surfaces closest to the observations.

ChrisAlbon



Steps in a Nutshell

The basic idea of a PCA can be broken into two steps

Locate the directions, or *components*, in a data set with high variance

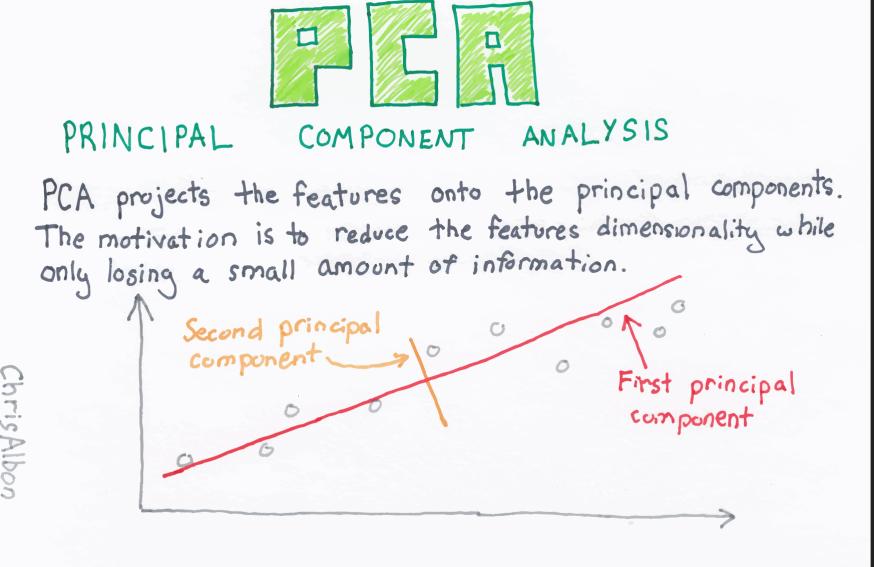
Find a limited number of components with high variance that in aggregate can explain most of the overall variance in the data

PRINCIPAL COMPONENTS

Principal components are the linear combination of features that have the maximum variance out of all linear combinations.

Alternative interpretation: Principal components are low dimensional linear surfaces closest to the observations.

ChrisAlbon



Reducing Complexity



Reducing Complexity



First an overview of some terms:



Reducing Complexity

First an overview of some terms:

Dimensionality - The number of input variables, or *features* in a dataset. In a spreadsheet, you can think of these as the column names.

id	read	write	math	science	socst
70	57	52	41	47	57
121	68	59	53	63	61
86	44	33	54	58	31
141	63	44	47	53	56
172	47	52	57	53	61
113	44	52	51	63	61
50	50	59	42	53	61
11	34	46	45	39	36
84	63	57	54	58	51
48	57	55	52	50	51
75	60	46	51	53	61
60	57	65	51	63	61
95	73	60	71	61	71

Reducing Complexity

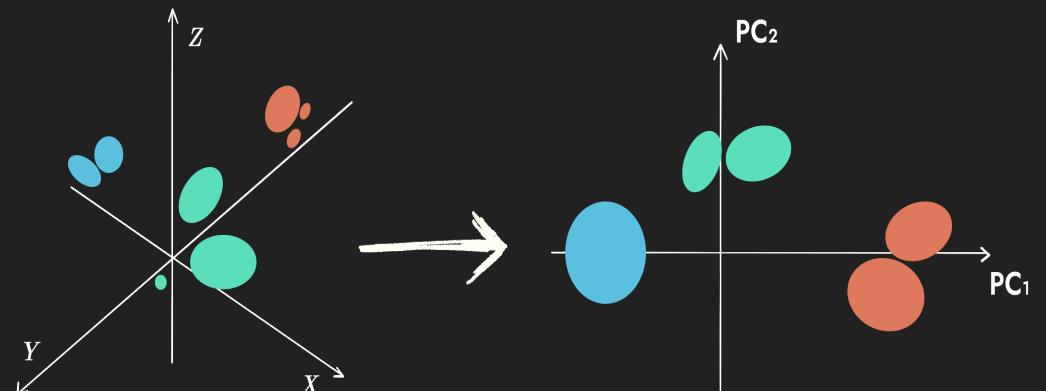


First an overview of some terms:

Dimensionality - The number of input variables, or *features* in a dataset. In a spreadsheet, you can think of these as the column names.

id	read	write	math	science	socst
70	57	52	41	47	57
121	68	59	53	63	61
86	44	33	54	58	31
141	63	44	47	53	56
172	47	52	57	53	61
113	44	52	51	63	61
50	50	59	42	53	61
11	34	46	45	39	36
84	63	57	54	58	51
48	57	55	52	50	51
75	60	46	51	53	61
60	57	65	51	63	61
95	73	60	71	61	71

Dimensionality Reduction - Statistical techniques used to reduce the number of input variables.





The Problem with Dimensions

Curse of Dimensionality - In brief terms, this refers to a few aspects



The Problem with Dimensions

Curse of Dimensionality - In brief terms, this refers to a few aspects

- *statistical.* the error rate increases as the number of features increases



The Problem with Dimensions

Curse of Dimensionality - In brief terms, this refers to a few aspects

- *statistical*. the error rate increases as the number of features increases
- *computational*. algorithms are harder to design and exponentially take more time to run in high dimensions



The Problem with Dimensions

Curse of Dimensionality - In brief terms, this refers to a few aspects

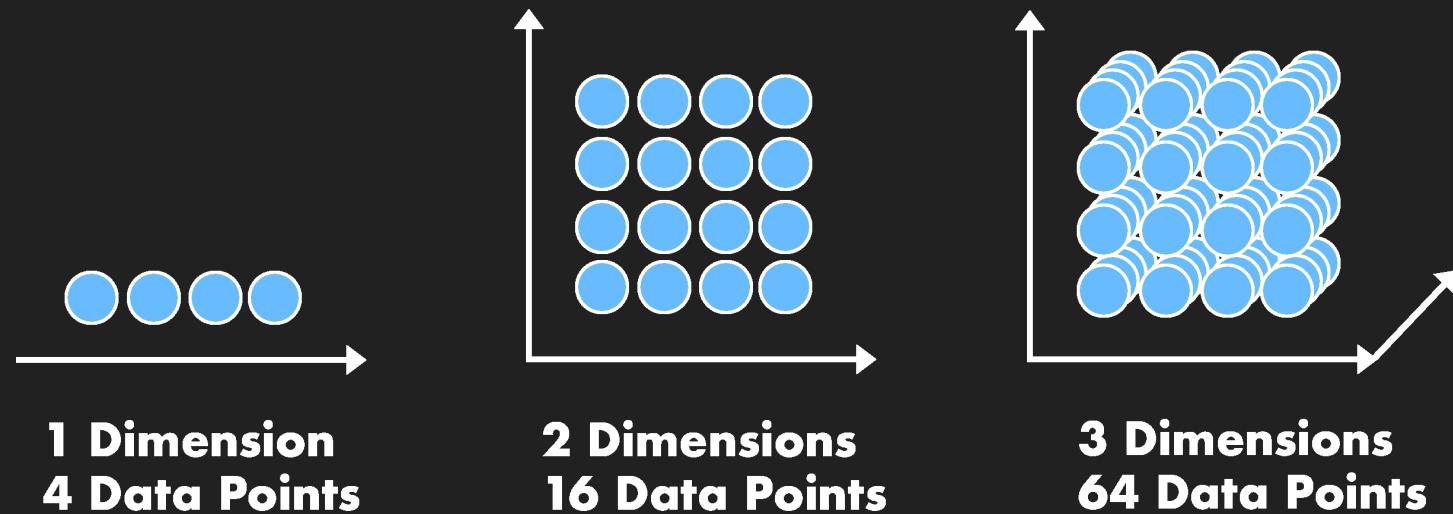
- *statistical*. the error rate increases as the number of features increases
- *computational*. algorithms are harder to design and exponentially take more time to run in high dimensions
- *practical*. higher number of dimensions theoretically allow more information to be stored, but in reality it rarely helps due to the higher possibility of noise and redundancy in real-world data



The Problem with Dimensions

Curse of Dimensionality - In brief terms, this refers to a few aspects

- *statistical*. the error rate increases as the number of features increases
- *computational*. algorithms are harder to design and exponentially take more time to run in high dimensions
- *practical*. higher number of dimensions theoretically allow more information to be stored, but in reality it rarely helps due to the higher possibility of noise and redundancy in real-world data





Fundamentals of What PCAs Can and Cannot Do

PCAs are one of the most traditional methods used for dimension reduction.



Fundamentals of What PCAs Can and Cannot Do

PCAs are one of the most traditional methods used for dimension reduction.

Primary benefit

It transforms the data into the most informative space, thereby allowing the use of lesser dimensions which retain needed information from the data while shedding much of the noise



Fundamentals of What PCAs Can and Cannot Do

PCAs are one of the most traditional methods used for dimension reduction.

Primary benefit

It transforms the data into the most informative space, thereby allowing the use of lesser dimensions which retain needed information from the data while shedding much of the noise

Primary drawback

It assumes linearity so any nonlinear relationship in a given data set is lost possibly causing loss in accuracy and the ability to estimate the likelihood of causality.



Fundamentals of What PCAs Can and Cannot Do

PCAs are one of the most traditional methods used for dimension reduction.

Primary benefit

It transforms the data into the most informative space, thereby allowing the use of lesser dimensions which retain needed information from the data while shedding much of the noise

Primary drawback

It assumes linearity so any nonlinear relationship in a given data set is lost possibly causing loss in accuracy and the ability to estimate the likelihood of causality.

Note as with most other procedures: *what you gain in efficiency, you lose in precision*. In a nutshell, there is no known perfect method that can both get rid of all of the noise and leave only relevant information. However with an ever growing machine learning library of approaches, we could get pretty close well within your lifetime!

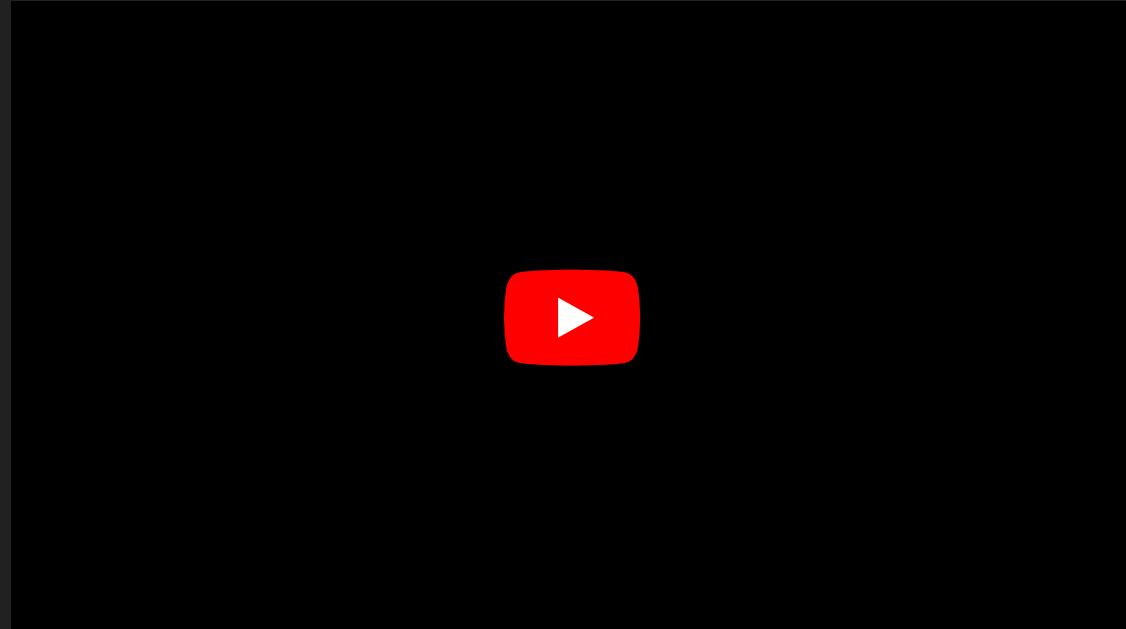
How Do PCAs Work?



How Do PCAs Work?



Before moving on please note that this is a nutshell explanation of the steps and avoids the mathematics¹. If you are interested in a more nuanced introduction coupled with the mathematics, watch this amazing lecture by Josh Starmer from StatQuest².



OK Now Really How Do PCAs Work?



OK Now Really How Do PCAs Work?



Let's look at a data set with 205 points randomly scattered in three-dimensions. Keep in mind that as you move along, the *PCA is carving out new dimensions which you will be able to see and interact with*.

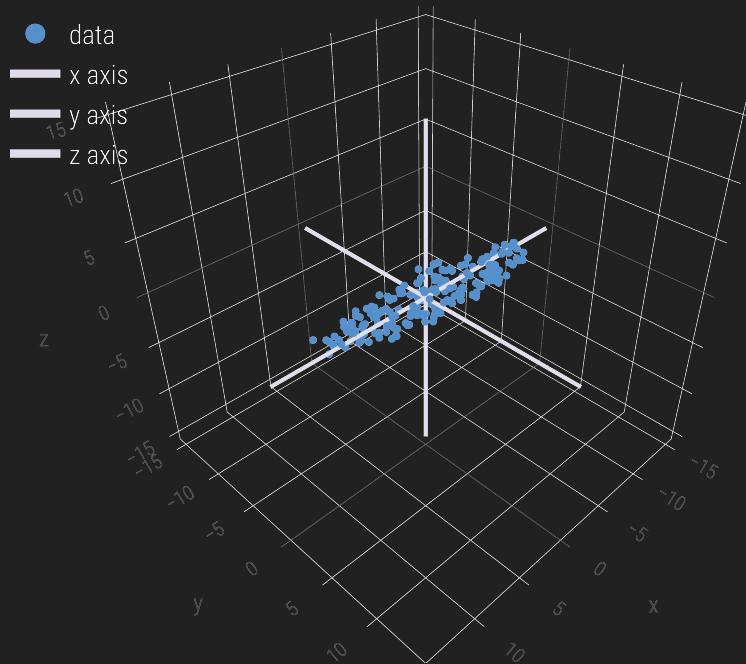
When applying a PCA, it locates the...

1. center point of data in multi-dimensional space



Look

Interact



1. center point of data in multi-dimensional space



Look

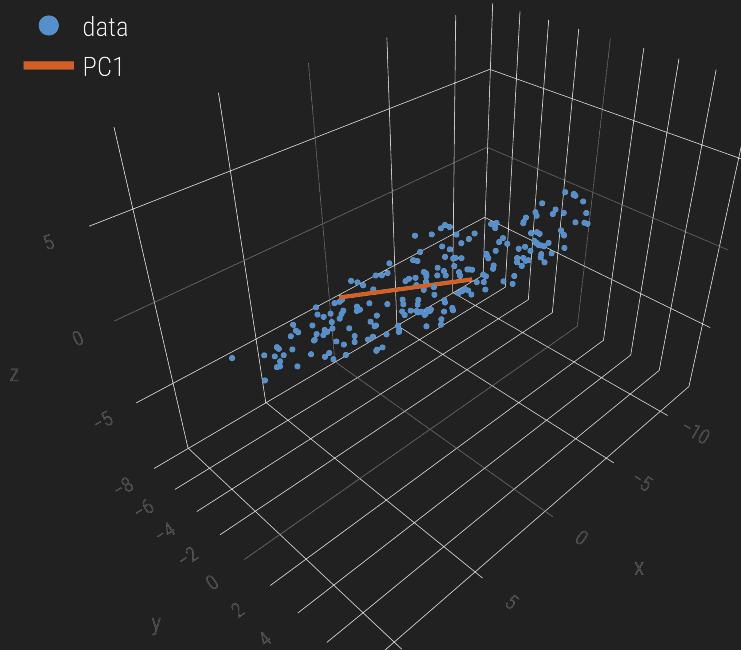
Interact

- data
- x axis
- y axis
- z axis

2. direction with the greatest variance. This is called the **1st component**



Look
Interact



2. direction with the greatest variance. This is called the **1st component**



Look

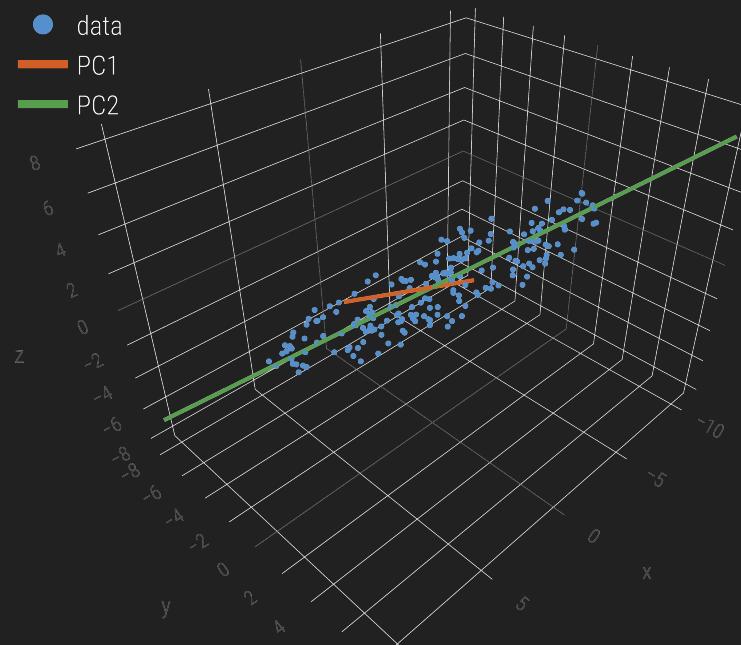
Interact

- data
- PC1

3. direction that is perpendicular, or *orthogonal* to the 1st component with the greatest variance. This is called the **2nd component**.



Look
Interact



3. direction that is perpendicular, or *orthogonal* to the 1st component with the greatest variance. This is called the **2nd component**.

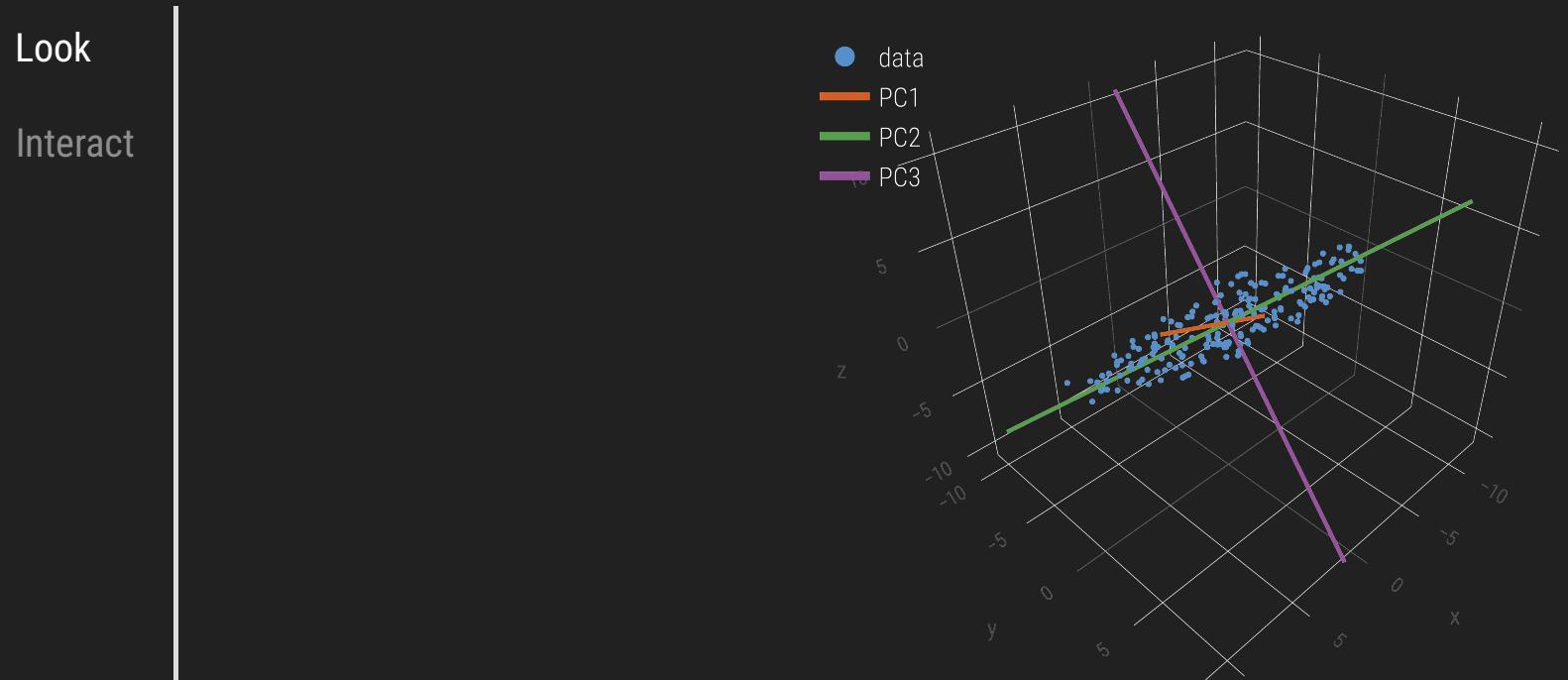


Look

Interact

- data
- PC1
- PC2

4. direction that is perpendicular, or *orthogonal* to the 1st and 2nd component with the greatest variance. This is called the **3rd component**.



4. direction that is perpendicular, or *orthogonal* to the 1st and 2nd component with the greatest variance. This is called the **3rd component**.



Look
Interact

- data
- PC1
- PC2
- PC3

| and it keeps going like this for as many dimensions as we have in a data set...



and it keeps going like this for as many dimensions as we have in a data set...

so you can probably imagine that big data sets with hundreds or thousands of columns and rows can take quite a bit of time...



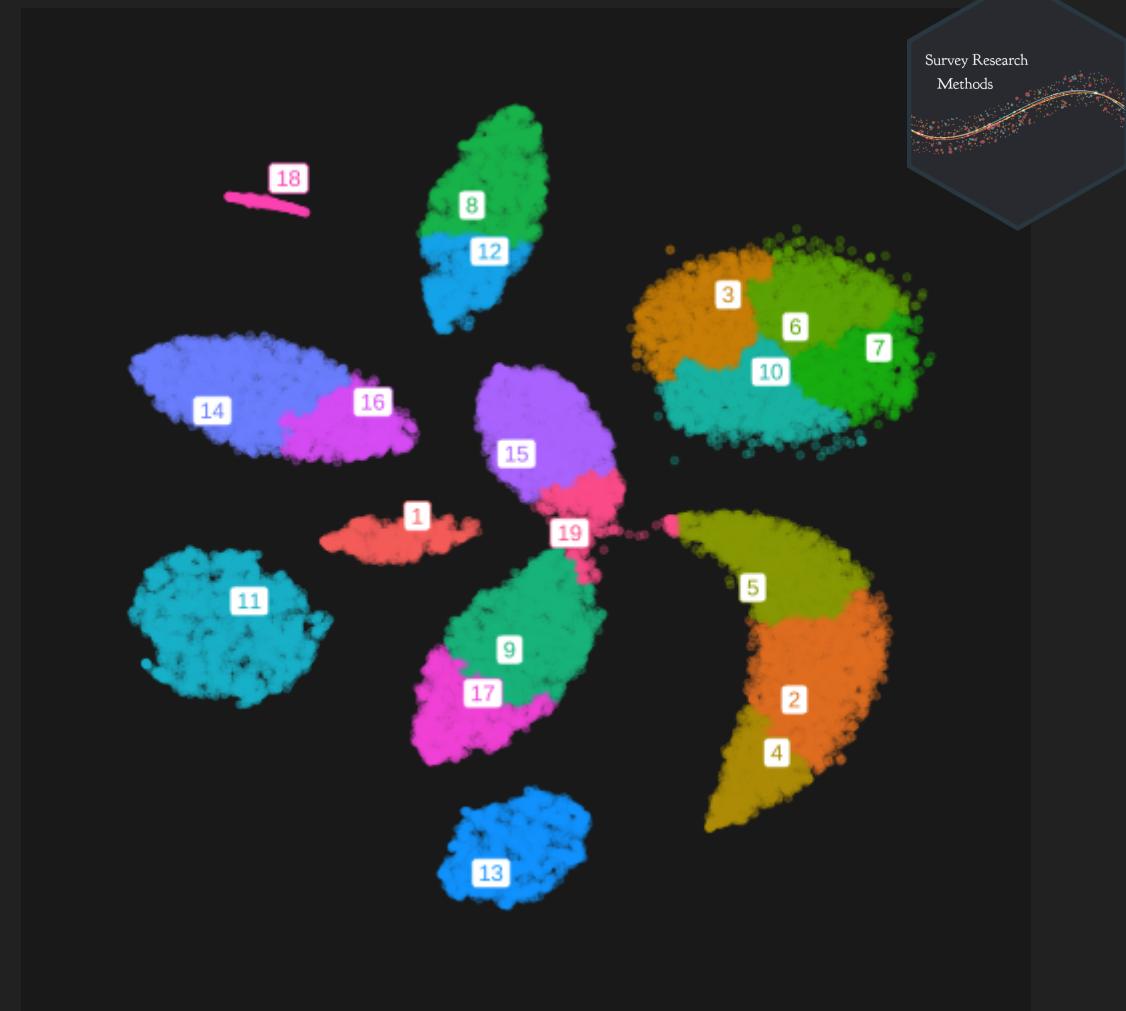
and it keeps going like this for as many dimensions as we have in a data set...

so you can probably imagine that big data sets with hundreds or thousands of columns and rows can take quite a bit of time...

but there are many other methods of reducing dimensions like

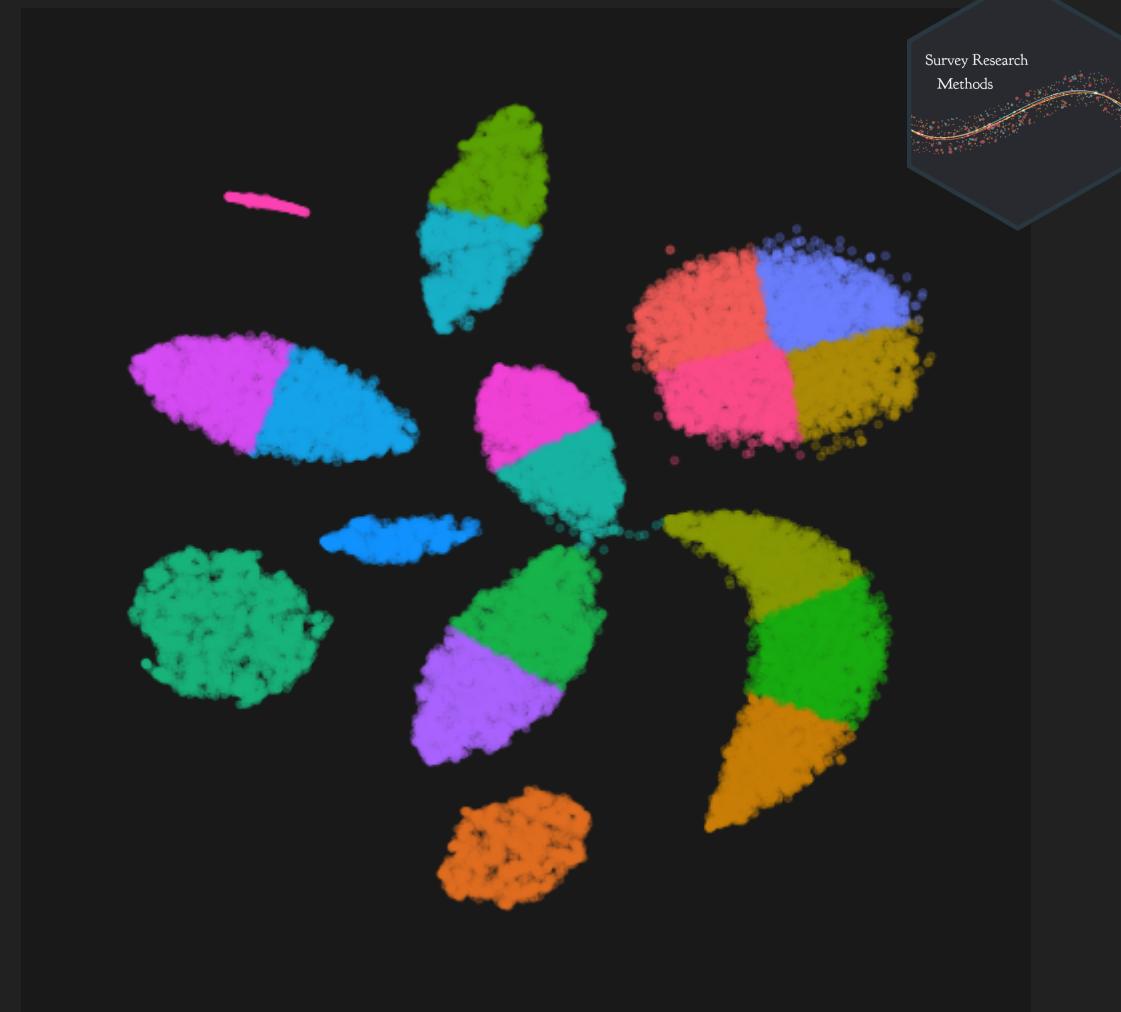


Hierarchical Clustering³



[3]

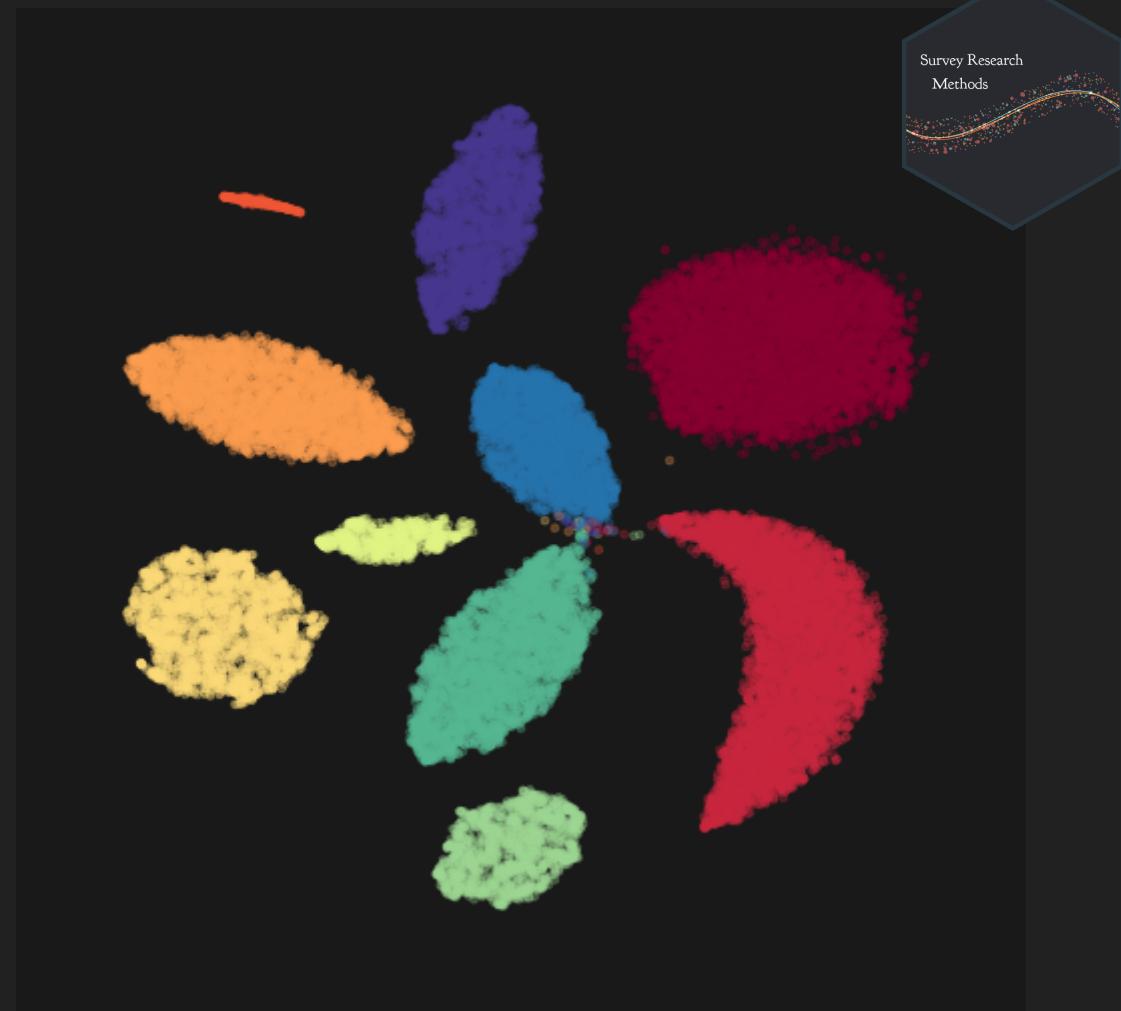
K-means Clustering⁴



[4]

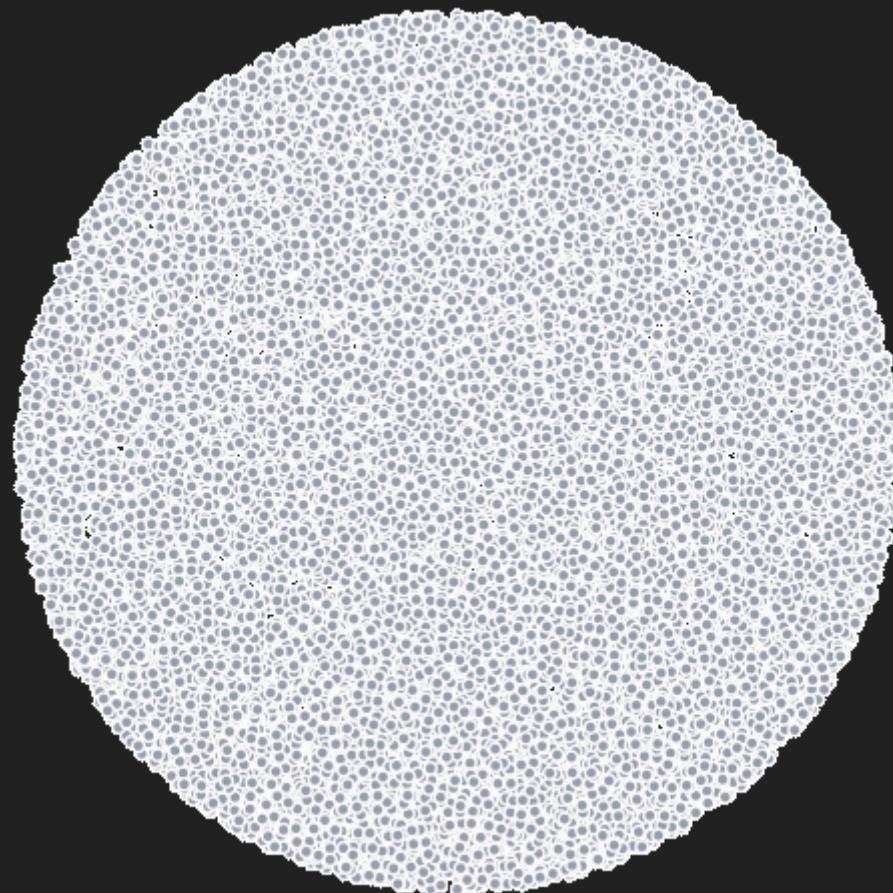
Survey Research
Methods

t-Distributed Stochastic Neighbor Embedding (t-SNE)⁵



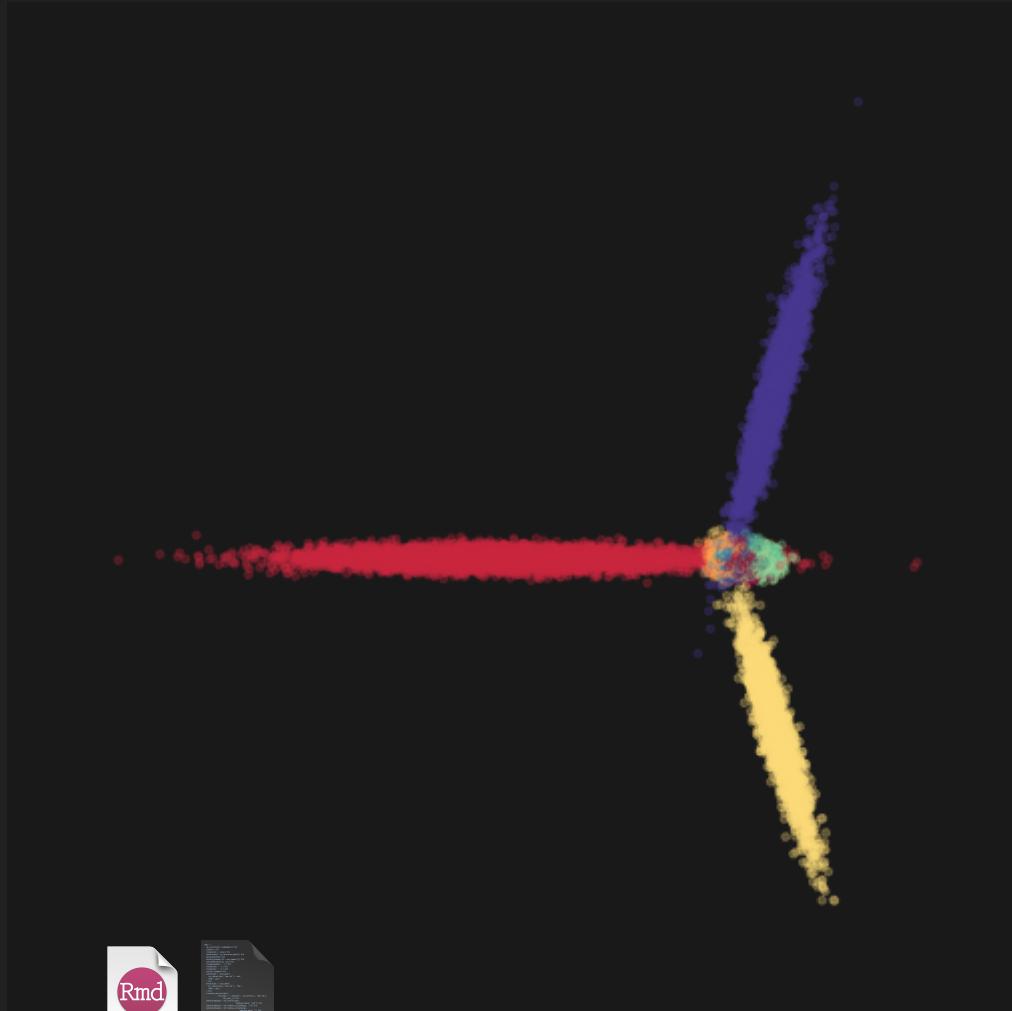
[5]

Below is an animation of the t-SNE process which shows a complex data set, data reduction and then clustering⁶

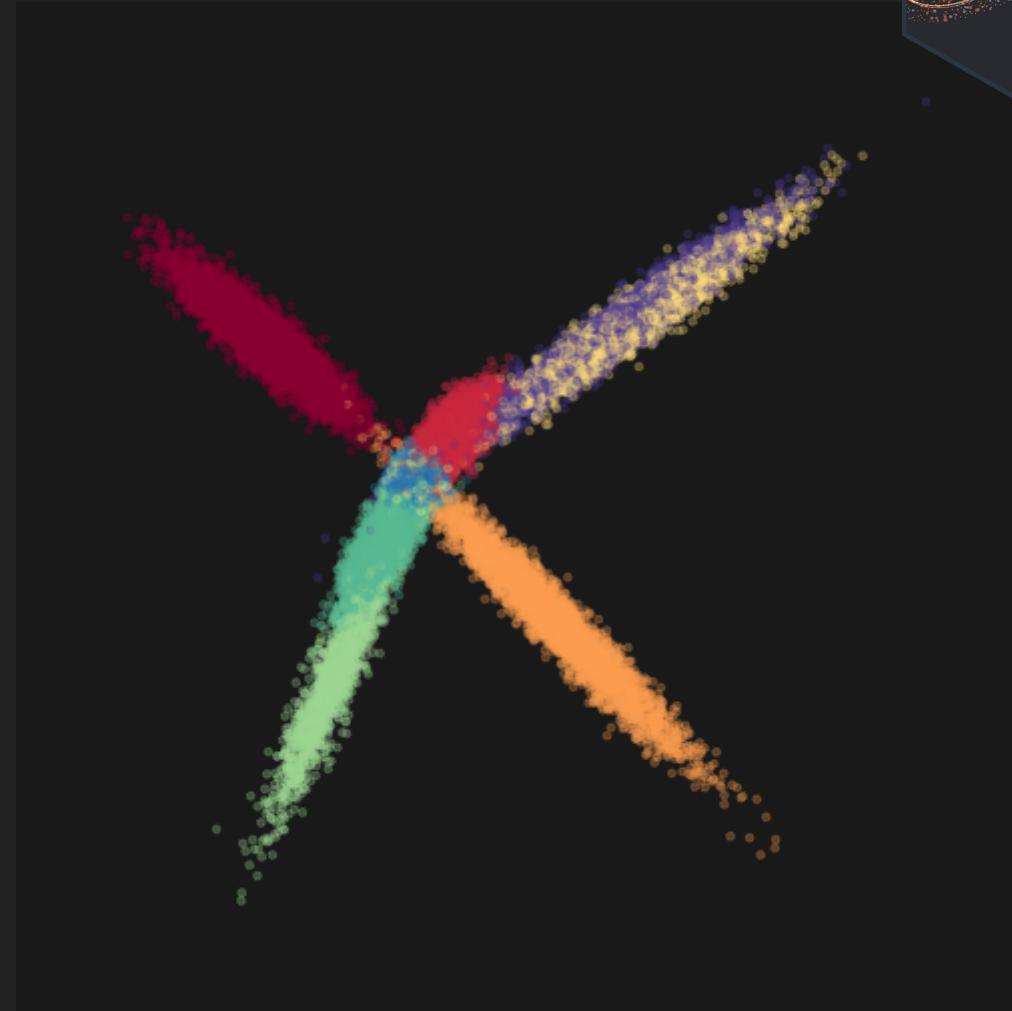


[6]

And just for fun here are two pca rotations of the example data set⁷



[7]



Surveys and PCAs



In general, using a PCA in survey data analysis helps you to understand

Surveys and PCAs



In general, using a PCA in survey data analysis helps you to understand

- how each item is similar to all others and the strength of that relationship

Surveys and PCAs



In general, using a PCA in survey data analysis helps you to understand

- how each item is similar to all others and the strength of that relationship
- which items are should likely be kept or removed



But Wait There's More!

Again this is just the tip of the iceberg. To really see the power of PCAs, take a look at machine learning. This is just one of many ways to deal with classification and dimensionality. Here are a couple resources. At this time, its good just to ignore the coding and to simply get a basic idea of each.

[8] If you cannot see the entire page, please load the site in a private window. Directions on how to do this are provided for





But Wait There's More!

Again this is just the tip of the iceberg. To really see the power of PCAs, take a look at machine learning. This is just one of many ways to deal with classification and dimensionality. Here are a couple resources. At this time, its good just to ignore the coding and to simply get a basic idea of each.

- 11 Dimensionality reduction techniques you should know in 2021⁸

[8] If you cannot see the entire page, please load the site in a private window. Directions on how to do this are provided for





But Wait There's More!

Again this is just the tip of the iceberg. To really see the power of PCAs, take a look at machine learning. This is just one of many ways to deal with classification and dimensionality. Here are a couple resources. At this time, its good just to ignore the coding and to simply get a basic idea of each.

- 11 Dimensionality reduction techniques you should know in 2021⁸
- Understanding Dimension Reduction and Principal Component Analysis in R for Data Science⁸

[8] If you cannot see the entire page, please load the site in a private window. Directions on how to do this are provided for





But Wait There's More!

Again this is just the tip of the iceberg. To really see the power of PCAs, take a look at machine learning. This is just one of many ways to deal with classification and dimensionality. Here are a couple resources. At this time, its good just to ignore the coding and to simply get a basic idea of each.

- 11 Dimensionality reduction techniques you should know in 2021⁸
- Understanding Dimension Reduction and Principal Component Analysis in R for Data Science⁸
- Workshop: Dimension reduction with R

[8] If you cannot see the entire page, please load the site in a private window. Directions on how to do this are provided for



Thats it!

If you have any questions, please reach out



Thats it!

If you have any questions, please reach out



This work is licensed under a
Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License