# An Introduction to Principal Component Analysis

# Outline

- LEARNING OBJECTIVES:
  - You will know why we use dimensionality reduction.
  - You will know the PCA implementation process.
  - You will know at a high level what the PCA formulas are doing in a small sample setting.

- YOU WILL NOT:
  - Know the linear algebra that makes PCA work.
  - You will not be able to magically think in more than 3 dimensions.

## Everyone has been here…

- You have all the columns below in your dataset. Which ones are you going to use to explain what is happening?



## Key Topics:

- **What is dimensionality reduction?**

- *How do we use PCA for dimensionality reduction?*

- *How is are the PCA formulas working to do this?*

## Key Topics:

- **What is dimensionality reduction?**

**-** *How do we use PCA for dimensionality reduction?*

**-** *How is are the PCA formulas working to do this?*

## High-level definition:

- Dimensionality reduction is used to end up with fewer features in our data set that we can then analyze.

## Two ways to do it.

- 1. FEATURE SELECTION: Select which columns to use and ignore the rest.

- 2. FEATURE EXTRACTION: Generalizes the columns into a set of new columns. Some of the new columns might be more useful than their original setup and some columns might be less useful. So we drop the ones that are less useful.

## PCA is a type of feature extraction.

**Breaking down what dimensionality reduction means for PCA.**

- Looking at this dataset below… what are the "Dimensions"?

| Blue | Red | Yellow |
|------|-----|--------|
| 20% | 57% | 23% |
| 60% | 40% | 0% |
| 50% | 0% | 50% |

## Key Topics:

*- What is dimensionality reduction?*

**- How do we use PCA for dimensionality reduction?**

*- How is are the PCA formulas working to do this?*

**How can you reduce these dimensions with feature extraction? … enter PCA.**

- PCA is going to look at all our features and then create new ones.

- These new ones are going to be a little special though… the first one of the new features will explain more of the data than the BLUE feature we currently have.
  - *If you're trying to figure out HOW pca does this, we will have an example after this section.*

**The oh so important process when using PCA.**

- 1. You first need to take all your variables and "Standardize" them.

  - *This is so that when the PCA formula looks to make new features from all the ones we have it can make apples from apples.*

- 2. You then get to simply decide how many new variables to keep.

**Wait…    What's the catch?**

# The catch.

- Each of the new features… or principle components are going to look very different.

| Principle Component 1 | Principle Component 2 | Principle Component 3 |
|:---:|:---:|:---:|
| 70.1% | 25.2% | 4.5% |
| 20.7% | 32.4% | 2.1% |
| 45.3% | 29.4% | 3.4% |

## Key Topics:

- *What is dimensionality reduction?*

**- How do we use PCA for dimensionality reduction?**

**-** *How is are the PCA formulas working to do this?*

## So your boss asks…

- Great predictive results! What was the most important thing that determined your findings?

## Uh…...

# The OTHER catch.

- Remember we said that we are dropping features… so... If we don't know what these new features mean… how do we know we aren't losing something really important?

## PCA explained with spoons and my daughters toys.
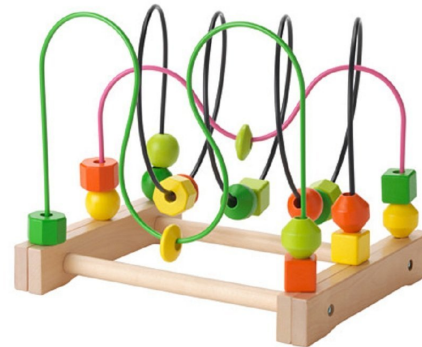
- So remember we have three "Dimensions".

  ***Blue***                    ***Red***                    ***Yellow***

- PCA will assign plot all of these data points in space that is equal to the number of features/dimensions we have. 3 features, so 3 dimensions.

# Cool so we have our data plotted.

- PCA will now make principle components… or... Spoons. One spoon per dimension/feature.

*Spoon 1*           *Spoon 2*           *Spoon 3*

## These spoons have rules…

- 1. The spoons have to be used to plot the data in the best way one at a time.

- 2. The spoons MUST be perpendicular to each other.

- *You can think of these as plotting a plane in our 3D space.*

## Key Topics:

*- What is dimensionality reduction?*

**-** *How do we use PCA for dimensionality reduction?*

**- How is are the PCA formulas working to do this?**

**Each spoon now says something about our data. In a new 3D plot.**

### *Spoon 1*

*Captures 75% of the data spread.*

### *Spoon 2*

*Captures an additional 22% of the data spread.*

### *Spoon 3*

*Captures an additional 3% of the data spread.*

## Key Topics:

*- What is dimensionality reduction?*

**-** *How do we use PCA for dimensionality reduction?*

**- How is are the PCA formulas working to do this?**

**So 3 features is nice… but what if I have 100 like with the beginning weather data?**

## Key Topics:

*- What is dimensionality reduction?*

**- How do we use PCA for** *dimensionality reduction?*

**- How is are the PCA formulas working to do this?**

**So in the end… why you would want to use PCA.**

1.     *If you **DO NOT** need to explain the data behind your results.*

2.     *Because PCA will give us the **BEST** features while using the least amount of columns.*

3.     *Because PCA makes cookies.*

That's all folks.

Questions?