

Implementing Factor Analysis and PCA in Python



Vitthal Srinivasan

CO-FOUNDER, LOONYCORN

www.loonycorn.com

Overview

Assemble a data set of returns from correlated stocks

Use Python to calculate principal components of the financial data

Eliminate low-value principal components using eigenvalues

Relate the principal components to underlying latent factors

PCA in Python

Explain Google's returns

Yahoo finance

Using returns of correlated stocks

Eigen Decomposition

Python library function

On covariance matrix

Principal Components

From eigen vectors

Uncorrelated components

Covariance and Correlation

Correlation matrix signals trouble

Multicollinearity problems

Scree Plot

Number of dimensions

Discard low-value dimensions

Interpret and Regress

Beta, bonds, sectors

Now regress Google

Demo

**Implement Eigen analysis and PCA in
Python**

Negative Indices in R

goog

DATE	GOOG. PRICE	NASDAQ. PRICE
2016-12-01	779	5550
2016-11-01	747	5324
2006-01-01	309	1900

Exclude

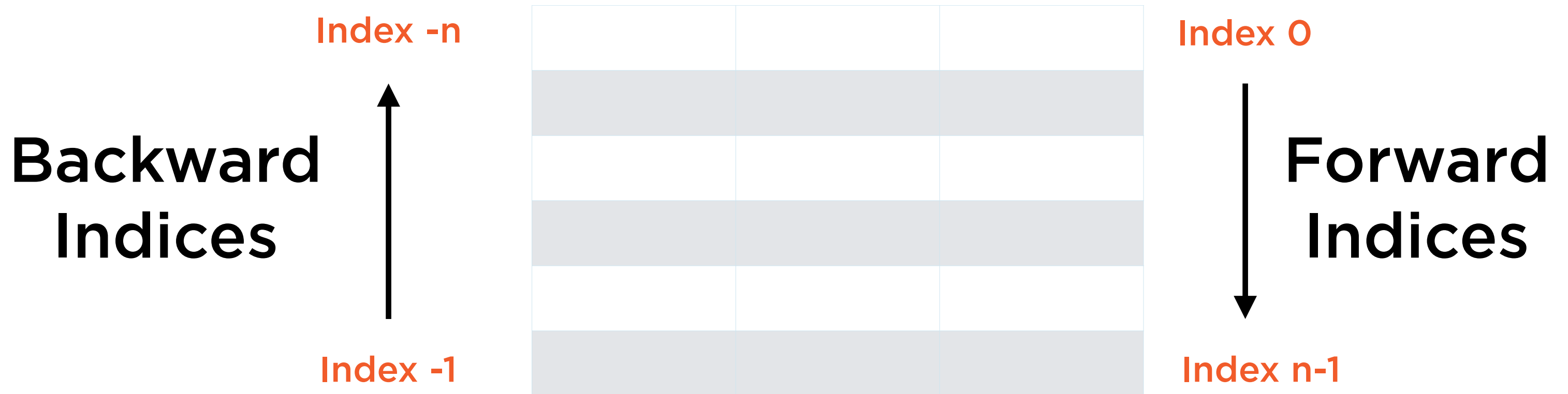
Row 1

Row nrow(goog)

Column 1

`goog[-nrow(goog),-1]`

Negative Indices In Python



PCA should always be applied on the
covariance matrix of standardised
vectors

Standardising Data

$$\begin{bmatrix} X_{11} & & X_{1k} \\ X_{21} & & X_{2k} \\ X_{31} & \dots & X_{3k} \\ \dots & & \dots \\ X_{n1} & & X_{nk} \end{bmatrix}$$

$\text{avg}(X_1) \quad \dots \quad \text{avg}(X_k)$

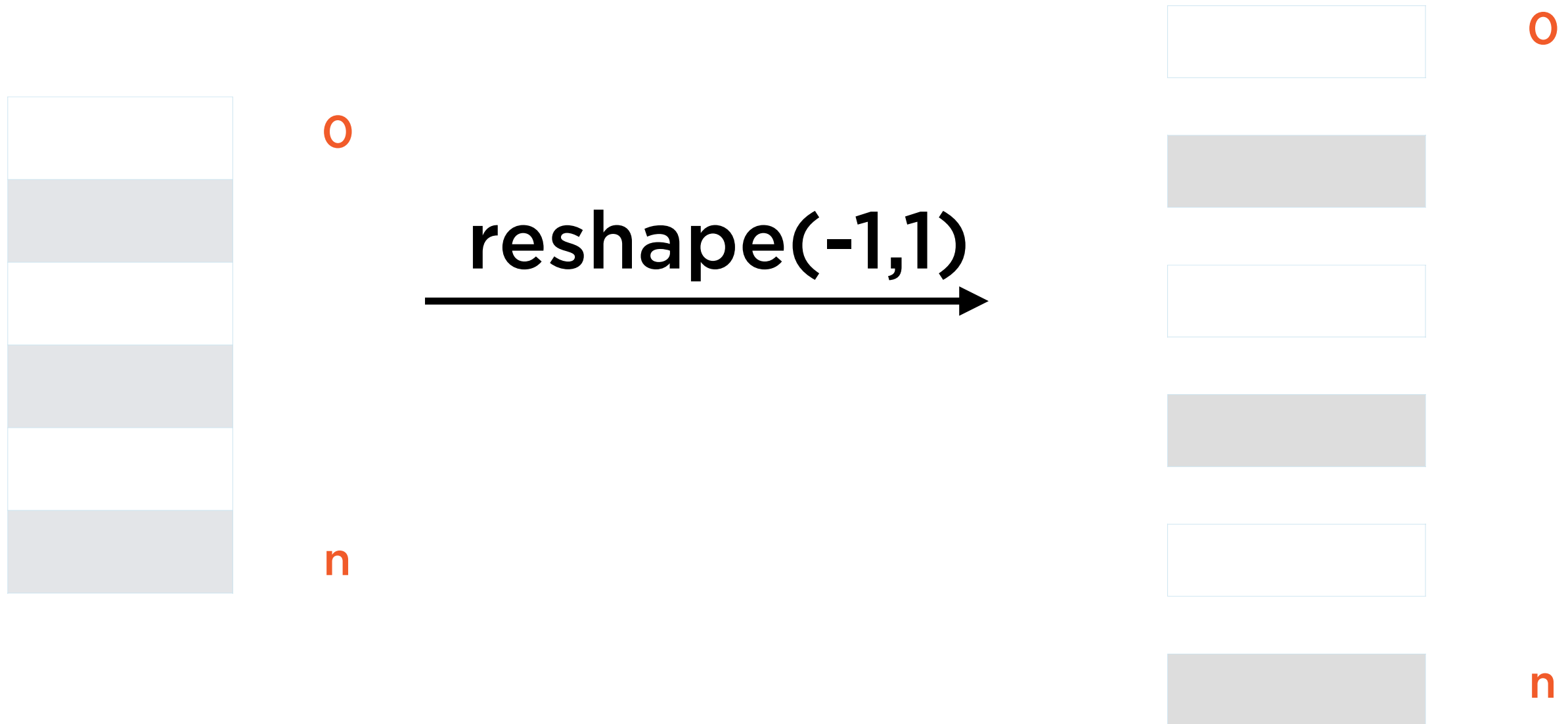
$\text{stdev}(X_1) \quad \dots \quad \text{stdev}(X_k)$

Standardising Data

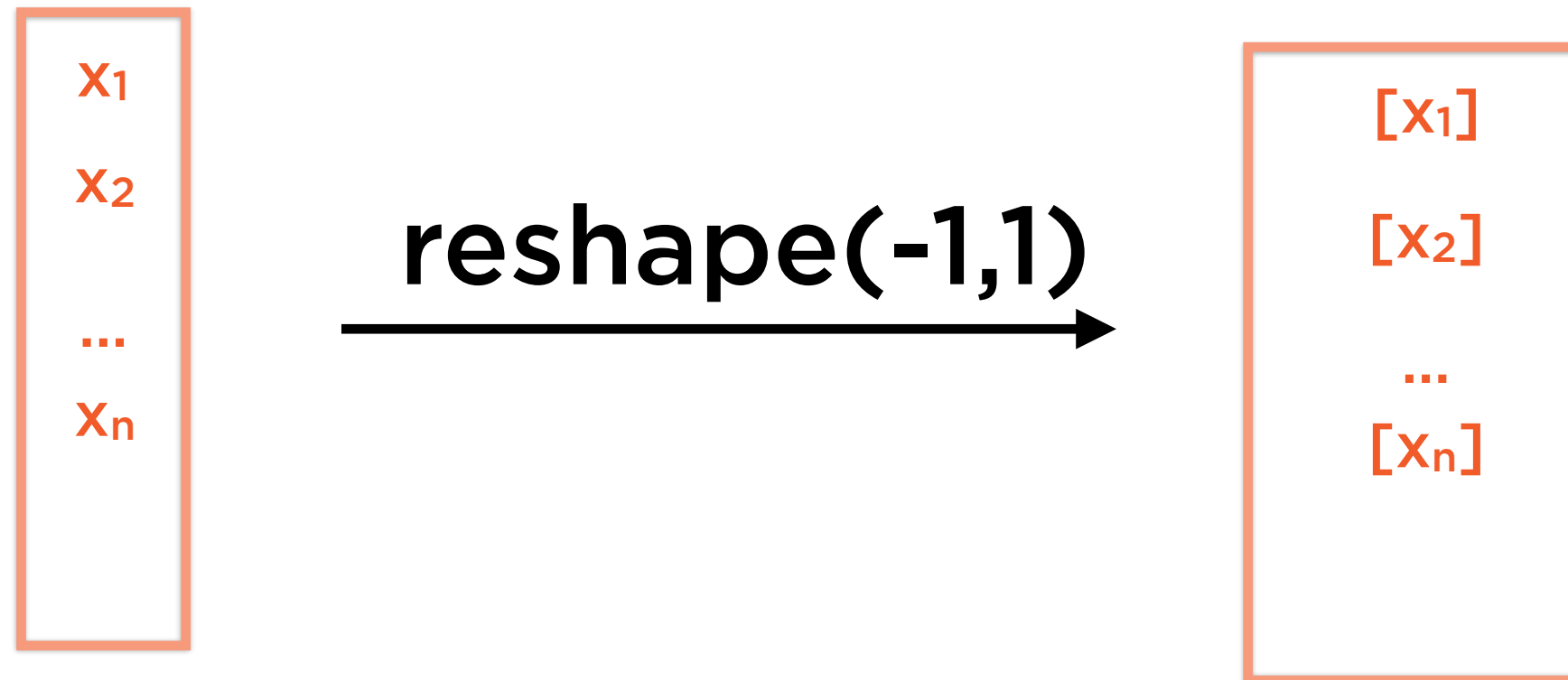
$$\begin{bmatrix} \frac{x_{11} - \text{avg}(X_1)}{\text{stdev}(X_1)} & \frac{x_{1k} - \text{avg}(X_k)}{\text{stdev}(X_k)} & \dots \\ \dots & \dots & \dots \\ \frac{x_{n1} - \text{avg}(X_1)}{\text{stdev}(X_1)} & \frac{x_{nk} - \text{avg}(X_k)}{\text{stdev}(X_k)} & \dots \end{bmatrix}$$

Each column of the standardised data has mean 0 and variance 1

Reshaping in NumPy




Reshaping in NumPy



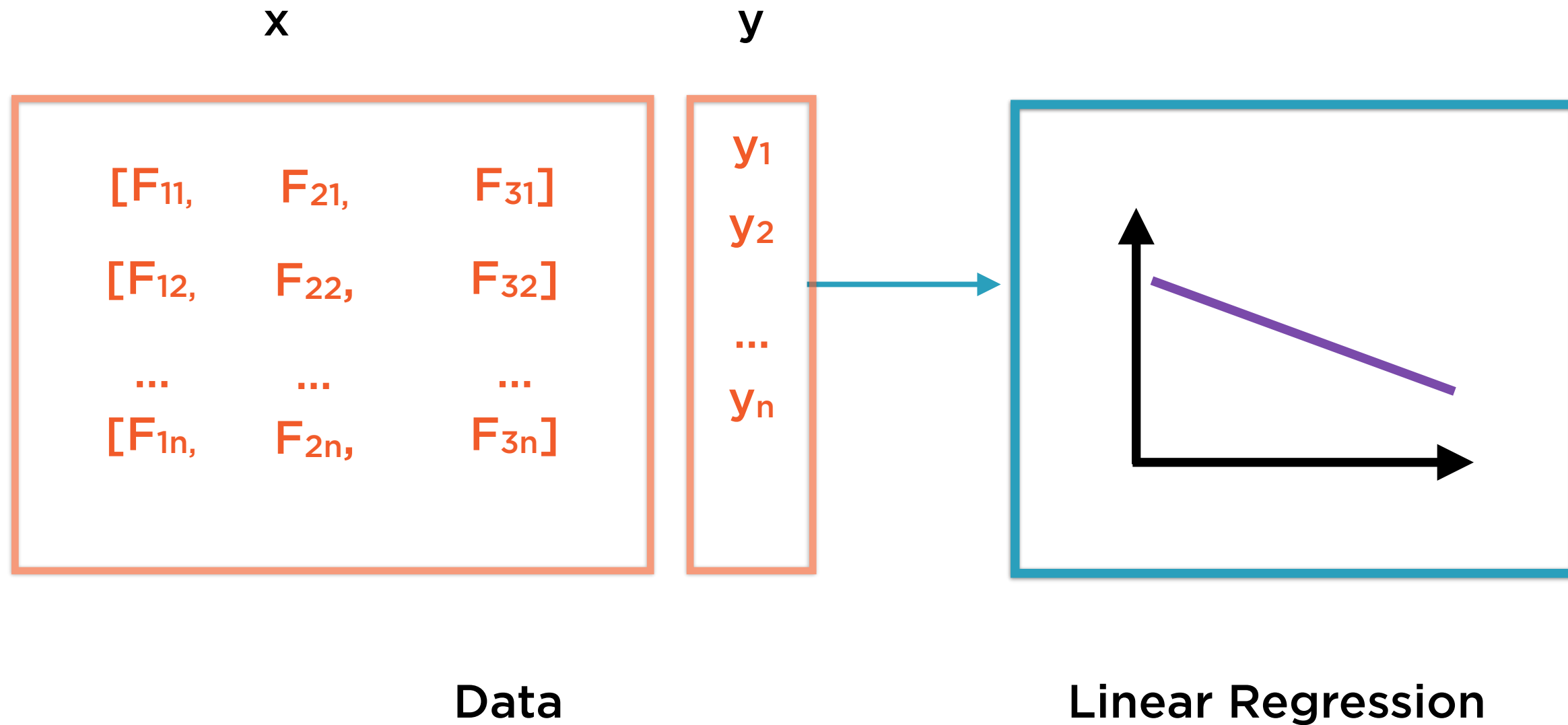
Rearranging Numbers in Python

PCA₁  [**F₁₁** **F₁₂** ...]

PCA₂  [**F₂₁** **F₂₂** ...]

PCA₃  [**F₃₁** **F₃₂** ...]

Rearranging Numbers in Python



Rearranging Numbers in Python

PCA₁.T



[F₁₁]
[F₁₂]
...
[F_{1n}]

PCA₂.T



[F₂₁]
[F₂₂]
...
[F_{2n}]

PCA₃.T



[F₃₁]
[F₃₂]
...
[F_{3n}]

Rearranging Numbers in Python

PCA₁.T



[F₁₁]
[F₁₂]
...
[F_{1n}]

PCA₂.T



[F₂₁]
[F₂₂]
...
[F_{2n}]

PCA₃.T



[F₃₁]
[F₃₂]
...
[F_{3n}]

Rearranging Numbers in Python

zip(PCA₁.T, PCA₂.T, PCA₃.T)



([F ₁₁]	[F ₂₁]	[F ₃₁])
([F ₁₂]	[F ₂₂]	[F ₃₂])
	
([F _{1n}]	[F _{2n}]	[F _{3n}])

Rearranging Numbers in Python

zip(PCA₁.T, PCA₂.T, PCA₃.T **)**



([F ₁₁]	[F ₂₁]	[F ₃₁])
([F ₁₂]	[F ₂₂]	[F ₃₂])
	
([F _{1n}]	[F _{2n}]	[F _{3n}])

Rearranging Numbers in Python

`zip(PCA1.T , PCA2.T, PCA3.T).reshape(-1,3)`



$[F_{11},$	$F_{21},$	$F_{31}]$
$[F_{12},$	$F_{22},$	$F_{32}]$
\dots	\dots	\dots
$[F_{1n},$	$F_{2n},$	$F_{3n}]$

Rearranging Numbers in Python

`zip(PCA1.T , PCA2.T, PCA3.T)` **`.reshape(-1,3)`**



$[F_{11},$

$F_{21},$

$F_{31}]$

$[F_{12},$

$F_{22},$

$F_{32}]$

...

...

...

$[F_{1n},$

$F_{2n},$

$F_{3n}]$

Principal Components Analysis

$[X_1 \ X_2 \ X_3 \ \dots \ X_k]$



Eigenvalue
Decomposition



Principal Components:

$[F_1 \ F_2 \ F_3 \ \dots \ F_k]$

\leftarrow \rightarrow

k columns

\updownarrow
n rows

Eigenvectors:

$[V_1 \ V_2 \ V_3 \ \dots \ V_k]$

\leftarrow \rightarrow

k columns

\updownarrow
k rows

Eigenvalues:

$[e_1 \ e_2 \ e_3 \ \dots \ e_k]$

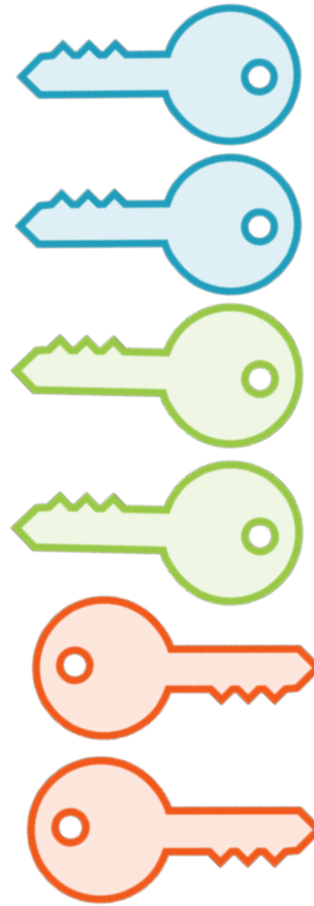
\leftarrow \rightarrow

k columns

\updownarrow
1 row

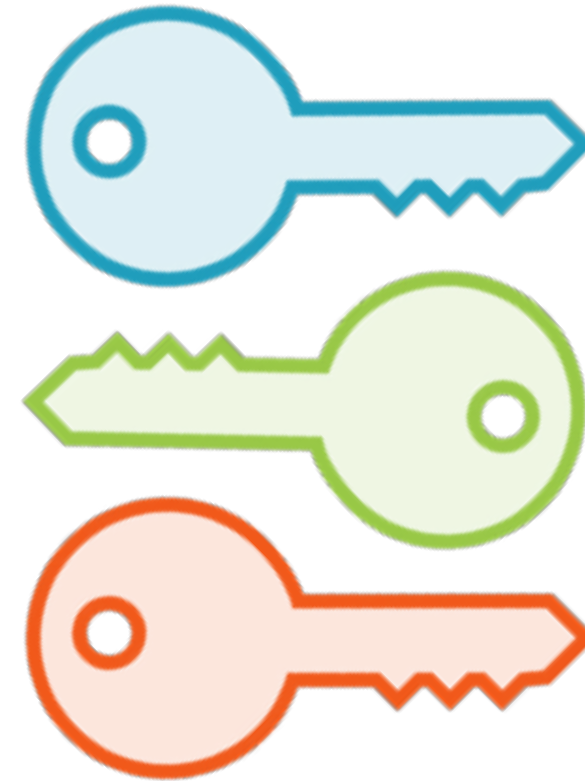
Keeping things simple is quite complicated

Similar, yet Different



Regression

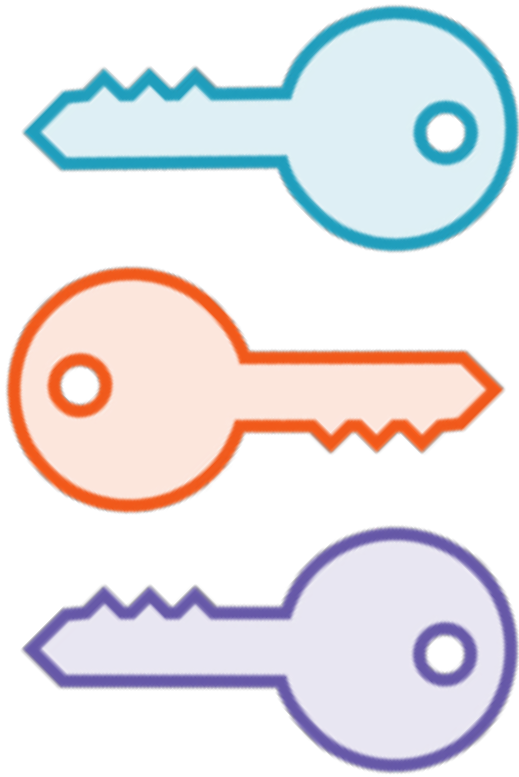
Connect the dots



Factor Analysis

Cut through the clutter

Regression



Causes

Independent variables



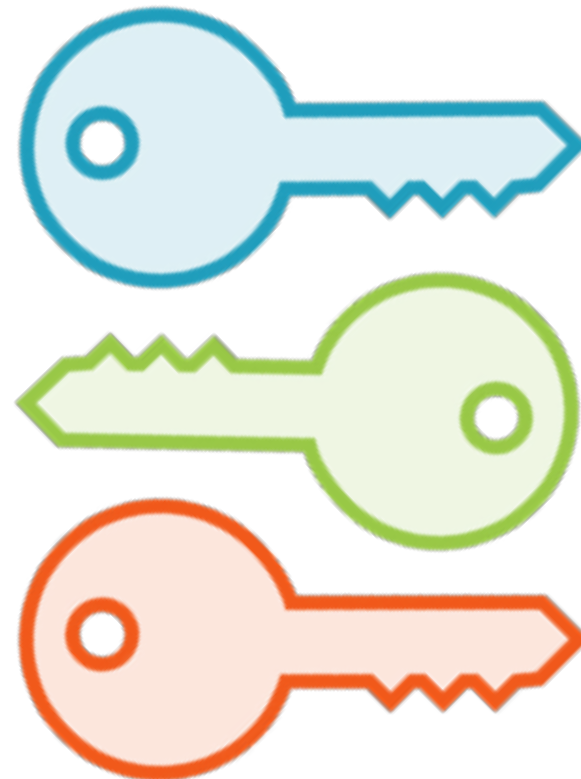
Effect

Dependent variable

Factor Analysis



**Many Observed
Causes**

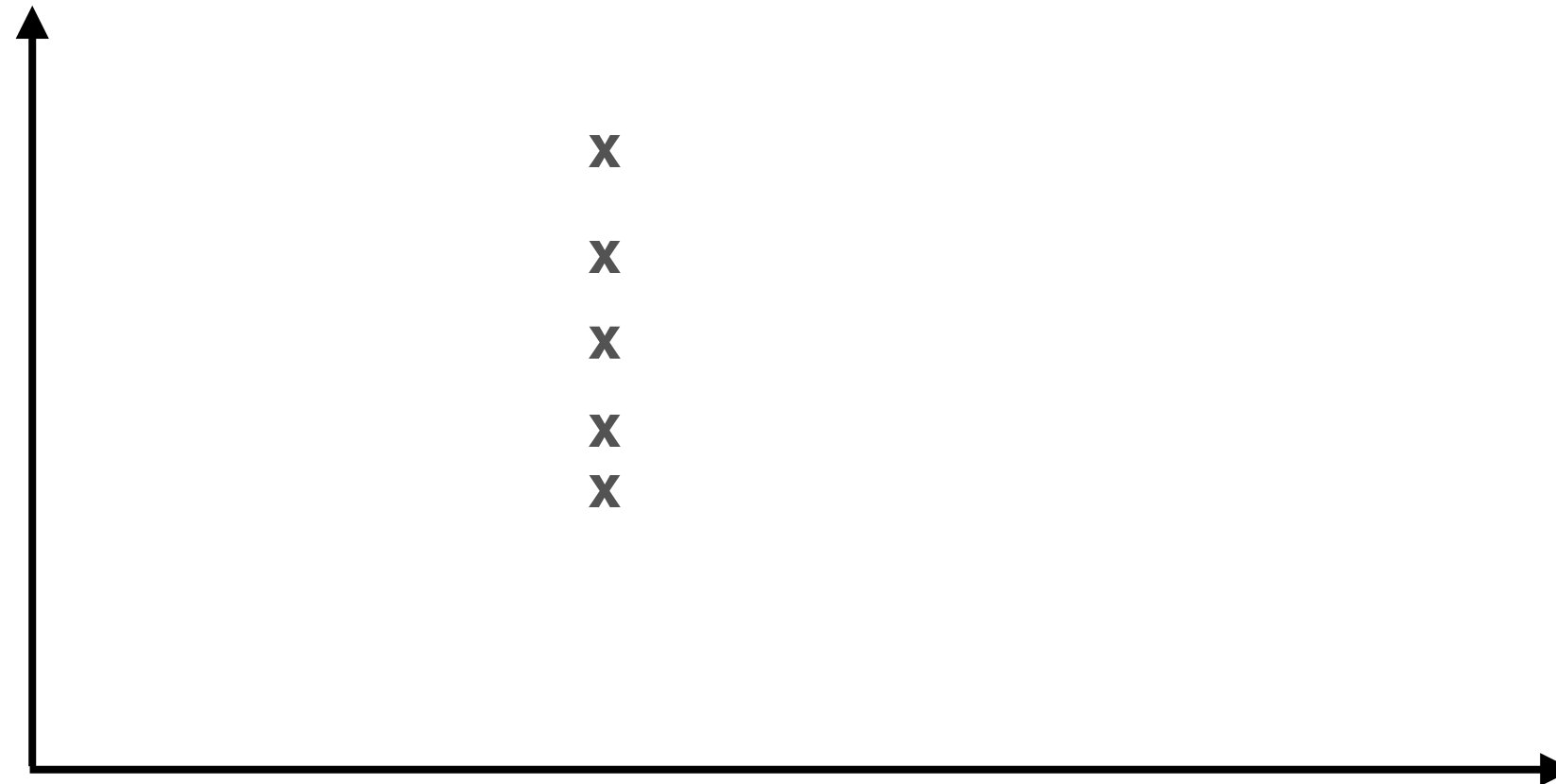


**Few Underlying
Causes**



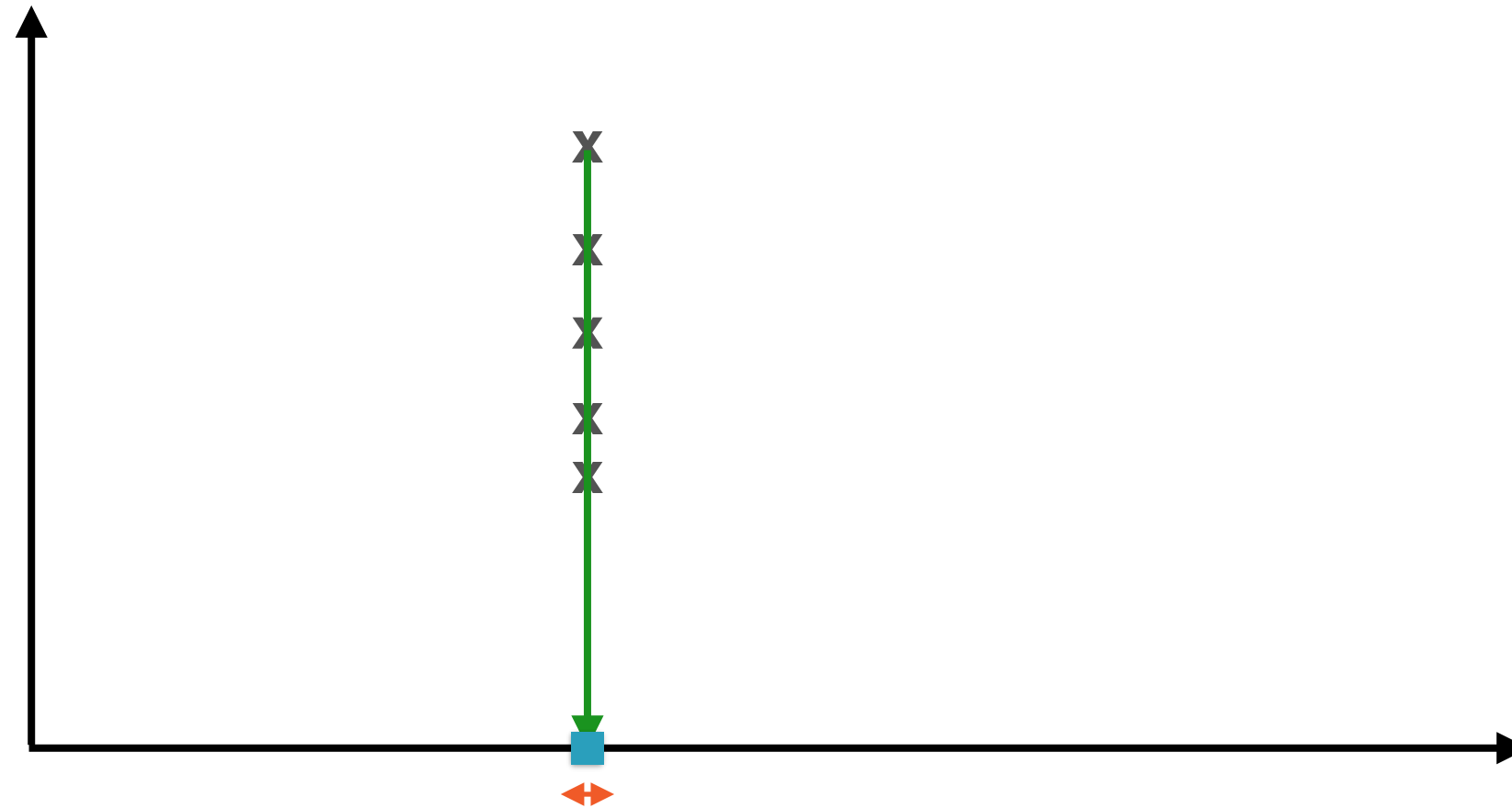
One Effect

A Question of Dimensionality



Pop quiz: Do we really need two dimensions to represent this data?

Bad Choice of Dimensions



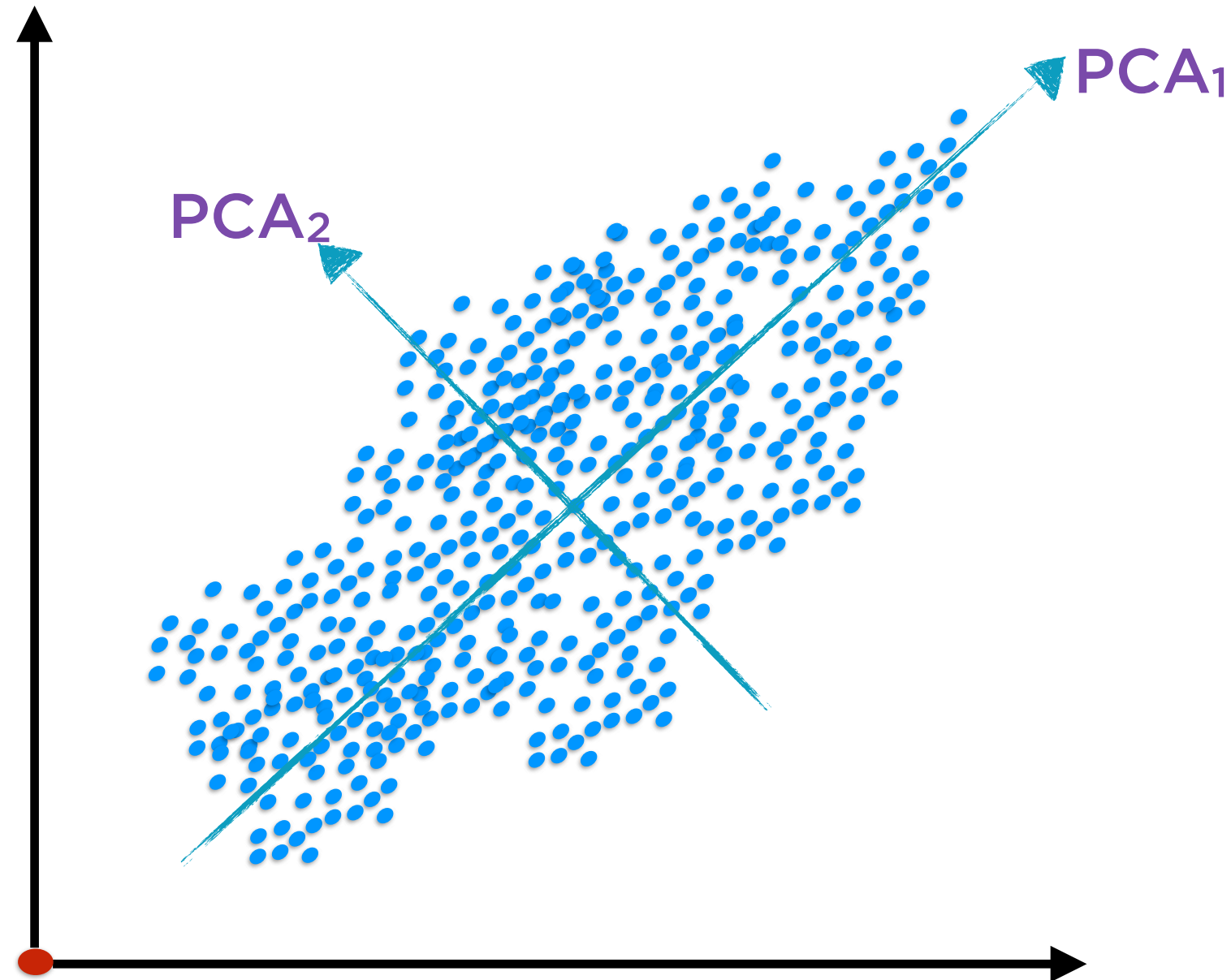
If we choose our axes (dimensions) poorly then we do need two dimensions

Good Choice of Dimensions



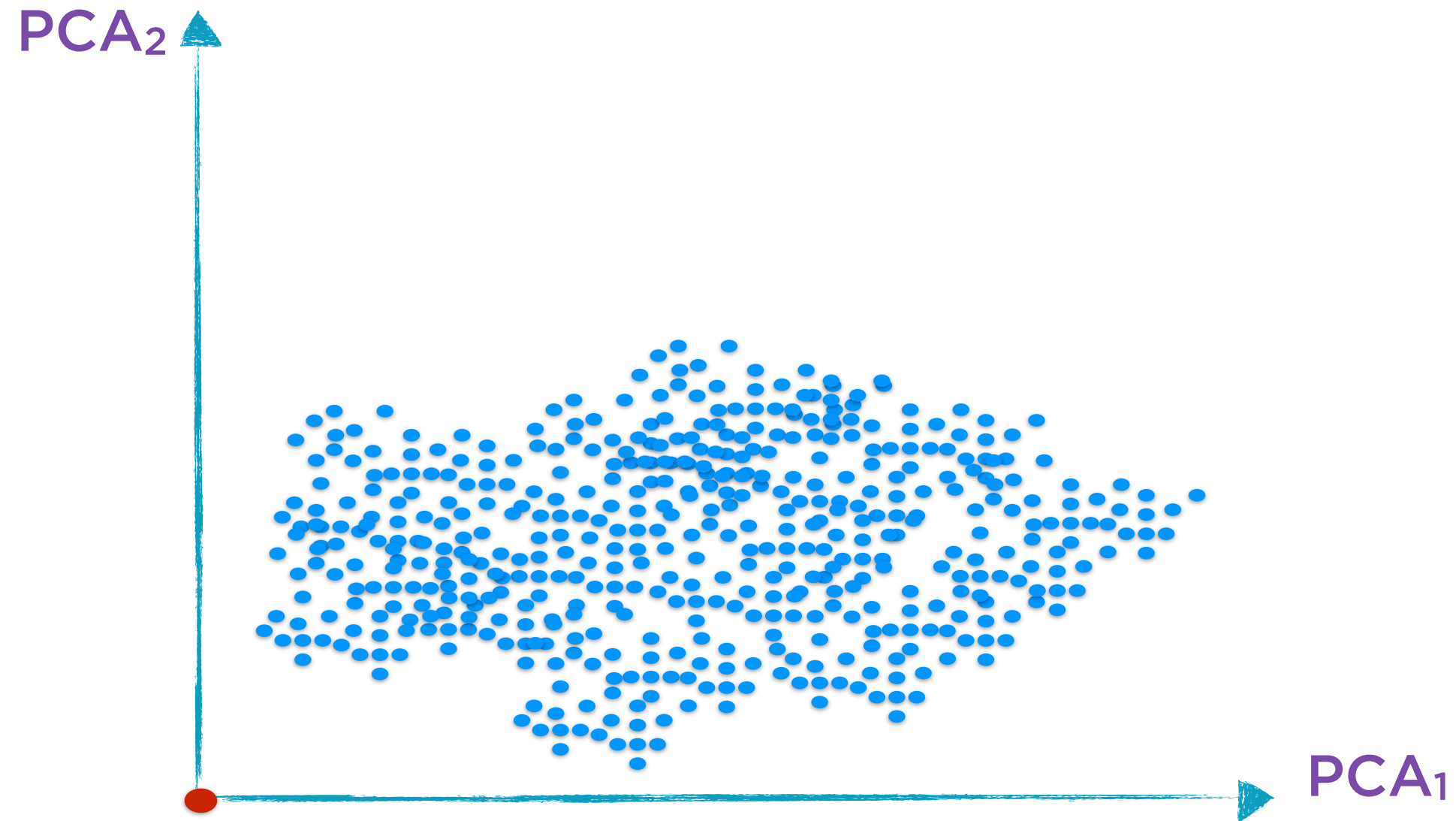
If we choose our axes (dimensions) well then one dimension is sufficient

Intuition Behind PCA



In general, there are as many principal components as there are dimensions in the original data

Intuition Behind PCA



Re-orient the data along these new axes

Summary

Python has powerful libraries for PCA and eigen analysis

PCA of equity returns reveals three important principal components