

# Implementing Factor Analysis and PCA in R

---



**Vitthal Srinivasan**

CO-FOUNDER, LOONYCORN

[www.loonycorn.com](http://www.loonycorn.com)

# Overview

**Assemble a data set of returns from correlated stocks**

**Use R to calculate principal components of the financial data**

**Eliminate low-value principal components using a Scree plot**

**Relate the principal components to underlying latent factors**

# PCA in R

## Explain Google's returns

Yahoo finance

Using returns of correlated stocks

## Eigen Decomposition

Built-in R 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 R**

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

# Data Frame: Data in Rows and Columns

Each row represents 1 observation	DATE	OPEN	...	ADJUSTED CLOSE	Each column represents 1 variable (a list or vector)
	2016-12-01	772	...	779	
	2016-11-01	758	...	747	
	2006-01-01	302	...	309	

# From File to Data Frame

DATE	OPEN	...	ADJUSTED CLOSE
2016-12-01	772	...	779
2016-11-01	758	...	747
2006-01-01	302	...	309

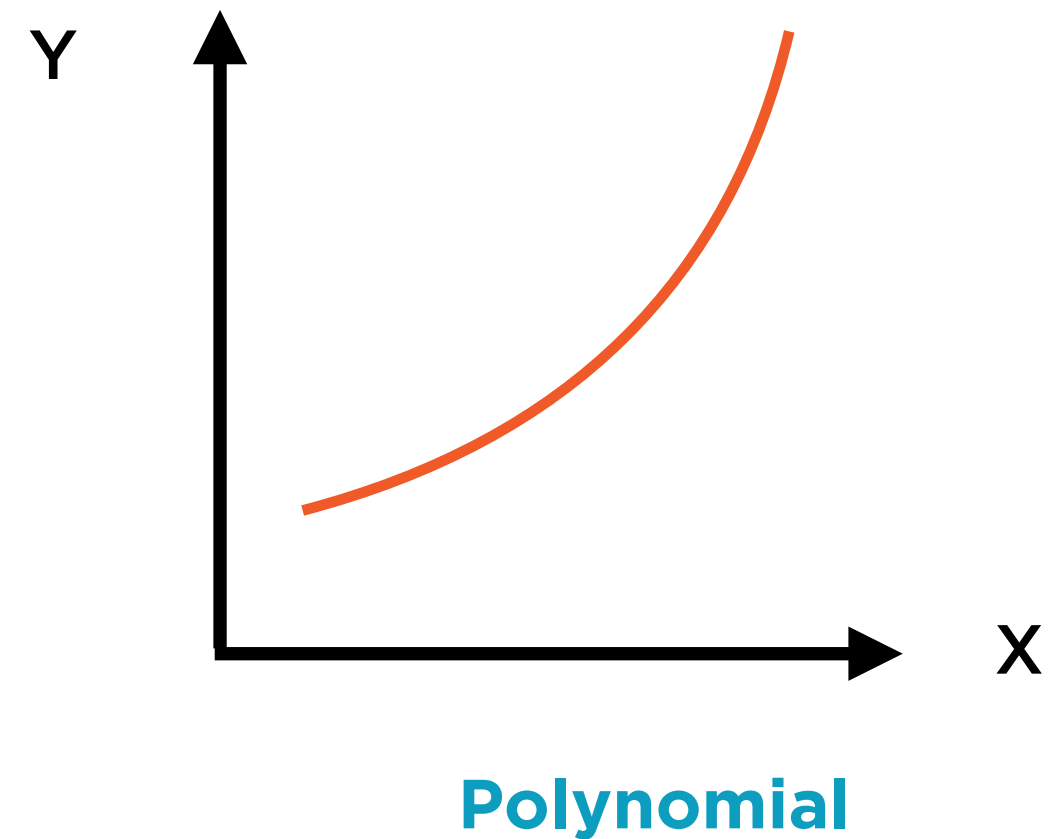
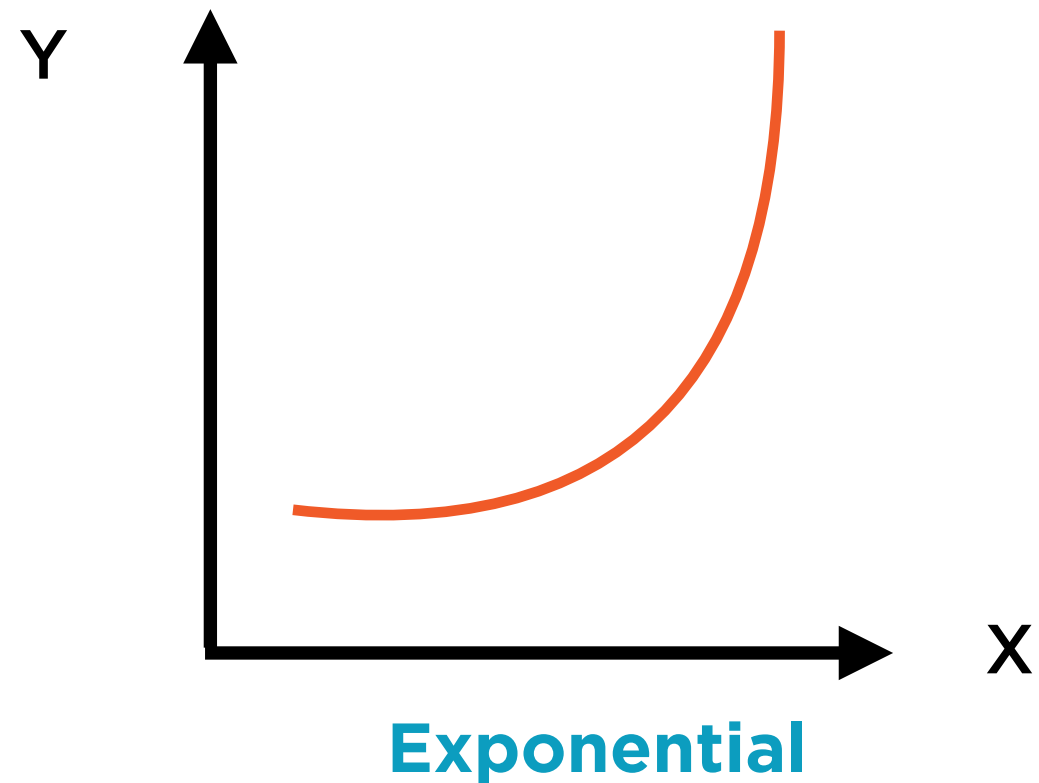
File

→  
`read.table`

DATE	OPEN	...	ADJUSTED CLOSE
2016-12-01	772	...	779
2016-11-01	758	...	747
2006-01-01	302	...	309

Data Frame

# Never Regress Non-Stationary Data



Smoothly trending data will lead to poor quality regression models



# First Differences

$$y'_{12} = \log y_2 - \log y_1$$

$$x'_{12} = \log x_2 - \log x_1$$

Regress  $y'$  and  $x'$

$$y'_{12} = (y_2 - y_1)/y_1$$

$$x'_{12} = (x_2 - x_1)/x_1$$

Regress  $y'$  and  $x'$

**Log Differences**

**Returns**

Take first differences of smooth data converting  
either to log differences or returns

# Negative Indices => Exclude Data

**goog**

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

Row 1

Row nrow(goog)

Column 1

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

# Negative Indices => Exclude Data

**goog**

	DATE	GOOG. PRICE	NASDAQ. PRICE	
	2016-12-01	779	5550	Row 1
	2016-11-01	747	5324	
<b>Exclude</b>	2006-01-01	309	1900	Row nrow(goog)

Column 1

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

# Negative Indices => Exclude Data

**goog**

**Exclude**

DATE	GOOG. PRICE	NASDAQ. PRICE	
2016-12-01	779	5550	Row 1
2016-11-01	747	5324	
2006-01-01	309	1900	Row nrow(goog)

Column 1

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

# Negative Indices => Exclude Data

**goog**

	DATE	GOOG. PRICE	NASDAQ. PRICE	
	2016-12-01	779	5550	Row 1
	2016-11-01	747	5324	
<b>Exclude</b>	2006-01-01	309	1900	Row nrow(goog)

Column 1

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

# Element-wise Operations

<b>779</b>	<b>5550</b>	/	<b>747</b>	<b>5324</b>	=	<b>779/747</b>	<b>5550/5324</b>
						...	...
						...	...

**goog[-nrow(goog),-1]/  
goog[-1,-1]**

# Prices to Returns

<b>779/747</b>	<b>5550/5324</b>		<b>1</b>	<b>1</b>		<b>779/747 - 1</b>	<b>5550/5324 - 1</b>
...	...		<b>1</b>	<b>1</b>		...	...
		-	<b>1</b>	<b>1</b>	=		
			<b>1</b>	<b>1</b>			
...	...		<b>1</b>	<b>1</b>		...	...

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

**This converts prices to returns**

# 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

# Principal Components Analysis

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



Eigenvalue  
Decomposition



Principal Components:

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



k columns



n rows

Eigenvectors:

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



k columns



k rows

Eigenvalues:

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



k columns



1 row

# Interpreting Eigenvalues

**[  $F_1$     $F_2$     $F_3$    ...    $F_k$  ]**



**$\text{var}(F_1) > \text{var}(F_2) > \text{var}(F_3) > \text{var}(F_k)$**

These vectors  $F_i$  are arranged in order of  
decreasing variance

The greater the variance of a principal  
component, the more important it is

# Interpreting Eigenvalues

**[  $F_1$     $F_2$     $F_3$    ...    $F_k$  ]**



**var( $F_1$ ) > var( $F_2$ ) > var( $F_3$ ) > var( $F_k$ )**



**Eigenvalue 1**

**Eigenvalue 2**

**Eigenvalue 3**

**Eigenvalue k**

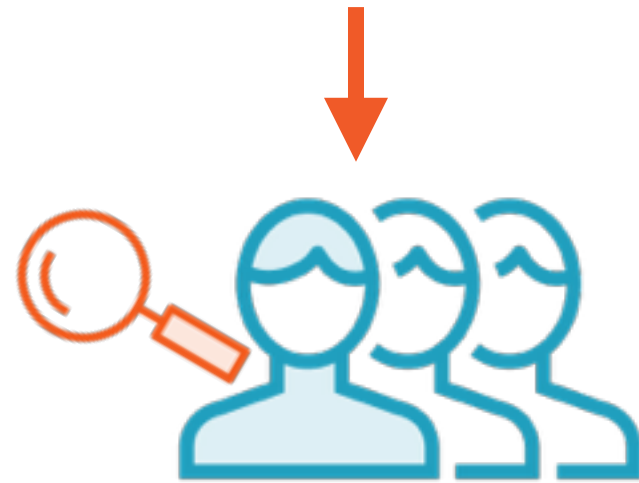
The greater the eigenvalue of a principal component, the more important it is

Use the Scree plot to determine how many principal components to discard

Exploratory Factor Analysis: Experts  
trace back principal components to  
observable factors

# PCA for Latent Factor Identification

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



$[ L_1 \quad L_2 \quad L_3 \quad \dots \quad L_k ]$



# 3 Latent Factors in Stock Returns

**Market Movements**

**Interest Rates**

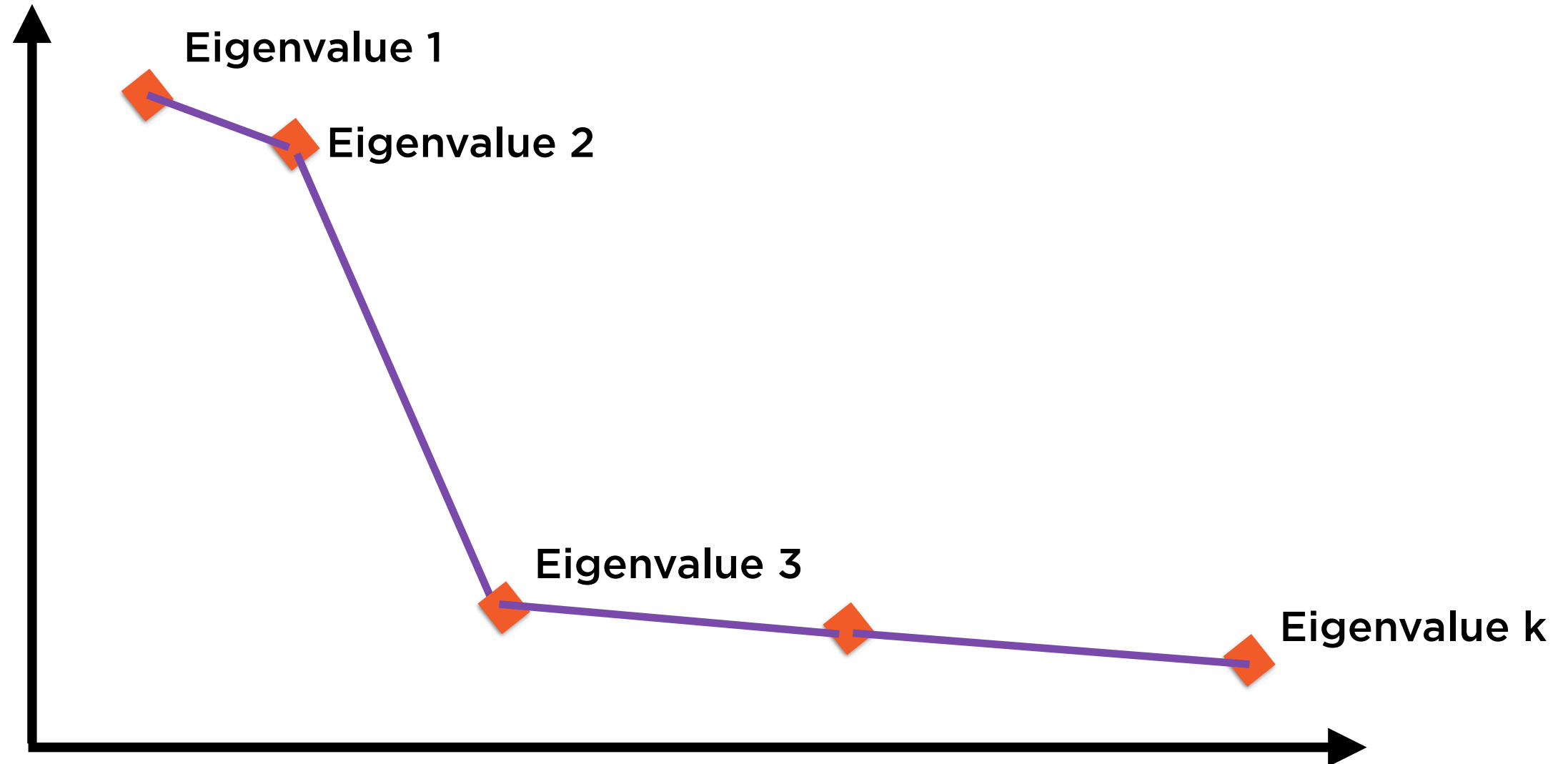
**Industry Sectors**

# Matrix Multiplication

$$\begin{array}{ccccc} \mathbf{F_i} & = & \mathbf{X} & & \mathbf{V_i} \\ \text{n rows,} & & \text{n rows,} & & \text{k rows,} \\ \text{1 column} & & \text{k columns} & & \text{1 column} \end{array}$$

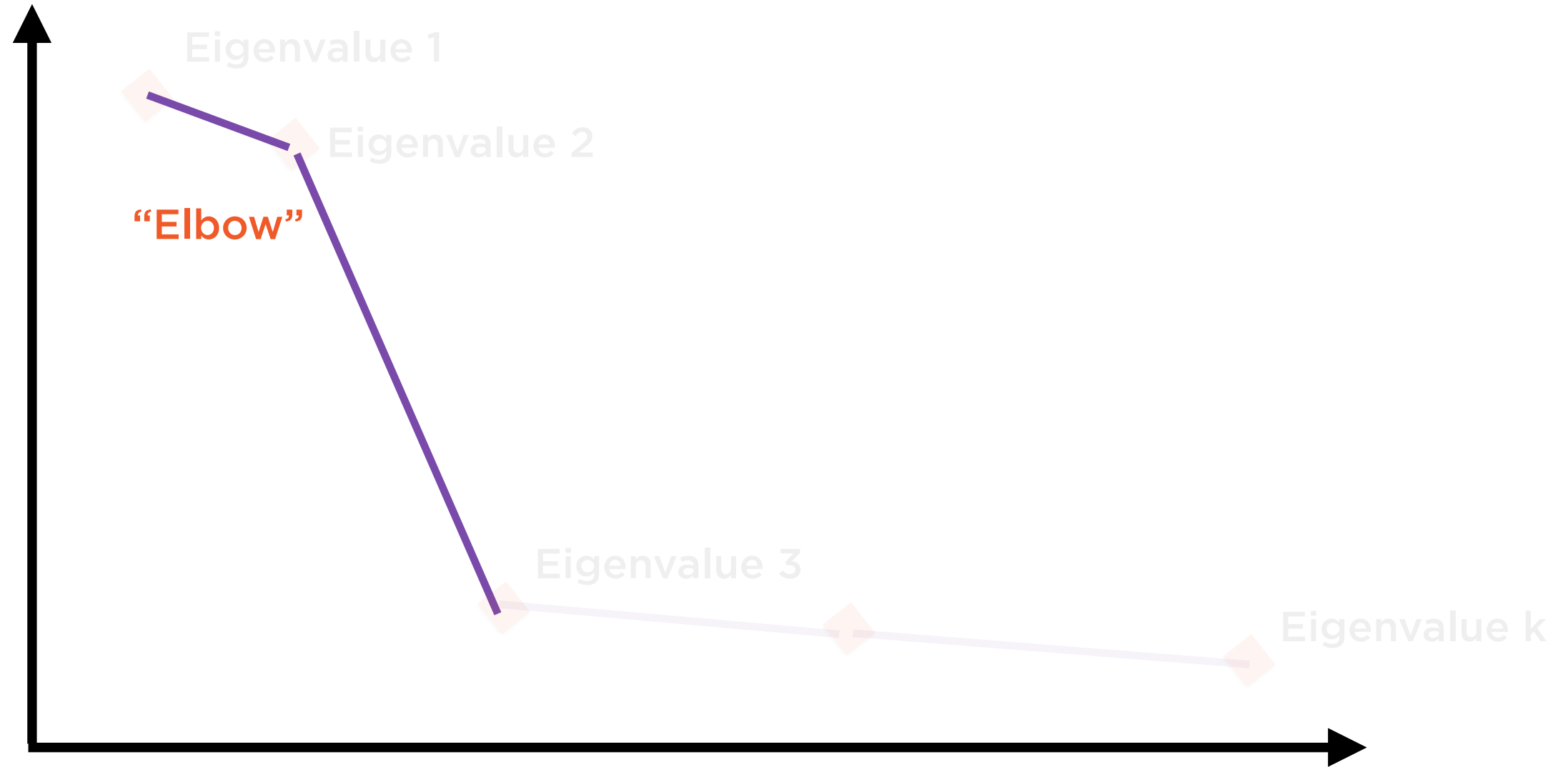
# Scree Plots

% of Total Variance  
Explained



# Scree Plots

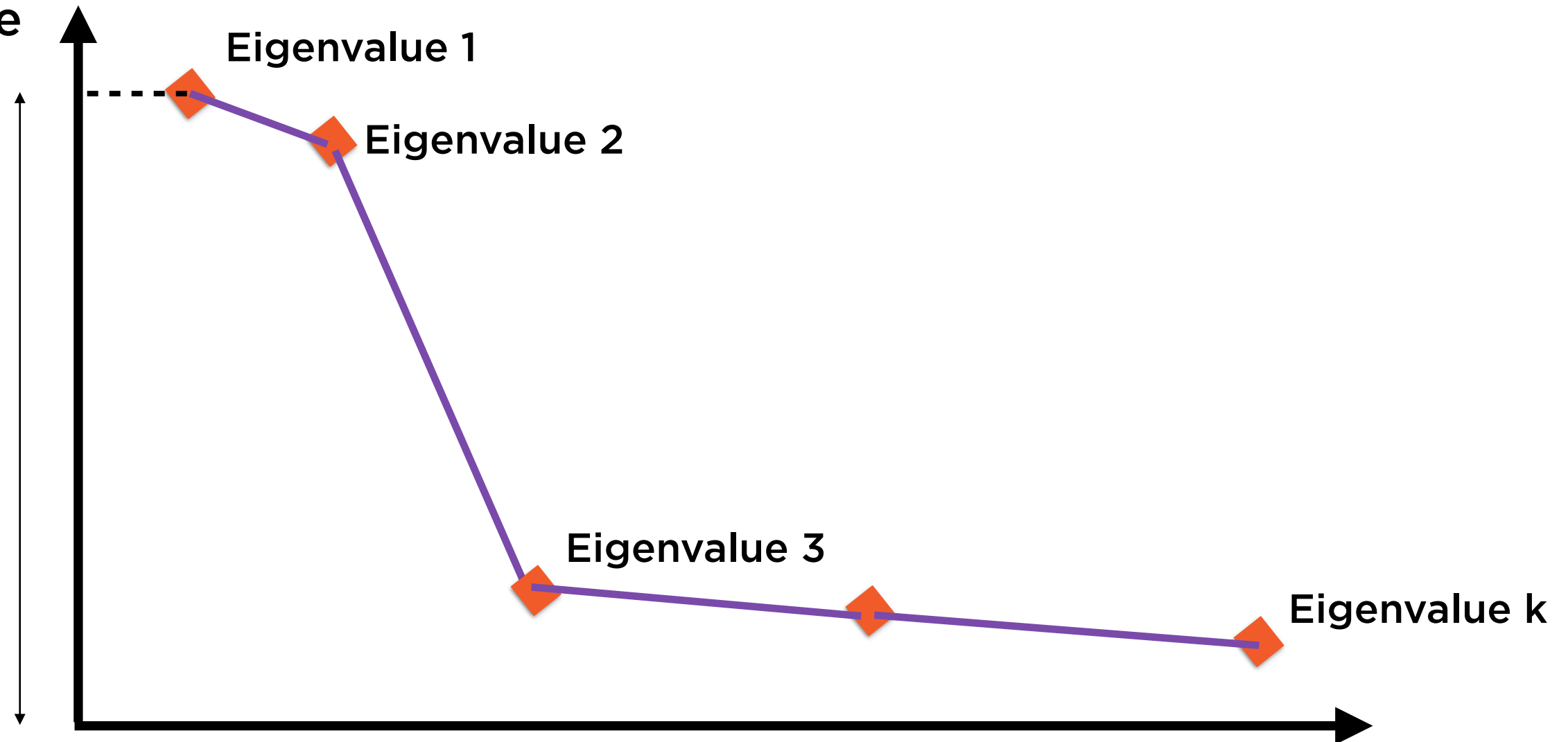
% of Total Variance  
Explained



# Scree Plots

% of Total Variance  
Explained

Proportion of  
variance explained  
by  $F_1$



# Summary

**R makes PCA very simple and easy-to-use**

**PCA of equity returns reveals three important principal components**

**These closely correlate with underlying economic factors**