

# Understanding Multiple Regression Models

---



**Vitthal Srinivasan**

CO-FOUNDER, LOONYCORN

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

# Overview

**Extend regression analysis to multiple explanatory variables**

**Interpret the results of a multiple regression**

**Mitigate the risks that accompany multiple regression**

**Apply multiple regression to include categorical variables**

# Introducing Multiple Regression

---

“A butterfly flapping its wings in  
Brazil can cause a tornado in Texas”

# Butterfly Effect

The butterfly effect is the concept that small causes can have large effects. In chaos theory, the butterfly effect is the sensitive dependence on initial conditions in which a small change in one state of a deterministic nonlinear system can result in large differences in a later state.

*Wikipedia*

# Simple Regression



**Cause**

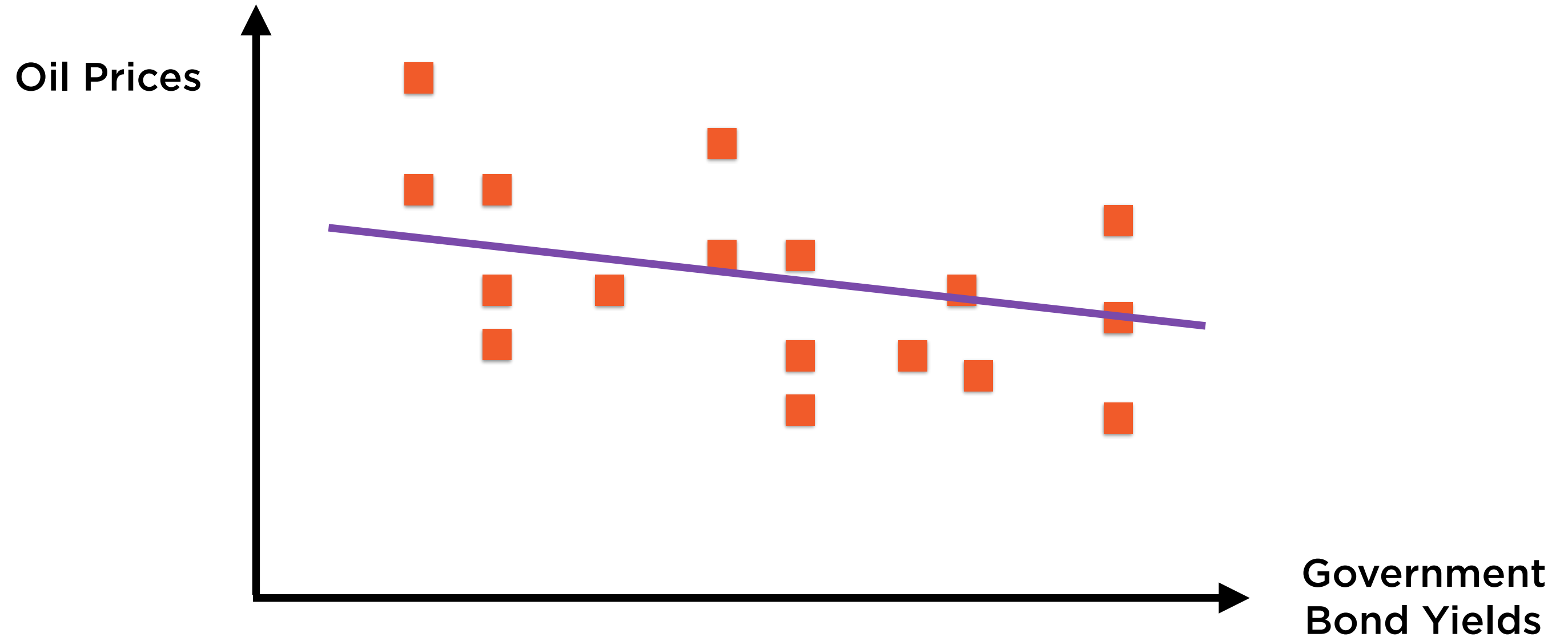
**Independent variable**



**Effect**

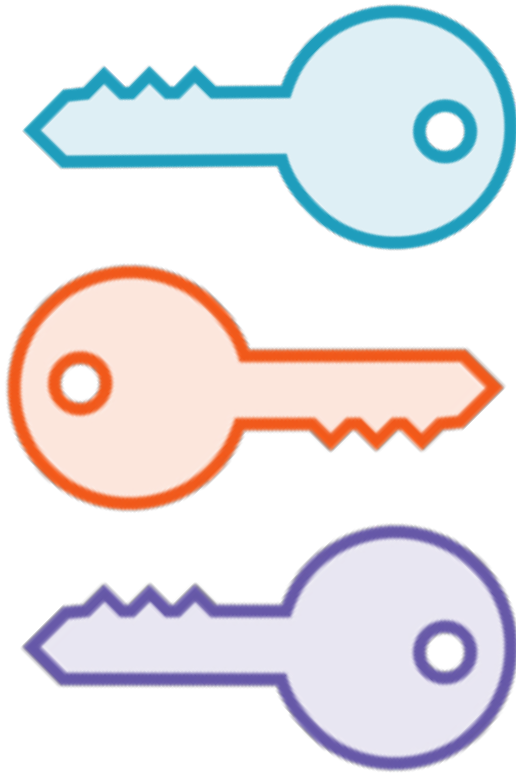
**Dependent variable**

# Simple Regression



One cause, one effect

# Multiple Regression



**Causes**

**Independent variables**

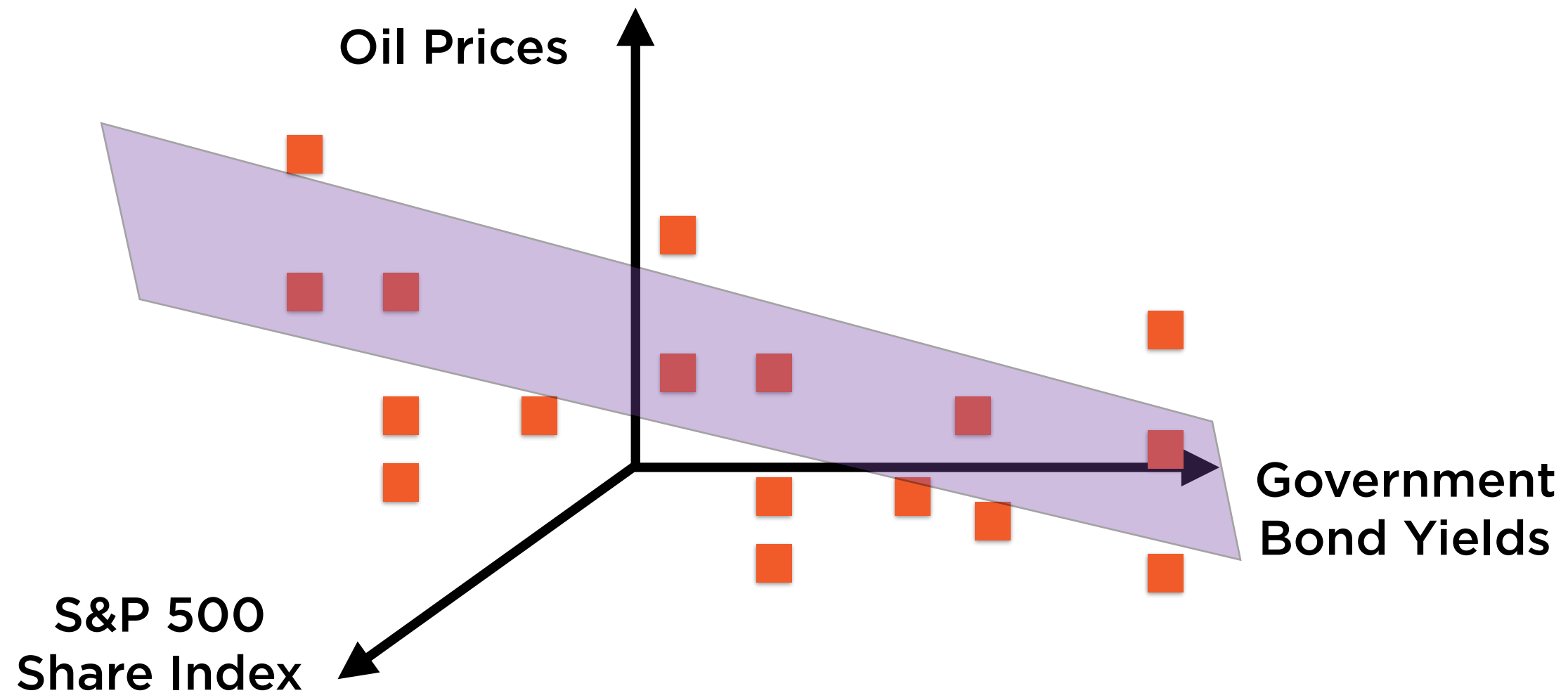


**Effect**

**Dependent variable**

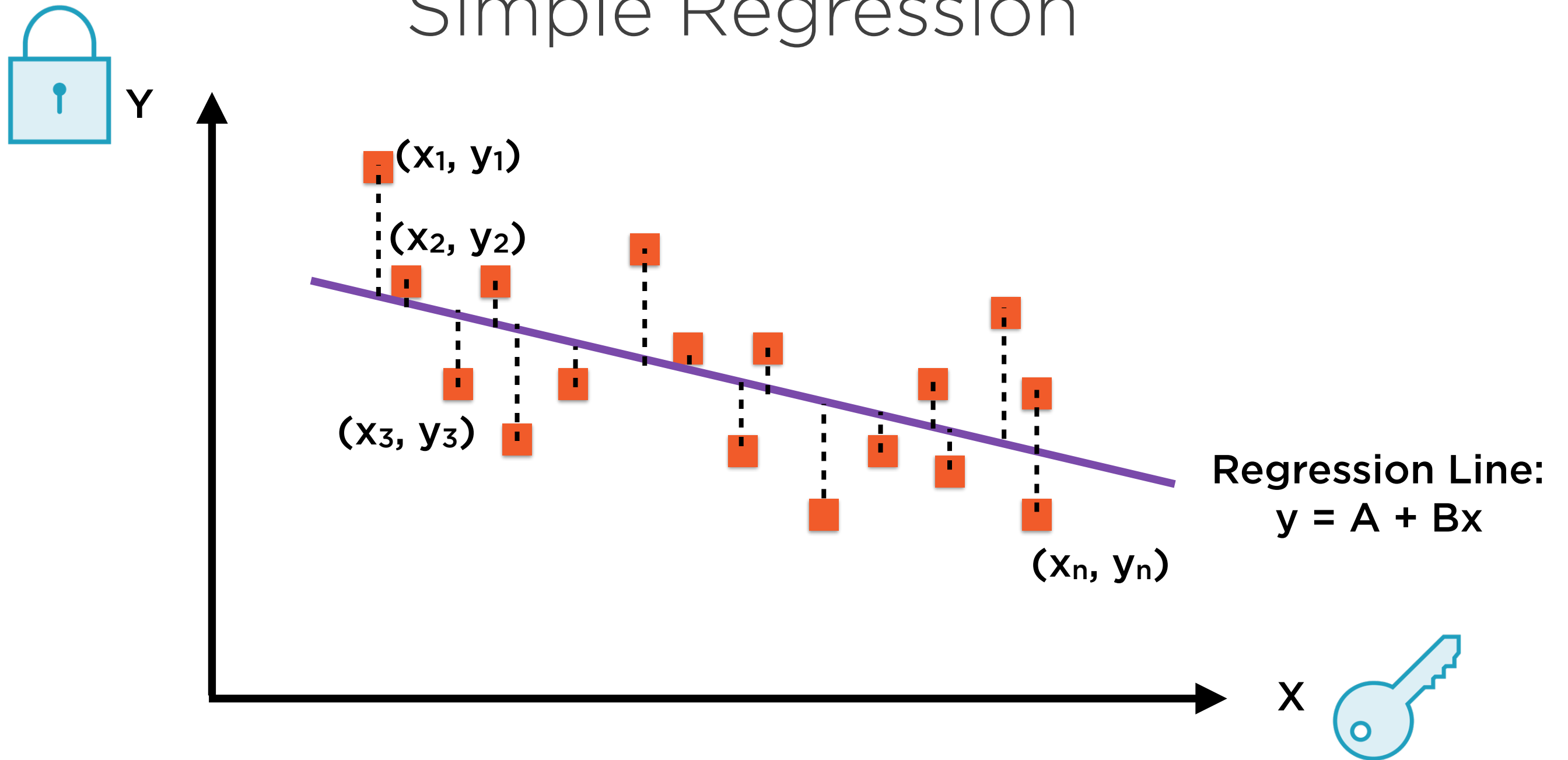


# Multiple Regression



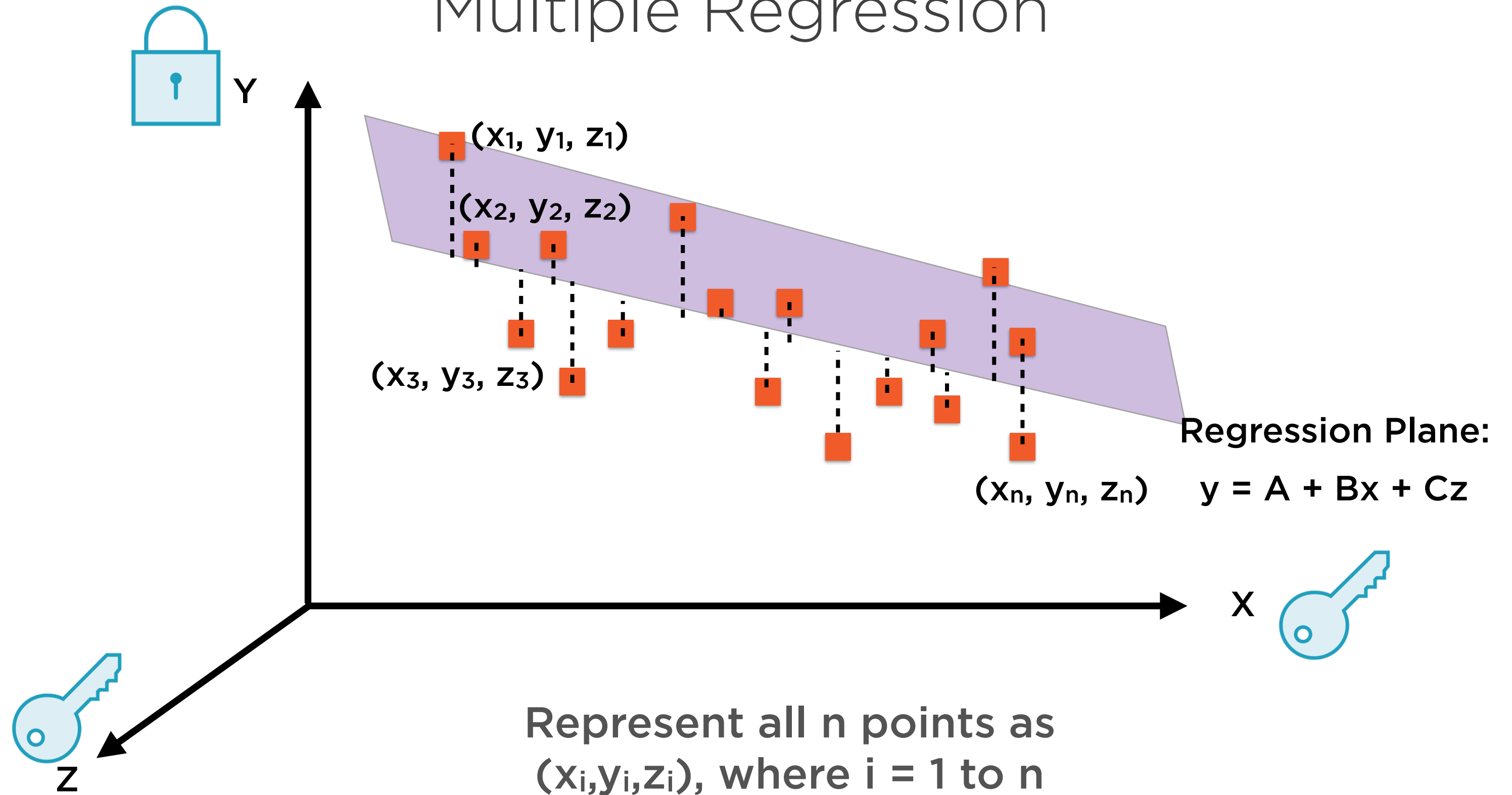
Many causes, one effect

# Simple Regression

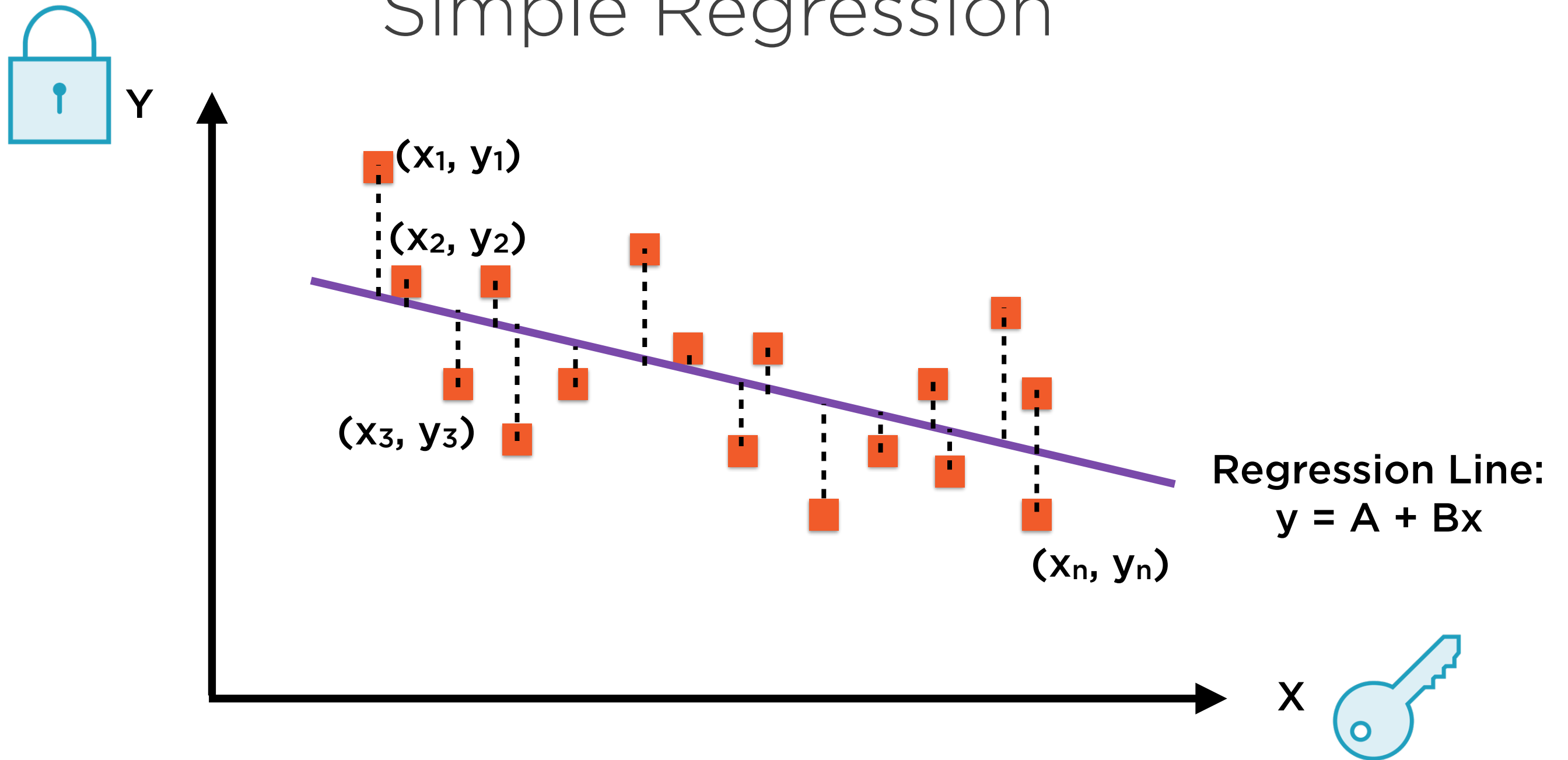


Represent all  $n$  points as  $(x_i, y_i)$ , where  $i = 1$  to  $n$

# Multiple Regression



# Simple Regression



Represent all  $n$  points as  
 $(x_i, y_i)$ , where  $i = 1$  to  $n$

# Simple Regression

**Regression Equation:**

$$y = A + Bx$$

$$y_1 = A + Bx_1$$

$$y_2 = A + Bx_2$$

$$y_3 = A + Bx_3$$

...

...

$$y_n = A + Bx_n$$

# Simple Regression

## Regression Equation:

$$y = A + Bx$$

$$y_1 = A + Bx_1 + e_1$$

$$y_2 = A + Bx_2 + e_2$$

$$y_3 = A + Bx_3 + e_3$$

...

...

$$y_n = A + Bx_n + e_n$$

# Simple Regression

**Regression Equation:**

$$y = A + Bx$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_n \end{bmatrix} = A \begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} + B \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \dots \\ e_n \end{bmatrix}$$

# Simple Regression

**Regression Equation:**

$$\text{EXXON}_t = A + B \text{ DOW}_t$$

$$\begin{bmatrix} E_1 \\ E_2 \\ E_3 \\ \dots \\ E_n \end{bmatrix} = A \begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} + B \begin{bmatrix} D_1 \\ D_2 \\ D_3 \\ \dots \\ D_n \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \dots \\ e_n \end{bmatrix}$$

$E_i$  = % return  
on Exxon stock  
on day  $i$

$D_i$  = % return of  
Dow Jones  
index on day  $i$



# Simple Regression



**Cause**

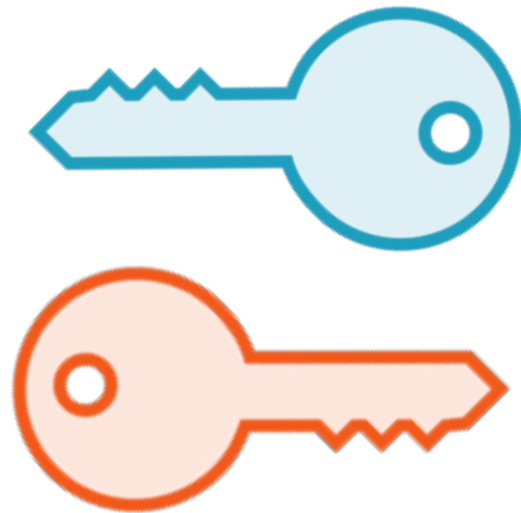
Changes in Dow Jones equity  
index



**Effect**

Changes in price of Exxon Stock

# Multiple Regression



## Causes

Dow Jones index,  
price of oil



## Effect

Exxon stock

# Multiple Regression

## Regression Equation:

$$\text{EXXON}_t = A + B \text{ DOW}_t + C \text{ OIL}_t$$

$$\begin{bmatrix} E_1 \\ E_2 \\ E_3 \\ \dots \\ E_n \end{bmatrix} = A \begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} + B \begin{bmatrix} D_1 \\ D_2 \\ D_3 \\ \dots \\ D_n \end{bmatrix} + C \begin{bmatrix} O_1 \\ O_2 \\ O_3 \\ \dots \\ O_n \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \dots \\ e_n \end{bmatrix}$$

$E_i$  = % return  
on Exxon stock  
on day  $i$

$D_i$  = % return of  
Dow Jones  
index on day  $i$

$O_i$  = % change  
in price of oil  
on day  $i$

# Multiple Regression

## Regression Equation:

$$y = A + Bx + Cz$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_n \end{bmatrix} = A \begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} + B \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ \dots \\ X_n \end{bmatrix} + C \begin{bmatrix} Z_1 \\ Z_2 \\ Z_3 \\ \dots \\ Z_n \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \dots \\ e_n \end{bmatrix}$$

# Multiple Regression

## Regression Equation:

$$y = A + Bx + Cz$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & z_1 \\ 1 & x_2 & z_2 \\ 1 & x_3 & z_3 \\ \dots & \dots & \dots \\ 1 & x_n & z_n \end{bmatrix} * \begin{bmatrix} A \\ B \\ C \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \dots \\ e_n \end{bmatrix}$$

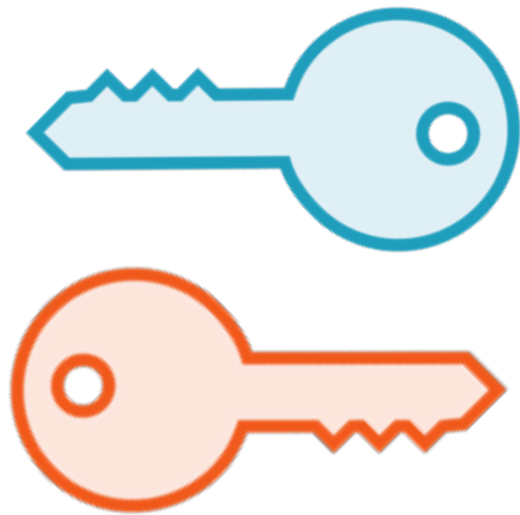
n Rows,  
1 Column

n Rows,  
3 Columns

3 Rows,  
1 Column

n Rows,  
1 Column

# Multiple Regression



**2 Causes**

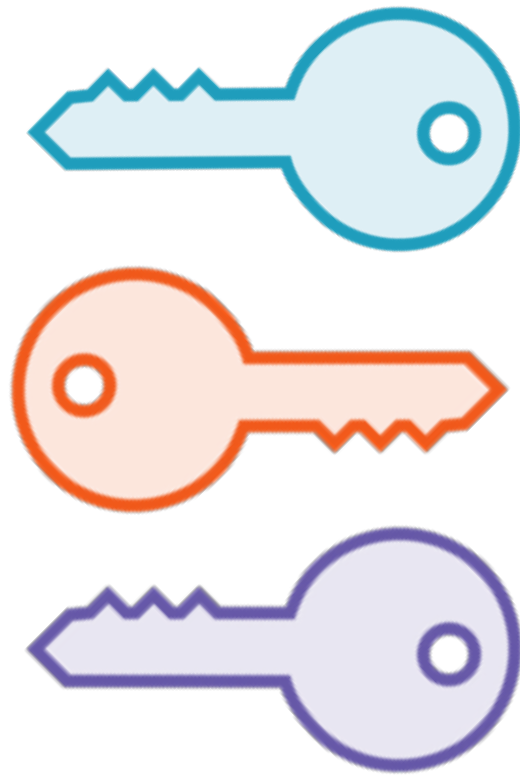
Dow Jones index,  
price of oil



**1 Effect**

Exxon stock

# Multiple Regression



**k Causes**

Dow Jones index,  
price of oil, bond yields...



**1 Effect**

Exxon stock

# Multiple Regression

## Regression Equation:

$$y = C_1 + C_2X_1 + \dots + C_{k+1}X_k$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_n \end{bmatrix} = C_1 \begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} + C_2 \begin{bmatrix} X_{11} \\ X_{21} \\ X_{31} \\ \dots \\ X_{n1} \end{bmatrix} + \dots + C_{k+1} \begin{bmatrix} X_{1k} \\ X_{2k} \\ X_{3k} \\ \dots \\ X_{nk} \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \dots \\ e_n \end{bmatrix}$$



# Multiple Regression

## Regression Equation:

$$y = C_1 + C_2X_1 + \dots + C_{k+1}X_k$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & \dots & X_{1k} \\ 1 & X_{21} & \dots & X_{2k} \\ 1 & X_{31} & \dots & X_{3k} \\ \dots & \dots & \dots & \dots \\ 1 & X_{n1} & \dots & X_{nk} \end{bmatrix} * \begin{bmatrix} C_1 \\ C_2 \\ \dots \\ C_{k+1} \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \dots \\ e_n \end{bmatrix}$$

n Rows,  
1 Column

n Rows,  
k Columns

k Rows,  
1 Column

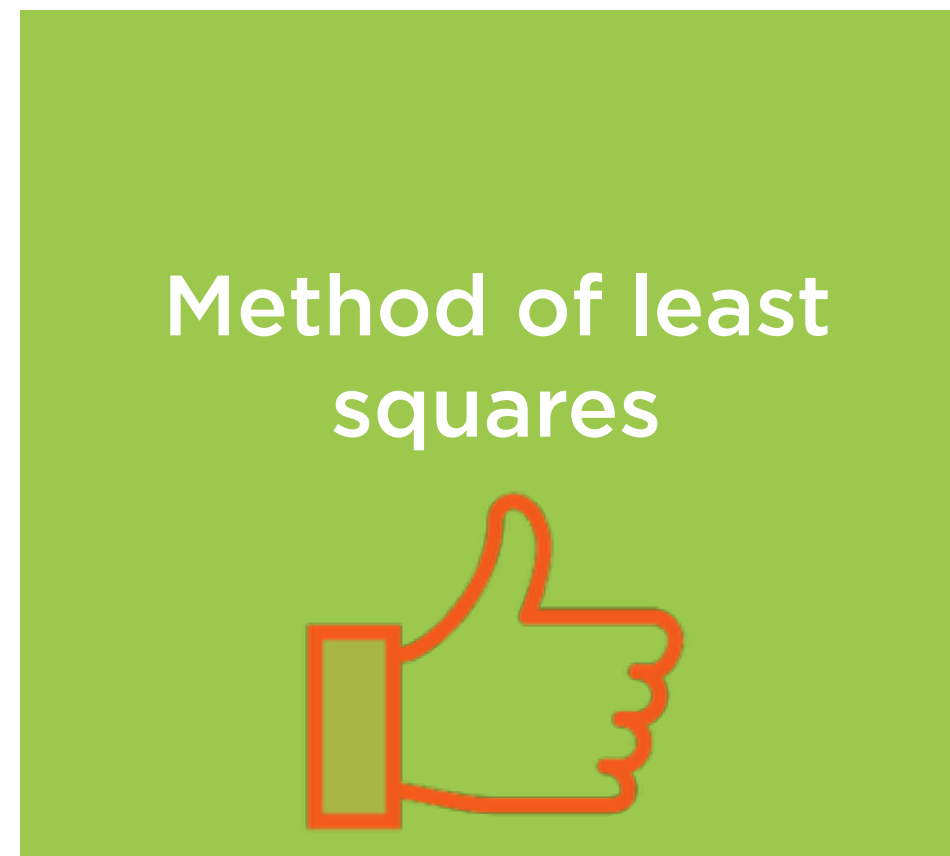
# Multiple Regression

## **Regression Equation:**

$$y = C_1 + C_2X_1 + \dots + C_{k+1}X_k$$

Multiple regression involves finding  $k+1$  coefficients,  $k$  for the explanatory variables, and 1 for the intercept

# Estimation Methods in Multiple Regression



**The method of least squares works for multiple regression too**

**Regression Equation:**

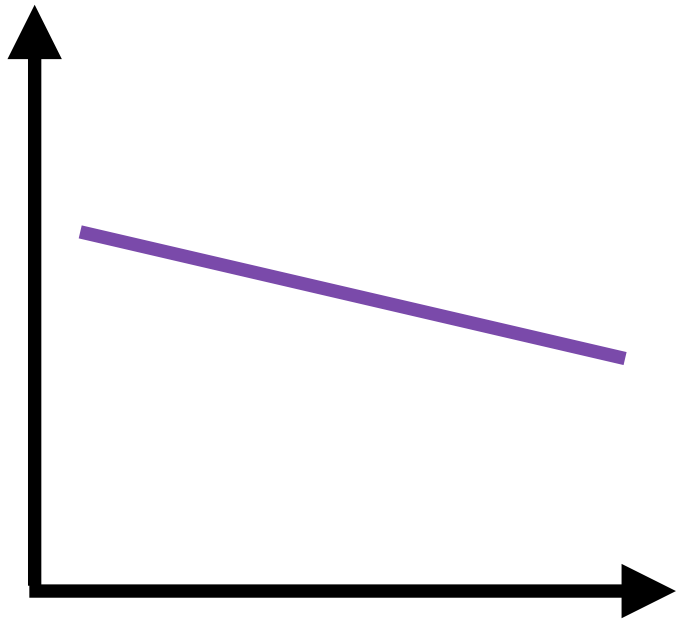
$$y = C_1 + C_2X_1 + \dots + C_{k+1}X_k$$

The “best fit” line is the one where the sum of the squares of the lengths of the errors is minimised

# Risks in Multiple Regression

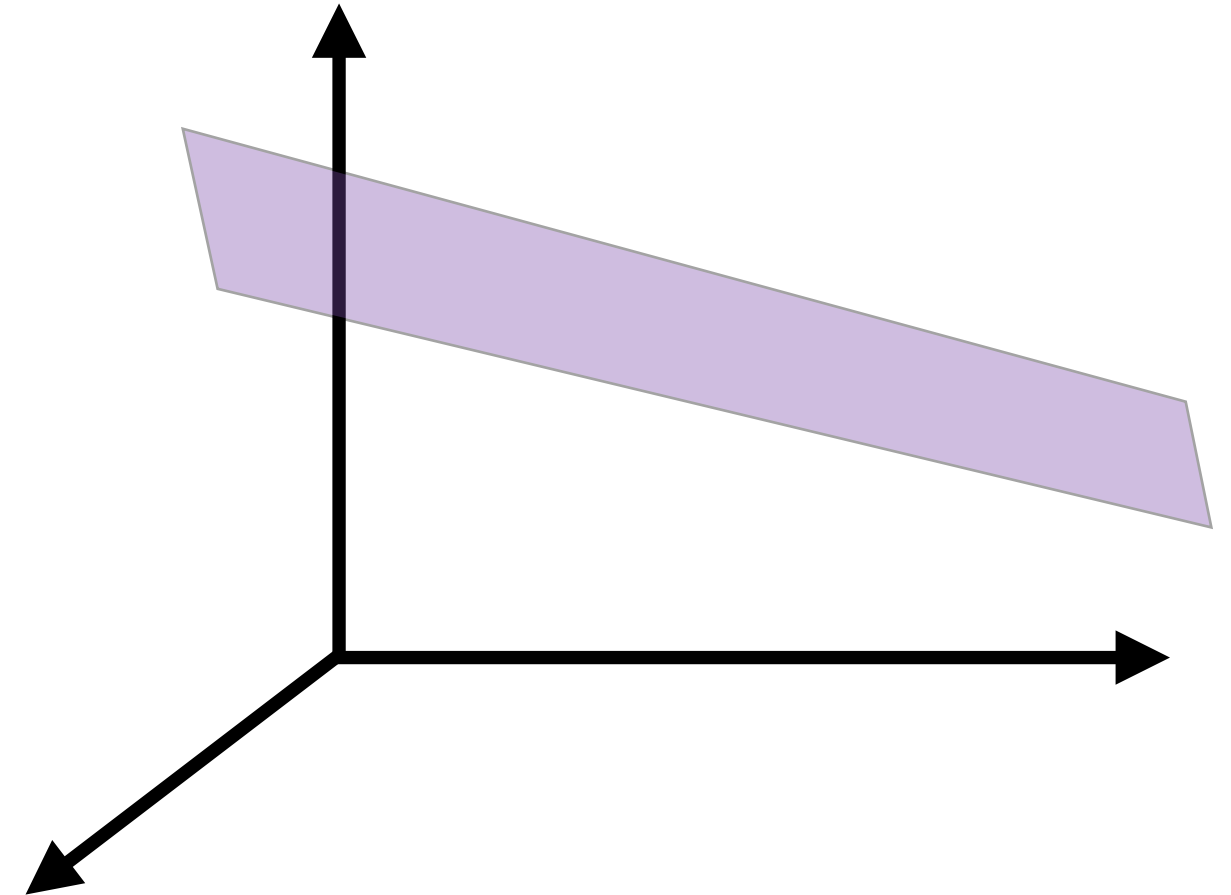
---

# Simple and Multiple Regression



**Simple Regression**

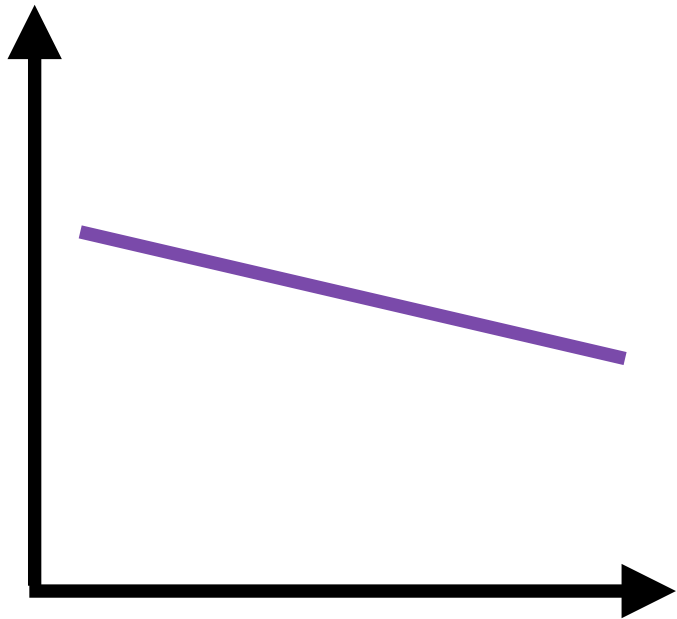
Data in 2 dimensions



**Multiple Regression**

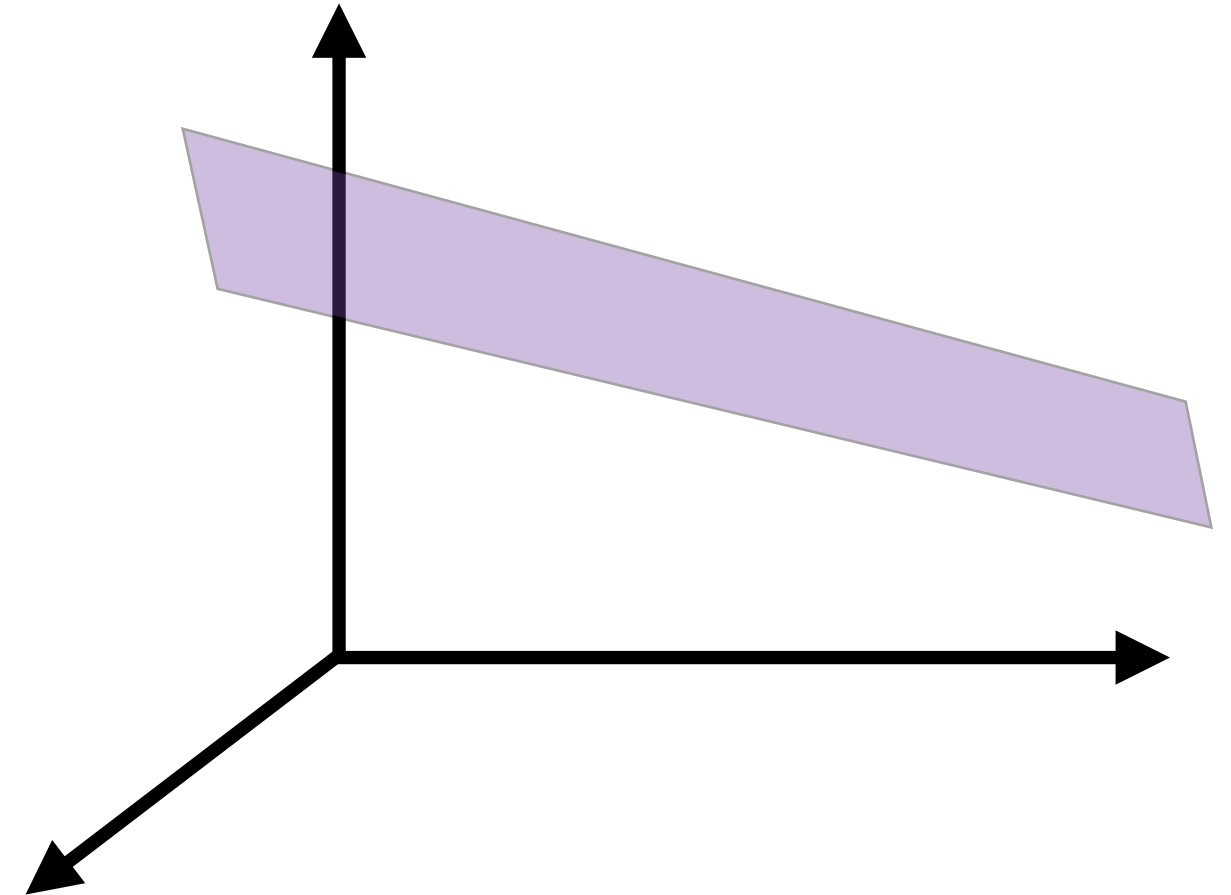
Data in  $> 2$  dimensions

# Simple and Multiple Regression



## Simple Regression

Risks exist, but can usually be mitigated analysing  $R^2$  and residuals



## Multiple Regression

Risks are more complicated, require interpreting regression statistics

# Risks in Simple Regression

**No cause-effect  
relationship**

Regression on completely  
unrelated data series

**Mis-specified  
relationship**

Non-linear (exponential  
or polynomial) fit

**Incomplete  
relationship**

Multiple causes exist, we  
have captured just one



# Diagnosing Risks in Simple Regression

**No cause-effect  
relationship**

low  $R^2$ , plot of  $X \sim Y$  has  
no pattern

**Mis-specified  
relationship**

high  $R^2$ , residuals are not  
independent of each  
other

**Incomplete  
relationship**

low  $R^2$ , residuals are not  
independent of  $x$

# Mitigating Risks in Simple Regression

## No cause-effect relationship

Wrong choice of X and Y  
- back to drawing board

## Mis-specified relationship

Transform X and Y -  
convert to logs or returns

## Incomplete relationship

Add X variables (move to  
multiple regression)

The big new risk with multiple regression is **multicollinearity**:  $X$  variables containing the same information

“If everyone is thinking alike, then  
somebody isn't thinking.”

**General Patton**

# Multiple Regression

## Regression Equation:

$$y = C_1 + C_2X_1 + \dots + C_kX_{k-1}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & X_{11} & & X_{1k-1} \\ 1 & X_{21} & & X_{2k-1} \\ 1 & X_{31} & \dots & X_{3k-1} \\ \dots & \dots & & \dots \\ 1 & X_{n1} & & X_{nk-1} \end{bmatrix}_{n \times k} * \begin{bmatrix} C_1 \\ C_2 \\ \dots \\ C_k \end{bmatrix}_{k \times 1}$$

n Rows,  
1 Column

n Rows,  
k Columns

k Rows,  
1 Column

# Multiple Regression

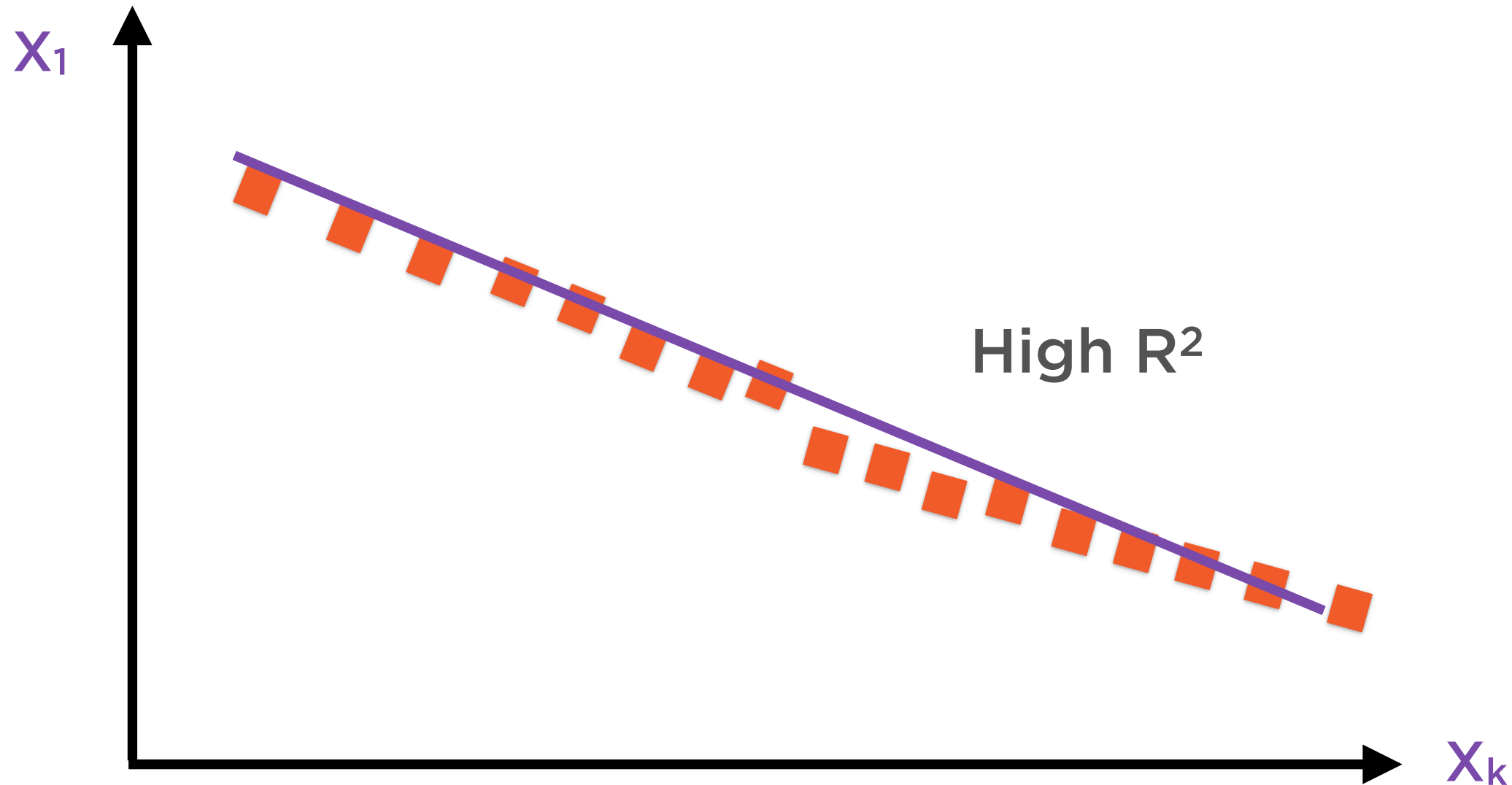
## Regression Equation:

$$y = C_1 + C_2X_1 + \dots + C_kX_{k-1}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & \boxed{\begin{matrix} X_{11} \\ X_{21} \\ X_{31} \\ \dots \\ X_{n1} \end{matrix}} & \dots & \boxed{\begin{matrix} X_{1k-1} \\ X_{2k-1} \\ X_{3k-1} \\ \dots \\ X_{nk-1} \end{matrix}} \\ \dots & & & \\ 1 & & & \end{bmatrix} * \begin{bmatrix} C_1 \\ C_2 \\ \dots \\ C_k \end{bmatrix}$$

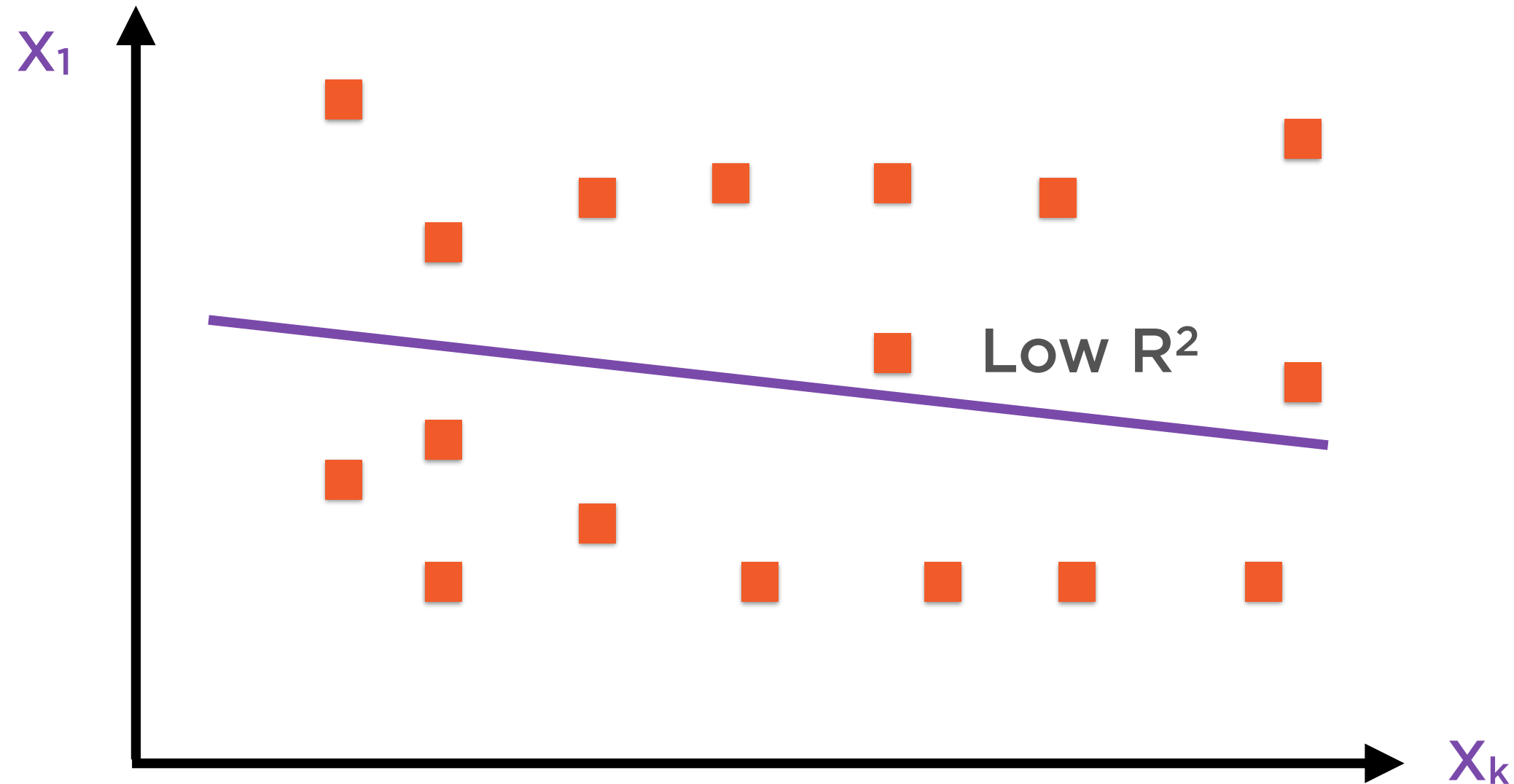
$X_1 \qquad X_k$

# Bad News: Multicollinearity Detected



Highly correlated explanatory variables

# Good News: No Multicollinearity Detected



Uncorrelated explanatory variables



# Multiple Regression

## Regression Equation:

$$\text{EXXON}_t = A + B \text{ DOW}_t + C \text{ OIL}_t$$

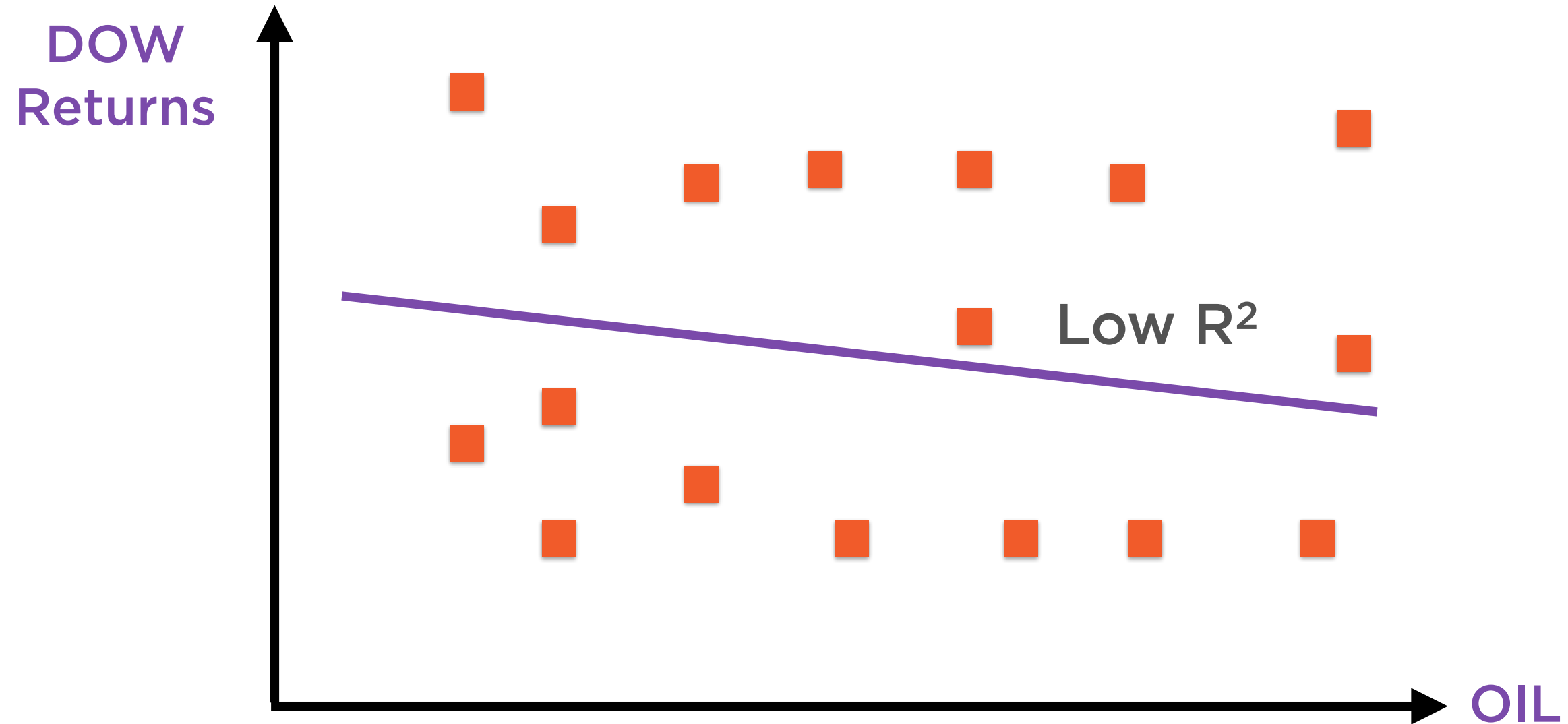
$$\begin{bmatrix} E_1 \\ E_2 \\ E_3 \\ \dots \\ E_n \end{bmatrix} = A \begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} + B \begin{bmatrix} D_1 \\ D_2 \\ D_3 \\ \dots \\ D_n \end{bmatrix} + C \begin{bmatrix} O_1 \\ O_2 \\ O_3 \\ \dots \\ O_n \end{bmatrix}$$

$E_i$  = % return  
on Exxon stock  
on day  $i$

$D_i$  = % return of  
Dow Jones  
index on day  $i$

$O_i$  = % change  
in price of oil  
on day  $i$

# Good News: No Multicollinearity Detected



Uncorrelated explanatory variables

# Multiple Regression

## Regression Equation:

$$\text{EXXON}_t = A + B \text{ DOW}_t + C \text{ NASDAQ}_t$$

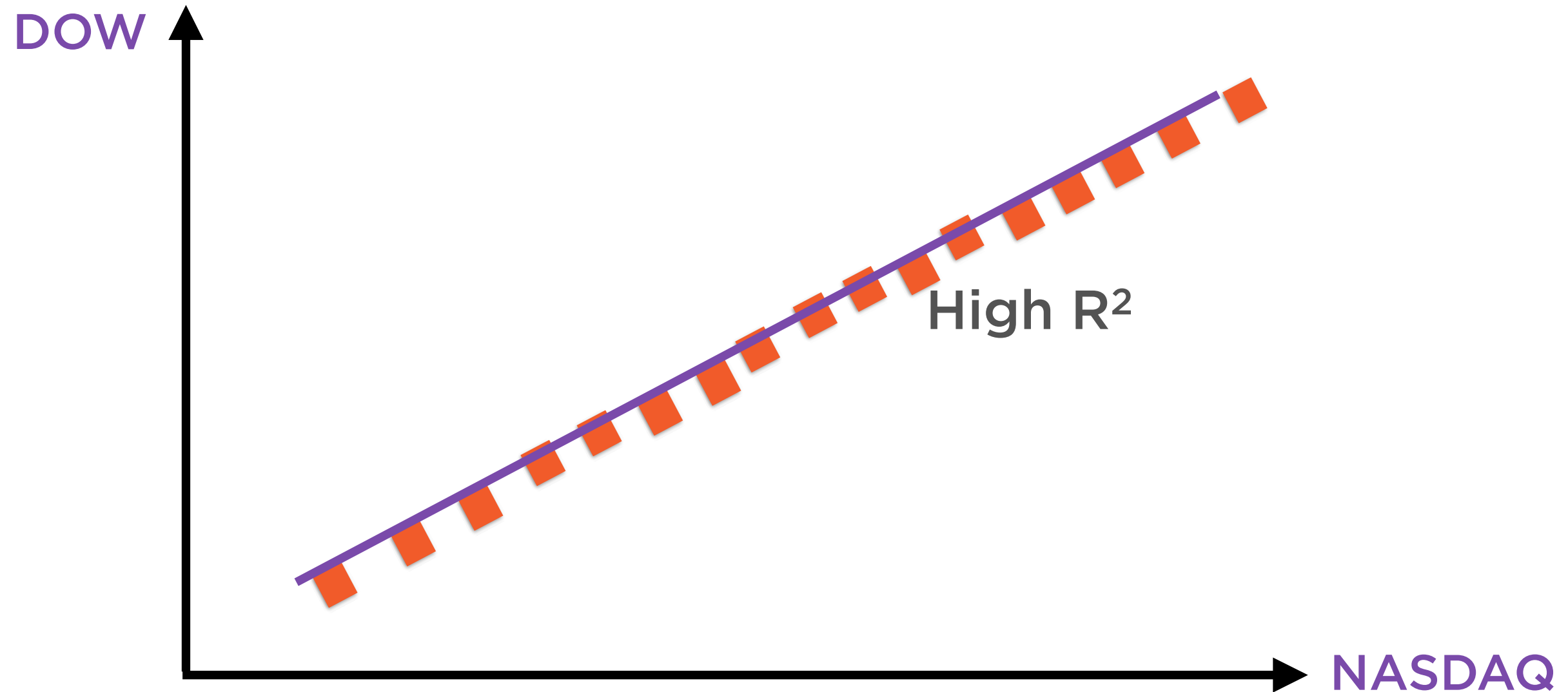
$$\begin{bmatrix} E_1 \\ E_2 \\ E_3 \\ \dots \\ E_n \end{bmatrix} = A \begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} + B \begin{bmatrix} D_1 \\ D_2 \\ D_3 \\ \dots \\ D_n \end{bmatrix} + C \begin{bmatrix} N_1 \\ N_2 \\ N_3 \\ \dots \\ N_n \end{bmatrix}$$

$E_i$  = % return  
on Exxon stock  
on day  $i$

$D_i$  = % return of  
Dow Jones  
index on day  $i$

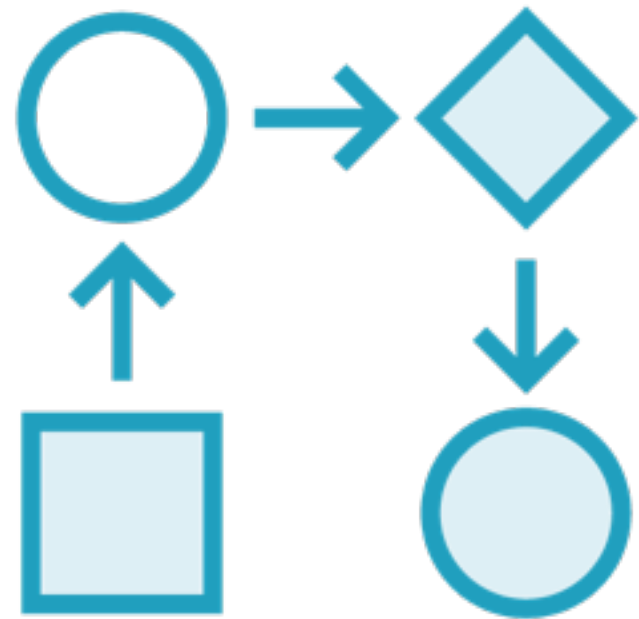
$N_i$  = % return of  
NASDAQ index  
on day  $i$

# Bad News: Multicollinearity Detected



Highly correlated explanatory variables

# Multicollinearity Kills Regression's Usefulness



## Explaining Variance

The  $R^2$  as well as the regression coefficients are not very reliable



## Making Predictions

The regression model will perform poorly with out-of-sample data

# Multicollinearity: Prevention and Cure



## **Common Sense**

Big-picture  
understanding of the  
data



## **Nuts and Bolts**

Setting up data right



## **Heavy Lifting**

Factor analysis,  
Principal components  
analysis (PCA)



Common  
Sense

**Think deeply about each x variable**

**Eliminate closely related ones**

**Dig down to underlying causes**

# Multiple Regression

**Proposed Regression Equation:**

$$\text{EXXON}_t = A + B \text{ DOW}_t + C \text{ NASDAQ}_t + D \text{ OIL}_t$$

% return on  
Exxon stock on  
day i

% return of  
Dow Jones  
index on day i

% return of  
NASDAQ index  
on day i

% change in  
price of oil on  
day i



# Common Sense

## Proposed Regression Equation:

$$\text{EXXON}_t = A + B \text{DOW}_t + C \text{NASDAQ}_t + D \text{OIL}_t$$

Dow Jones  
Industrial Average

30 Large-cap US stocks

NASDAQ 100 Index

100 large tech stocks

Oil Prices

Price of barrel of oil

# Common Sense

## Proposed Regression Equation:

$$\text{EXXON}_t = A + B \text{DOW}_t + C \text{NASDAQ}_t + D \text{OIL}_t$$

Dow Jones  
Industrial Average  
30 Large-cap US stocks

NASDAQ 100 Index  
100 large tech stocks

Oil Prices  
Price of barrel of oil

Do we really need both Dow and NASDAQ returns as  
explanatory variables?

# Common Sense

## Proposed Regression Equation:

$$\text{EXXON}_t = A + B \text{DOW}_t + C \text{NASDAQ}_t + D \text{OIL}_t$$

Dow Jones  
Industrial Average

30 Large-cap US stocks

NASDAQ 100 Index

100 large tech stocks

Oil Prices

Price of barrel of oil

If yes - consider keeping one, and constructing a  
new explanatory variable of their difference

Common Sense

**Proposed Regression Equation:**

$$\text{EXXON}_t = A + B \text{ DOW}_t + C \text{ OIL}_t$$

**Dow Jones  
Industrial Average**

30 Large-cap US stocks

**Oil Prices**

Price of barrel of oil

What underlying factors drive both US large-cap stocks and the price of oil?

# Common Sense

**GDP growth**

**Interest rates**

**US dollar strength**

**Seasonality**

**What underlying factors drive both US large-cap  
stocks and the price of oil?**

# Common Sense

**GDP growth**

**Interest rates**

**US dollar strength**

**Seasonality**

**What underlying factors drive both US large-cap  
stocks and the price of oil?**

# Multiple Regression

Original Regression Equation:

$$\text{EXXON}_t = A + B \text{DOW}_t + C \text{NASDAQ}_t + D \text{OIL}_t$$

**Revised Regression Equation:**

$$\text{EXXON}_t = A + B \text{DOW}_t + C \text{INTEREST}_t + D \text{GDP}_t$$



## Nuts and Bolts

**‘Standardise’ the variables**

**Rely on adjusted- $R^2$ , not plain  $R^2$**

**Set up dummy variables right**

**Distribute lags**





Nuts and Bolts

**‘Stepwise regression’ - use with care**

**Automated selection of x variables**

**Variants**

- Backward elimination
- Forward selection
- Iterative Elimination



Heavy Lifting

**Find underlying factors that drive the correlated x variables**

**Principal Component Analysis (PCA) is a great tool**

# Multiple Regression

**Proposed Regression Equation:**

$$\text{HOME}_t = A + B \text{ 5-yr}_t + C \text{ 10-year}_t + D \text{ 2-year}_t + \\ E \text{ 1-year}_t + F \text{ 3-month}_t + G \text{ 1-day}_t + \dots$$

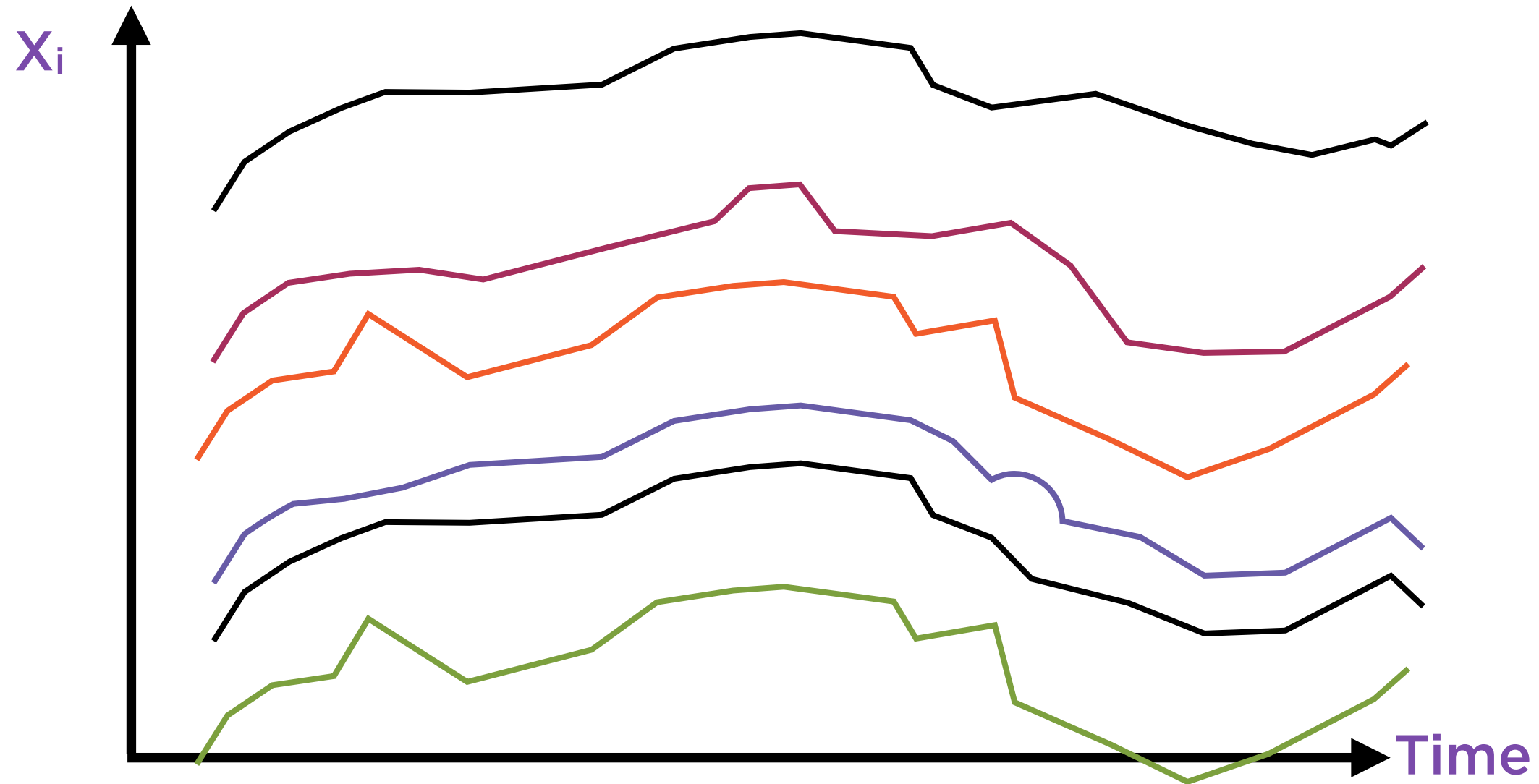
% change in  
home prices in  
month i

yield on 5-  
year bond in  
month i

yield on 10-  
year bond in  
month i

...

# Bad News: Multicollinearity Detected



Highly correlated explanatory variables

# Factor Analysis

**3-month  
government bonds**

**1-year government  
bonds**

**5-year government  
bonds**

**1-day (overnight)  
money market**

**30-year government  
bonds**

**5-year swap rate  
(inter-bank)**

**Interest rates on a wide variety of fixed-income  
instruments**

# Factor Analysis on Interest Rates

**Level**

How high are interest rates?

**Slope**

How steep is the yield curve?

**Twist**

How convex is the yield curve?

**Three uncorrelated factors explain most variation in all interest rates**



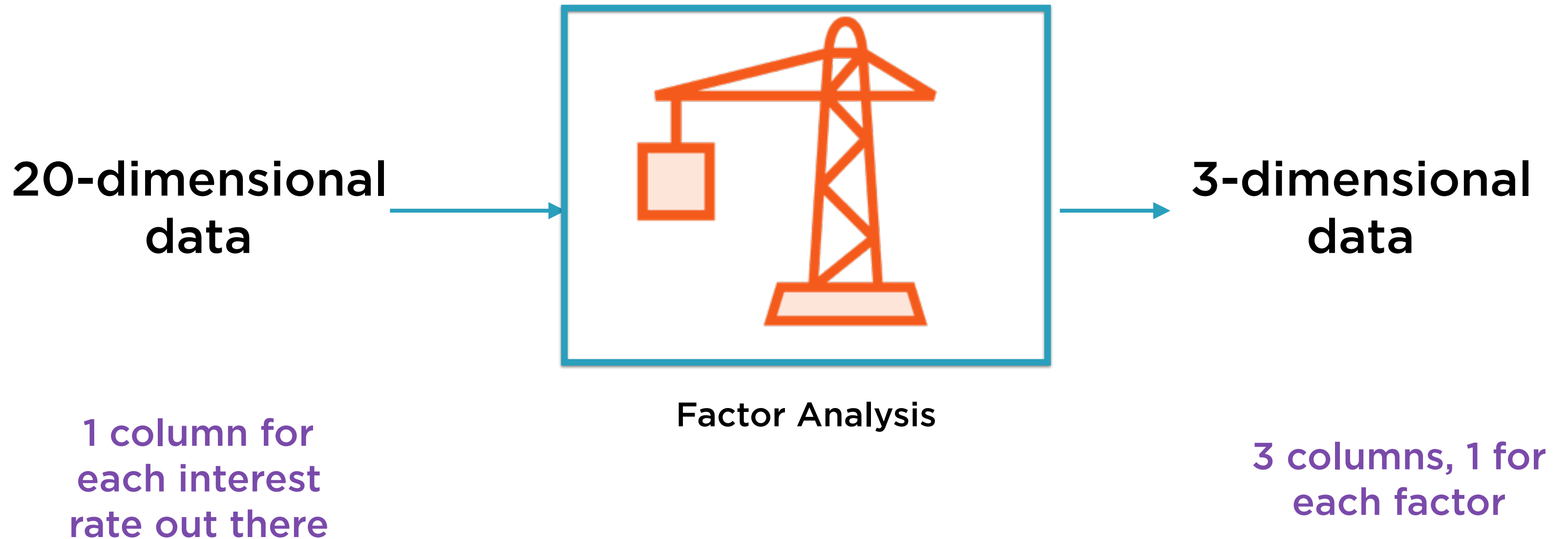
## Factor Analysis

**The factors identified are guaranteed to be uncorrelated**

**However, they may not have an intuitive interpretation**

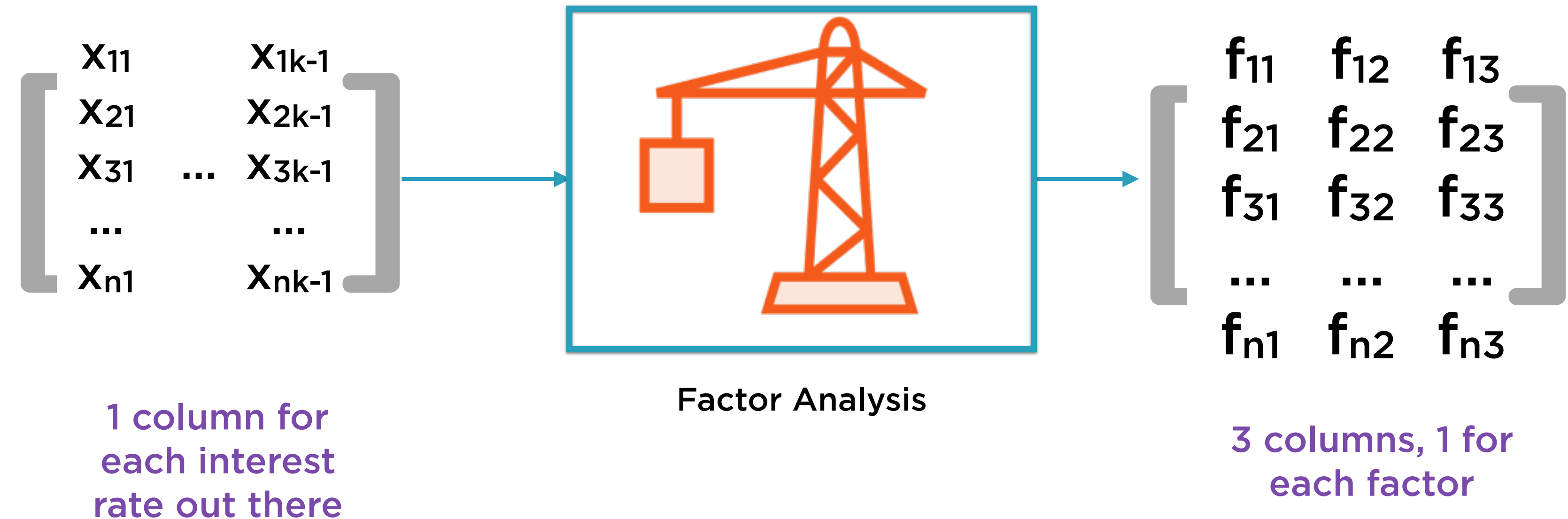
**Principal Component Analysis is one procedure for factor analysis**

# Dimensionality Reduction via Factor Analysis





# Dimensionality Reduction via Factor Analysis

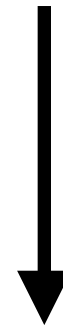


Factor Analysis is a dimensionality-reduction technique to identify a few underlying causes in data

# Multiple Regression

Proposed Regression Equation:

$$\text{HOME}_t = A + B \text{ 5-yr}_t + C \text{ 10-year}_t + D \text{ 2-year}_t + \\ E \text{ 1-year}_t + F \text{ 3-month}_t + G \text{ 1-day}_t + \dots$$



Principal  
Component  
Analysis

Revised Regression Equation:

$$\text{HOME}_t = A + B \text{ LEVEL}_t + C \text{ SLOPE}_t + D \text{ TWIST}_t$$

# Benefits of Multiple Regression

---

# Simple Regression Is a Great Tool

## Powerful

Perfectly suited to two  
common use-cases

## Versatile

Easily extended to non-  
linear relationships

## Deep

The first “crossover hit”  
from Machine Learning

# Multiple Regression Is Even Better

## Powerful

Also controls for effects  
different causes

## Versatile

Also works with  
categorical data

## Deep

Especially if combined  
with factor analysis

# Multiple Regression Is Even Better

## Powerful

Also controls for effects  
different causes

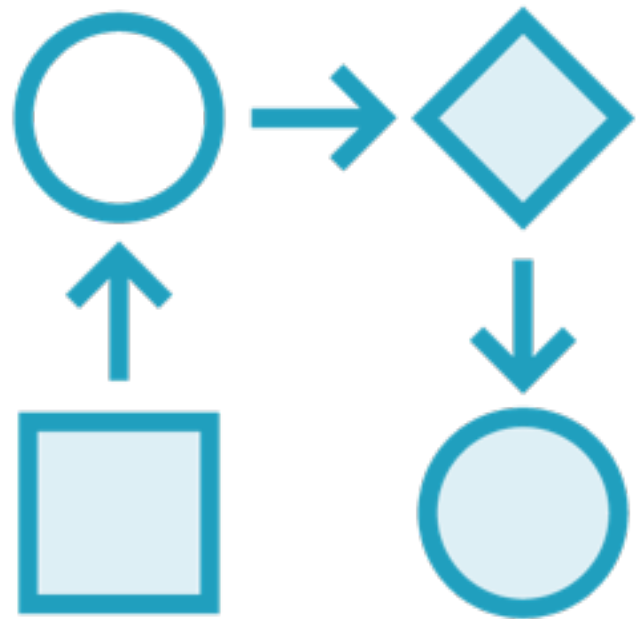
## Versatile

Also works with  
categorical data

## Deep

Especially if combined  
with factor analysis

# Two Common Applications of Regression



## Explaining Variance

How much variation in one data series is caused by another?



## Making Predictions

How much does a move in one series impact another?



# Controlling For Different Causes

**Proposed Regression Equation:**

$$\text{EXXON}_t = A + B \text{ DOW}_t + C \text{ OIL}_t$$

% return on  
Exxon stock on  
day i

% return of  
Dow Jones  
index on day i

% change in  
price of oil on  
day i

**All else being equal, how much will Exxon stock move by if oil prices increase by 1%?**

# Controlling For Different Causes

**Proposed Regression Equation:**

$$\text{EXXON}_t = A + B \text{DOW}_t + C \text{OIL}_t$$

**All else being equal**, how much will Exxon stock move by if oil prices increase by 1%?

# Controlling For Different Causes

$$\text{EXXON}_t = A + B \text{ DOW}_t + C \text{ OIL}_t$$

$$\text{EXXON}_? = A + B \text{ DOW}_t + C (\text{OIL}_t + 1\%)$$

**All else being equal, how much will Exxon stock move by if oil prices increase by 1%?**

# Controlling For Different Causes

$$\text{EXXON}_t = A + B \text{ DOW}_t + C \text{ OIL}_t$$

$$\text{EXXON}_? = A + B \text{ DOW}_t + C (\text{OIL}_t + 1\%)$$

$$\text{Change in EXXON} = \text{EXXON}_t - \text{EXXON}_?$$

All else being equal, how much will Exxon stock move by if oil prices increase by 1%?

# Controlling For Different Causes

$$\begin{aligned} \text{EXXON}_t &= A + B \text{DOW}_t + C \text{OIL}_t \\ - \quad \text{EXXON}_? &= A + B \text{DOW}_t + C (\text{OIL}_t + 1\%) \end{aligned}$$

---

All else being equal, how much will Exxon stock move by if oil prices increase by 1%?

# Controlling For Different Causes

$$\text{EXXON}_t = A + B \text{DOW}_t + C \text{OIL}_t$$

$$- \text{EXXON}_? = A + B \text{DOW}_t + C (\text{OIL}_t + 1\%)$$

---

$$\text{Change in EXXON} = C$$

All else being equal, how much will Exxon stock move by if oil prices increase by 1%?

Regression coefficients tell how much  $y$  changes for a unit change in each predictor, **all others being held constant**

# Multiple Regression Is Even Better

## Powerful

Also controls for effects  
different causes

## Versatile

Also works with  
categorical data

## Deep

Especially if combined  
with factor analysis



# Multiple Regression Is Even Better

## Powerful

Also controls for effects  
different causes

## Versatile

Also works with  
categorical data

## Deep

Especially if combined  
with factor analysis

# Multiple Regression Is Even Better

## Powerful

Also controls for effects  
different causes

## Versatile

Also works with  
categorical data

## Deep

Especially if combined  
with factor analysis

# Interpreting the Results of a Regression Analysis

---

# Interpreting Results of a Simple Regression

**$R^2$**

Measures overall quality of fit - the higher the better (up to a point)

**Residuals**

Check if regression assumptions are violated

Standard errors of individual coefficients are usually of little significance

# Interpreting Results of a Multiple Regression

**Adjusted  $R^2$**

**Residuals**

**F-statistic**

**Standard Errors  
of coefficients**

**$R^2$**

$$e = y - y'$$

$$\Rightarrow y = y' + e$$

$$\Rightarrow \text{Variance}(y) = \text{Variance}(y' + e)$$

$$\Rightarrow \text{Variance}(y) = \text{Variance}(y') + \text{Variance}(e) + \text{Covariance}(y', e)$$

---

## A Not-Very-Important Intermediate Step

Variance of the dependent variable can be decomposed into variance of the regression fitted values, and that of the residuals

$$e = y' - y$$

$$\Rightarrow y = y' + e$$

$$\Rightarrow \text{Variance}(y) = \text{Variance}(y' + e)$$

$$\Rightarrow \text{Variance}(y) = \text{Variance}(y') + \text{Variance}(e) + \text{Covariance}(y', e)$$

Always = 0

---

## A Leap of Faith

This is important - more on why in a bit

$$\text{Variance}(y) = \text{Variance}(y') + \text{Variance}(e)$$

---

## Variance Explained

Variance of the dependent variable can be decomposed into variance of the regression fitted values, and that of the residuals



$$\text{Variance}(y) = \text{Variance}(y') + \text{Variance}(e)$$

---

Total Variance (*TSS*)

A measure of how volatile the dependent variable is, and of much it moves around

$$\text{TSS} = \text{Variance}(y') + \text{Variance}(e)$$

---

Explained Variance (*ESS*)

A measure of how volatile the fitted values are - these come from the regression line

$$\text{TSS} = \text{Variance}(y)$$

$$\text{TSS} = \text{ESS} + \text{Variance}(e)$$

---

## Residual Variance ( $RSS$ )

This is the variance in the dependent variable that can not be explained by the regression

$$\text{TSS} = \text{Variance}(y) \quad \text{ESS} = \text{Variance}(y')$$

$$\text{TSS} = \text{ESS} + \text{RSS}$$

---

## Variance Explained

Variance of the dependent variable can be decomposed into variance of the regression fitted values, and that of the residuals

$$\text{TSS} = \text{Variance}(y) \quad \text{ESS} = \text{Variance}(y') \quad \text{RSS} = \text{Variance}(e)$$

$$R^2 = ESS / TSS$$

---

$R^2$

The percentage of total variance explained by the regression. Usually, the higher the  $R^2$ , the better the quality of the regression (upper bound is 100%)

$$R^2 = ESS / TSS$$

---

$R^2$

**In multiple regression, adding explanatory variables always increases  $R^2$ , even if those variables are irrelevant and increase danger of multicollinearity**

**Adjusted-R<sup>2</sup> = R<sup>2</sup> x (Penalty for adding irrelevant variables)**

---

Adjusted-R<sup>2</sup>

**Increases if irrelevant\* variables are deleted**

**(\*irrelevant variables = any group whose F-ratio < 1)**



Nuts and Bolts

**‘Stepwise regression’ - use with care**

**Automated selection of x variables**

**Variants**

- Backward elimination
- Forward selection
- Iterative Elimination



# Extending Multiple Regression to Categorical Variables

---

# A Simple Regression

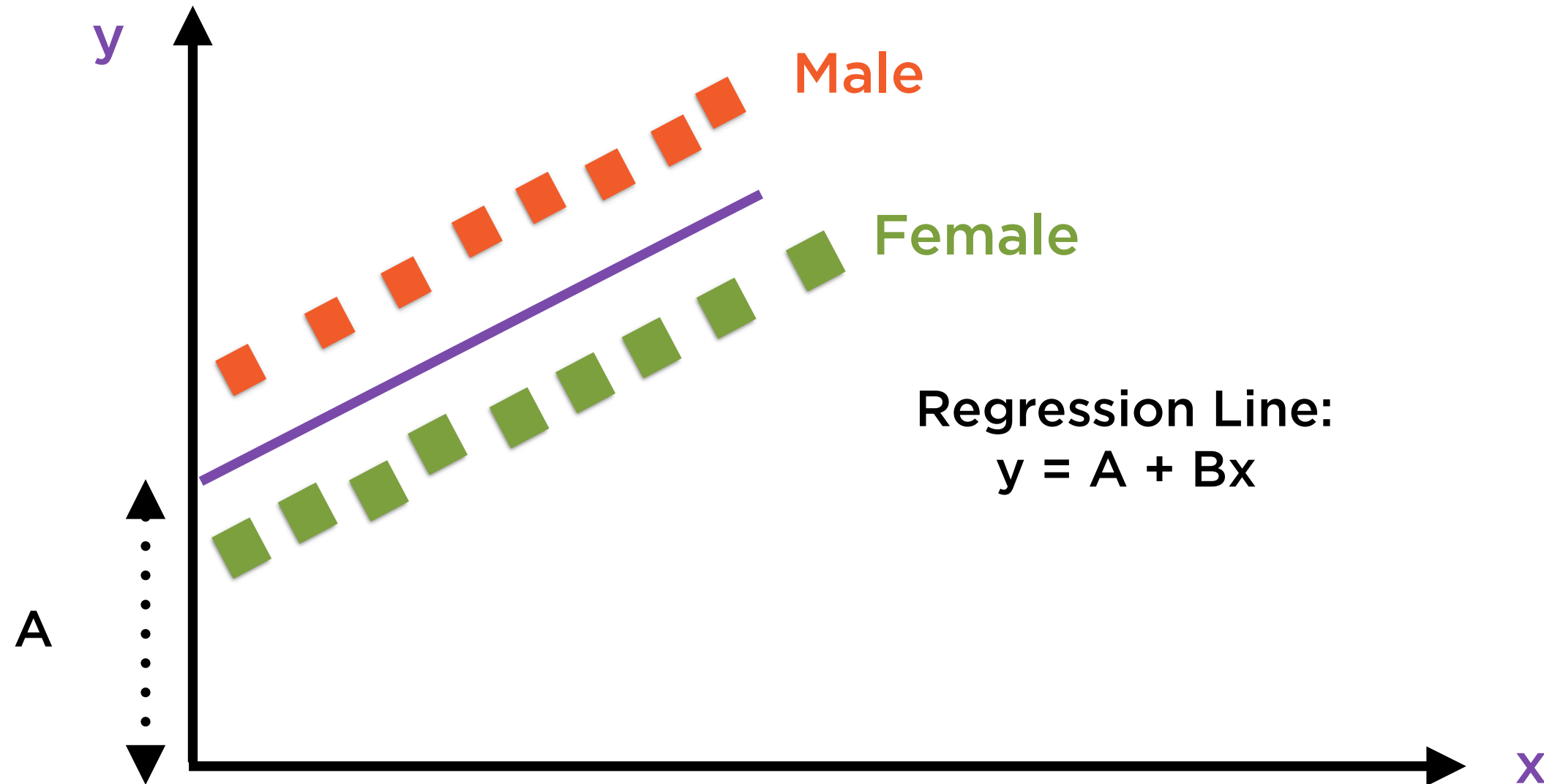
**Proposed Regression Equation:**

$$y = A + Bx$$

Height of  
individual

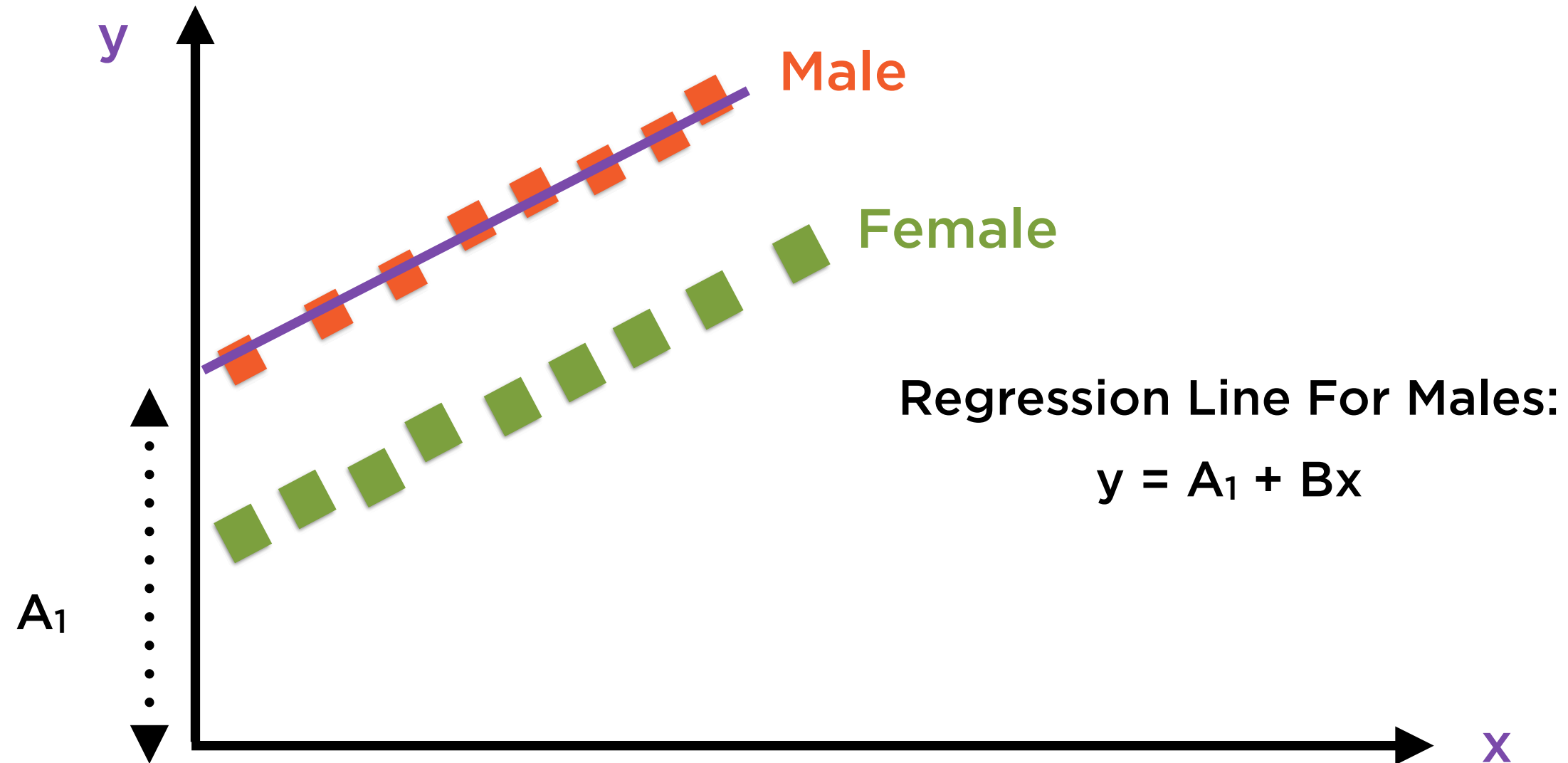
Average height  
of parents

# A Simple Regression



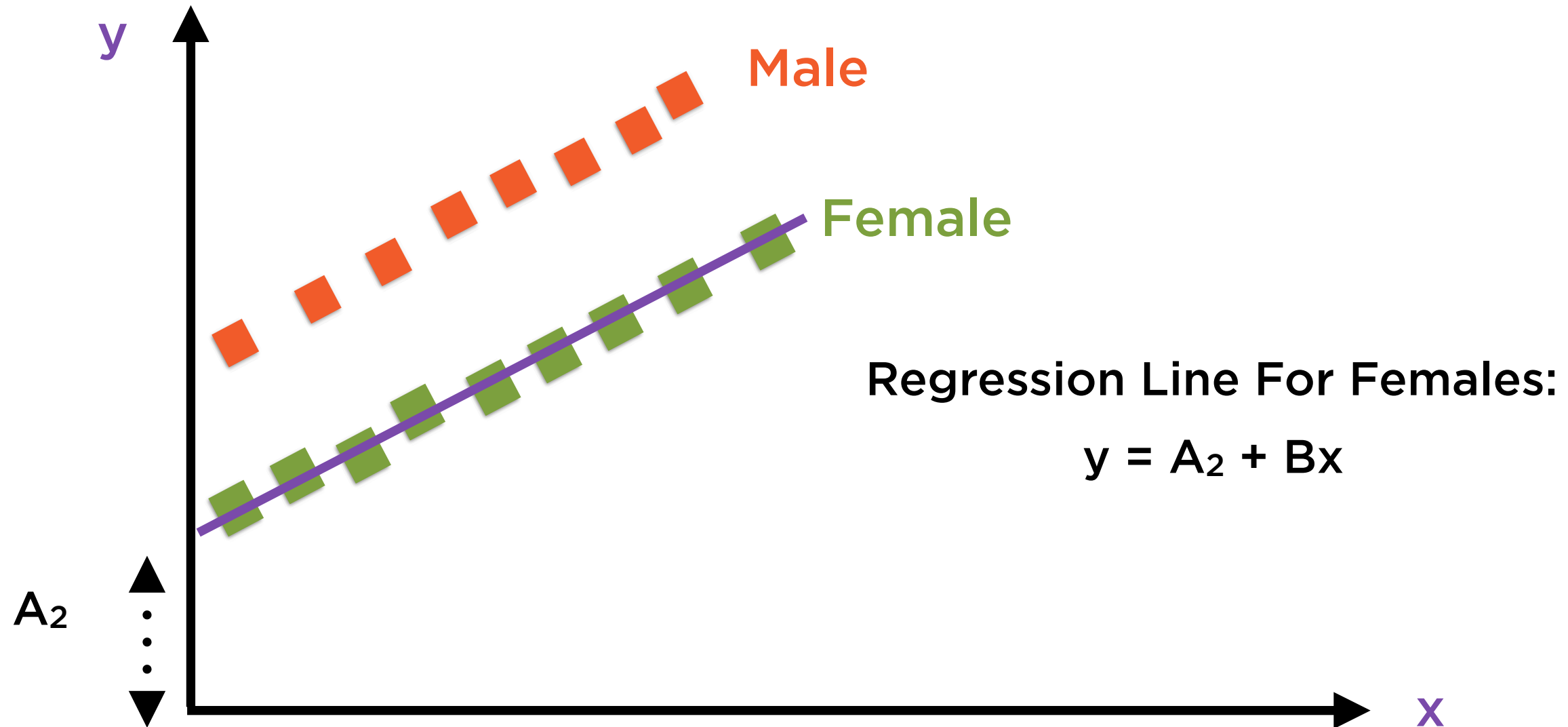
Not a great fit - regression line is far from all points!

# A Simple Regression



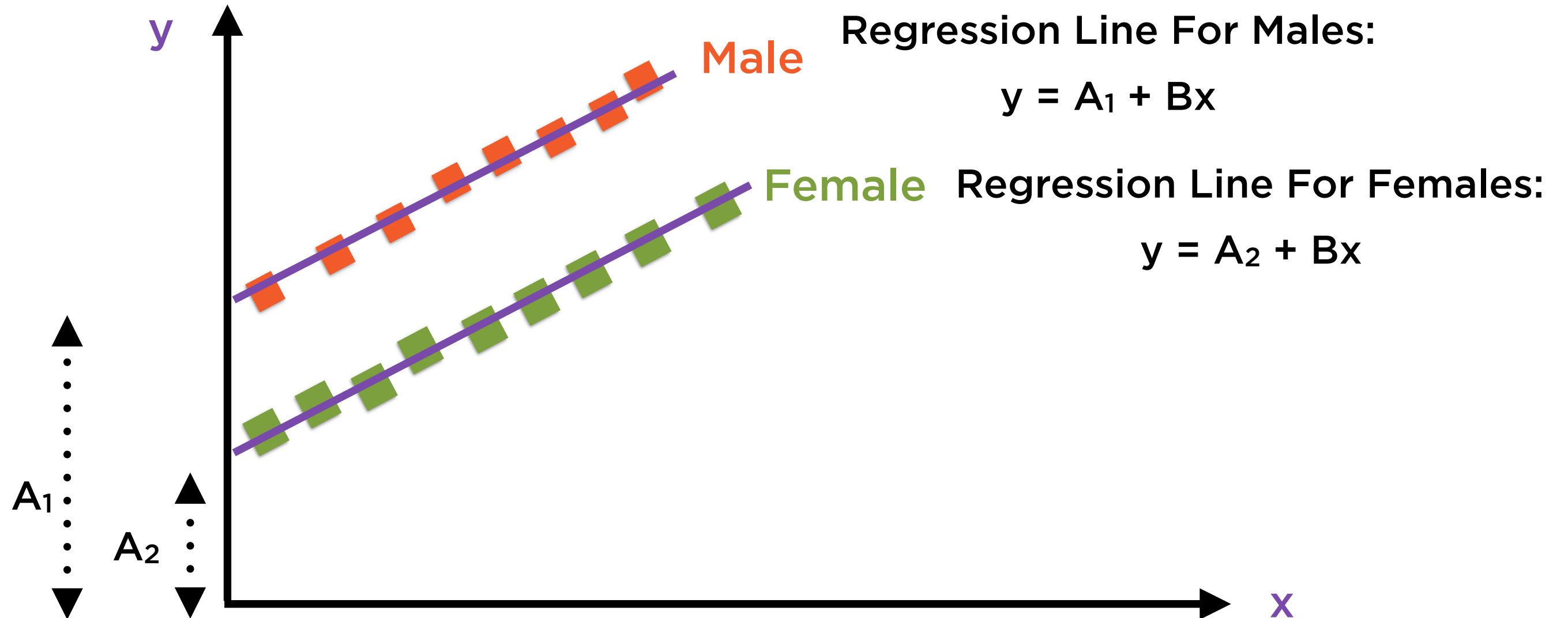
We can easily plot a great fit for males...

# A Simple Regression



...and another great fit for females

# A Simple Regression



Two lines - same slope, different intercepts

# Adding A Dummy Variable

Regression Line For Males:

$$y = A_1 + Bx$$

Regression Line For Females:

$$y = A_2 + Bx$$

**Combined Regression Line:**

$$y = A_1 + (A_2 - A_1)D + Bx$$

$D = 0$  for males

$= 1$  for females

# Adding A Dummy Variable

Regression Line For Males:

$$y = A_1 + Bx$$

Regression Line For Females:

$$y = A_2 + Bx$$

**Combined Regression Line:**

$$y = A_1 + (A_2 - A_1)D + Bx$$

**D = 0** for males

$$y = A_1 + \cancel{(A_2 - A_1)D} + Bx$$

$$= A_1 + Bx$$



# Adding A Dummy Variable

Regression Line For Males:

$$y = A_1 + Bx$$

Regression Line For Females:

$$y = A_2 + Bx$$

**Combined Regression Line:**

$$y = A_1 + (A_2 - A_1)D + Bx$$

$D = 1$  for females

$$y = \cancel{A_1} + (A_2 - \cancel{A_1}) + Bx$$

$$= A_2 + Bx$$

# Adding A Dummy Variable

Original Regression Equation:

$$y = A + Bx$$

Height of  
individual

Average height  
of parents

**Combined Regression Line:**

$$y = A_1 + (A_2 - A_1)D + Bx$$

$D = 0$     for males

$= 1$     for females

# Adding A Dummy Variable

**Combined Regression Line:**

$$y = A_1 + (A_2 - A_1)D + Bx$$

$$\begin{aligned} D &= 0 && \text{for males} \\ &= 1 && \text{for females} \end{aligned}$$

**The data contained 2 groups, so we added 1 dummy variable**

Given data with  $k$  groups, set up  $k-1$   
dummy variables, else  
multicollinearity occurs

# Dummy and Other Categorical Variables

## Dummy Variables

Binary - 0 or 1

## Categorical Variables

Finite set of values - e.g. days of week, months of year...

**To include non-binary categorical variables, simply add more dummies**

# Testing for Seasonality

**Proposed Regression Equation:**

$$y = A + BQ_1 + CQ_2 + DQ_3$$

Average stock  
returns

Quarter of the  
year

**The data contains 4 groups, so we  
added 3 dummy variables**

# Testing for Seasonality

$$y = A + BQ_1 + CQ_2 + DQ_3$$

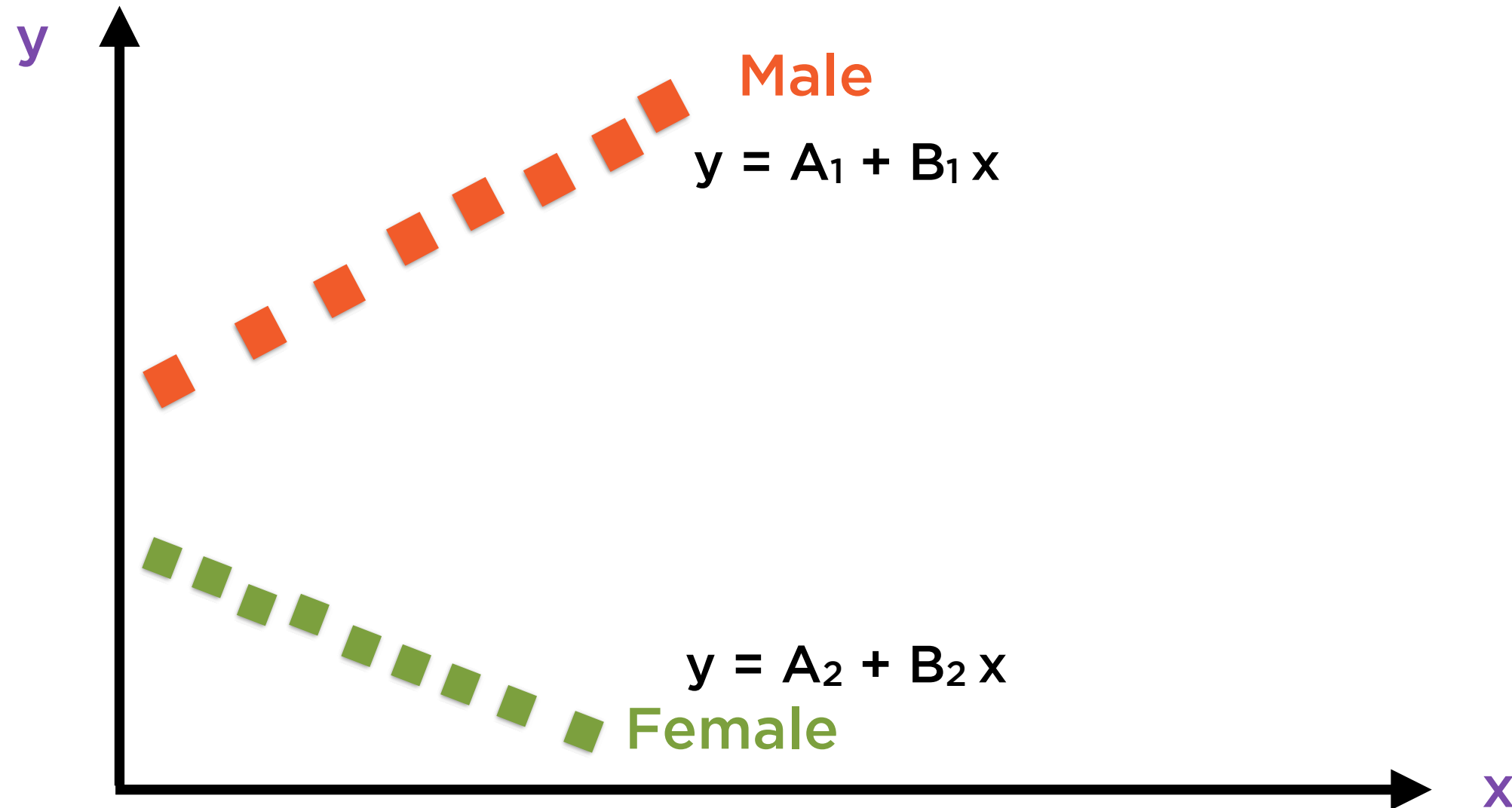
**The data contains 4 groups, so we added 3 dummy variables**

$Q_1 = 1$  for Jan, Feb, Mar  
 $= 0$  for other quarters

$Q_2 = 1$  for Apr, May, Jun  
 $= 0$  for other quarters

$Q_3 = 1$  for July, Aug, Sep  
 $= 0$  for other quarters

# Different Groups, Different Slopes



Dummy variables can also be extended for use where groups have different slopes



# Adding A Dummy Variable

Regression Line For Males:

$$y = A_1 + B_1 x$$

Regression Line For Females:

$$y = A_2 + B_2 x$$

**Combined Regression Line:**

$$y = A_1 + (A_2 - A_1)D_1 + B_1x + (B_2 - B_1)D_2$$

$D_1 = 0$  for males  
 $= 1$  for females

$D_2 = 0$  for males  
 $= x$  for females

# Adding A Dummy Variable

Regression Line For Males:

$$y = A_1 + B_1 x$$

Regression Line For Females:

$$y = A_2 + B_2 x$$

**For males:**

$$\begin{aligned} y &= A_1 + (A_2 - A_1) \boxed{D_1} + \\ &\quad B_1 x + (B_2 - B_1) \boxed{D_2} \\ &= A_1 + B_1 x \end{aligned}$$

$D_1 = 0$   
 $D_2 = 0$

# Adding A Dummy Variable

Regression Line For Males:

$$y = A_1 + B_1 x$$

Regression Line For Females:

$$y = A_2 + B_2 x$$

**For females:**

$$y = A_1 + (A_2 - A_1)(1) +$$

$$D_1 = 1$$

$$B_1 x + (B_2 - B_1)x$$

$$D_2 = x$$

$$= A_1 + (A_2 - A_1) +$$

$$B_1 x + (B_2 - B_1)x$$

$$= A_2 + B_2 x$$

# Dummy Variables

**X**

Linear regression

**Y**

Logistic regression

# Summary

**Understood the formidable benefits of multiple regression**

**Mitigated the significant risks that come with those benefits**

**Understood the utility of Adjusted- $R^2$**

**Used multiple regression to work with categorical variables**