

A hiker with a large blue backpack is seen from behind, walking up a rocky mountain trail. The hiker is wearing brown pants and a backpack with a blue top and orange bottom. The trail is made of light-colored rocks and is surrounded by dry grass and moss. The background is a steep, light-colored rock face.

# The Curious

DataRockie Basecamp



A hiker with a large blue backpack is seen from behind, climbing a steep, rocky mountain trail. The hiker is wearing brown pants and is positioned in the center of the frame. The background is a massive, light-colored rock face with some patches of moss and small plants. The foreground is filled with large, grey rocks and some dry, yellowish grass. The overall scene is rugged and adventurous.

# Correlation and Regression

Slide Version 1.0 | 31 Dec 2023

# Welcome

สวัสดีครับทุกคน 🥰

คอร์สนี้เรามาเรียนสถิติพื้นฐานที่ใช้กัน**เยอะสุด**ในงาน  
data analyst เวลาหาความสัมพันธ์ของตัวแปร

ลืมทุกอย่างที่เคยเรียนมา แล้วมาเรีมใหม่ในคอร์สนี้ 555+



# The Contents

- Welcome
- Correlation
- Linear Regression
- Error Calculation
- Model Interpretation
- Hands-On Tutorials



# Your Instructor

แอดทอย คนดี คนเดิม

Bachelor of Economics, TH  
Master of Economics, UK  
1 Learning Hour Every Day

เรียน เขียน แชร 100



# Data

ข้อมูลในทางสถิติหลักๆจะแบ่งออกเป็น 2 แบบ

1. เชิงปริมาณ เช่น ยอดขาย กำไร
2. เชิงคุณภาพ เช่น จังหวัด ประเทศ ประเภทสินค้า

# Ad Data

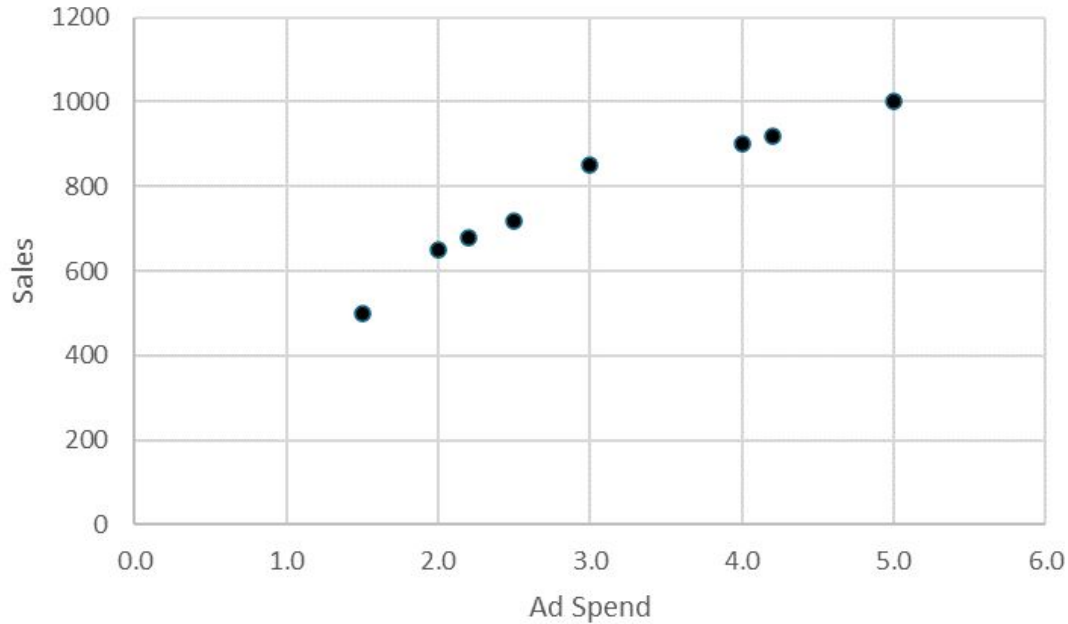
ข้อมูลตัวอย่างในคอร์สนี้เป็น  
แบบตัวเลขเชิงปริมาณ

ตัวแปรต้นคือ **\$ Ad Spend**  
ตัวแปรตามคือ **\$ Sales**

$Sales = f(Ad\ Spend)$

Ad	Sales
1.5	500
2.0	650
2.2	680
2.5	720
3.0	850
4.0	900
4.2	920
5.0	1000

# Basic Scatter Plot





A hiker with a large blue backpack is seen from behind, walking on a rocky mountain trail. The hiker is wearing brown pants and a yellow belt. The trail is made of light-colored rocks and is surrounded by dry grass and moss. The background is a steep, light-colored rock face. The text "Before Correlation" is overlaid in the center of the image.

**Before Correlation**

# Covariance

ยุคแรก นักสถิติใช้ค่า COV เพื่อหาความสัมพันธ์  
เชิงเส้นตรงของตัวแปรสองตัว

$$cov_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N - 1}$$

$cov_{x,y}$  = covariance between variable x and y

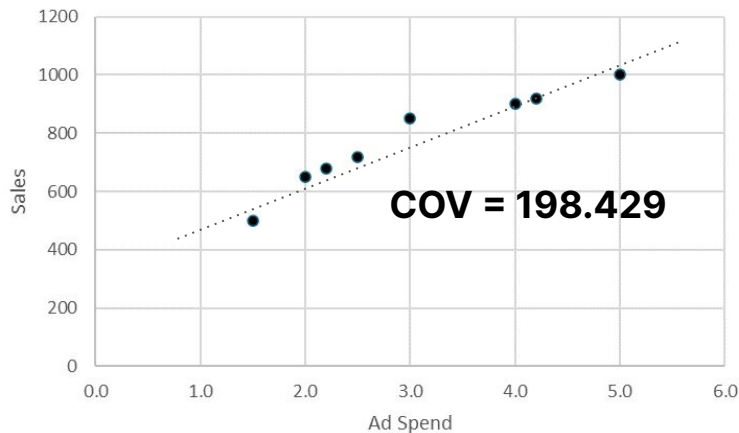
$x_i$  = data value of x

$y_i$  = data value of y

$\bar{x}$  = mean of x

$\bar{y}$  = mean of y

$N$  = number of data values



# Limitation of Covariance

แต่ข้อจำกัดของ COV คืออธิบายผลยาก เพราะมีค่า  
ได้ตั้งแต่ **-infinity** ถึง **+infinity**





A hiker with a large blue backpack is seen from behind, walking up a steep, rocky mountain trail. The hiker is wearing brown pants and a backpack with a blue top and orange bottom. The trail is composed of light-colored rocks and patches of dry grass. The background is a massive, light-colored rock face with some moss or lichen. The sky is overcast and grey.

# Correlation

# The Birth of Correlation

Karl Pearson จึงได้คิดวิธีการปรับสูตร  
Covariance ให้อ่านค่าได้ง่ายขึ้น

กำเนิดเป็น **Pearson Correlation** ที่มี  
ค่าอยู่ระหว่าง **-1** ถึง **+1**



Karl Pearson (1857-1936)

# Correlation Formula

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

simple formula

$$\frac{\text{COV}(x,y)}{\text{sd.x} * \text{sd.y}}$$

$r$  = correlation coefficient

$x_i$  = values of the x-variable in a sample

$\bar{x}$  = mean of the values of the x-variable

$y_i$  = values of the y-variable in a sample

$\bar{y}$  = mean of the values of the y-variable



# Simple Functions



ทุกวันนี้เราไม่ต้องคำนวณมือเองเหมือนสมัยก่อนแล้ว  
แค่เรียกใช้ **function** ได้เลย

**=COVARIANCE.S()**

$$cov_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N - 1}$$

$cov_{x,y}$  = covariance between variable x and y

$x_i$  = data value of x

$y_i$  = data value of y

$\bar{x}$  = mean of x

$\bar{y}$  = mean of y

$N$  = number of data values

**=CORREL()**

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

$r$  = correlation coefficient

$x_i$  = values of the x-variable in a sample

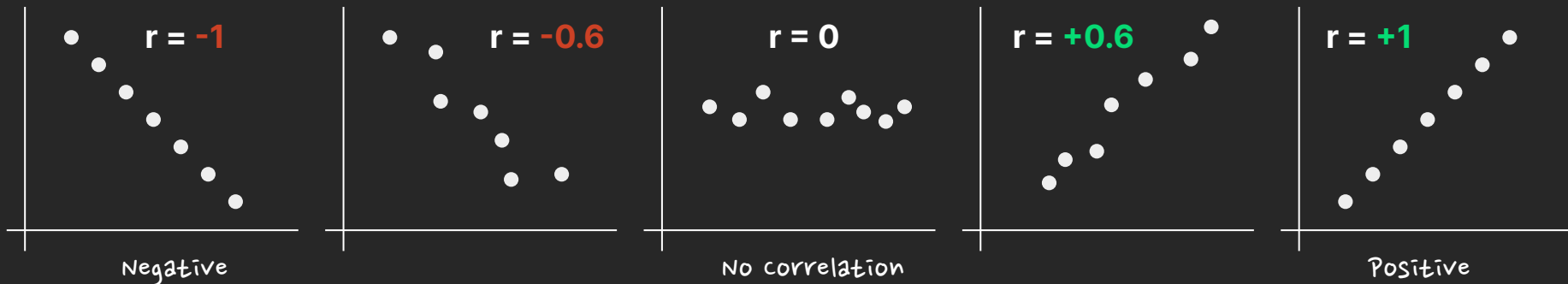
$\bar{x}$  = mean of the values of the x-variable

$y_i$  = values of the y-variable in a sample

$\bar{y}$  = mean of the values of the y-variable

# Interpretation

ค่า correlation หรือ  $r$  บอกความสัมพันธ์ของตัวแปรแบบ  
ตัวเลขสองตัว (เหมือน covariance)



# Use Case

Correlation ใช้ตอบคำถามต่อไปนี่ว่าตัวแปรสองตัวมี ...

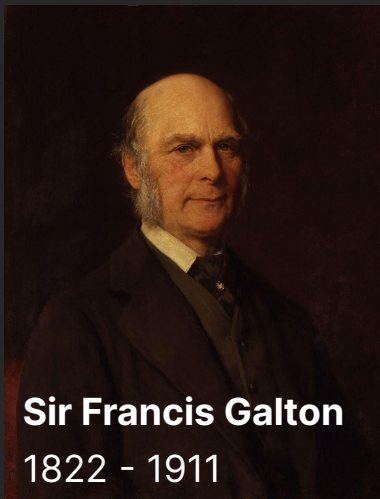
- ความสัมพันธ์เชิงเส้นตรงกันหรือเปล่า
- ความสัมพันธ์เป็นเชิง + หรือเชิง -
- ความสัมพันธ์นั้นเข้มแค่ไหน strong หรือ weak



# Regression

A hiker with a large blue backpack is seen from behind, climbing a steep, rocky mountain trail. The hiker is wearing brown pants and a yellow backpack. The trail is narrow and rocky, with patches of dry grass and moss. The background is a steep, light-colored rock face with some green moss. The sky is overcast.

# The Inventors



ปี **1809** Gauss พัฒนาเทคนิค **Least Squares Method**

ปี **1885** Galton ได้นำเสนอไอเดียเรื่อง **Regression** ให้กับโลกสถิติ

“Regression Towards The Mean”

# Regression

นักสถิติคิดโมเดล Regression เพื่อใช้ตอบคำถามที่ Correlation ตอบไม่ได้

ตัวอย่างเช่น  $\text{Sales} = f(\text{Ad Spend})$

✓ ถ้า Ad Spend เปลี่ยน 1 หน่วย แล้ว Sales จะเปลี่ยนเท่าไร? ปัจจัยอื่นๆคงที่



# A Primer to Linear Regression

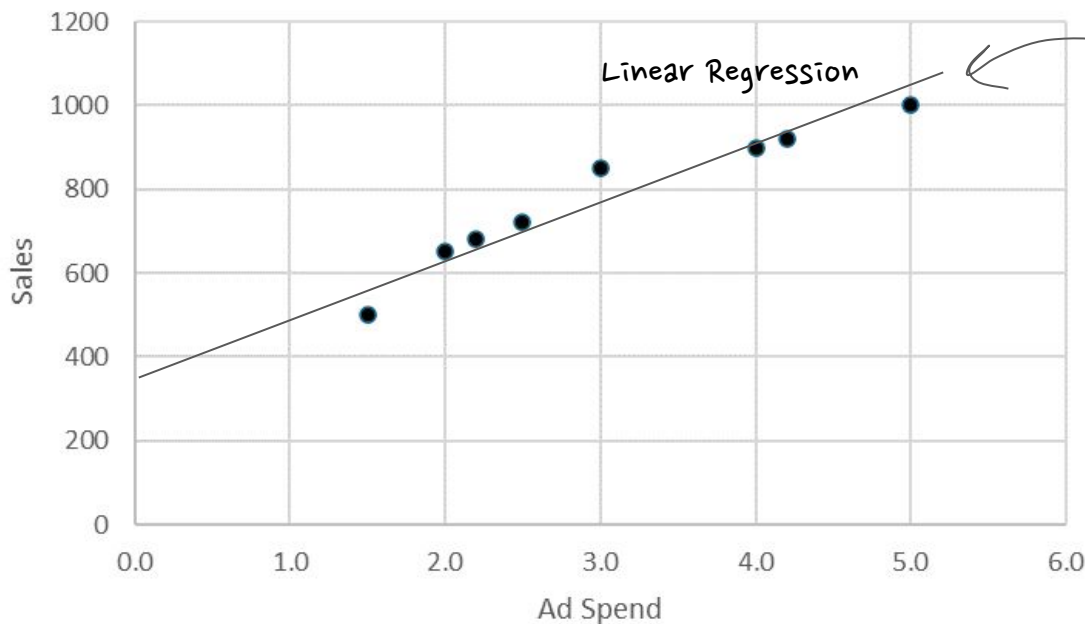


และโมเดล Regression รูปแบบที่ง่ายที่สุด

ที่สอนในคลาสสถิติพื้นฐานทั่วโลกคือ **Linear Regression**



# It's Just A Straight Line

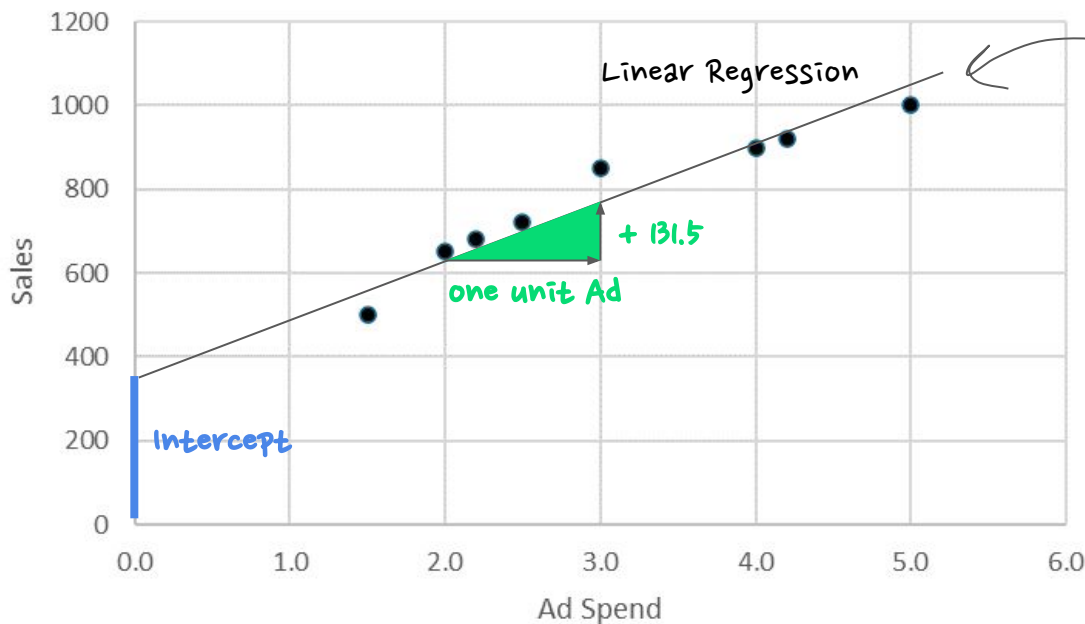


$$\text{Sales} = f(\text{Ad})$$

$$\text{Sales} = \text{intercept} + \text{slope} * \text{Ad}$$

$$\text{Sales} = 376 + 131.5 * \text{Ad}$$

# A Simple Idea

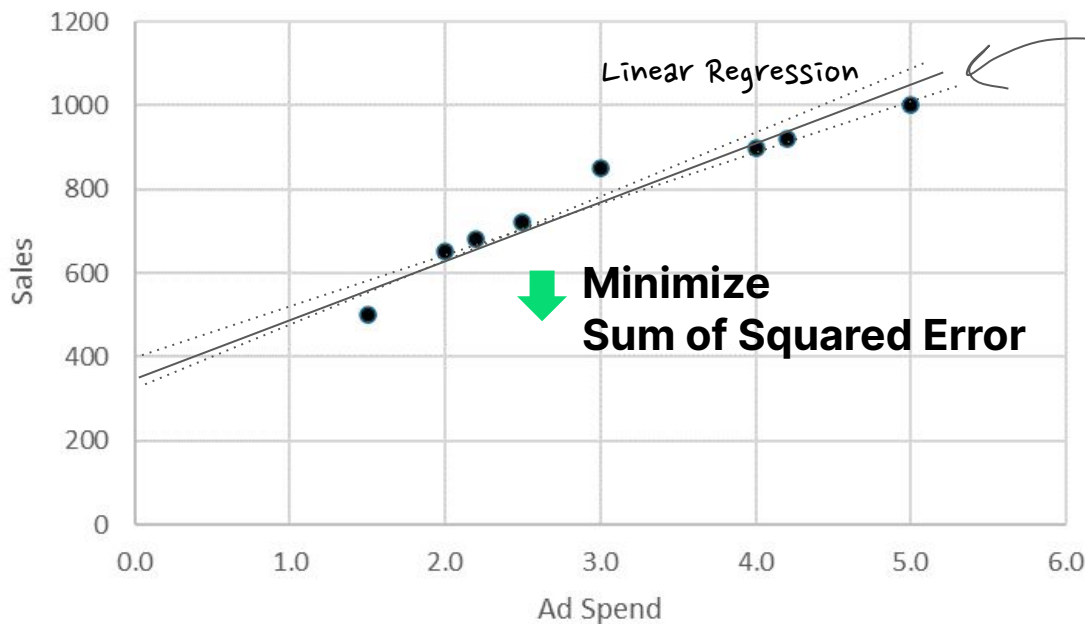


$$\text{Sales} = f(\text{Ad})$$

$$\text{Sales} = \text{intercept} + \text{slope} * \text{Ad}$$

$$\text{Sales} = 376 + 131.5 * \text{Ad}$$

# The Best Line



$$\text{Sales} = f(\text{Ad})$$

$$\text{Sales} = \text{intercept} + \text{slope} * \text{Ad}$$

$$\text{Sales} = 376 + 131.5 * \text{Ad}$$

ทำให้ค่า **Total Error** ของโมเดล  
มีค่าต่ำที่สุด (minimize)

# Error Calculation

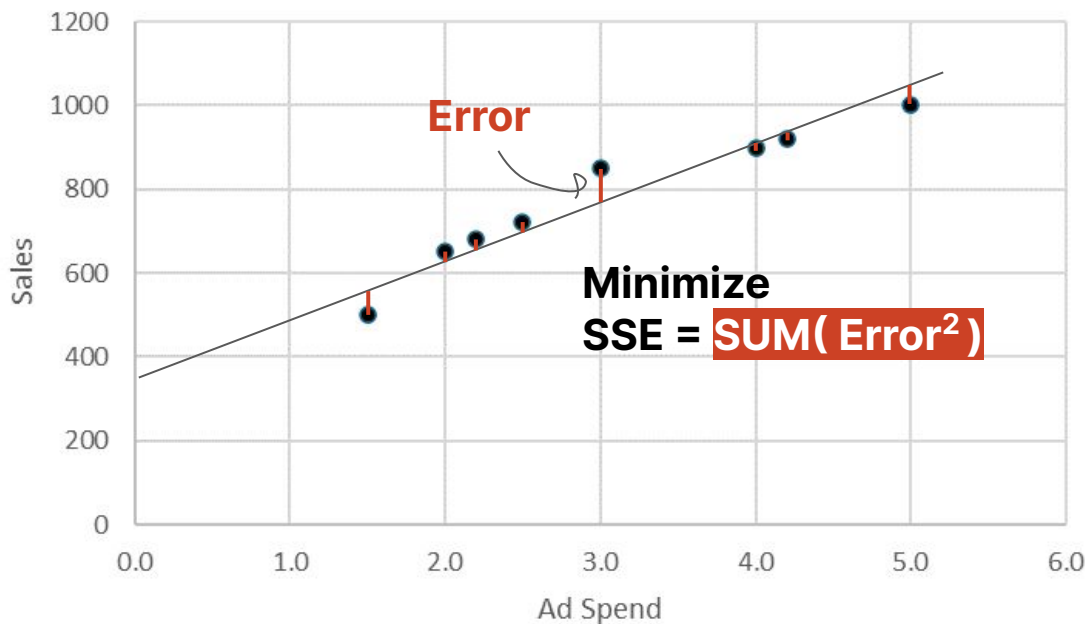
Error คือ **ความแตกต่าง** ระหว่างค่าจริง กับสิ่งที่โมเดลทำนาย

Error = Actual - Prediction

✓ Error โดยรวมของโมเดล Linear Regression เรียกว่า **SSE** (sum of squared error)



# Error Calculation



$$\text{Sales} = f(\text{Ad})$$

$$\text{Sales} = \text{intercept} + \text{slope} * \text{Ad}$$

$$\text{Sales} = 376 + 131.5 * \text{Ad}$$

ทำให้ค่า **Total Error** ของโมเดลมีค่าต่ำที่สุด (minimize)

# Our Model Tries To ..

หาค่า intercept และ slope ที่ทำให้ค่า sum of squared error ของโมเดลที่ค่าต่ำที่สุด

Find {**intercept**, **slope**}  
That **Minimize SSE**



ง่ายจนงง 555+

# R-Squared

คือค่าที่ใช้วัดสิ่งที่โมเดล linear regression อธิบายได้ โดยมีค่าวิ่งอยู่ระหว่าง  $[0 - 1]$

ค่ายิ่งเข้าใกล้ 1 แปลว่าโมเดลเราทำงานได้ดี -> ตัวแปรต้น X อธิบายตัวแปร Y ได้ดี

PS.  $R^2$  มีชื่อเต็มๆว่า Coefficient of Determination

# R-Squared

เราสามารถคำนวณค่า R-Squared ได้ด้วยสูตร

$$\text{R-Squared} = SS_m / SS_t$$

หรือ  $1 - SS_r / SS_t$



ตัวย่อที่เราใช้ในสูตร

m = model

r = residual

t = total variance

$$\text{และ } SS_m + SS_r = SS_t$$



# R-Squared

เราสามารถคำนวณค่า R-Squared ได้ด้วยสูตร

$$\text{R-Squared} = 4 / 10 = 40\%$$

$$\text{หรือ } 1 - 6 / 10 = 40\%$$



ตัวย่อที่เราใช้ในสูตร

m = model

r = residual

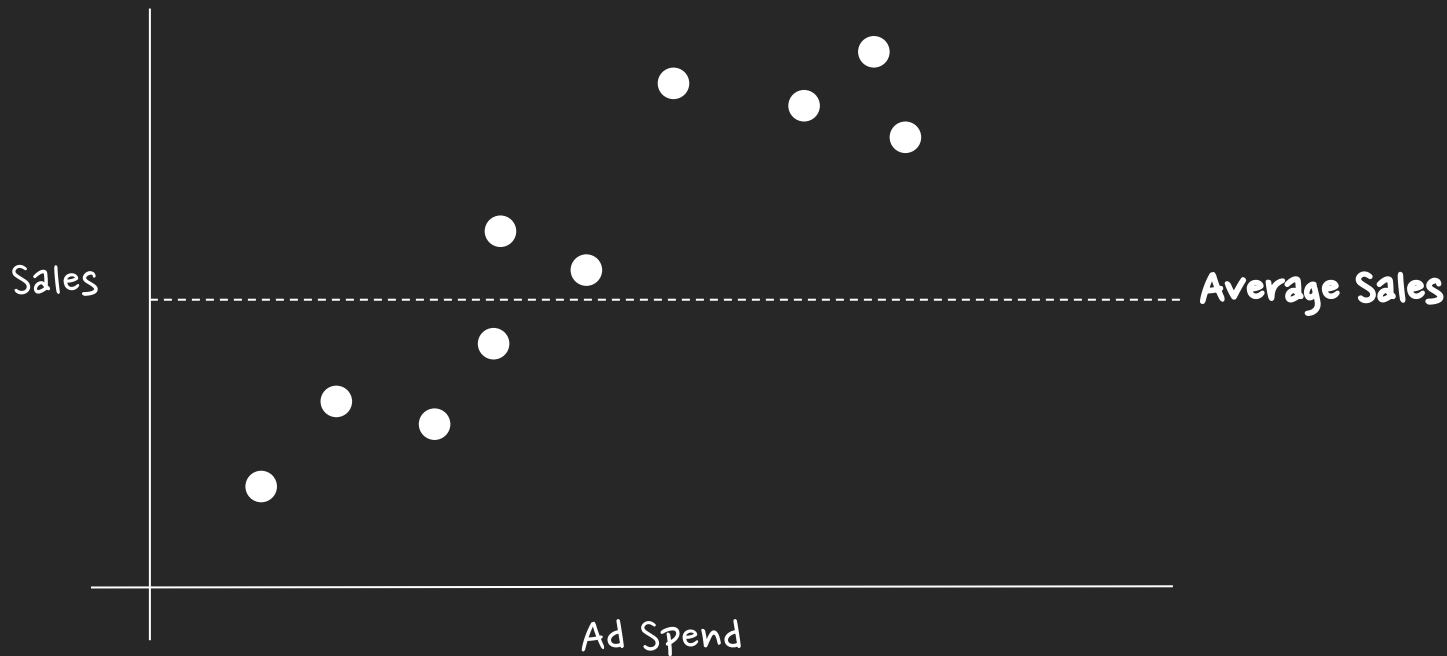
t = total variance

$$\text{และ } SS_m + SS_r = SS_t$$

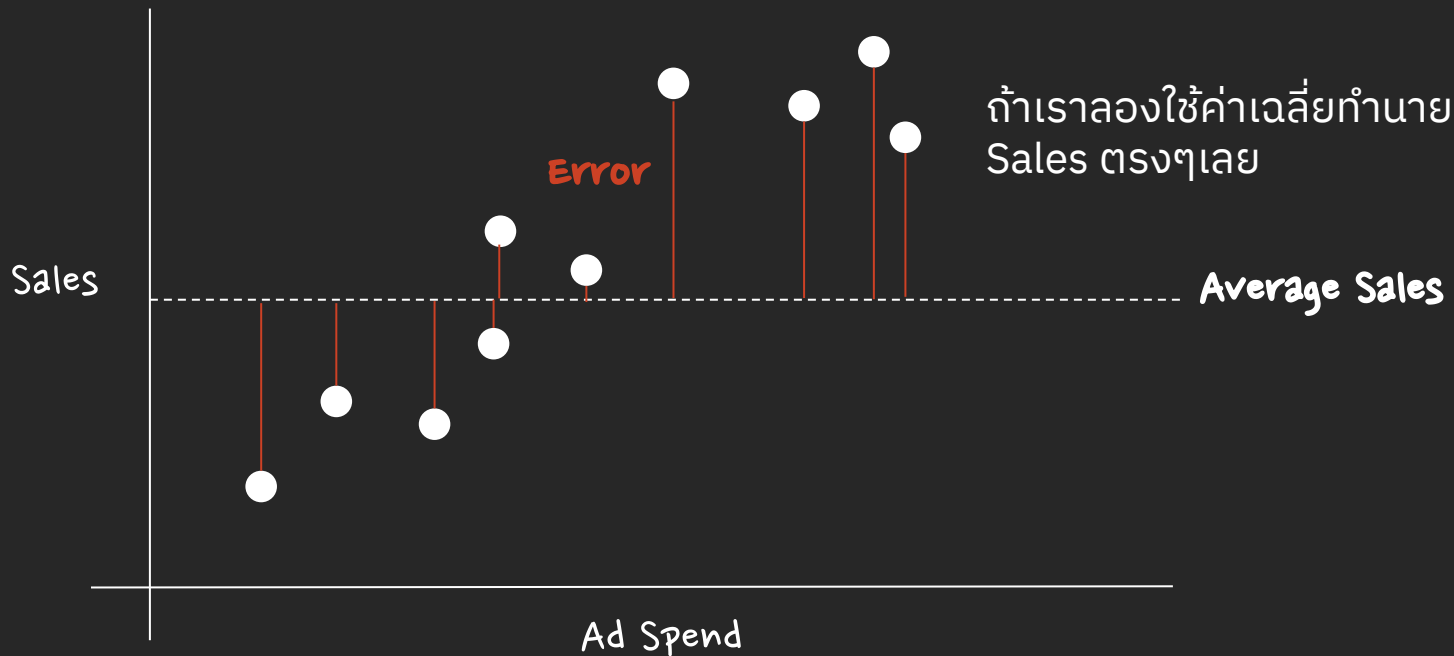
A hiker with a large blue backpack is seen from behind, climbing a steep, rocky mountain trail. The hiker is wearing brown pants and a dark shirt. The trail is narrow and rocky, with patches of dry grass and moss. The background shows a steep, light-colored rock face. The text "Geek Mode" is overlaid in the center of the image.

# Geek Mode

# R-Squared Visualization

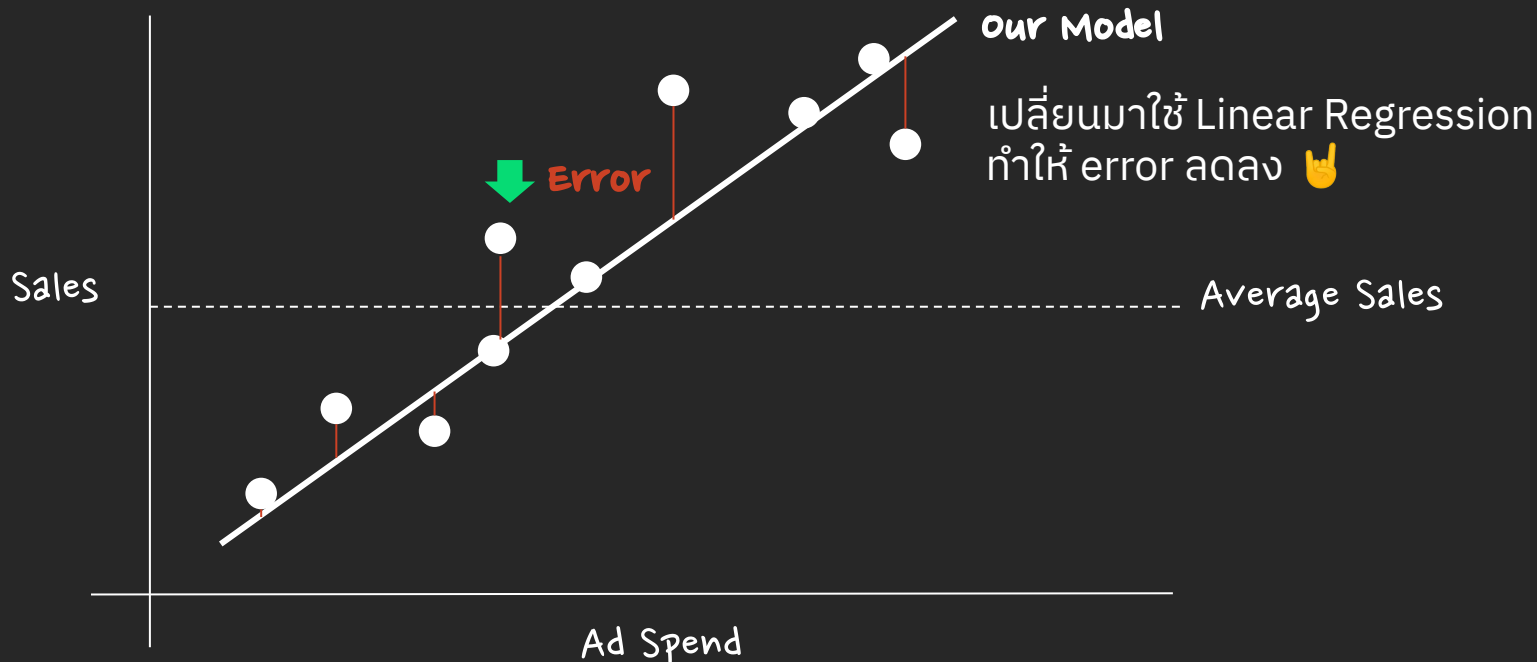


# R-Squared Visualization

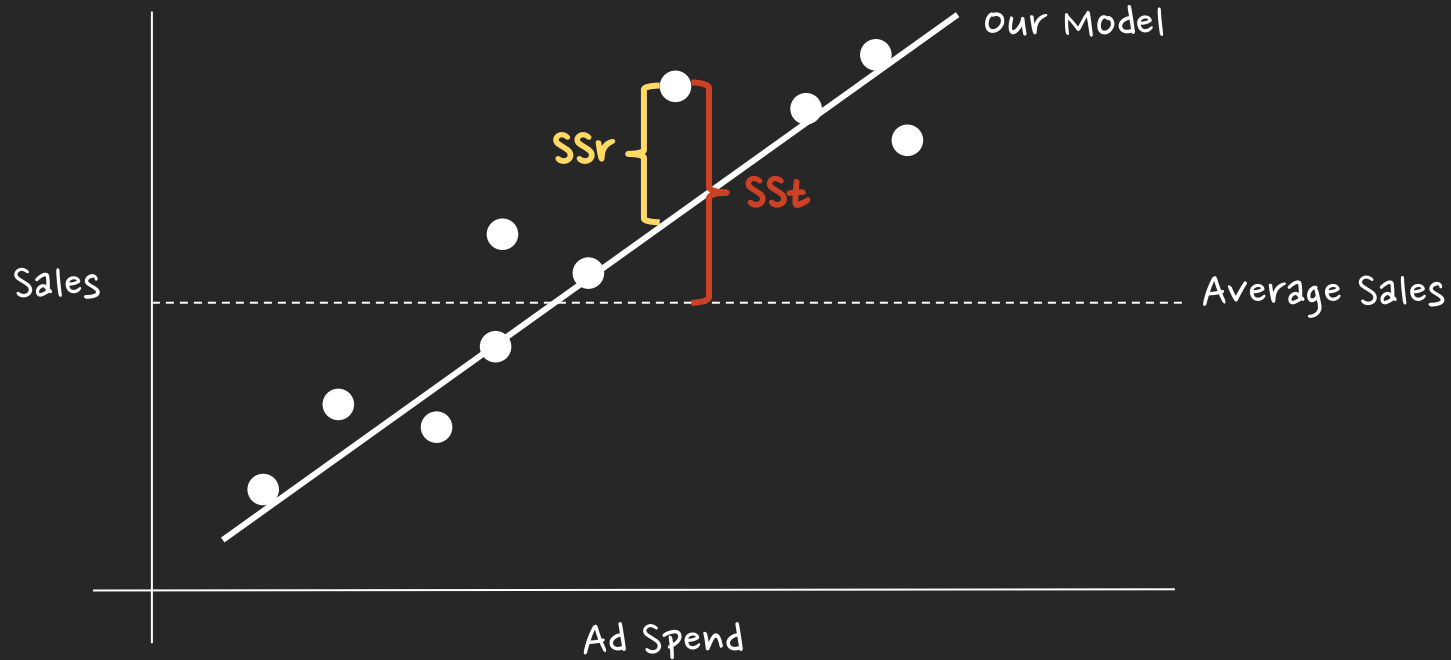




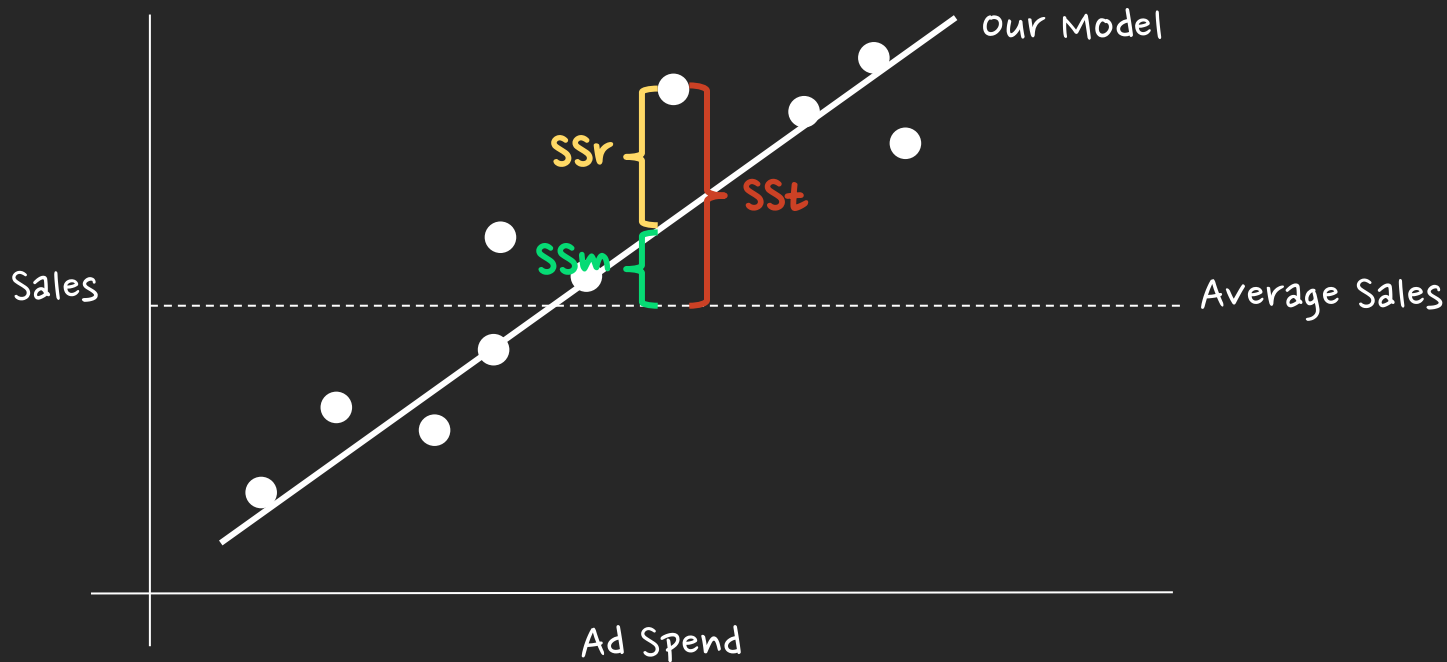
# R-Squared Visualization



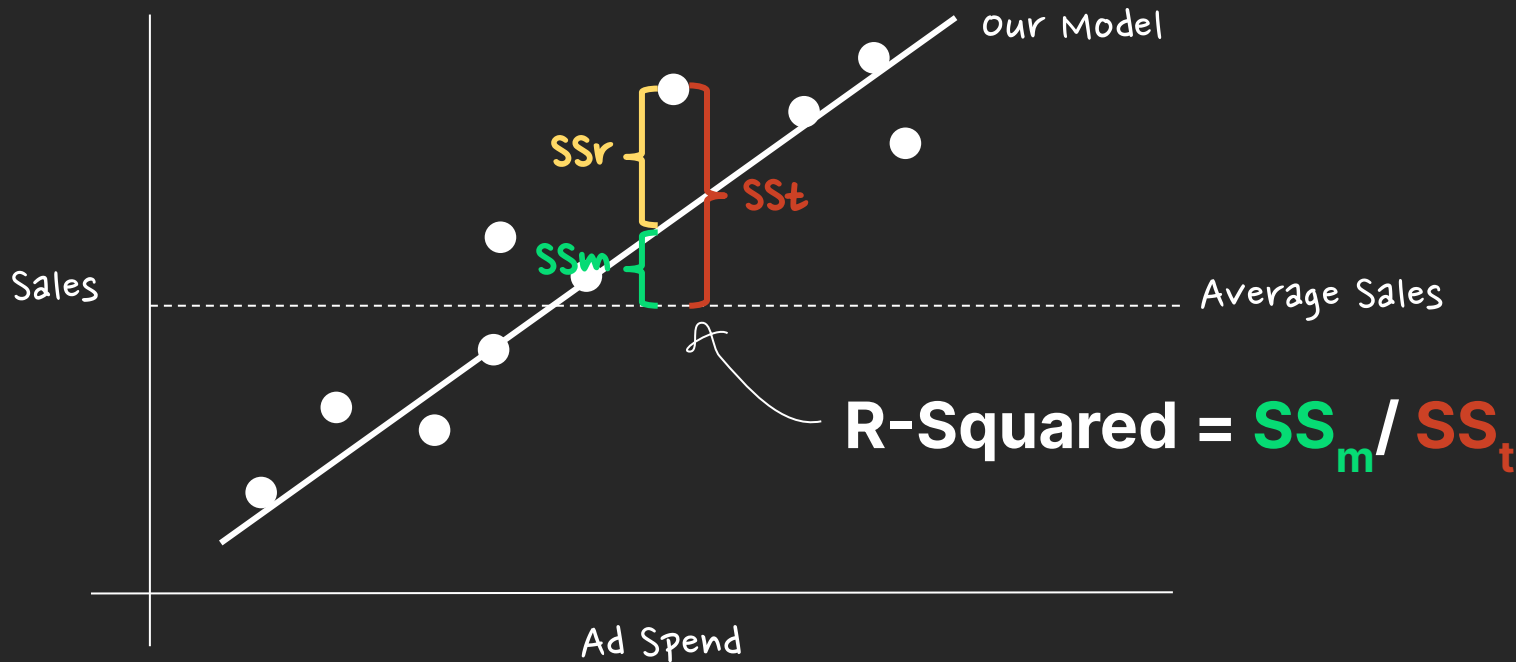
# R-Squared Visualization



# R-Squared Visualization

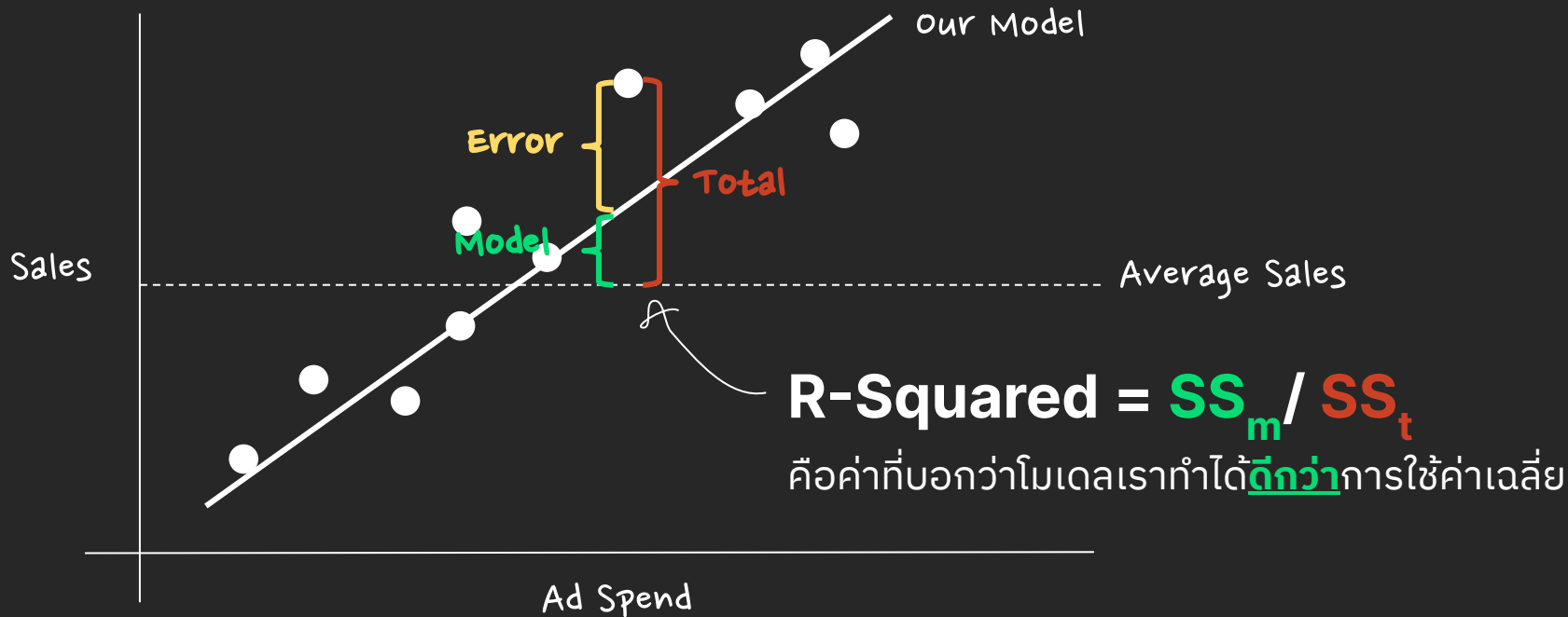


# R-Squared Visualization

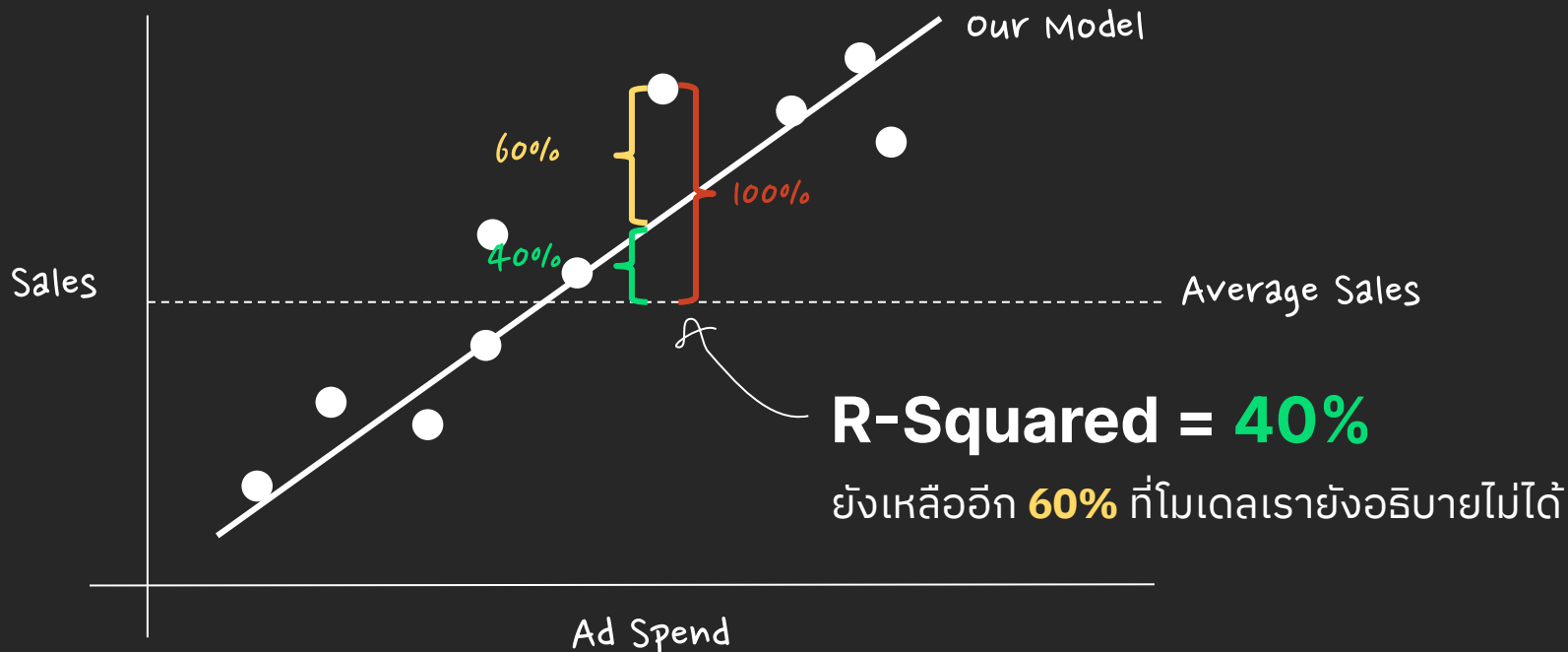




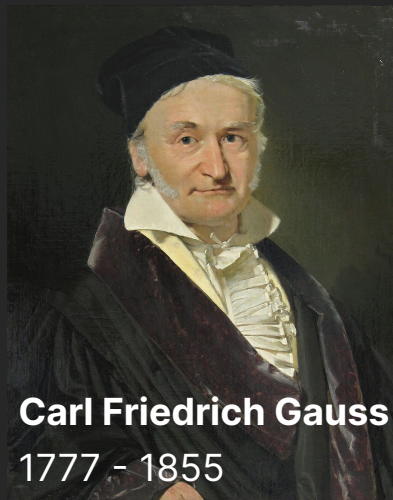
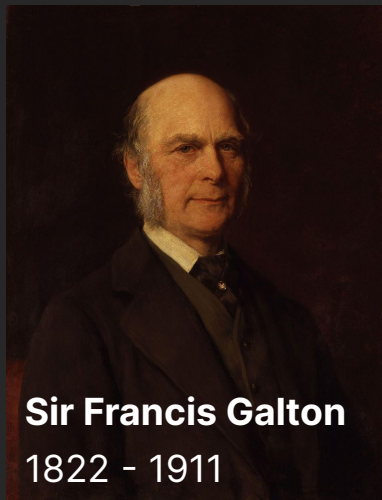
# R-Squared Visualization



# R-Squared Visualization



# Recap: The Inventors



ปี **1809** Gauss พัฒนาเทคนิค **Least Squares Method**

ปี **1885** Galton ได้นำเสนอไอเดียเรื่อง **Regression** ให้กับโลกสถิติ

“**Regression** Towards **The Mean**”

# Mind Blown

ง่ายจนง 555+





A hiker with a large blue and orange backpack is seen from behind, walking up a rocky mountain trail. The trail is narrow and composed of light-colored rocks and patches of dry grass. The background is a steep, light-colored rock face with some moss. The sky is overcast and grey.

**Simple vs. Multiple**

# Two Types of LR

[1] **Simple** Linear Regression

Sales =  $f(\text{Ad})$

[2] **Multiple** Linear Regression

Sales =  $f(\text{Ad}, \text{TV}, \text{Radio})$

# The Model

## [1] Simple Linear Regression

$$\text{Sales} = b_0 + b_1 * \text{Ad}$$

## [2] Multiple Linear Regression

$$\text{Sales} = b_0 + b_1 * \text{Ad} + b_2 * \text{TV} + b_3 * \text{Radio}$$

โดยที่  $b_0$  คือ **intercept** และ  $b_1, b_2, b_3, ..$  คือ **slope**



A hiker with a large blue backpack is seen from behind, walking on a rocky mountain trail. The hiker is wearing brown pants and a yellow backpack. The trail is made of light-colored rocks and is surrounded by dry grass and moss. The background is a steep, light-colored rock face. The word "Prediction" is overlaid in white text in the center of the image.

Prediction

# Predict New Data

## Simple Linear Regression

$$\text{Sales} = f(\text{Ad})$$

$$\text{Sales} = 100 + 50 * \text{Ad}$$

✓ ถ้าเราใช้เงินโฆษณา **\$5 Million USD** เราจะได้ยอดขายกลับมาเท่าไร?

$$\text{Sales} = 100 + 50 * 5 = \boxed{350}$$

Prediction



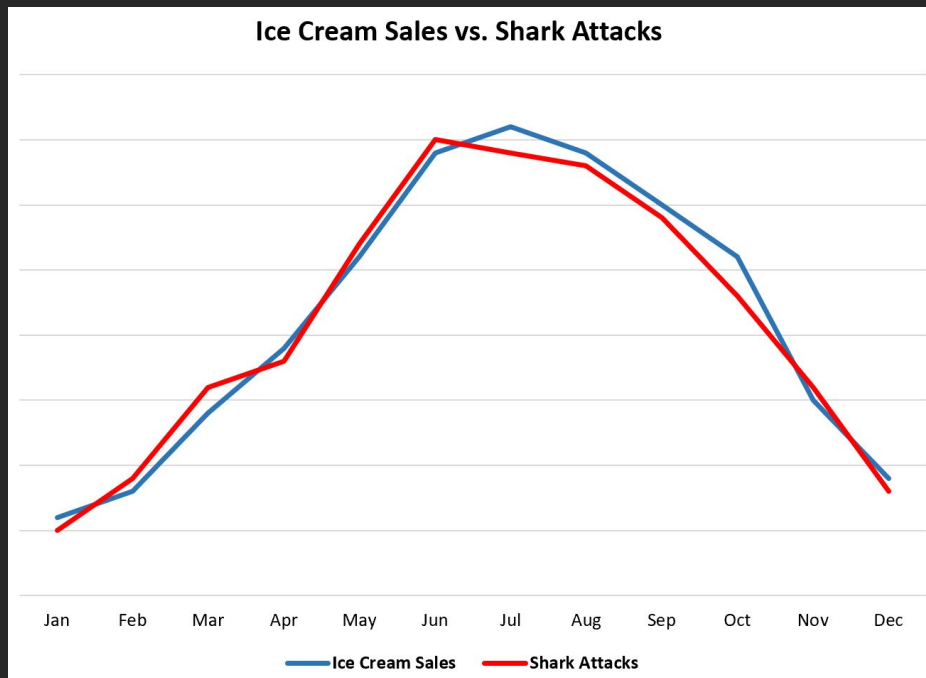
A hiker with a large blue backpack is seen from behind, walking on a rocky mountain trail. The hiker is wearing brown pants and a backpack with a blue top and orange bottom. The trail is narrow and rocky, with patches of dry grass and moss. The background is a steep, light-colored rock face. The word "Caution" is overlaid in white text in the center of the image.

**Caution**

# Correlation Does Not Imply Causation

ตัวแปรสองตัวมีความสัมพันธ์กัน (Correlation)  
ไม่ได้แปลว่า **X ทำให้เกิด Y** (Causation)

# Ice Cream vs. Shark

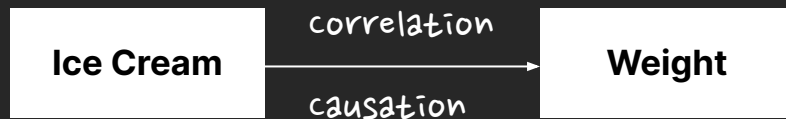


ชอบกินไอติม ไม่ได้ แปล  
ว่าจะถูกลาถกั๊ด

# Ice Cream vs. Weight



กินไอติมเยอะ น้ำหนักขึ้น?



# Remember This

ทุกความสัมพันธ์แบบ Causation ต้องมี Correlation เสมอ

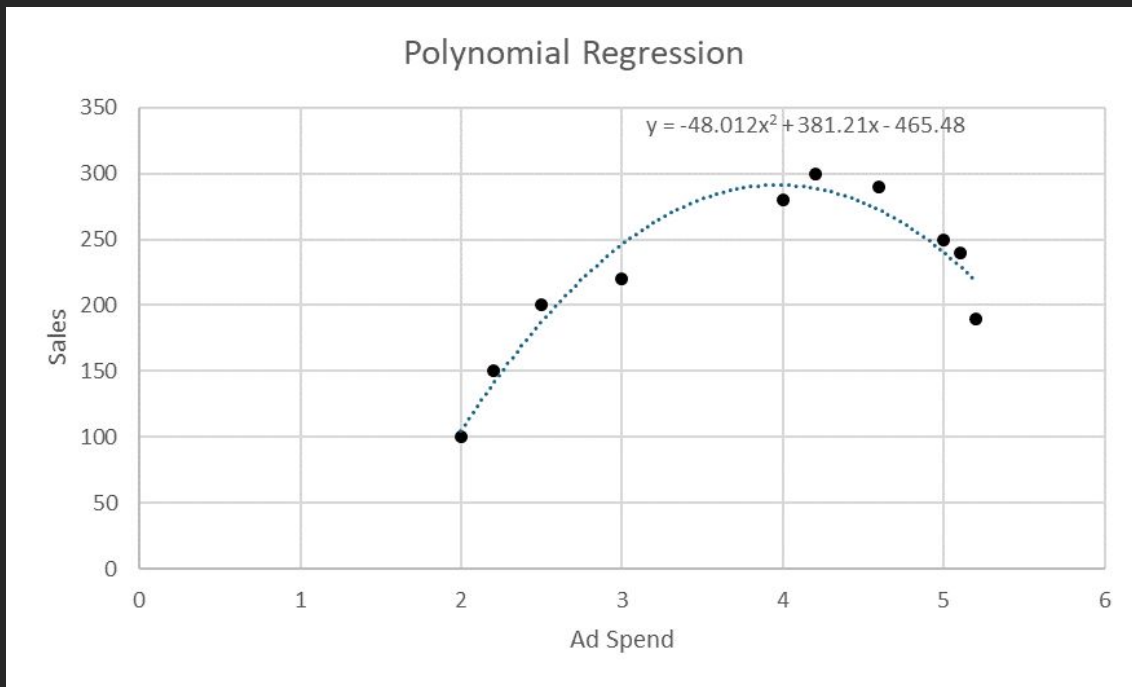
แต่ไม่ใช่ทุกความสัมพันธ์แบบ Correlation จะเป็น Causation



A hiker with a large blue backpack is seen from behind, climbing a steep, rocky mountain trail. The hiker is wearing a dark shirt and light-colored pants. The trail is narrow and rocky, with patches of dry grass and moss. The background shows a steep, light-colored rock face with some vegetation. The sky is overcast.

# Other Types of Regression

# Polynomial Regression



เวลาเจอกับ **Non-Linear**  
Relationship

# Logistic Regression



# Take-Home Cheat Sheets



นักศึกษาใช้ Correlation และ Regression ในการ  
โมเดลความสัมพันธ์ตัวแปรเชิงปริมาณ

Note - ถ้าเรียนต่อไปจะรู้ว่า Regression สามารถ  
รับตัวแปรประเภทอื่นได้ด้วย เช่น dummy 0,1



Correlation หรือ r มีค่าอยู่ระหว่าง [-1, +1]

Linear Regression ตัวพื้นฐานมี 2 แบบคือ

- Simple มีตัวแปรต้นหนึ่งตัว
- Multiple มีตัวแปรต้นมากกว่าหนึ่งตัว

นักศึกษาใช้ Regression เพื่ออธิบายว่าถ้า x เปลี่ยน  
1 หน่วย y จะเปลี่ยนเท่าไร ปัจจัยอื่นคงที่



Functions ที่ต้องใช้ให้เป็นใน Excel/ Sheets

- COV.S()
- CORREL()
- INTERCEPT()
- SLOPE()
- LINEST()



วิธีวัดความแม่นยำของโมเดล Linear Regression  
ทำได้หลายวิธี

- R-Squared ยิ่งเข้าใกล้ 1 ยิ่งดี
- MAE ยิ่งเข้าใกล้ 0 ยิ่งดี

สูตรของ R-Squared คือ  $SS_{\text{model}} / SS_{\text{total}}$   
หรือ  $1 - SS_{\text{residual}} / SS_{\text{total}}$





A hiker with a large blue backpack is seen from behind, walking up a rocky mountain trail. The hiker is wearing brown pants and a backpack with a blue top and orange bottom. The trail is made of light-colored rocks and is surrounded by dry grass and moss. The background is a steep, light-colored rock face.

# The Curious

DataRockie Basecamp