

**YOuRClass**

# **Basic Data Science**



# Afifa Ayu Widhiyanthi

Chairperson at Data Science Indonesia

## Experiences

- Data Analyst at INA Digital
- B.S in Computer Science Bina Nusantara University, Jakarta
- B.S in Statistics Bina Nusantara University, Jakarta
- President, Bina Nusantara Statistics Student Association 2019
- ex Tokopedia, ex Trevo



Afifa Widhiyanthi



@afifadayu

**YOuRClass**

# **Data Science 101**

**Week 1**

## Konten

**1**

Pengenalan Data Science

---

**2**

Pengenalan Python

---

**3**

Jenis-jenis Data

---

**4**

Statistika Dasar

---

# **Pengenalan Data Science**

**1**



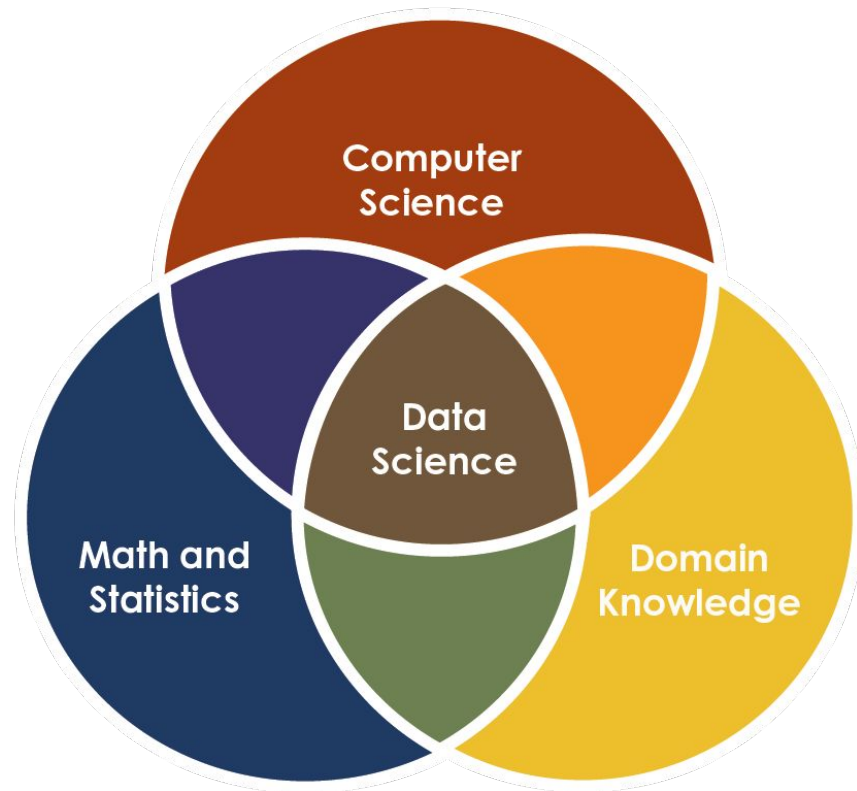
Kumpulan fakta  
berisi angka, tulisan,  
atau pengukuran



Proses memahami  
dunia melalui  
observasi,  
eksperimen, dan  
logika



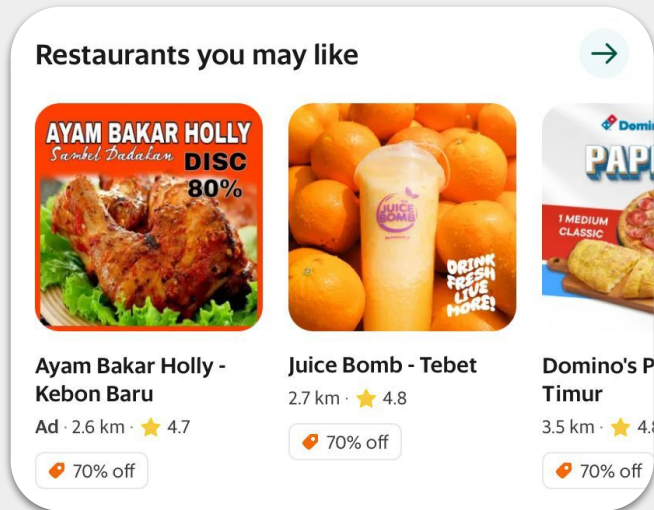
Proses dalam  
memahami sebuah  
masalah  
menggunakan data.



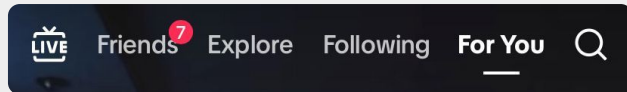
## Dimana ada data?

- Sistem rekomendasi dalam *tech company (e-commerce/ride-hailing)*, biasanya kamu lihat tulisan seperti ini “Anda mungkin suka”
- Pada media sosial, konten yang kamu **sukai** atau **komentari** rekomendasinya akan lebih sering muncul di “For You”

Rekomendasi pada aplikasi *online*



Rekomendasi pada *social media*





Data tersebut menghasilkan  
**Big Data**

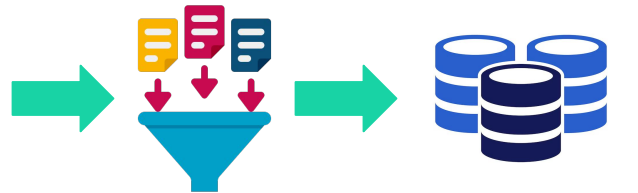
## Big Data

Ketika **volumenya sangat besar** yang diproses oleh mesin, alat, atau manusia. Sehingga, tidak bisa ditangani oleh penyimpanan tradisional.

### Use case,



Menonton Netflix



Data dikumpulkan

Data disimpan



**Maksimum 17 miliar sel**

Baris = 1.048.576

Kolom = 6.384 kolom (XFD)

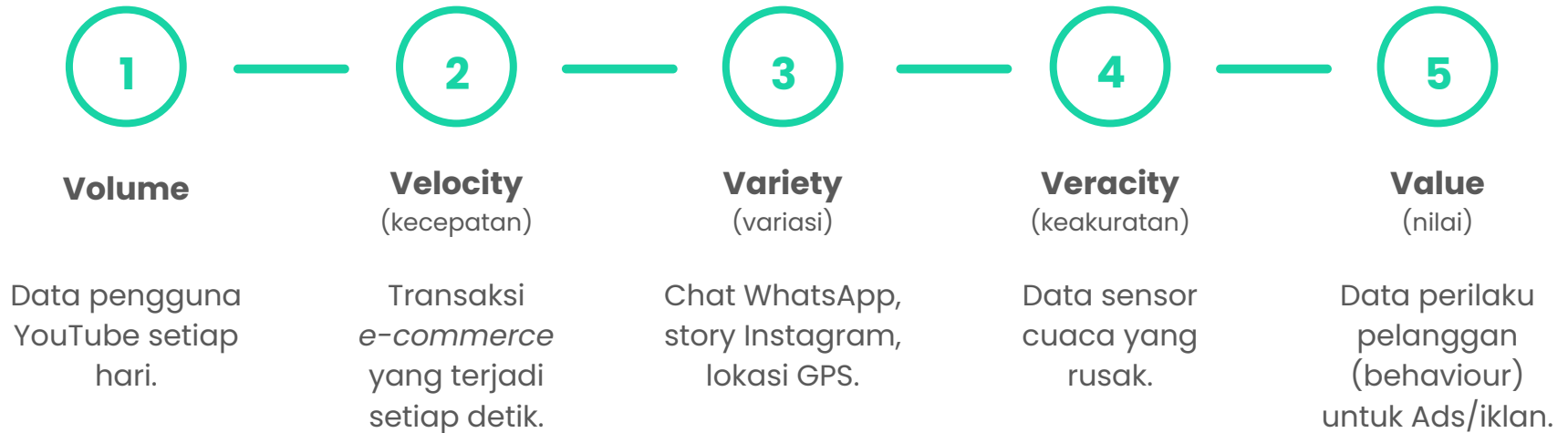


**Maksimum 10 juta sel**

Kolom = 18.278 kolom/tab

# Ciri-ciri Big Data

disebut sebagai **5V**



## Langkah lanjutan

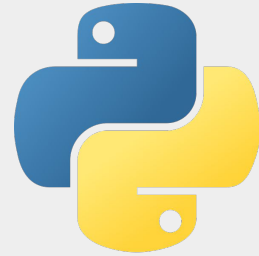
### Menjadi **Storyteller**

Tujuannya untuk **memberikan informasi** yang terdapat pada **data tersebut kepada orang lain**, sehingga seseorang dapat memahami dan mengetahui informasi yang kita ketahui.

## Tools

Tools/alat tersebut membantu kita melihat dan memahami tren dan pola dalam data.

1. Python
2. R
3. SQL (Google BigQuery)
4. Microsoft Excel
5. Tableau
6. Google Looker



Google  
Big Query



+ a b l e a u



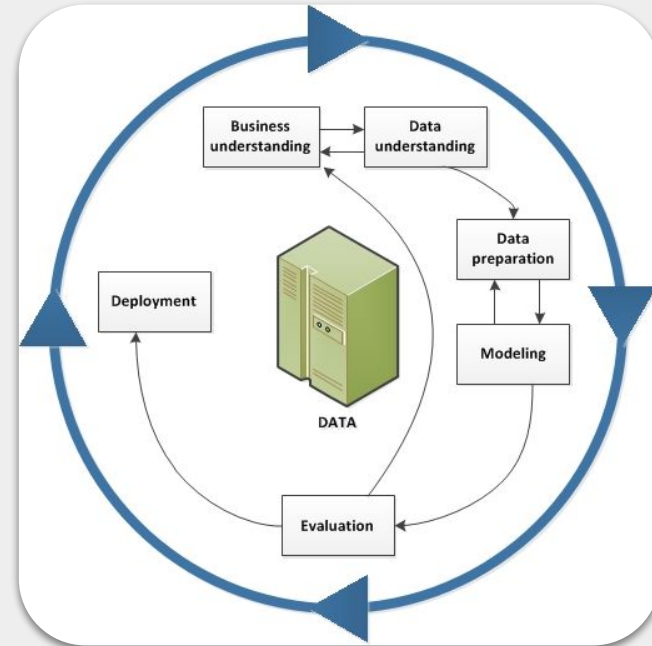
Looker

# Framework Data Science

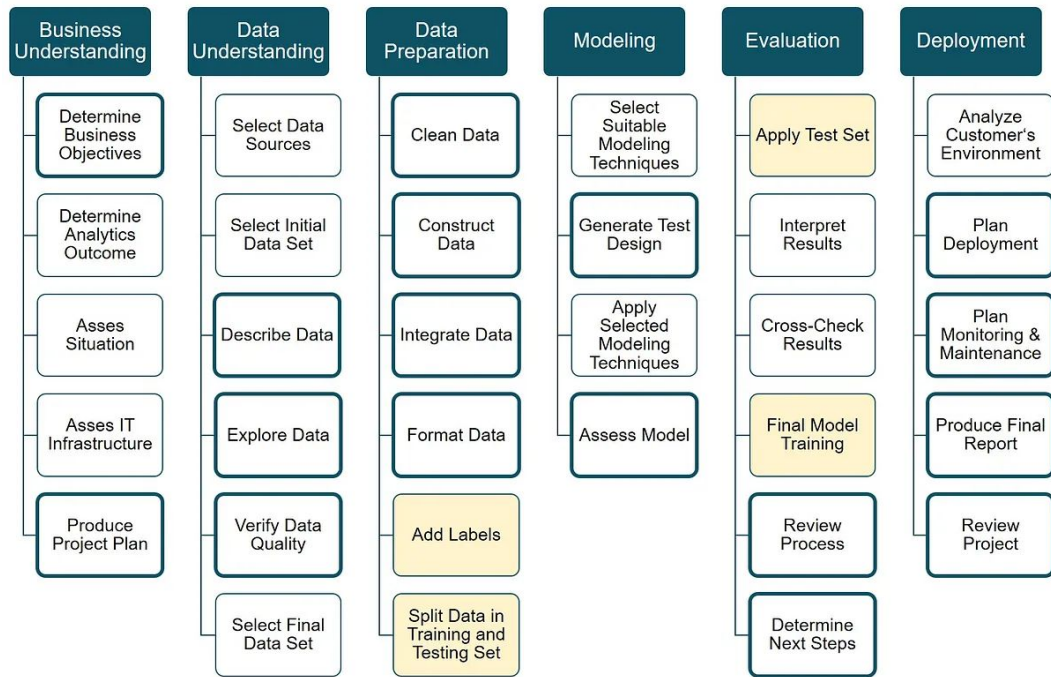
Framework/rangka Data Science bertujuan agar pekerjaan menjadi lebih terstruktur, efisien, dan hasilnya relevan dengan kebutuhan **bisnis atau pengguna**.

## Cross-Industry Standard Process for Data Mining [Proses Standar Lintas Industri untuk Data Mining]

- Sebuah proses yang memiliki enam fase untuk menggambarkan Data Science Life Cycle.



## CRISP-DM

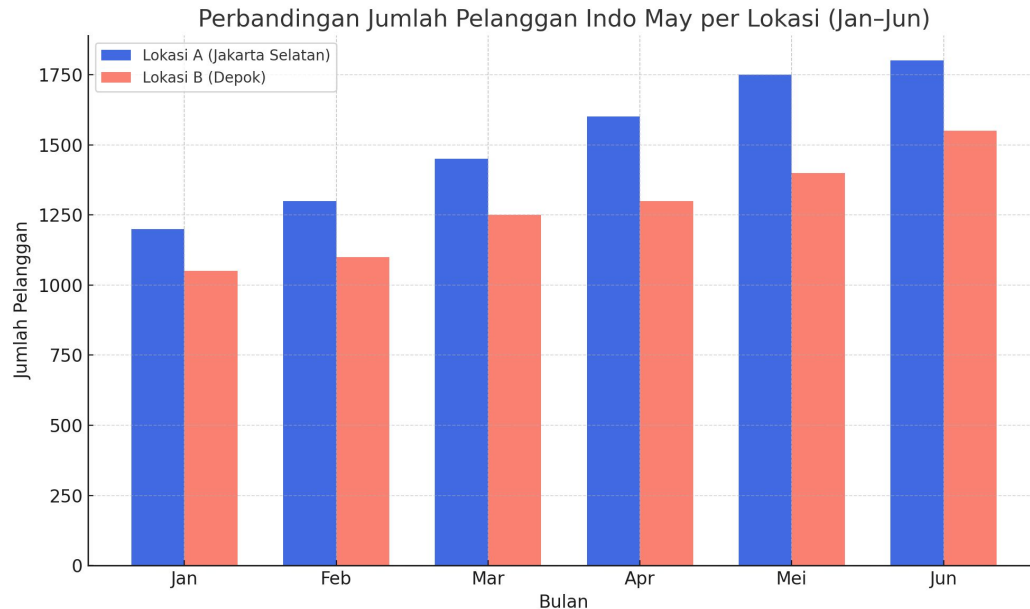


Jumlah pelanggan Indo May pada 2 lokasi berbeda selama 6 bulan terakhir

Bulan	Lokasi A (Jakarta Selatan)	Lokasi B (Depok)
Januari	1.200	1.050
Februari	1.300	1.100
Maret	1.450	1.250
April	1.600	1.300
Mei	1.750	1.400
Juni	1.800	1.550



# Jumlah pelanggan Indo May pada 2 lokasi berbeda selama 6 bulan terakhir



Data Science bukan tentang model  
paling keren, tapi tentang  
menyelesaikan **masalah nyata**  
**dengan cara paling bijak**

# Pengenalan Python

2

# Python

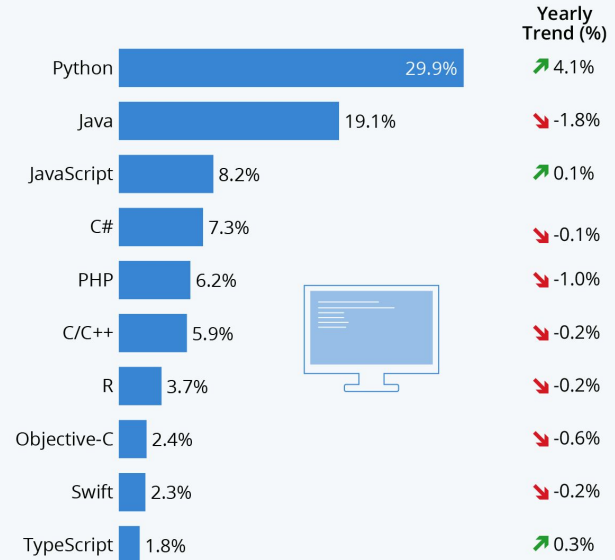
Bahasa pemrograman tingkat tinggi, dengan pengaplikasian dalam pemrograman web, penulisan skrip, komputasi ilmiah, dan kecerdasan buatan (AI).

## Kenapa Python?

- Mudah dipelajari
- Komunitas yang besar
- Memiliki **library** untuk berbagai kebutuhan
- Kode yang mudah dibaca dan maintain.

## Python Remains Most Popular Programming Language

Popularity of each programming language based on share of tutorial searches in Google



Yearly trend compares percent change from Feb 2019 to Feb 2020  
Sources: GitHub, Google Trends



# Library

Potongan kode yang dapat digunakan kembali. Library Python berisi kumpulan modul dan package terkait.



# Pandas

**Pandas** berasal dari **Panel Data**. Panel Data terdiri dari pengamatan selama beberapa periode waktu untuk individu yang sama. Diimplementasikan pada tahun 2008 oleh Wes McKinney.

Pandas dapat melakukan, sebagai berikut:



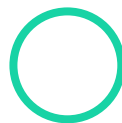
Import file  
csv/excel



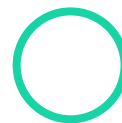
Menambah  
kolom



Filter data



Penanganan  
data hilang



Penggabungan  
data



Membuat grafik

# Fungsi dan Parameter

Fungsi adalah kode yang hanya berjalan saat dipanggil. Fungsi dapat menerima input saat dipanggil untuk menentukan sebuah perintah yang akan dijalankan yang disebut parameter.

## Contoh:

Input : `Fungsi(Parameter)`  
`print("hello world")`

Output : **hello world**



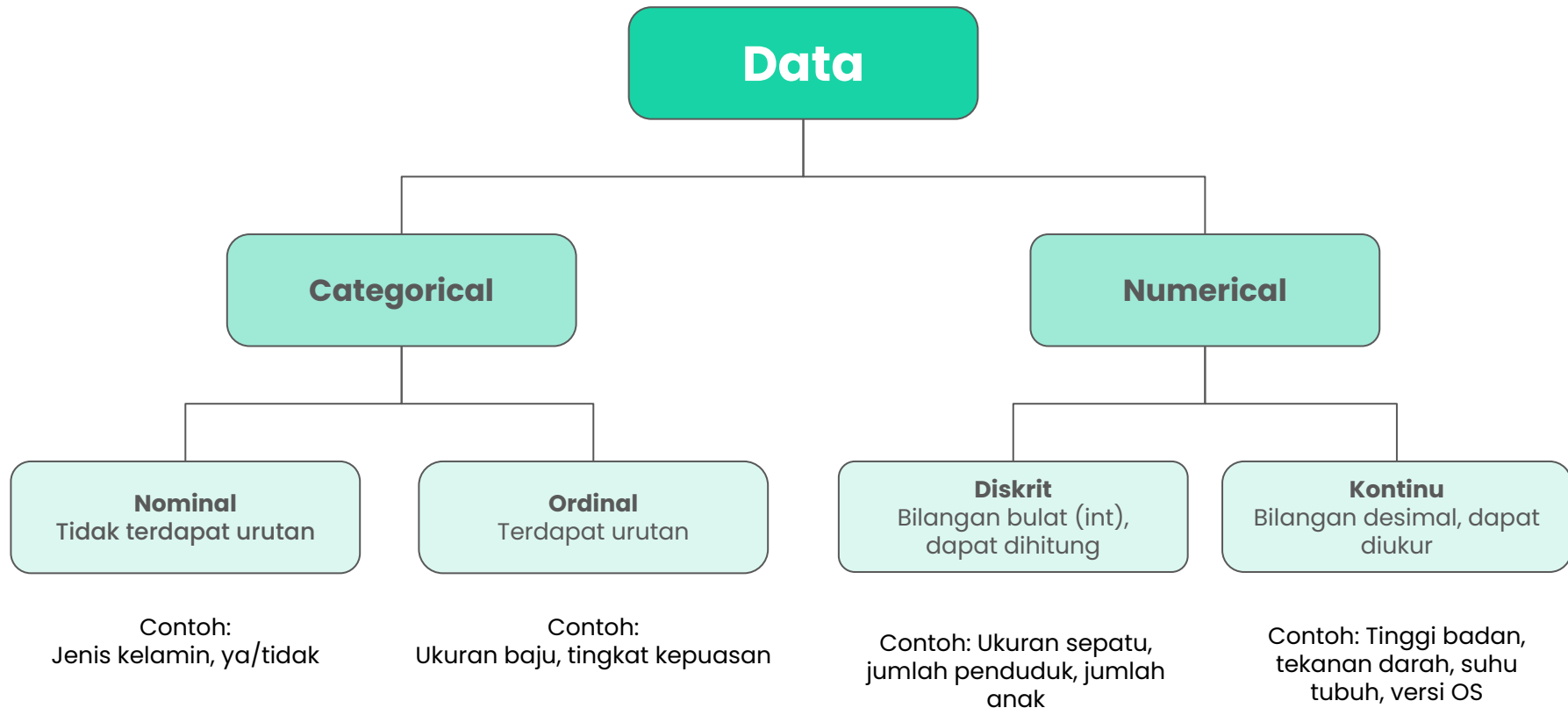
=Fungsi(Parameter)

F2							<code>=SUMIF(C2:C12, "&lt;"&amp;F1, B2:B12)</code>
	A	B	C	D	E	F	
1	Item	Amount	Delivery date		Before	4/11/2024	
2	Apples	\$250	4/10/2024		Total	\$700	
3	Bananas	\$450	4/10/2024				
4	Oranges	\$250	4/11/2024				



# Jenis-jenis Data

3



# Latihan

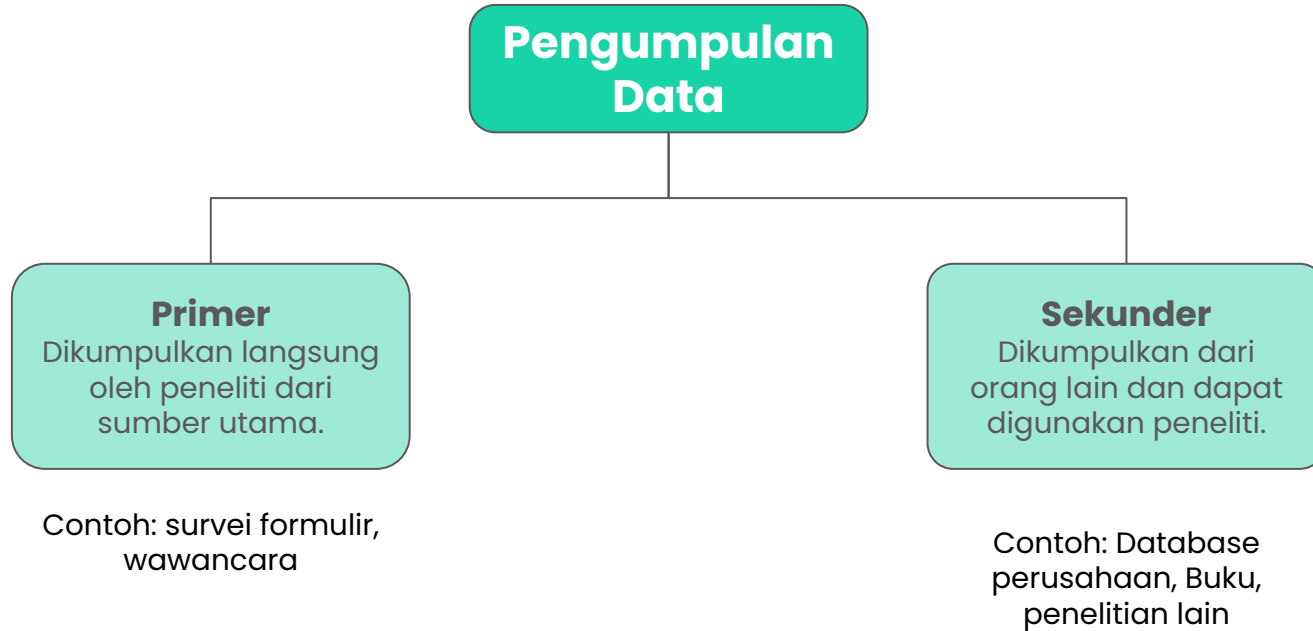
<b>Nama</b>	<b>Jumlah Buku Favorit</b>	<b>Tinggi Badan (cm)</b>	<b>Golongan Darah</b>	<b>Tingkat Minat Membaca</b>	<b>Waktu Tempuh ke Sekolah</b>
Nina	3	158.2	A	Tinggi	25 menit
Bima	5	165.0	O	Sedang	15 menit
Sari	2	149.7	B	Rendah	40 menit

# Latihan

Nama	Jumlah Buku Favorit	Tinggi Badan (cm)	Golongan Darah	Tingkat Minat Membaca	Waktu Tempuh ke Sekolah
Nina	3	158.2	A	Tinggi	25 menit
Bima	5	165.0	O	Sedang	15 menit
Sari	2	149.7	B	Rendah	40 menit

- Kontinu → Tinggi Badan (cm), Waktu Tempuh ke Sekolah
- Diskrit → Jumlah Buku Favorit
- Nominal → Nama, Golongan Darah
- Ordinal → Tingkat Minat Membaca

# Teknik pengumpulan data



# **Statistika Dasar**

**3**

# Statistika

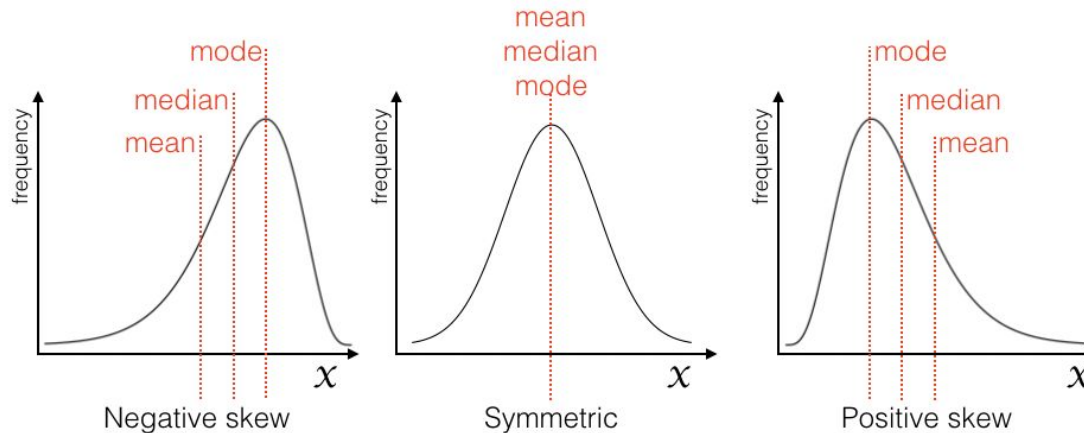
Ilmu yang mempelajari **cara mengumpulkan, mengolah, menganalisis, menyajikan, dan menginterpretasi data** untuk membuat keputusan.

$\Sigma$	Summation	$X$	An individual value, an observation
$S$	The standard deviation of sample data	$X_1$	A particular (1 <sup>st</sup> ) individual value
$\sigma$	The standard deviation of population data	$X_i$	For each, all, individual values
$S^2$	The variance of sample data	$\bar{X}$	The mean, average of sample data
$\sigma^2$	The variance of population data	$\bar{\bar{X}}$	The grand mean, grand average
$R$	The range of data	$\mu$	The mean of population data
$\bar{R}$	The average range of data	$p$	A proportion of sample data
$k$	Multi-purpose notation, i.e. # of subgroups, # of classes	$P$	A proportion of population data
$ y $	The absolute value of some term	$n$	Sample size
$>, <$	Greater than, less than	$N$	Population size
$\geq, \leq$	Greater than or equal to, less than or equal to		

# Measure of Central Tendency

Sebuah teori pemusatan data yang terdiri dari: Mean, Median, Modus

Mean vs median vs mode indicates skew





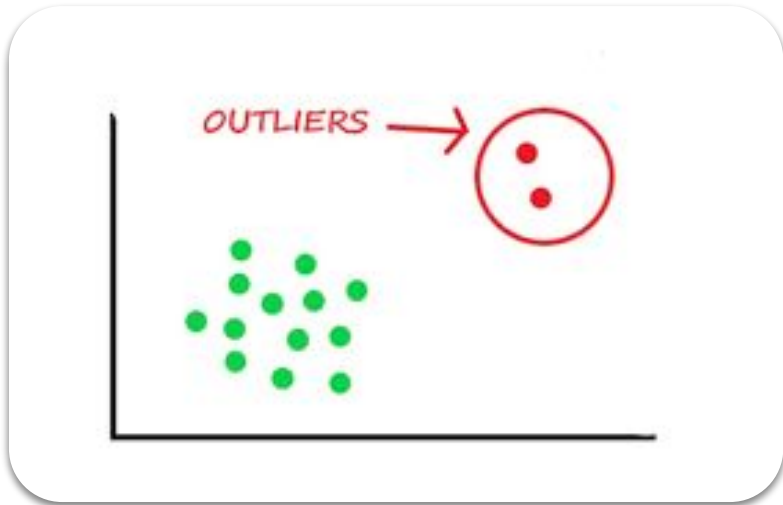
# Measure of Variability

Sebuah teori penyebaran data yang terdiri dari: Range, Variance, Standar Deviasi

1. Range → Jarak antara data terbesar dan terkecil
2. Variance → Mengukur rata-rata perbedaan dari setiap poin
3. Standard Deviation → Menunjukkan **seberapa tersebar** data terhadap nilai rata-rata
  - a. SD rendah : sebagian besar data berada di dekat nilai rata-rata
  - b. SD tinggi : data lebih tersebar

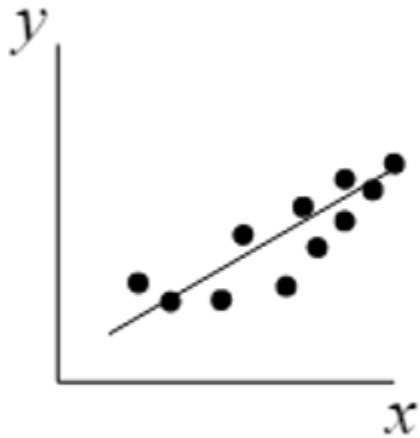
# Outlier

Data yang paling berbeda dari data lainnya, sehingga Outlier dapat memengaruhi nilai mean & standar deviasi.

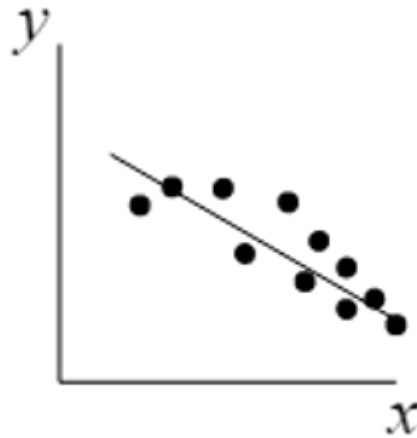


# Korelasi

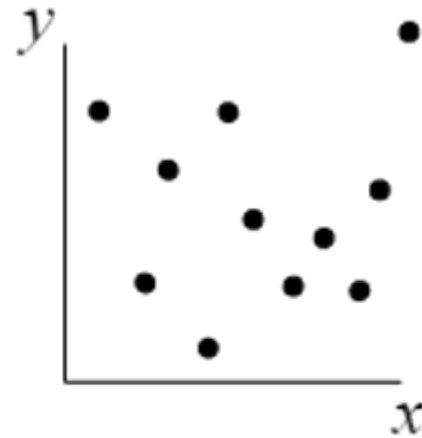
Mengukur hubungan antara dua variabel, Nilainya berisi rentang antara -1 hingga +1.



Positive



Negative



No correlation

**Hands On!**

**YOuRClass**

**Terima kasih**



Afifa Widhiyanthi



@afifadayu