

# Data Preparation

Informatics Research Group  
School of Electrical Engineering and Informatics  
Institut Teknologi Bandung

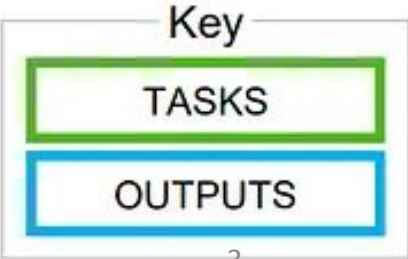
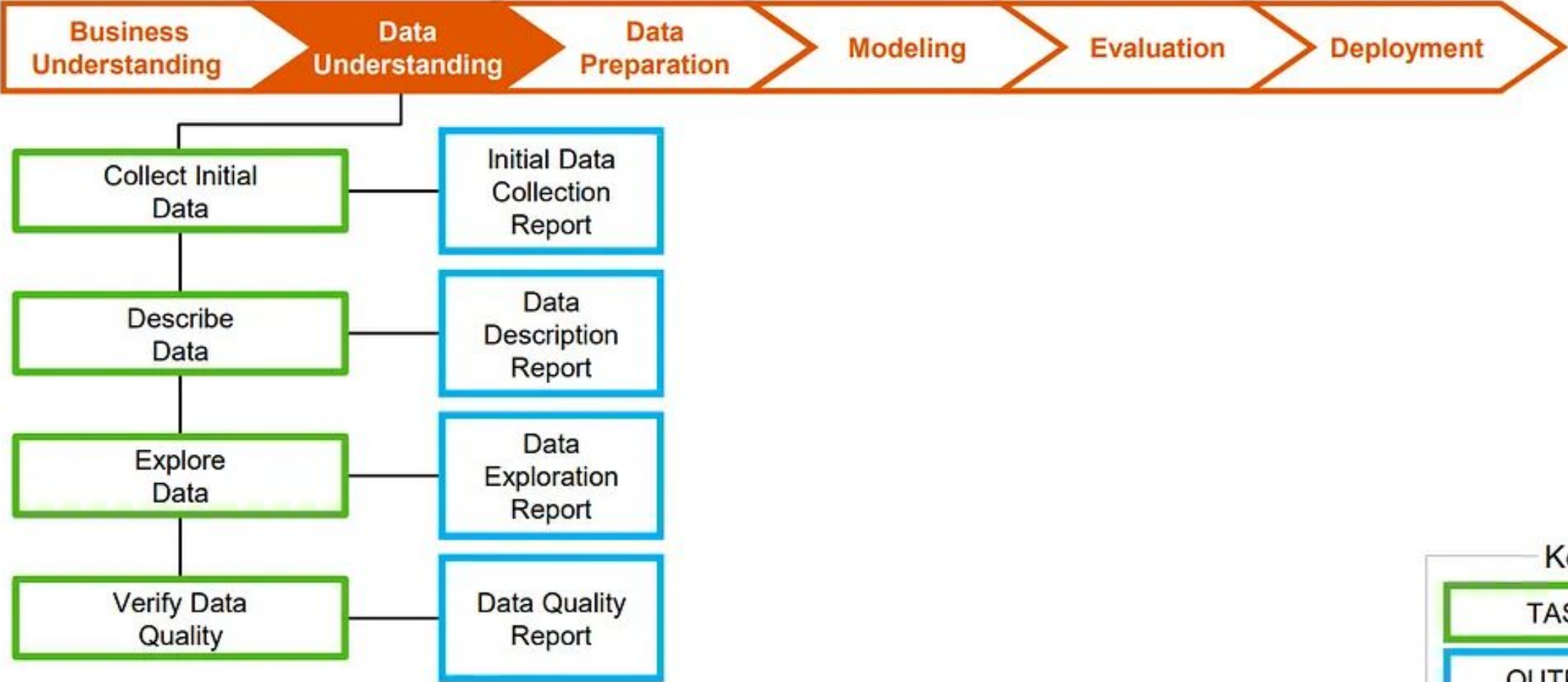
Sources:

Materi Pelatihan Associate Data Scientist – Pusat Artificial Intelligence ITB



Data Understanding Phase – Overview

**CRISP-DM – Phase 2: Data Understanding**



## Collect Initial Data

- **Collect Initial Data** or acquire the data and its access to the data listed in the projects resources. Collecting initial data also means you need to have a checklist of the dataset you have acquired, the dataset location, the methods to acquire the datasets, and record any problems encountered and any solutions to the problems for the other users or project members to be aware of.

## Describe Data

- **Describe Data** by examining the properties of the data acquired, provide a description report regarding the format of the data, quantity of data and even the records and fields in each table or datasets.

## Explore Data

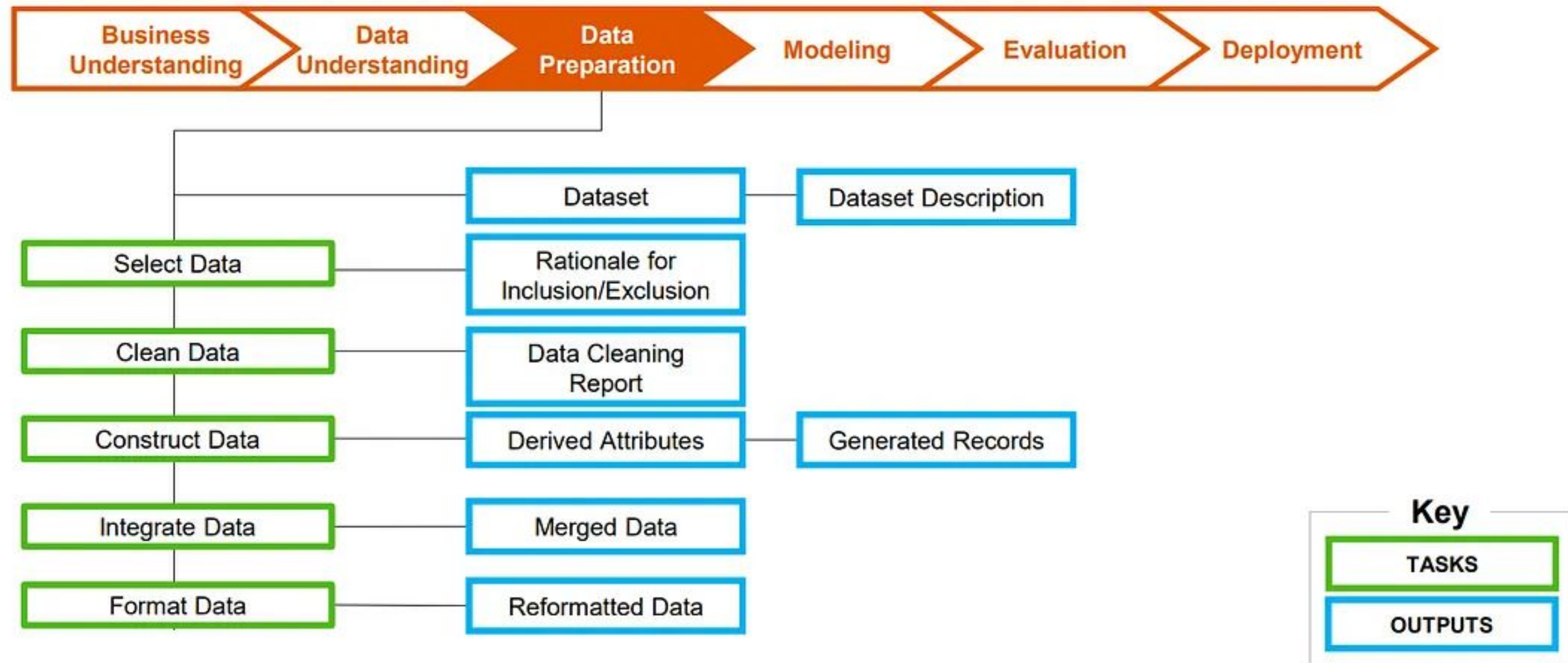
- **Explore Data** by using data science questions that can be quickly answered through querying, visualization, and reporting or summary report. In this stage, you will be able to find your first or initial hypothesis and their impact on the project.

## Verify Data Quality

- **Verify Data Quality** by examining if the data is complete. If the data has errors or are there missing values and if there is, what is the percentage of the missing values versus the overall data obtained.



## CRISP-DM – Phase 3: Data Preparation



## Persiapan Data

- Proses dimana data yang sesuai dikumpulkan, dipilih, dibersihkan, dan diorganisir sesuai dengan kebutuhan bisnis untuk digunakan pada tahap pemodelan.
- Dilakukan setelah tahap pemahaman data.
  - Laporan hasil pemahaman data digunakan sebagai dasar untuk menentukan aksi apa yang harus dilakukan pada tahap ini.





## Select Data

Select data or decide on the data to be used for analysis. One of the criteria in selecting the data is that it should be relevant to the data science goal that was identified in the business understanding phase. In selecting data, you also need to list the data to be excluded and included and the reasons for these decisions

## Clean data

Clean data by raising the data quality to the level required by the selected analysis techniques. Here, you also need to describe what decisions and actions were taken to address the data quality problems reported during the Verify Data Quality task of the Data Understanding Phase

## Construct data

Construct data by including derived attributes, entire new records, or transformed values for existing attributes. This may be conducting encoding methods especially for categorical variables or feature engineering

## Integrate data

Integrate data by combining from multiple tables or records to create records or values. SQL knowledge and skill is very important and would come in handy in this part

## Format data

Format data by transforming the data but not necessarily change its meaning but might be required by the modeling tool. An example would be transforming your data either by standardization or normalization



## Tujuan Persiapan Data

- Meningkatkan kualitas data
- Memudahkan pemodelan





# Berdasar Pengalaman dalam Persiapan Data

Data preparation is more than half of every data mining process

- Maxim of data mining: most of the effort in a data mining project is spent in data acquisition and preparation, and informal estimates vary from 50 to 80 percent



# Proses dalam Persiapan Data

## 1. Pemilihan Data

- a. *Record selection*
- b. *Feature selection*

## 2. Perbaikan Data

- a. Mengisi *missing values*
- b. Perbaikan error
- c. Penanganan *outlier*
- d. Penghapusan duplikasi

## 3. Konstruksi Data

- a. Reduksi data
- b. Mengubah representasi data
- c. Encoding

## 4. Integrasi Data

- a. *Data Join*
- b. *Append*





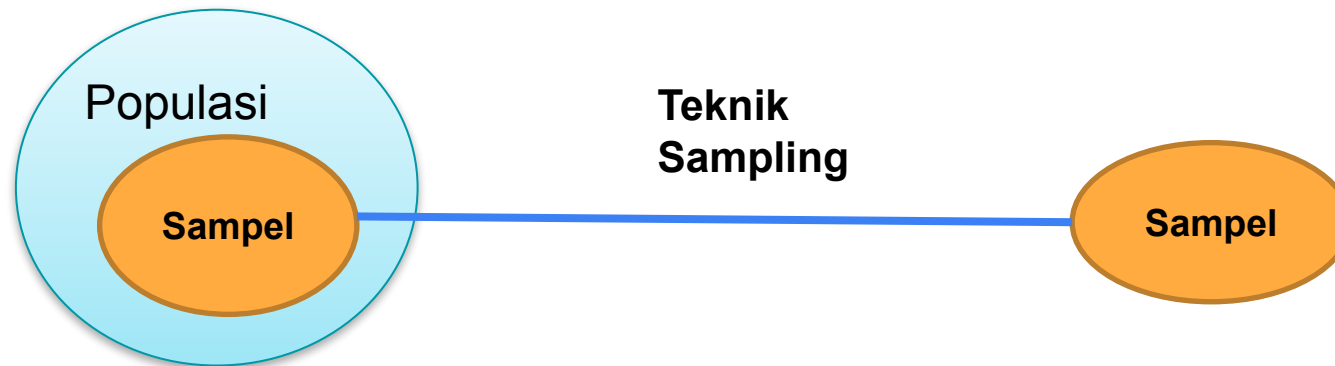
# Pemilihan Data

- a. Record selection
- b. Feature selection



## Record Selection (Sampling)

- **Sampling** adalah proses dalam analisis statistik dimana peneliti mengambil sejumlah **pengamatan** yang telah ditentukan sebelumnya dari **populasi** yang lebih **besar**.



- Pengambilan sampel memungkinkan peneliti untuk melakukan studi tentang kelompok besar dengan menggunakan sebagian **kecil** dari populasi.



# Kategori Metode Sampling

- **Probability Sampling:**

Teknik pengambilan sampel dimana sampel dari populasi yang lebih besar dipilih dengan menggunakan metode berdasarkan teori probabilitas.

- **Non-Probability Sampling**

Teknik pengambilan sampel dimana peneliti memilih sampel berdasarkan penilaian subjektif, bukan pemilihan acak.



# Data Preparation

Informatics Research Group  
School of Electrical Engineering and Informatics  
Institut Teknologi Bandung

Sources:

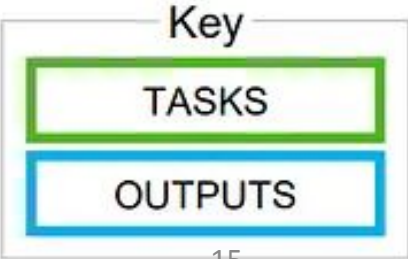
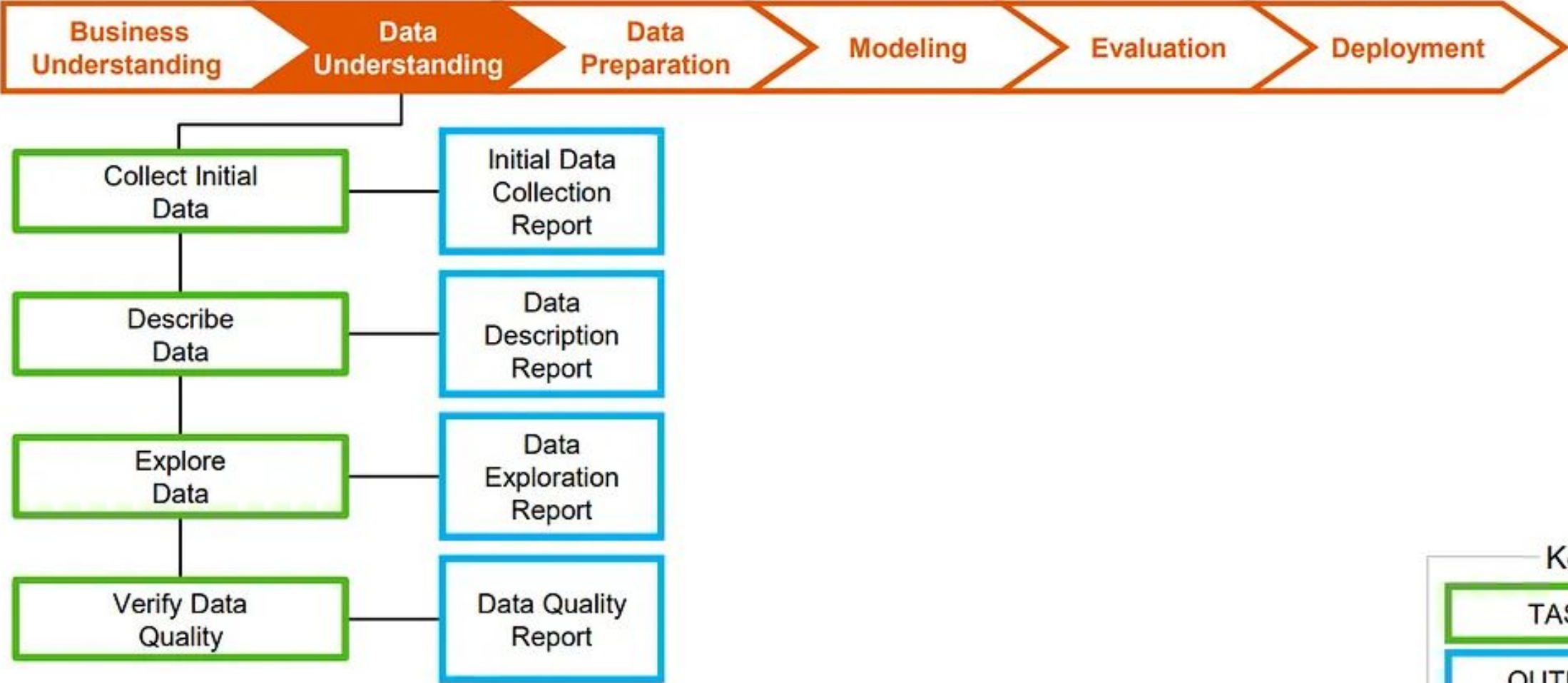
Materi Pelatihan Associate Data Scientist – Pusat Artificial Intelligence ITB





Data Understanding Phase – Overview

**CRISP-DM – Phase 2: Data Understanding**



## Collect Initial Data

- **Collect Initial Data** or acquire the data and its access to the data listed in the projects resources. Collecting initial data also means you need to have a checklist of the dataset you have acquired, the dataset location, the methods to acquire the datasets, and record any problems encountered and any solutions to the problems for the other users or project members to be aware of.

## Describe Data

- **Describe Data** by examining the properties of the data acquired, provide a description report regarding the format of the data, quantity of data and even the records and fields in each table or datasets.

## Explore Data

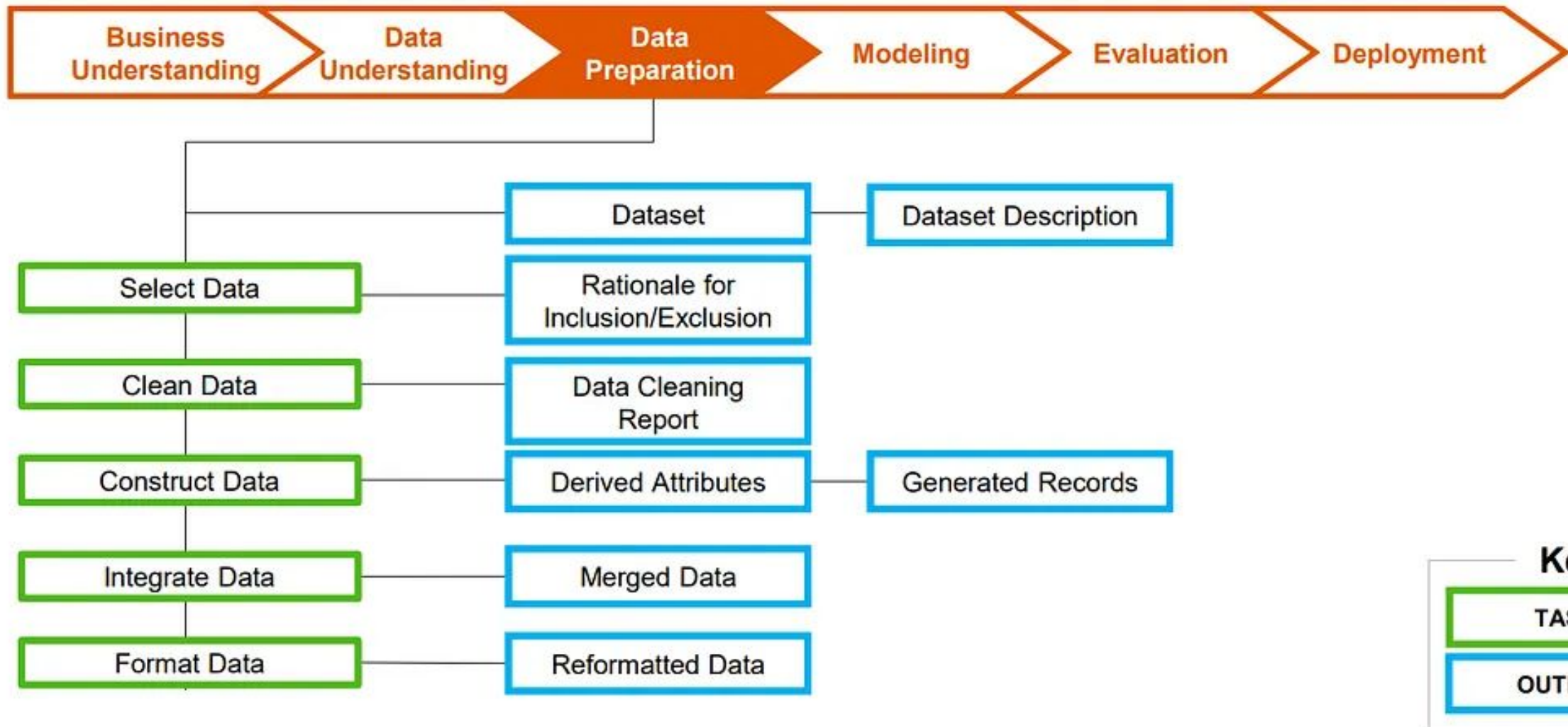
- **Explore Data** by using data science questions that can be quickly answered through querying, visualization, and reporting or summary report. In this stage, you will be able to find your first or initial hypothesis and their impact on the project.

## Verify Data Quality

- **Verify Data Quality** by examining if the data is complete. If the data has errors or are there missing values and if there is, what is the percentage of the missing values versus the overall data obtained.



## CRISP-DM – Phase 3: Data Preparation



## Persiapan Data

- Proses dimana data yang sesuai dikumpulkan, dipilih, dibersihkan, dan diorganisir sesuai dengan kebutuhan bisnis untuk digunakan pada tahap pemodelan.
- Dilakukan setelah tahap pemahaman data.
  - Laporan hasil pemahaman data digunakan sebagai dasar untuk menentukan aksi apa yang harus dilakukan pada tahap ini.



## Select Data

Select data or decide on the data to be used for analysis. One of the criteria in selecting the data is that it should be relevant to the data science goal that was identified in the business understanding phase. In selecting data, you also need to list the data to be excluded and included and the reasons for these decisions

## Clean data

Clean data by raising the data quality to the level required by the selected analysis techniques. Here, you also need to describe what decisions and actions were taken to address the data quality problems reported during the Verify Data Quality task of the Data Understanding Phase

## Construct data

Construct data by including derived attributes, entire new records, or transformed values for existing attributes. This may be conducting encoding methods especially for categorical variables or feature engineering

## Integrate data

Integrate data by combining from multiple tables or records to create records or values. SQL knowledge and skill is very important and would come in handy in this part

## Format data

Format data by transforming the data but not necessarily change its meaning but might be required by the modeling tool. An example would be transforming your data either by standardization or normalization





## Tujuan Persiapan Data

- Meningkatkan kualitas data
- Memudahkan pemodelan





# Hukum Persiapan Data

## Data Preparation Law (Data Mining Law 3)

Data preparation is more than half of every data mining process

- Maxim of data mining: most of the effort in a data mining project is spent in data acquisition and preparation, and informal estimates vary from 50 to 80 percent



# Proses dalam Persiapan Data

## 1. Pemilihan Data

- a. *Record selection*
- b. *Feature selection*

## 2. Perbaikan Data

- a. Mengisi *missing values*
- b. Perbaikan error
- c. Penanganan *outlier*
- d. Penghapusan duplikasi

## 3. Konstruksi Data

- a. Transformasi data
- b. *Encoding*
- c. *Dimensionality reduction*

## 4. Integrasi Data

- a. *Data Join*
- b. *Append*





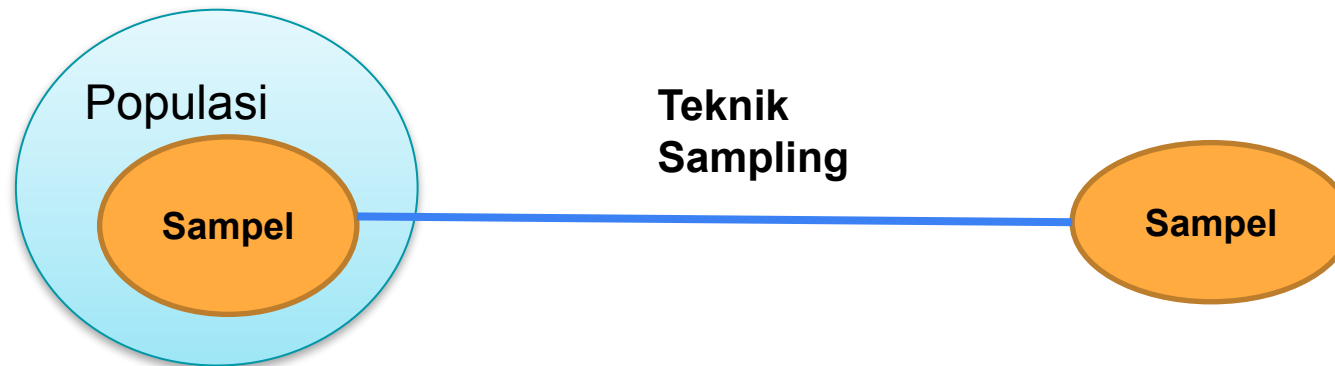
# Pemilihan Data

- a. Record selection
- b. Feature selection



## Record Selection (Sampling)

- **Sampling** adalah proses dalam analisis statistik dimana peneliti mengambil sejumlah **pengamatan** yang telah ditentukan sebelumnya dari **populasi** yang lebih **besar**.



- Pengambilan sampel memungkinkan peneliti untuk melakukan studi tentang kelompok besar dengan menggunakan sebagian **kecil** dari populasi.



# Kategori Metode Sampling

- **Probability Sampling:**

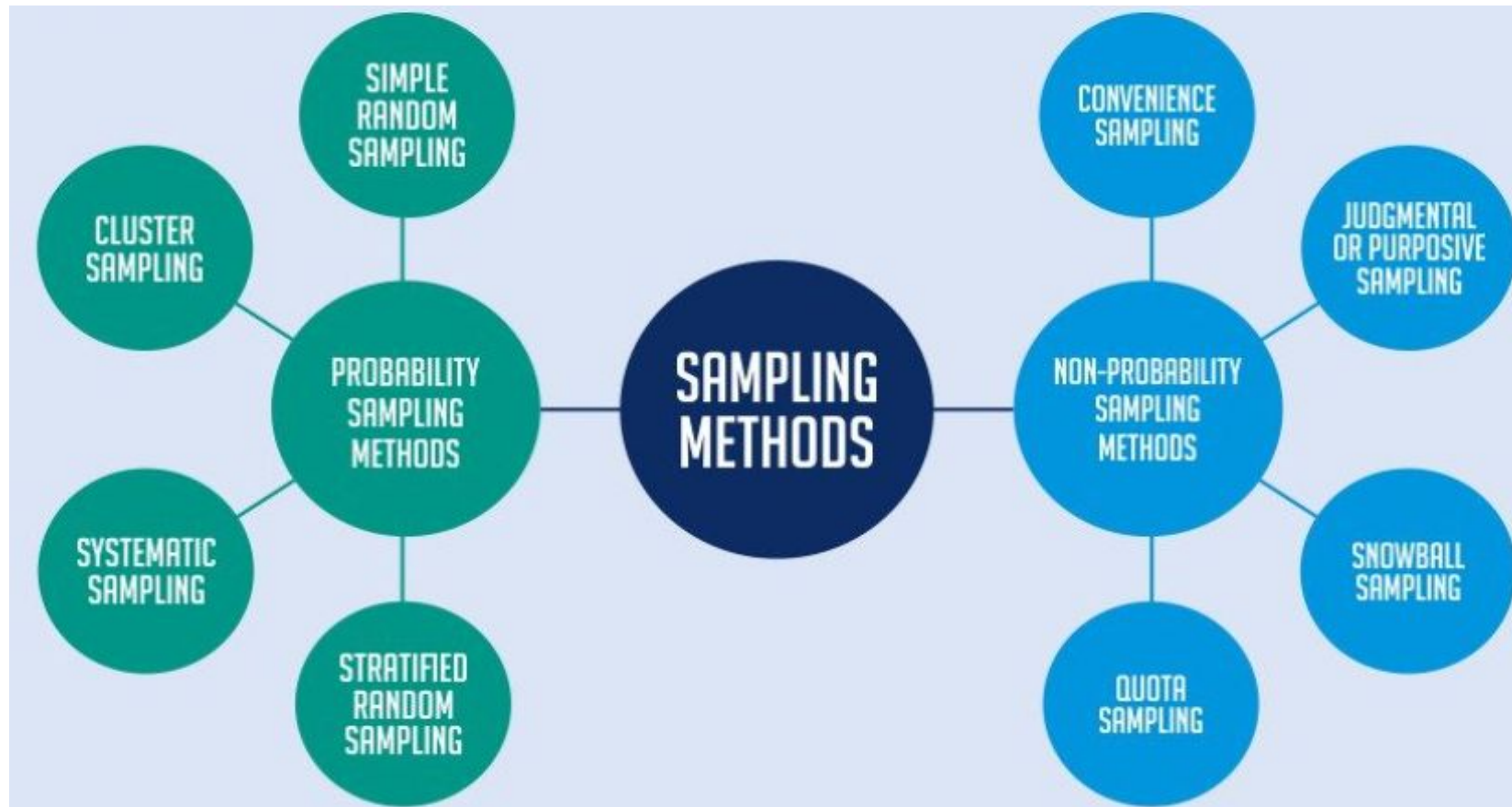
Teknik pengambilan sampel dimana sampel dari populasi yang lebih besar dipilih dengan menggunakan metode berdasarkan teori probabilitas.

- **Non-Probability Sampling**

Teknik pengambilan sampel dimana peneliti memilih sampel berdasarkan penilaian subjektif, bukan pemilihan acak.



# Kategori Metode Sampling





# Types Probability Sampling

## Types of probability sampling



### 1. Simple Random Sampling

Pengambilan sampel dari anggota populasi dengan menggunakan acak tanpa memperhatikan apapun.

### 2. Systematic Sampling

Pengambilan sampel secara sistematis dengan interval (jarak) tertentu antar sampel yang terpilih.

### 3. Stratified Random Sampling

Pengambilan sampel dengan cara membagi populasi ke dalam kelompok-kelompok yang homogen (disebut strata), dan dari tiap stratum tersebut diambil sampel secara acak.

### 4. Cluster Sampling

Pengambilan sampel dilakukan terhadap sampling unit, dimana sampling unitnya terdiri dari satu kelompok (cluster). Tiap item (individu) di dalam kelompok yang terpilih akan diambil sebagai sampel.



# When to use probability sampling?

1

## When you want to reduce the sampling bias

Probability sampling leads to higher quality findings because it provides an unbiased representation of the population.



2

## When the population is usually diverse

This sampling method will help pick samples from various socio-economic strata, background, etc. to represent the broader population.



3

## To create an accurate sample

Researchers use proven statistical methods to draw a precise sample size to obtain well-defined data.



Learn more:  
[www.questionpro.com/blog/probability-sampling/](https://www.questionpro.com/blog/probability-sampling/)

QuestionPro



# Tipe Non-Probability Sampling

- **Purposive Sampling**

Sample dipilih berdasarkan pertimbangan tertentu untuk memperoleh sampel dengan karakteristik yang dikehendaki.

- **Quota Sampling**

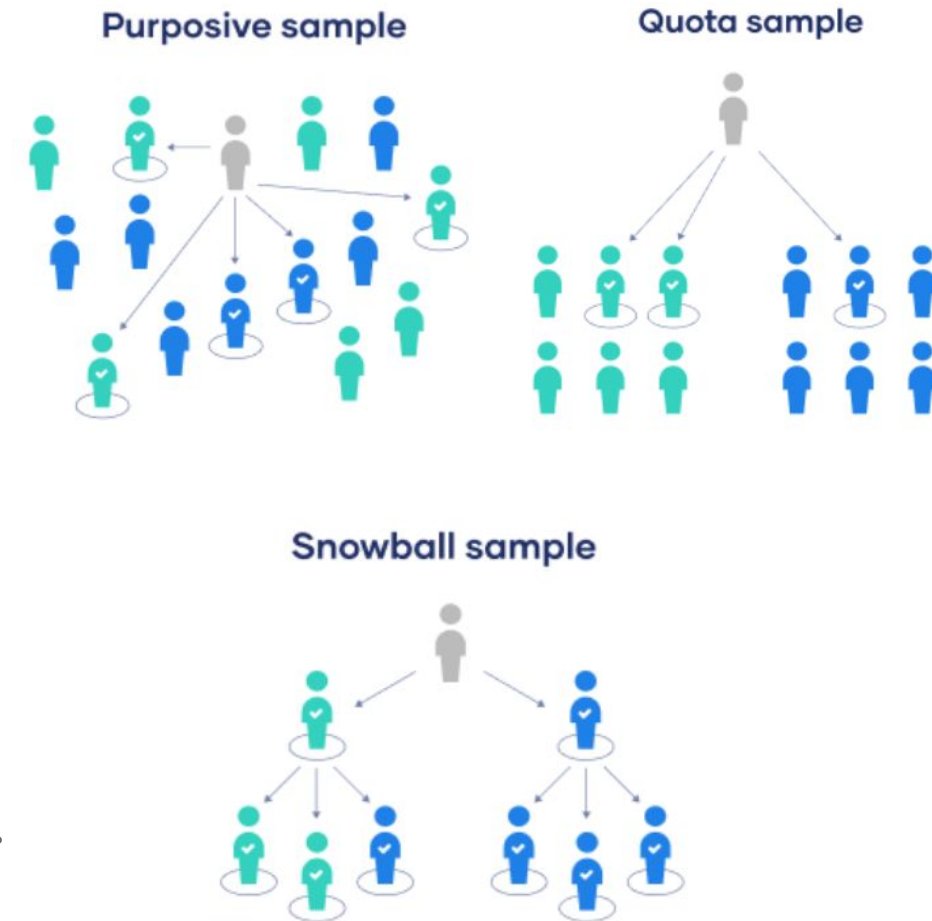
Pengambilan sampel hanya berdasarkan pertimbangan peneliti dengan besar dan kriteria sampel ditentukan lebih dahulu.

- **Saturation Sampling**

Semua anggota populasi digunakan sebagai sampel karena jumlah populasi tidak banyak, atau ingin membuat generalisasi dengan kesalahan sangat kecil.

- **Snowball Sampling**

Sampel diambil secara berantai, mulai dari ukuran sampel yang kecil semakin menjadi besar.



# Tahapan Melakukan Sampling

Step 1

Identify and define Target Population

Step 2

Select Sampling Frame

Step 3

Choose Sampling methods

Step 4

Determine Sample Size

Step 5

Collect the required Data







# Feature Selection



# Feature Selection

- Pemilihan fitur adalah proses pengurangan jumlah variabel masukan saat mengembangkan model machine learning.







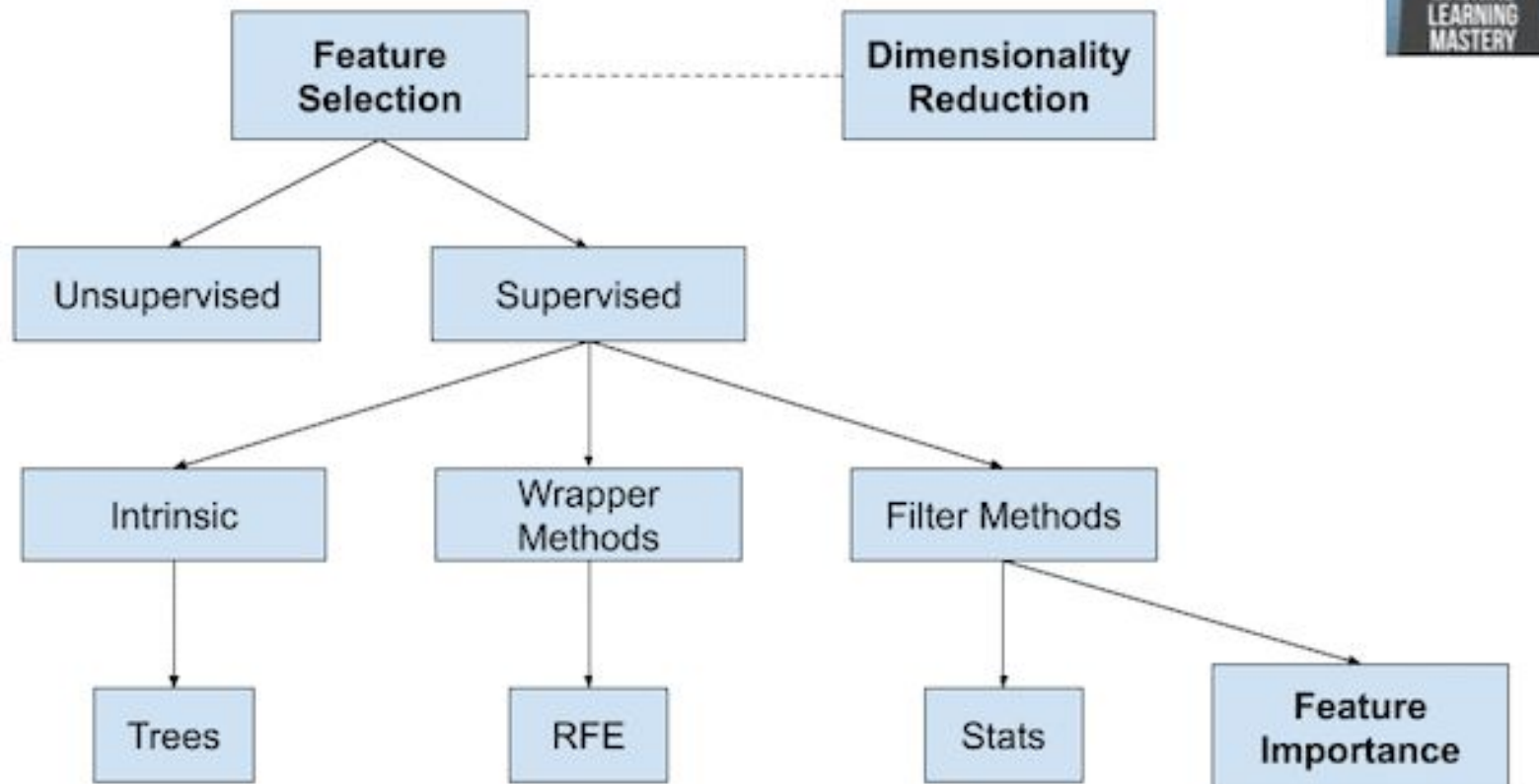
## Tujuan Feature Selection

- Mengurangi jumlah variabel masukan untuk mengurangi biaya komputasi pemodelan
- Dalam beberapa kasus, ditujukan untuk meningkatkan kinerja model.



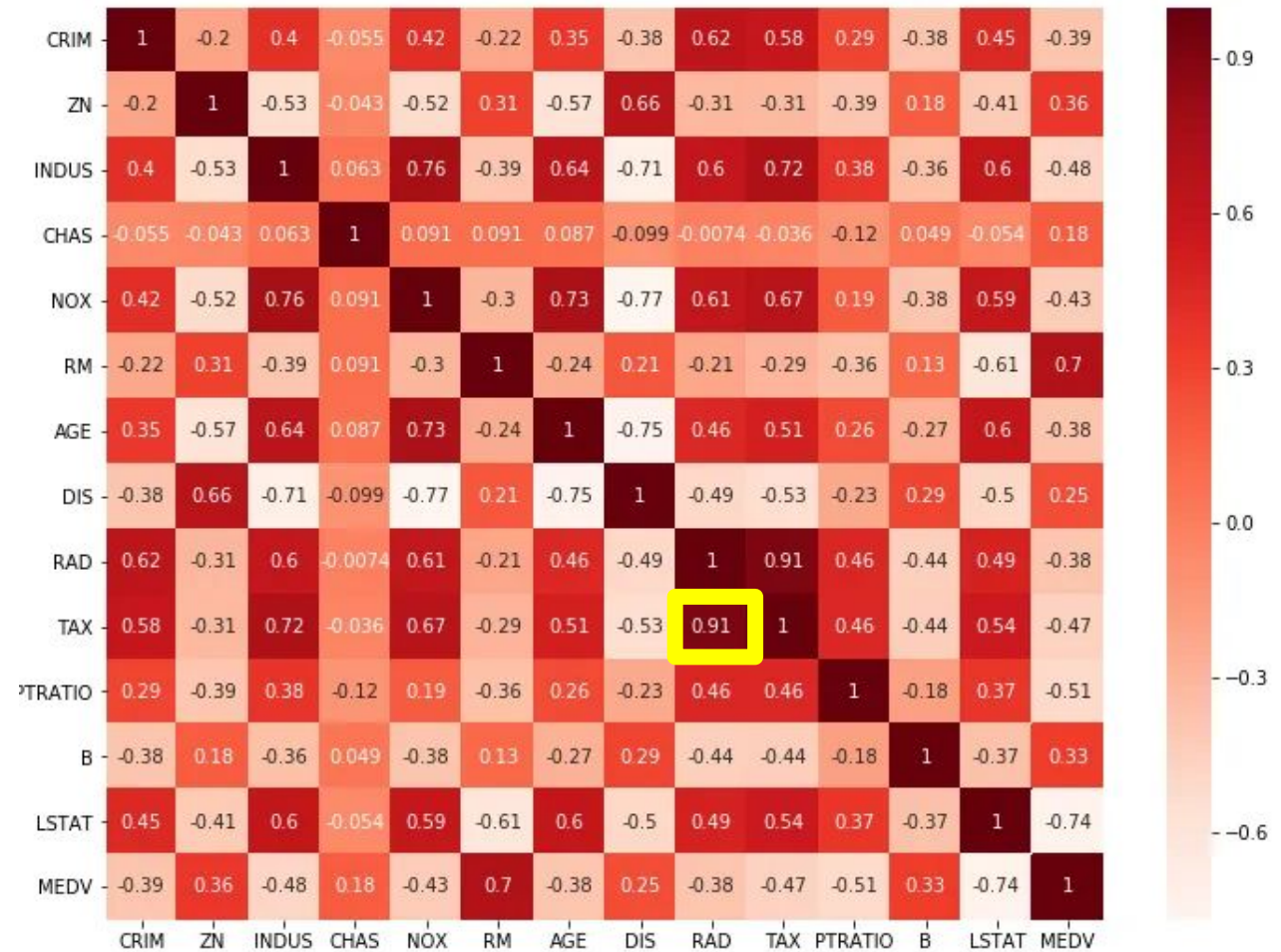
# Teknik Feature Selection

## Overview of Feature Selection Techniques



# Teknik Feature Selection berbasis Unsupervised

- Tidak melibatkan variable target
- Jika ada dua fitur yang memiliki korelasi yang kuat, hilangkan salah satu fitur



## Teknik Feature Selection berbasis Filter

- Mengevaluasi hubungan antara setiap variabel masukan dan variabel target (kelas) menggunakan:
  - statistik
  - memilih fitur yang penting (*feature importance*)
- Memilih variabel masukan yang memiliki hubungan paling kuat dengan variabel target.



## Teknik Statistik untuk Feature Selection berbasis Filter

- Memilih variabel masukan yang memiliki hubungan yang paling kuat dengan variabel target secara statistik.
- Pilihan ukuran statistik bergantung pada tipe data variabel masukan dan keluaran.

Tipe Variable Masukan	Tipe Variable Target	Teknik Statistik
Numerik	Numerik	<ul style="list-style-type: none"><li>● Pearson's correlation coefficient (linear)</li><li>● Spearman's rank coefficient (nonlinear)</li></ul>
Numerik	Categorical	<ul style="list-style-type: none"><li>● ANOVA correlation coefficient (linear)</li><li>● Kendall's rank coefficient (nonlinear)</li></ul>
Categorical	Numerik	<ul style="list-style-type: none"><li>● ANOVA correlation coefficient (nonlinear)</li><li>● Kendall's rank coefficient (linear)</li></ul>
Categorical	Categorical	<ul style="list-style-type: none"><li>● Chi-Squared test (contingency tables)</li><li>● Mutual Information</li></ul>



# Proses dalam Persiapan Data (review)

## 1. Pemilihan Data

- a. *Record selection*
- b. *Feature selection*

## 2. Perbaikan Data

- a. Mengisi *missing values*
- b. Perbaikan error
- c. Penanganan *outlier*
- d. Penghapusan duplikasi

## 3. Konstruksi Data

- a. Reduksi data
- b. Mengubah representasi data
- c. Encoding

## 4. Integrasi Data

- a. *Data Join*
- b. *Append*







# Proses Pembersihan Data

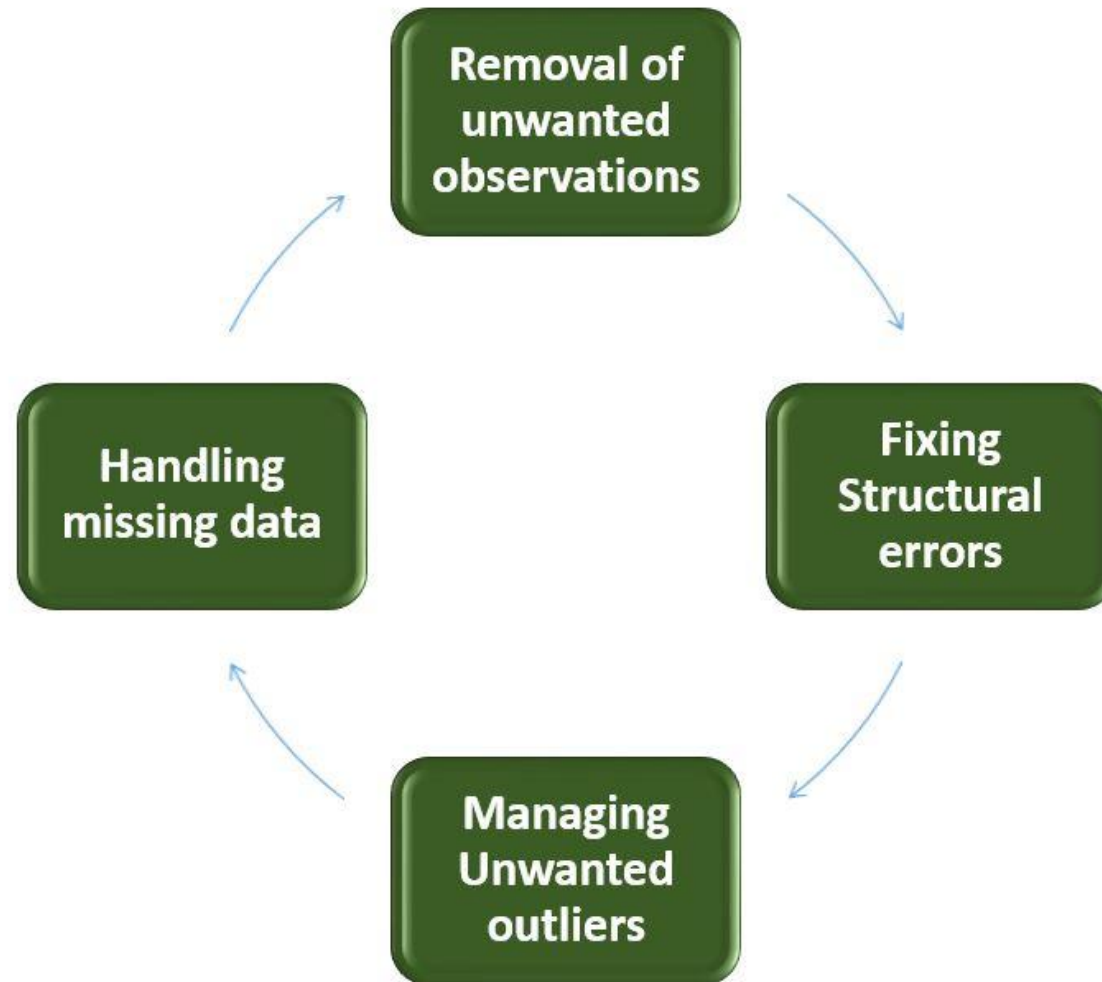


Associate Data Science Online Training

EDUNEX ITB



# Langkah-langkah Pembersihan Data



# Jenis Kesalahan Data dan Alternatif Cara Mengatasinya

## 1. Kesalahan nilai fitur di dalam sebuah dataset (1)

### Jenis Error

- Kesalahan selama proses data entry
- Nilai fitur yang meragukan/tidak mungkin (*impossible values*)
- Pengulangan white space (*unreadable or undetected characters*)

### Tindakan Mengatasinya

- Meningkatkan kapasitas staf data entry
- Menggunakan dukungan software untuk memvalidasi data.
- Perbaiki Data
- Menggunakan software untuk menghilangkan unreadable atau undetected characters dari data input.



# Jenis Kesalahan Data dan Alternatif Cara Mengatasinya

## 1. Kesalahan nilai fitur di dalam sebuah dataset (2)

### Jenis Error

- Tidak ada nilai fitur (*missing value*)
- Pencilan data (*outlier*)

### Tindakan Mengatasinya

- Perbaiki Data
- Menghapus sampel
- Perbaiki Data
- Menghapus sampel



# Jenis Kesalahan Data dan Alternatif Cara Mengatasinya

## 2. Ketidakkonsistenan nilai fitur di dalam sebuah dataset (2)

### Jenis Error

- Deviasi dari nilai fitur yang standar
- Perbedaan unit pengukuran  
(contoh: centimeter dengan meter)
- Perbedaan level agregasi (Contoh: akumulasi per hari dengan per minggu)

### Tindakan Mengatasinya

- Meningkatkan kapasitas staf data entry
- Menggunakan dukungan software untuk memvalidasi data
- Perbaiki data
- Menghitung ulang.
- Menyamakan tingkat pengukuran menggunakan teknik agregasi atau ekstrapolasi



## Contoh Data Kotor: Data Mengandung Missing Values

Row No.	age	education	balance	duration	campaign	y
1	58	tertiary	2143	261	1	no
2	44	secondary	29	151	1	no
3	33	secondary	?	76	1	no
4	47	unknown	1506	92	1	no
5	33	unknown	1	198	1	no
6	35	tertiary	231	139	1	no
7	28	tertiary	?	217	1	no
8	42	tertiary	2	380	1	no
9	58	primary	121	50	1	no
10	43	secondary	593	55	1	no
11	41	secondary	270	222	1	no
12	29	secondary	?	137	1	no
13	53	secondary	?	517	1	no
14	58	unknown	71	71	1	no
15	57	secondary	162	174	1	no
16	51	primary	229	353	1	no
17	45	unknown	13	98	1	no
18	57	primary	52	38	1	no

ExampleSet (45,211 examples, 0 special attributes, 6 regular attributes)





## Contoh Data Kotor: Data Tidak Konsisten

Row No.	nama_nasa...	jenis_kelamin	umur	jml_pinjaman	jkw
1	x1	P	40	345000	1
2	x2	L	31	350000	7
3	x3	L	29	649926	6
4	x4	P	2	459168	19
5	x5	WANITA	34	3055499	8
6	x6	L	49	2000000	19
7	x7	L	29	8333334	10
8	x8	L	27	4435001	8
9	x9	L	29	560000	19
10	x10	LAKI-LAKI	49	1443750	15
11	x11	LAKI-LAKI	42	3066000	10
12	x12	PRIA	26	4071669	20
13	x13	L	29	228655000	19
14	x14	L	55	840000	4
15	x15	L	38	3000000	24
16	x16	WANITA	29	1640000	19
17	x17	L	41	930000.010	4



# Pembersihan Data Kuantitatif

- Data kuantitatif: bilangan bulat atau bilangan floating point dalam berbagai bentuk (set, tensor, deret waktu)
- Tantangan: konversi unit (terutama untuk unit yang mudah berubah seperti mata uang)
- Teknik pembersihan: Normalisasi data
  - perbaikan data missing value
  - perbaikan data outlier
  - perbaikan data salah
  - perbaikan data tidak konsisten



# Pembersihan Data Kategori (Kualitatif)


- Data kategori: nama atau kode untuk menetapkan data ke dalam grup, tidak ada urutan atau jarak yang ditentukan
- Masalah umum: salah mengeja saat entri data
- Dasar teknik pembersihan: Normalisasi data



# Pembersihan Data Text

- Tantangan utama:
  - Duplikasi data
  - Salah ketik
  - Karakter yang salah
- Cara pembersihan data:
  - Penghilangan duplikat
  - Perbaikan salah ketik
  - Penghilangan karakter yang error





# Teknik Pembersihan Data: Duplikat Data

Menghapus duplikat data



## Teknik Pembersihan Data: Missing value, tidak konsisten, outlier

### Urutan pembersihan berdasarkan prioritas:

1. Data yang salah/kosong diisi dengan nilai sebenarnya
  - Lakukan validasi data dan data diisi/diganti dengan nilai sebenarnya
  - Validasi bisa dilakukan secara manual atau menggunakan kode program
  - Contoh menggunakan kode program:
    - Jenis kelamin laki-laki, atribut “hamil/tidak” pasti berisi “Tidak”
    - Usia bisa dihitung ulang dari tanggal lahir





# Teknik Pembersihan Data: Missing value, tidak konsisten, outlier

## Urutan pembersihan berdasarkan prioritas:

2. Jika tidak diketahui nilai seharusnya:

- Diisi dengan nilai yang paling mungkin
  - Nilai yang sama untuk label yang sama
  - Data non time series:

Nilai tengah

- mean : data numerik jika outlier sudah dihilangkan
- median : data numerik jika outlier belum dihilangkan
- modus : data kategorikal

Dari:

- Kelompok data yang sama, misal nilai gaji dari level pekerjaan yang sama
- Data dengan label yang sama
- Keseluruhan data untuk data non time series
- Data time series:
  - Nilai data sebelumnya atau setelahnya



## Teknik Pembersihan Data: Missing value, tidak konsisten, outlier

**Urutan pembersihan berdasarkan prioritas:**

3. Dihapus, jika tidak memungkinkan diperbaiki (jumlah data masih cukup banyak)

**Jika dihapus data menjadi sangat sedikit** (tidak cukup digunakan sebagai data latih untuk membangun model):

Harus mengumpulkan atau mencari lagi data lain



# Proses dalam Persiapan Data (review)

## 1. Pemilihan Data

- a. *Record selection*
- b. *Feature selection*

## 2. Perbaikan Data

- a. Mengisi *missing values*
- b. Perbaikan error
- c. Penanganan *outlier*
- d. Penghapusan duplikasi

## 3. Konstruksi Data

- a. Reduksi data
- b. Mengubah representasi data
- c. Encoding

## 4. Integrasi Data

- a. *Data Join*
- b. *Append*





# Konstruksi Data: Reduksi Data



# Reduksi Data

- Reduksi data dilakukan untuk memperoleh dataset yang lebih sedikit dari sisi volume, namun tetap menghasilkan analisis yang sama
- Mengapa perlu dilakukan:
  - Jumlah data yang dikumpulkan (dari basis data atau data warehouse) sangat besar (orde terabytes)
  - Analisis data terhadap data yang sangat besar akan membutuhkan waktu yang sangat lama



# Metode Reduksi Data

- **Reduksi Dimensi (*Dimensionality reduction*)**  
Mengurangi jumlah kolom/atribut data
- **Pengurangan Data (*Numerosity reduction*)**  
Mengurangi jumlah instan atau sample data





# Teknik Reduksi Dimensi (Dimensionality reduction)

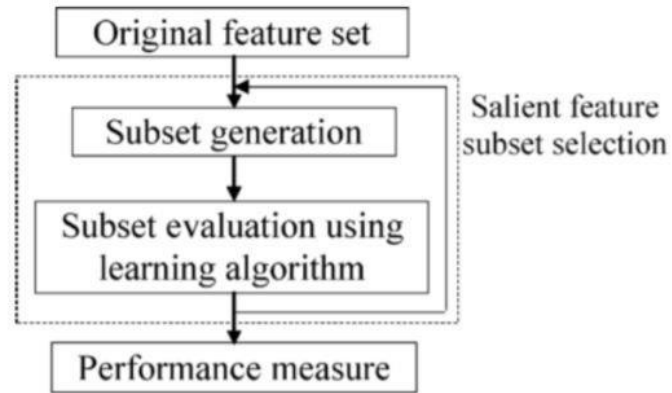
1. Feature Extraction
2. Feature Selection
  - a. Filter Approach
  - b. Wrapper Approach
  - c. Embedded Approach



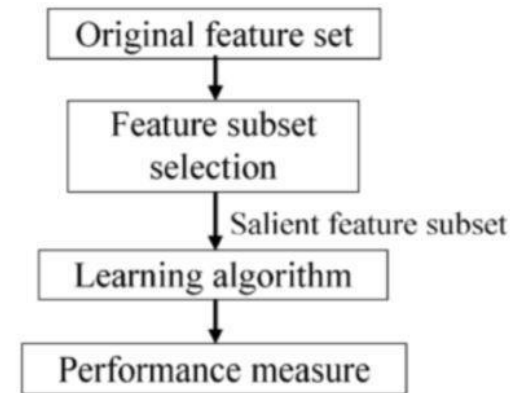
## Reduksi Dimensi dengan Ekstraksi Fitur: **Principal Component Analysis (PCA)**

1. Normalisasi data input: Setiap atribut berada dalam kisaran yang sama
2. Hitung vektor ortonormal (unit), yaitu Komponen utama
3. Setiap data input (vektor) adalah kombinasi linear dari nilai  $k$  vektor komponen utama.
4. Komponen utama diurutkan dalam urutan decreasing (menurun) “Signifikansi” atau kekuatan
5. Karena komponen diurutkan, ukuran data dapat dikurangi dengan menghilangkan komponen yang lemah, yaitu komponen-komponen dengan varian rendah.

# Reduksi Dimensi dengan Seleksi Fitur



Wrapper Approach



Filter Approach

1. Dalam pendekatan wrapper, fitur-fitur digunakan untuk melatih model pembelajaran yang telah ditentukan. Fitur dikurangi secara bertahap dengan melihat kinerja model menaik atau menurun ketika fitur tersebut dihilangkan. Dapat menggunakan forward selection, backward elimination, randomized hill climbing, dll.
2. Dalam pendekatan filter, analisis statistik dari set fitur diperlukan, tanpa menggunakan model pembelajaran apapun. Dapat menggunakan information gain, chi square, log likelihood ratio, dll. (sudah dibahas sebelumnya)
3. Pendekatan yang embedded memanfaatkan kekuatan pelengkap pendekatan wrapper dan filter. Dapat menggunakan decision tree, weighted naïve bayes, dll.

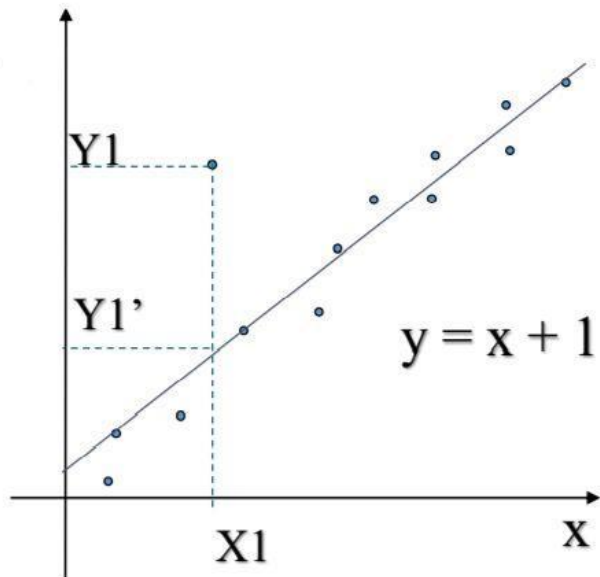
# Teknik Pengurangan Data (Numerosity Reduction)

1. Parametrik:
  - a. Regresi
  - b. log-linear model, dll
2. Non-Parametrik
  - a. Histogram
  - b. Clustering
  - c. Sampling



# Numerosity Reduction dengan Pendekatan Parametrik

## Regresi:



Regresi linier memodelkan hubungan antara dua atribut dengan memodelkan persamaan linier ke kumpulan data.

Misalkan kita perlu memodelkan fungsi linier antara dua atribut.

$$y = wx + b$$

y adalah atribut respons

x adalah atribut prediktor.

Jika kita membahas dari segi data mining, atribut x dan atribut y adalah atribut numerik database, sedangkan w dan b adalah koefisien regresi.

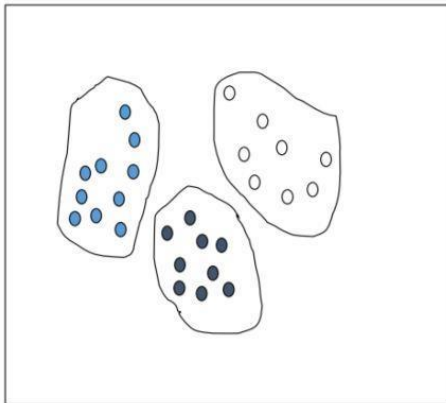
Nilai w dan b dijadikan data baru untuk menggantikan n data pada kelompok data yang bisa diwakili oleh model regresi



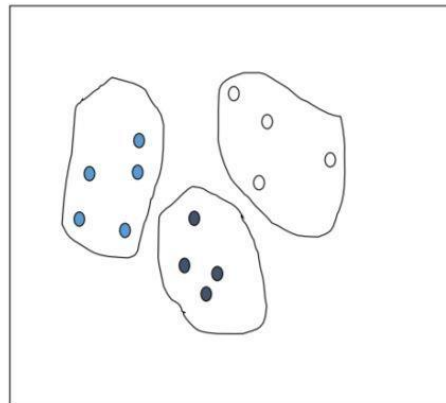
# Numerosity Reduction dengan Pendekatan Non-Parametrik

## Clustering:

Raw Data



Cluster/Stratified Sample



Teknik clustering mengelompokkan objek-objek yang mirip, sehingga objek-objek dalam satu cluster akan mirip satu sama lain, tetapi berbeda dengan objek-objek di cluster lain.

Seberapa mirip objek di dalam cluster dapat dihitung menggunakan fungsi jarak.

**Centroid dari cluster digunakan sebagai data baru untuk mewakili data-data lain pada cluster yang sama**





# Konstruksi Data: Mengubah Representasi Data (Transformasi & Encoding)



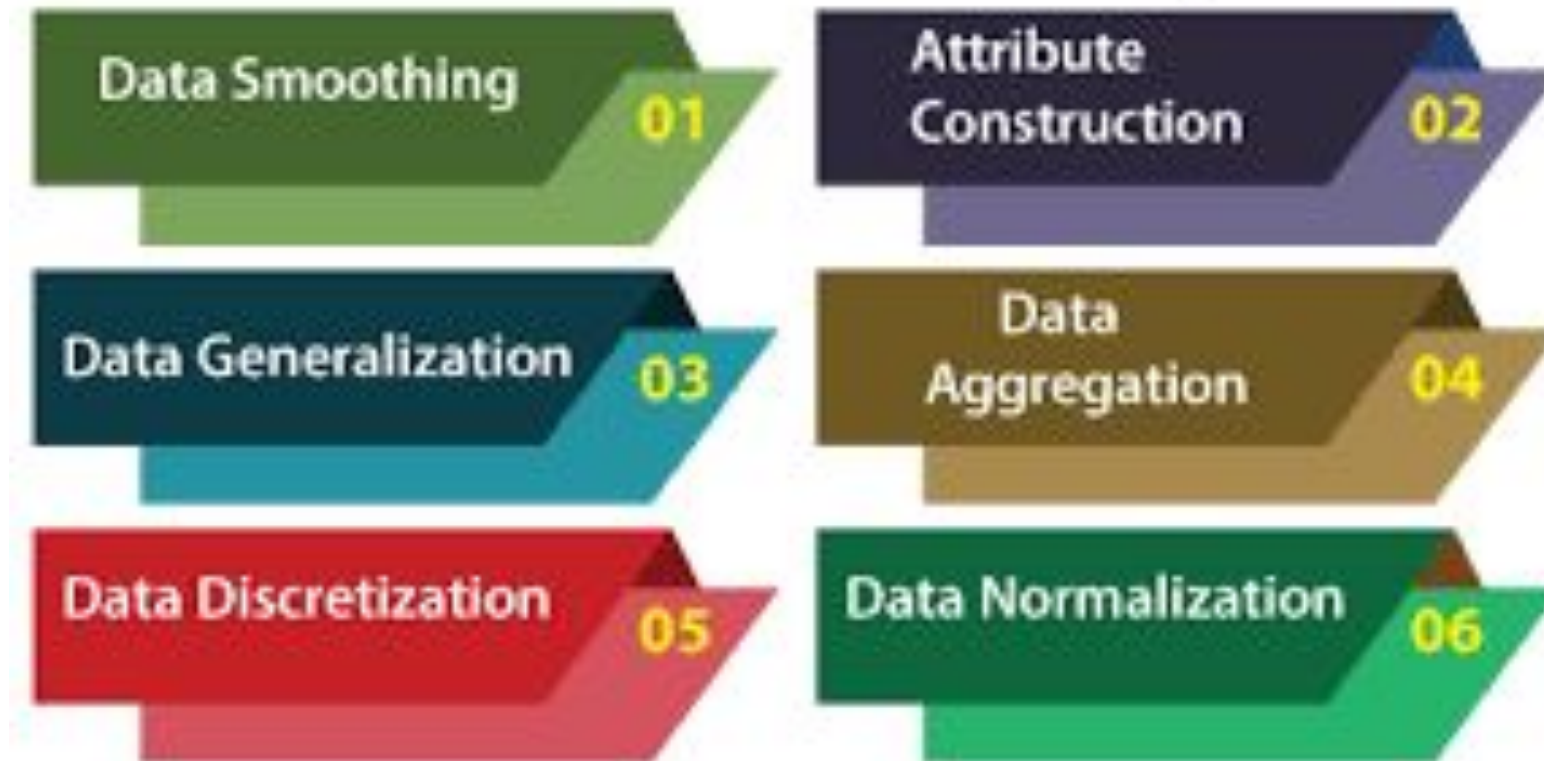
# Transformasi Data

Prosedur untuk:

- mengubah
- memformat
- menskalakan
- membersihkan data mentah dalam format tertentu yang diperlukan, baik itu untuk:
  - aplikasi
  - sistem
  - algoritme
  - atau model pembelajaran mesin.



# Teknik Transformasi Data



# 1. Data Smoothing (Binning)

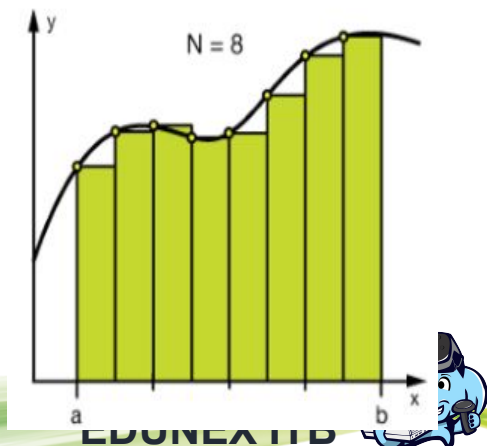
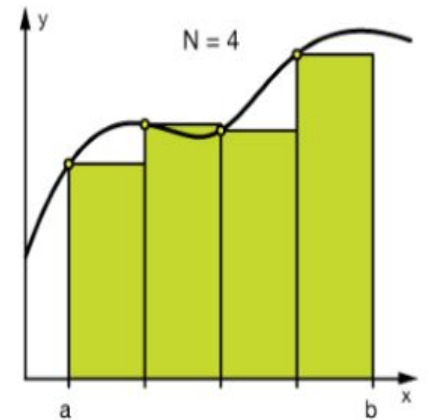
Data binning adalah proses mengelompokkan nilai-nilai data kontinu menjadi interval atau "bin" yang lebih besar, yang mewakili rentang nilai tertentu.

- **Pro:**

- Dapat diterapkan pada data kategorik dan numerik.
- Model lebih robust dan mencegah *overfitting*.

- **Kontra:**

- Meningkatnya biaya kinerja perhitungan.
- Mengorbankan informasi.
- Untuk kolom data numerik, dapat menyebabkan redundansi untuk beberapa algoritma.
- Untuk kolom data kategorik, label dengan frekuensi rendah berdampak negatif pada robustness model statistik.
- Untuk ukuran data dengan 100 ribu baris, disarankan menggabungkan label/kolom dengan record yang  $< 100$  menjadi kategori baru, misal "Lain-lain".



## 2. Attribute Construction

- Atribut baru dibuat untuk membantu proses data mining dari atribut yang sudah ada.
- Contoh: Membuat atribut baru 'area' dari atribut 'tinggi' dan 'lebar'.

### 3. Data Generalization

- Mengubah atribut data tingkat rendah menjadi atribut data tingkat tinggi menggunakan hierarki konsep.
- Kegunaan:
  - Mendapatkan gambaran data yang lebih jelas.
  - Mengurangi kedetailan data
- Contoh:
  - Data umur dapat berupa (10, 20, 30, 40) dalam sebuah dataset dapat ditransformasikan ke tingkat konseptual yang lebih tinggi menjadi nilai kategoris anak-anak, remaja, dewasa, tua.
  - Menggantikan alamat rumah individu dengan nama kota atau negara bagian.



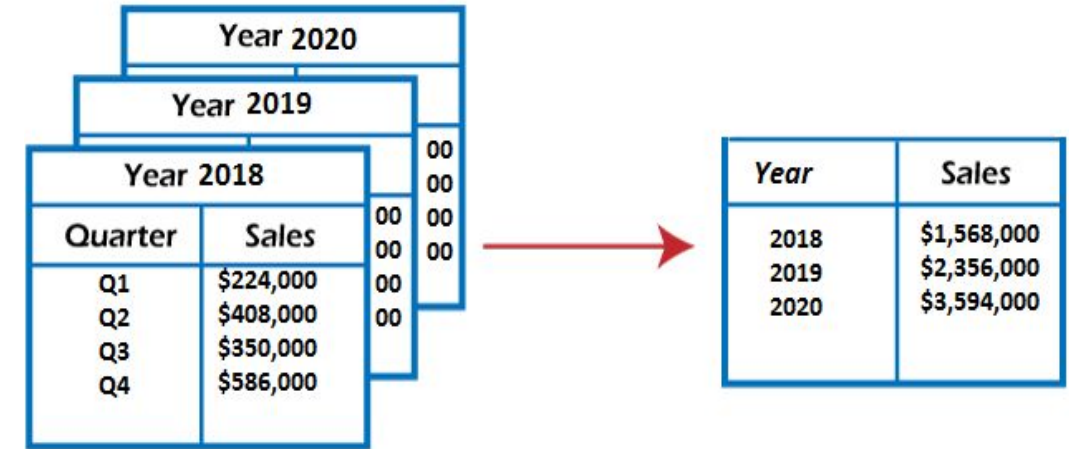


## 4. Data Aggregation

- Metode menyimpan dan menyajikan data dalam format ringkasan.

- Contoh:

Kumpulan data laporan penjualan suatu perusahaan yang memiliki data penjualan triwulanan setiap tahun dapat di-agregasi untuk mendapatkan laporan penjualan tahunan perusahaan.

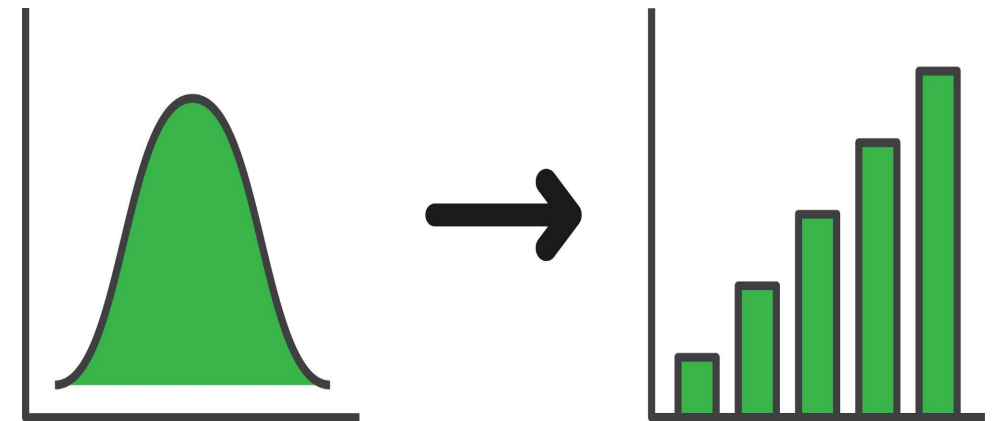


Aggregated Data



## 5. Data Discretization

- Proses mengubah fungsi, model, dan variabel kontinu menjadi diskrit.
- Contoh:
  - misal berat badan  $< 65$  kg (ringan);  $65 - 80$  kg (mid);  $> 80$  kg (berat).
  - pembulatan sebuah nilai riil ke nilai terdekat



Discretization Process

## 6. Data Normalization

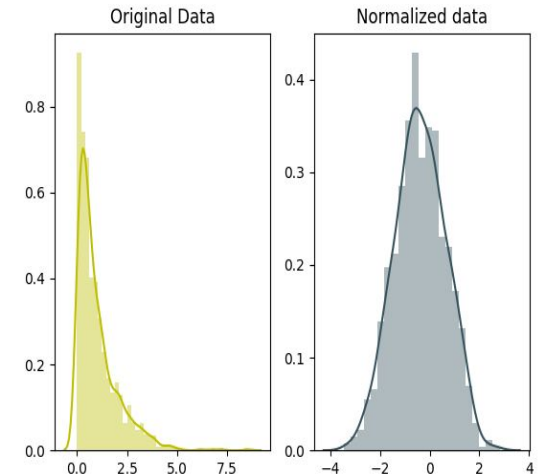
- Definisi: Teknik penskalaan di mana nilai-nilai digeser dan diubah skalanya sehingga nilainya berkisar antara 0 dan 1 ( rentang  $[0,1]$  ).

- Min-Max Normalization:

Misal  $X_{max}$  dan  $X_{min}$  masing-masing adalah nilai maksimum dan minimum dari fitur.

- Ketika nilai  $X$  adalah nilai minimum dalam kolom, pembilangnya adalah 0, dan karenanya  $X'$  adalah 0.
- Sebaliknya, ketika nilai  $X$  adalah nilai maksimum dalam kolom, pembilangnya sama dengan penyebutnya sehingga nilai  $X'$  adalah 1.
- Jika nilai  $X$  berada di antara nilai minimum dan maksimum, maka nilai  $X'$  berada di antara 0 dan 1.

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$



# Encoding

- Ada algoritma Machine Learning yang tidak dapat menerima input berupa data kategorikal, harus berupa data numerik

- Contoh Teknik:

- ordinal encoding

	boro	boro_ordinal	salary
0	Manhattan	2	103
1	Queens	3	89

- one-hot encoding

	boro	boro_Bronx	boro_Brooklyn	boro_Manhattan	boro_Queens	salary
0	Manhattan	0	0	1	0	103
1	Queens	0	0	0	1	89
2	Manhattan	0	0	1	0	142

	boro	salary	vegan
0	Manhattan	103	0
1	Queens	89	0
2	Manhattan	142	0
3	Brooklyn	54	1
4	Brooklyn	63	1
5	Bronx	219	0





# Integrasi Data

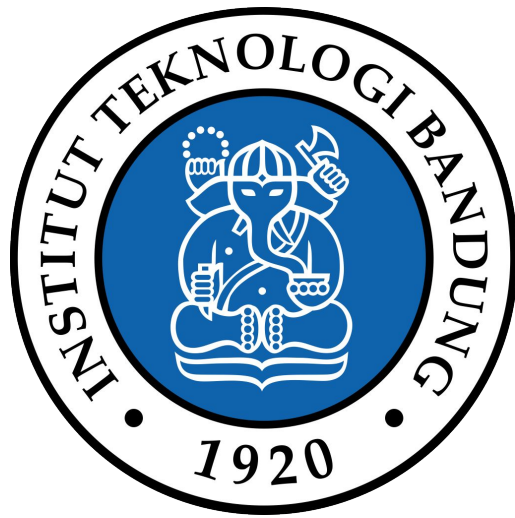
- **Join**

Menggabungkan dua data yang memiliki satu fitur sama.

- **Append**

Menambah instans dari data yang persis sama fitur-fiturnya.





Salam  
Semoga Bermanfaat

