**Modul 3**

# Introduction to Data Science

**Data Science Program**

**Purwadhika**
Startup and Coding School

# Outline

- Data Science Challenges
- Data Science Workflow
- Data Science Roles
- Group Assignment

**Purwadhika**
Startup and Coding School

# WHAT IS DATA SCIENCE

What is it?

- Is it a Role or Position?
- Is it a Process?
- Is it a problem / challenge ?

Correlations to this term :

- Big Data
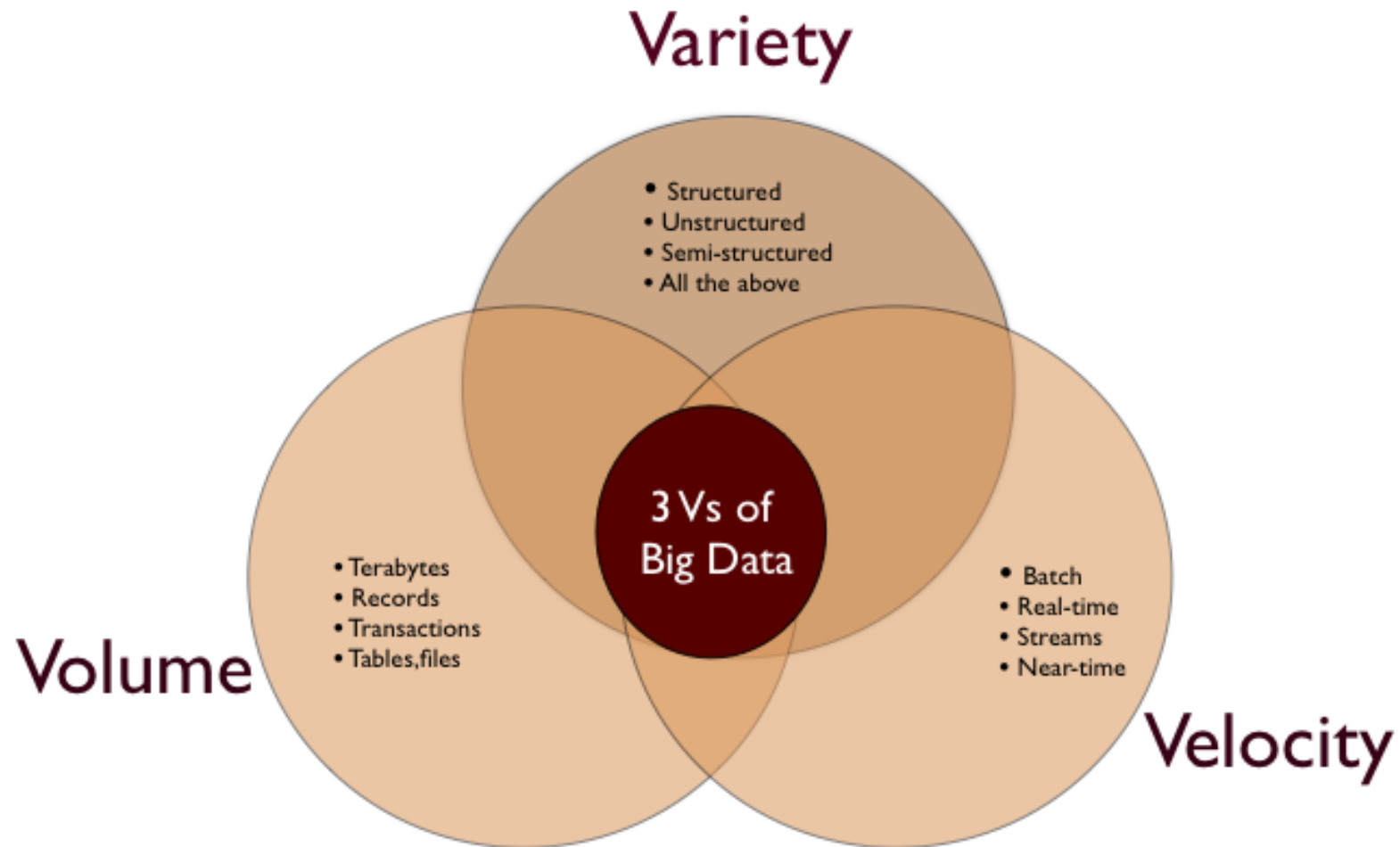- Data Driven
- Machine Learning
- AI
- Distributed computing

**Purwadhika**
Startup and Coding School

# THE RISE OF INTERNET

DIGITIZATION OF EVERYTHING

THEN

NOW

Purwadhika
Startup and Coding School

# THE RISE OF INTERNET

# BIG DATA : 3V

# MULTI DISCIPLINARY



**BANKING AND FINANCIAL SERVICES**

**PUBLIC SECTOR SERVICES**

**BANKING AND FINANCIAL SERVICES**

**SPORTS AND EDUCATION**

**OIL & GAS AND MANUFACTURING**

**HEALTHCARE AND MEDICAL SERVICES**

**TRANSPORTATION AND LOGISTICS**

**TELECOM AND ICT SERVICES**

**MEDIA AND ENTERTAINMENT**

**TRAVEL AND HOSPITALITY**

**CONSTRUCTION AND REAL ESTATE**

We need to understand the PROBLEM

1. How the management think
2. How the customer think
3. How the market shifts

**Purwadhika**
Startup and Coding School

# THE QUESTIONS

*"Kami mau pasang iklan, tapi tidak tahu channel mana yang paling efektif"*

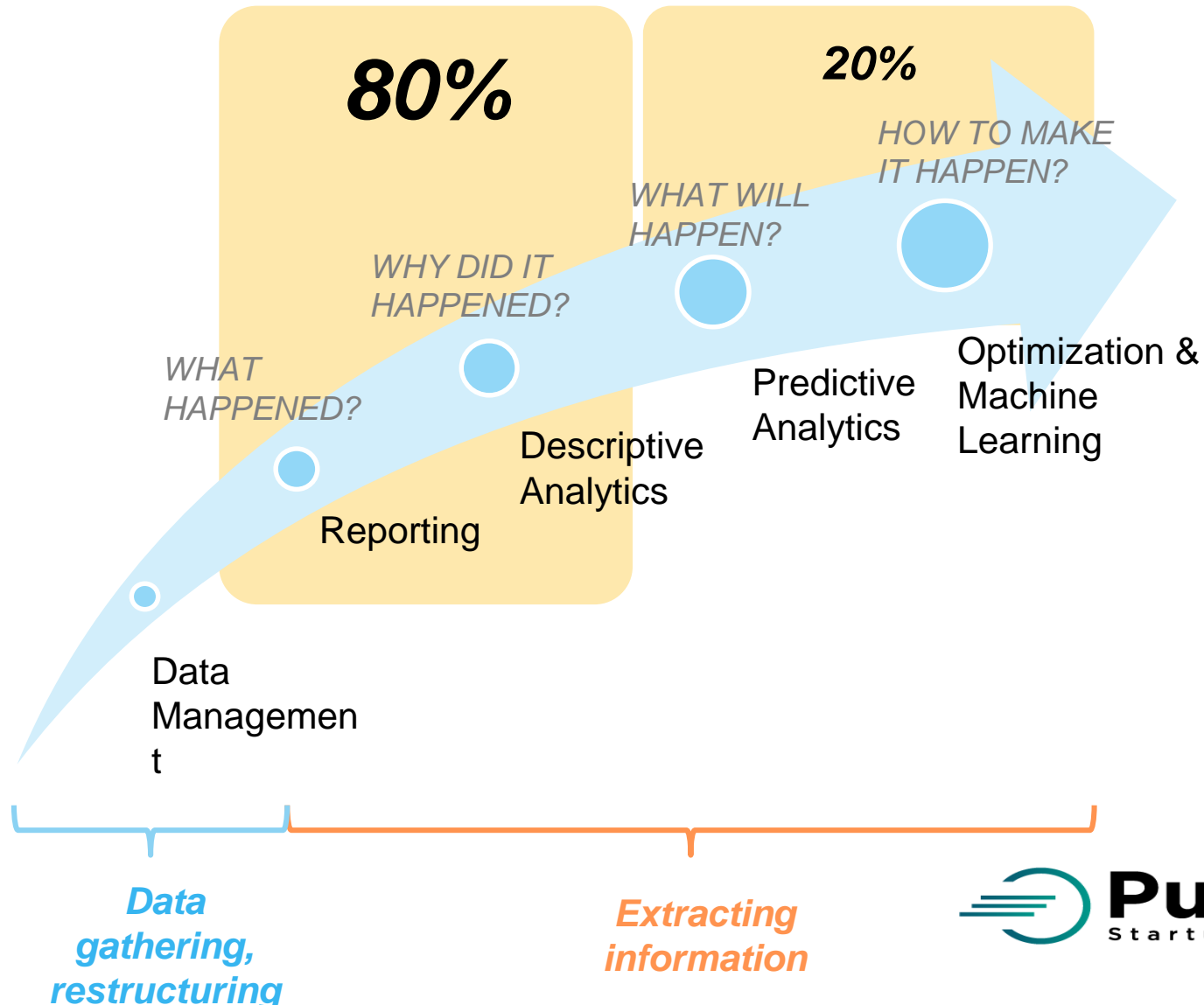*"Ada beberapa produk kami yang tidak laku, walau review sangat bagus"*

*"Kredit nasabah kami banyak yang macet"*

*"Stock barang selalu habis/terlalu banyak"*

*"Kami tidak tahu seberapa efisien sales person kami"*
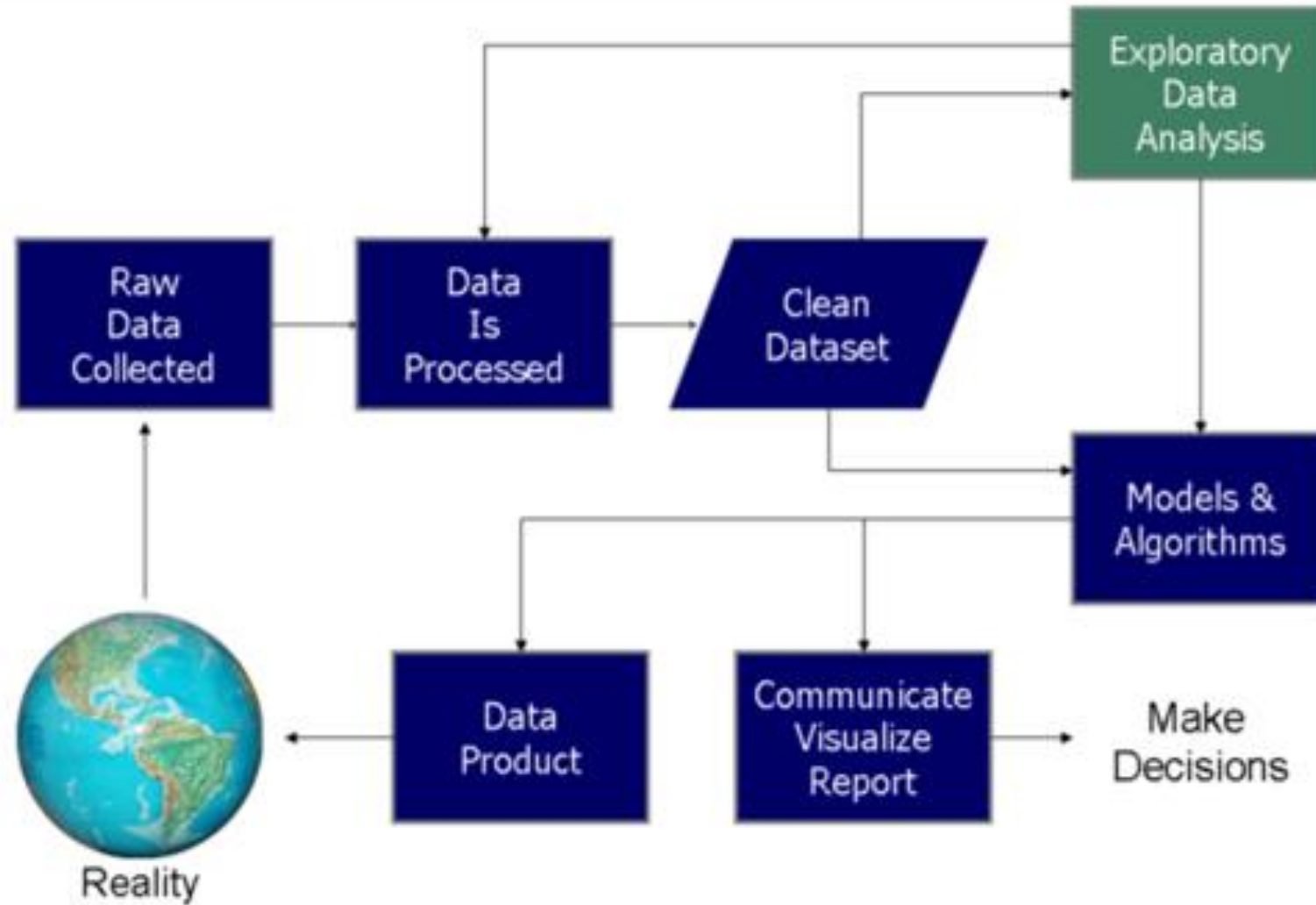
# DATA SCIENCE CHALLENGES

**80%**          **20%**

*HOW TO MAKE IT HAPPEN?*

*WHAT WILL HAPPEN?*

*WHY DID IT HAPPENED?*

Optimization & Machine Learning

*WHAT HAPPENED?*

Predictive Analytics

Descriptive Analytics

Reporting

Data Management

**Data gathering, restructuring**

**Extracting information**

**Purwadhika**
Startup and Coding School

# DATA SCIENCE WORKFLOW

# DATA SCIENCE WORKFLOW

# DATA SCIENCE WORKFLOW

Ask Questions

- Who are the customers?
- Why are they buying our product?
- How do we predict if a customer is going to buy our product?
- What is different from segments who are performing well and those that are performing below expectations?
- How much money will we lose if we don't actively sell the product to these groups?

**Purwadhika**
Startup and Coding School

# DATA SCIENCE WORKFLOW

What needs to be considered :
- Data Sources
- Data Location
- Data Format
- Data Types
- Acquisition Methods
- Data Privacy

**Purwadhika**
Startup and Coding School

# DATA SCIENCE WORKFLOW

Data Sources :
- Users Profile
- Users Activity/Transaction
- Enterprise ressources
- World trends/activity

**Purwadhika**
Startup and Coding School

# DATA SCIENCE WORKFLOW

Data Location
- Inter Department
- Across Department
- External Data
- Public Data

# DATA SCIENCE WORKFLOW

Data Format
- Hard copy
- Digital Documents
- Database
- Streams

**Purwadhika**
Startup and Coding School

# DATA SCIENCE WORKFLOW

## Data Types

- Numerical
- Text
- Image
- Audio
- Video

**Purwadhika**
Startup and Coding School

# DATA SCIENCE WORKFLOW

Data Access
- Data warehousing
- REST API
- Web scraping

# DATA SCIENCE WORKFLOW

Data Privacy
- User Consent : User needs to give consent for any usage purposes
- Data Privacy Law :
  - EU General Data Protection Regulation
  - RUU Perlindungan Data Pribadi

# DATA SCIENCE WORKFLOW

Data preparation

- Data cleansing
  - Format normalization
  - Typing inconsistency
- Handling NULL values
- Handling outliers
- Feature Selection/Engineering

# DATA SCIENCE WORKFLOW : DATA ANALYSIS AND VISUALIZATION

DATA ANALYSIS
- Always aim to answer the problem definition
- Identify
  - Variations
  - Correlations
  - Trends
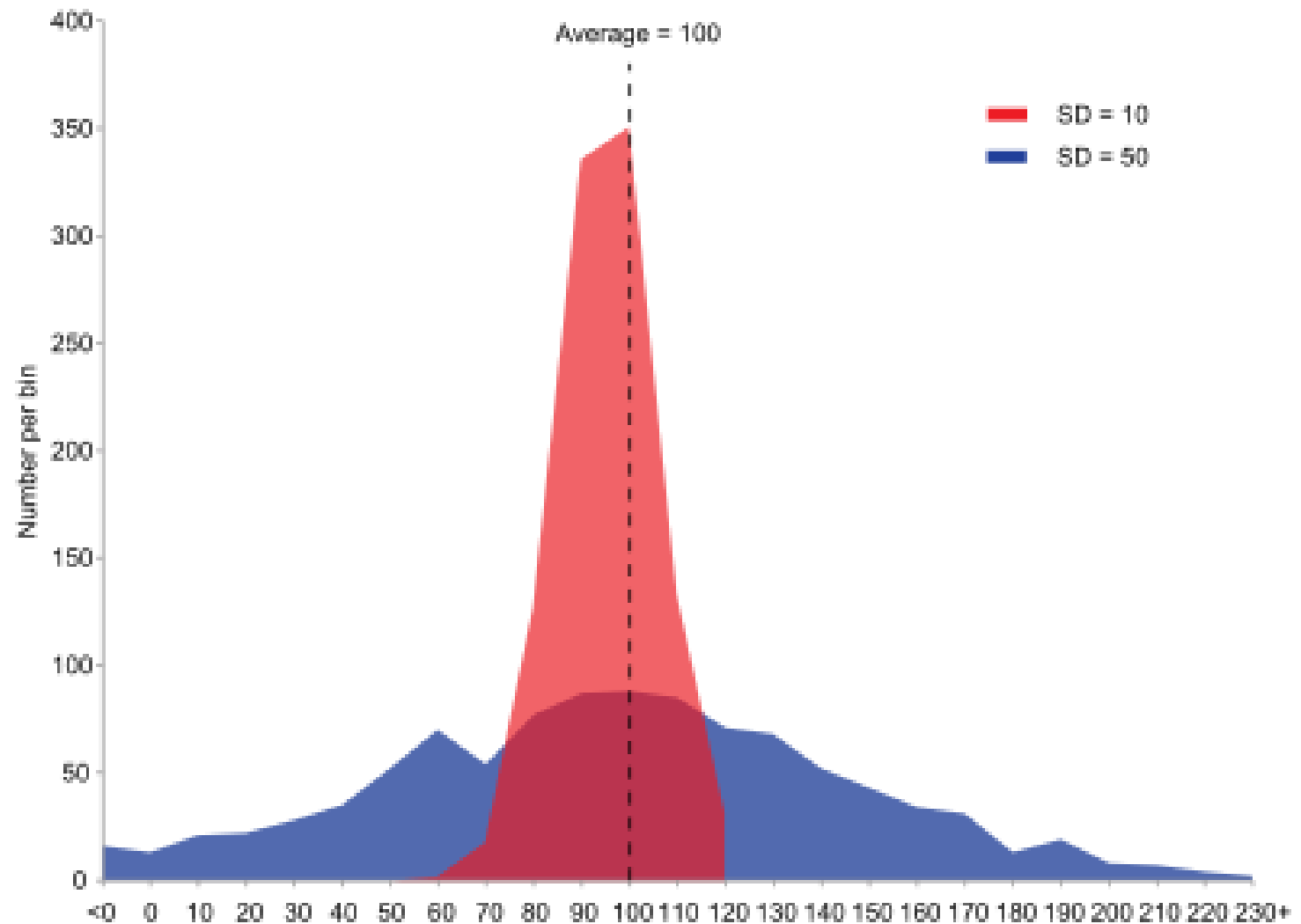  - Outliers

**Purwadhika**
Startup and Coding School

# Data Analysis Term

- **Variation** is a measure of how far a set of random numbers are spread out from their mean value. If one is familiar with Standard Deviation, variance could be calculated by squaring the Standard Deviation.

- **Correlation** is a degree of relationship between two or more numerical variable. In statistical term, it could be define as how pair of variable are linearly related.

- **Trend** is a general pattern of how condition, output, process, or general tendency of series data point move in a certain direction over course of time.

- **Outlier** is a data point that differ significantly from the other data observation
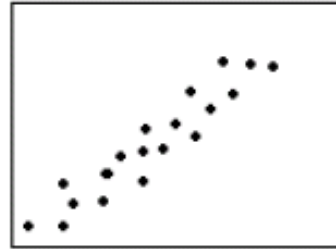
**Purwadhika**
Startup and Coding School

# Variation

Example of samples from two populations with the same mean but different variances. The red population has mean 100 and variance 100 (SD=10) while the blue population has mean 100 and variance 2500 (SD=50).
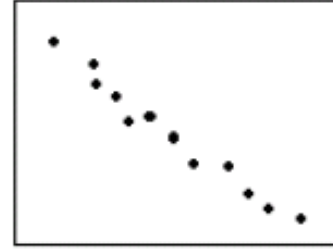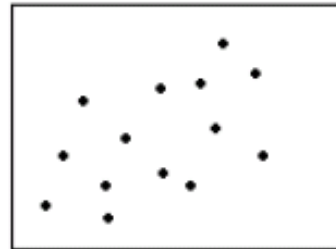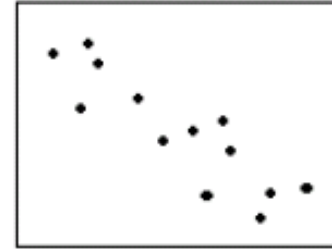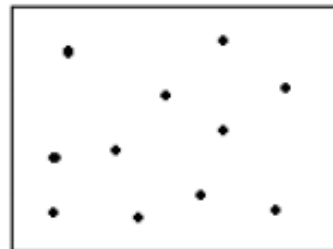
# Correlation
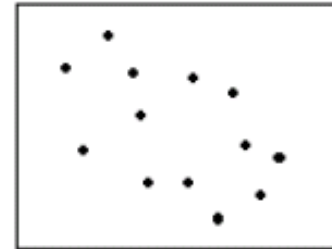


Degree of Correlation

Strong Positive

Strong Negative
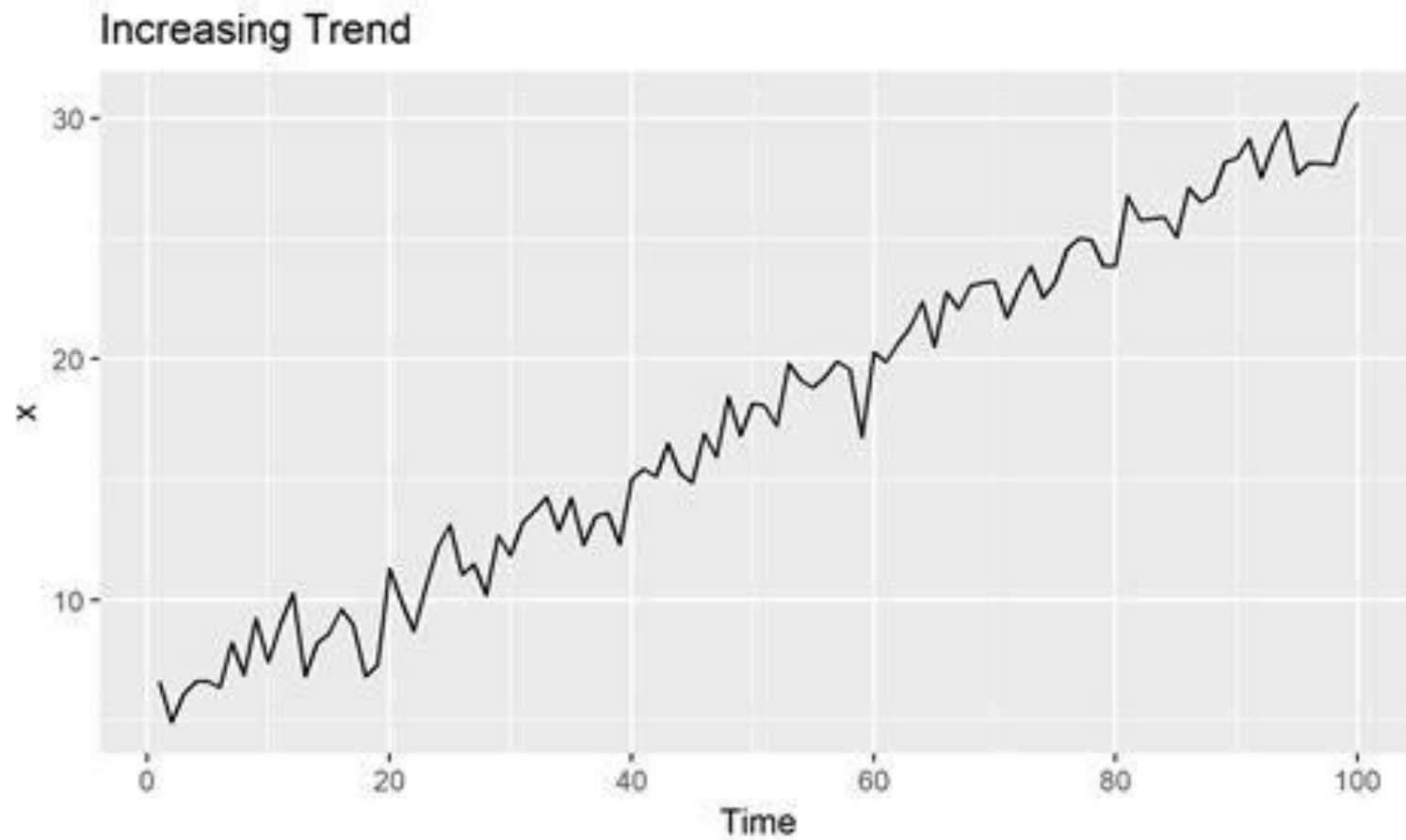
Weak Positive

Moderate Negative

None

Weak Negative

Purwadhika
Startup and Coding School

# Trend



Increasing Trend

# Outlier



Outlier

Outlier

Copyright 2014. Laerd Statistics.

**Purwadhika**
Startup and Coding School

# DATA SCIENCE WORKFLOW : DATA ANALYSIS AND VISUALIZATION

DATA Visualization

- Know the audience
- Visualization is all about perception

**Purwadhika**
Startup and Coding School

# DATA SCIENCE WORKFLOW : DATA ANALYSIS AND VISUALIZATION

# DATA SCIENCE WORKFLOW : DATA ANALYSIS AND VISUALIZATION



Average Annual Global Temperature in Fahrenheit 1880-2015



Average global temperature, 1880 to 2014

Purwadhika
Startup and Coding School

# DATA SCIENCE ROLES

## Data Scientist

**Activities**
- Data cleansing and Preparation
- Evaluating statistical models
- Build ML Model

**Tools**
- R
- Python
- Matlab
- Stata
- SQL
- Spark

**Skills and Talents**
- Statistical theories and methodologies
- Database systems
- Programming skills

**Purwadhika**
Startup and Coding School

# DATA SCIENCE ROLES

## Data Engineer

**Activities**
- Data Integration
- Product Development (Dashboard, API)
- Scalability and Automation

**Skills and Talents**
- Programming skills
- Database system and modelling
- IT Infrastructure and Cloud environment

**Tools**
- Database systems: SQL, NoSQL
- Python, Node
- Google Cloud Platform, Amazon AWS
- Distributed System

**Purwadhika**
Startup and Coding School

# DATA SCIENCE ROLES

## Business Analyst

### Activities
- Framing the problem
- Data Exploration
- Presenting Analysis insights

### Skills and Talents
- Business and Domain knowledge
- Communication
- Database query language

### Tools
- Dashboard
- Visualization tools :Tableau, QlikView
- Open Refine
- Powerpoint and Excel

**Purwadhika**
Startup and Coding School

# DATA SCIENCE ROLES

## Domain Expert

**Activities**
- Framing the problem
- Provides Consultation to the real world problems

**Skills and Talents**
- Business and Domain knowledge
- Communication

**Tools**
- (depends on the field)

**Purwadhika**
Startup and Coding School

# DATA SCIENCE ROLES

## Other roles

- Database Admin : Query/Prepare data to be processed/analyse
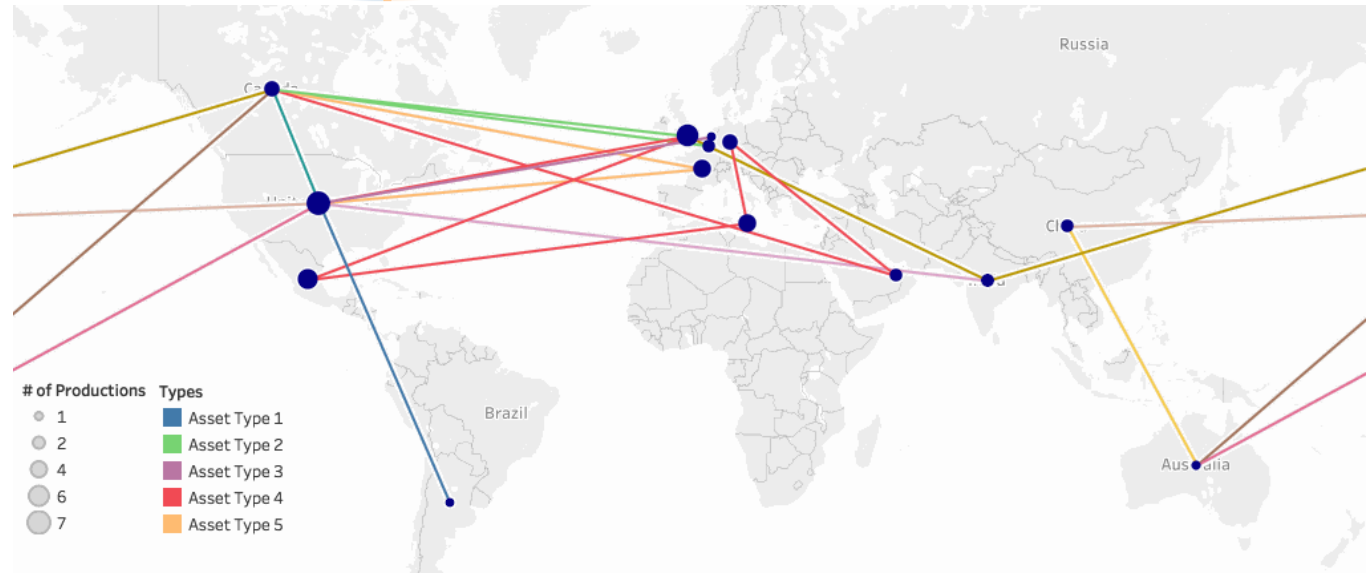- Data Architect : Design information archtect
- Statistician :
- Developer

**Purwadhika**
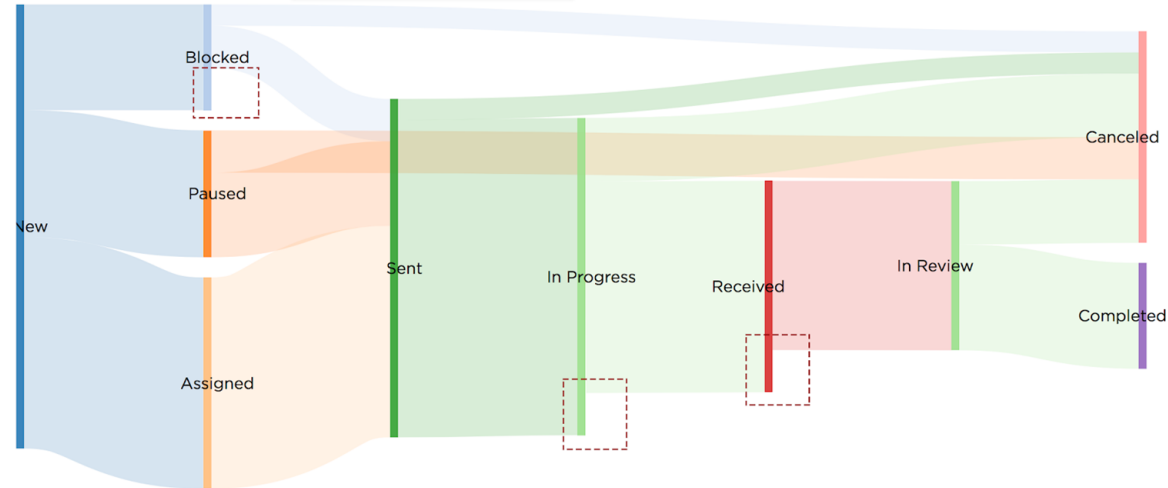Startup and Coding School

# DATA SCIENCE ROLES

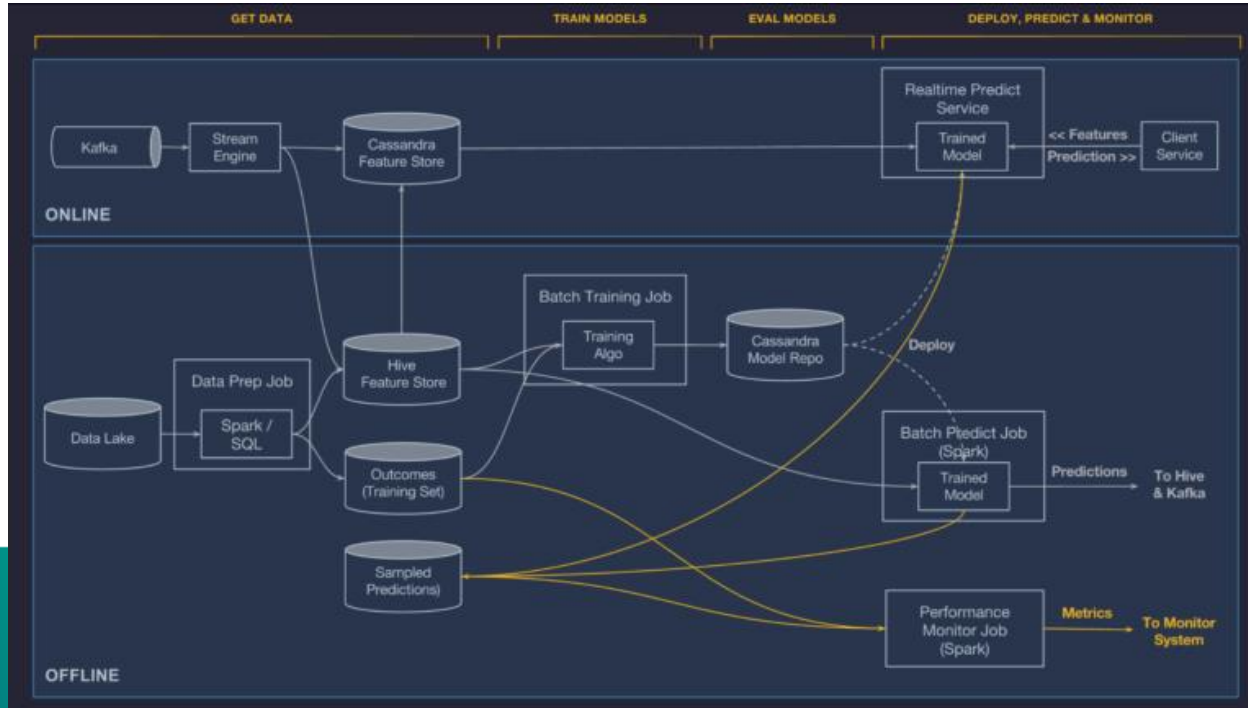# Keep it lean, grow as you go

# USECASE : NETFLIX

- **Pre-production cost estimation**
  - Location
  - Crews
  - Schedule
- **Shooting schedules**
- **Post-Production assets progression**
- **Prioritazion of location**

# Study Case : **UberEATS** estimated time of delivery model



**Michelangelo:** Uber's Machine Learning Platform

**Machine Learning Workflow:**
- Manage data
- Train models
- Evaluate models
- Deploy models
- Make predictions
- Monitor predictions

# ASSIGNMENT

Challenge :

Unilever akan mengeluarkan varian shampo baru. Direksi meminta bantuan kepada tim Data Science untuk memberikan rekomendasi spesifikasi varian tersebut.

**Purwadhika**
Startup and Coding School

# ASSIGNMENT

PROBLEM IDENTIFICATION :

Define the problem, identify the questions

- What is the problem ?
- Who is having the problem ?
- When is it happening ?
- Where is it happening ?
- What are the expected output?
- What have happened in the past?

**Purwadhika**
Startup and Coding School

# ASSIGNMENT

Plan the data driven Process!

- **Data Acquisition :**
  What data do I need, and how to access them?

- **Data Preparation :**
  Define the ideal data format, and ways to prepare them

Purwadhika
Startup and Coding School

# ASSIGNMENT

Plan the data driven Process!

- **Data Acquisition :**
What data do I need, and how to access them?

- **Data Preparation :**
Define the ideal data format, and ways to prepare them

- **Data Analysis :**
What insigths do you need, and how to analyse them?

- **Data Visualization:**
How and to whom do you share your insights

**Purwadhika**
Startup and Coding School