# HANDS-ON WITH DATA VISUALIZATION

## Darren Chiu

# DARREN CHIU



**Technologist**

**+**

**Educator**

Academic Background: #Computer Engineering #Business #Psychology
Work: #Software Engineer #Entrepreneur #Tech Consultant
Industry: #Investment Banking #Education #Media #Digital Marketing

# Data Science

- Data Science = Analysis + Prediction
- EDA (Exploratory Data Analysis): Find useful insights thru analysing past data
- Prediction: Usually done by Machine Learning technology
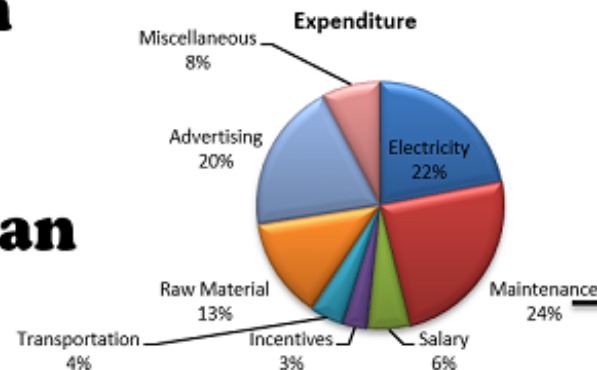
# Exploratory Data Analysis... How?

- "Dimension" as a framework to guide your EDA journey

- 1D : descriptive statistics

- 2D : correlational study

- 3D+ : multi dimensional correlational study
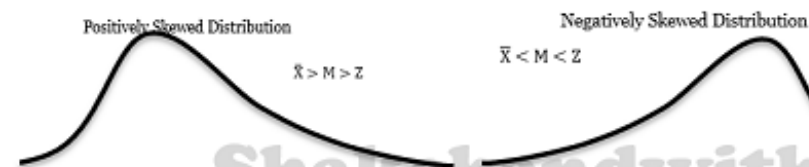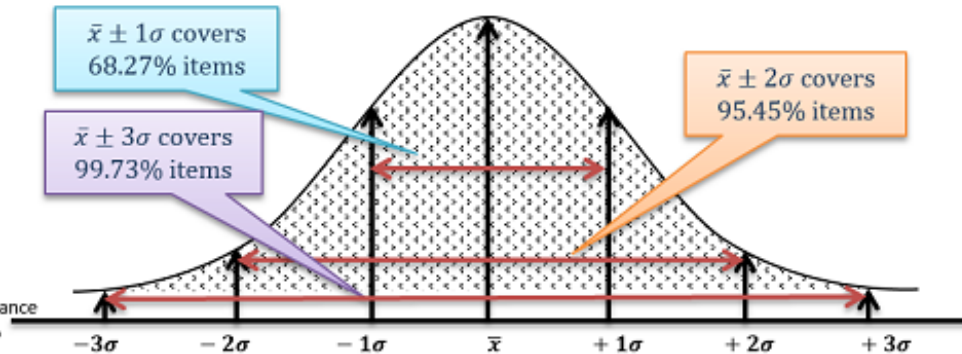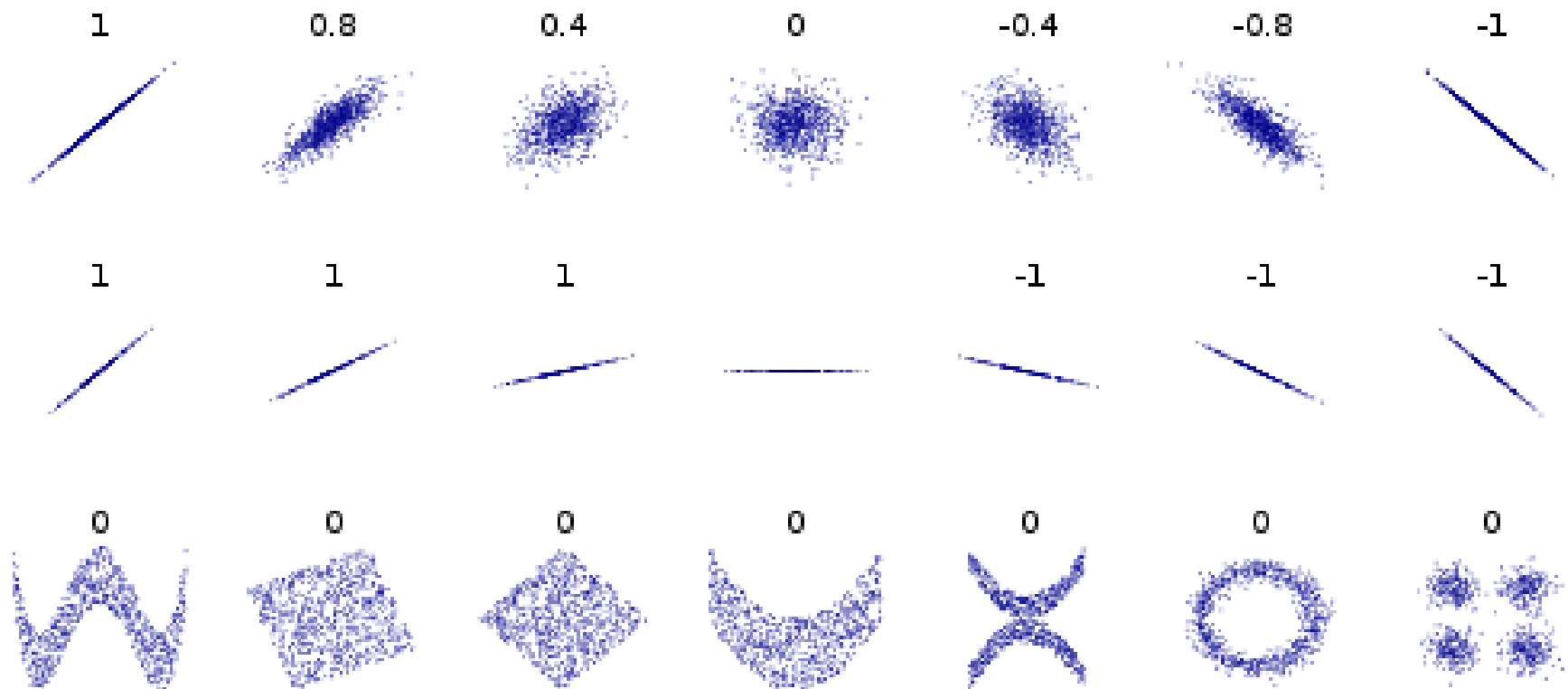
# 1D: Descriptive Statistics

**Mean**

**Median**

**Mode**

$$Std.\,Dev.\,\sigma = \sqrt{\dfrac{\Sigma(x - \overline{x})^2}{n}}$$

**Expenditure**

- Miscellaneous 8%
- Advertising 20%
- Electricity 22%
- Raw Material 13%
- Transportation 4%
- Incentives 3%
- Salary 6%
- Maintenance 24%

$\overline{x} \pm 1\sigma$ covers 68.27% items

$\overline{x} \pm 3\sigma$ covers 99.73% items

$\overline{x} \pm 2\sigma$ covers 95.45% items

$-3\sigma$   $-2\sigma$   $-1\sigma$   $\overline{x}$   $+1\sigma$   $+2\sigma$   $+3\sigma$

Positively Skewed Distribution

$\overline{X} > M > Z$

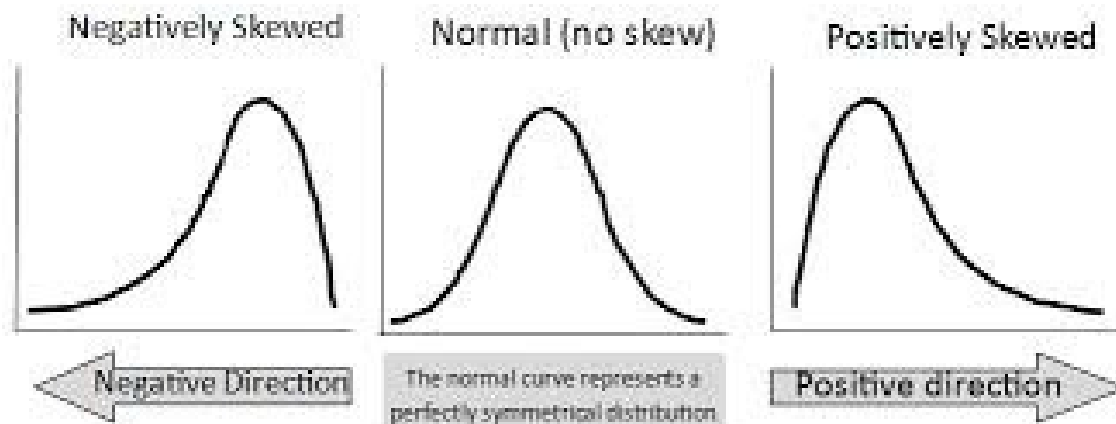Negatively Skewed Distribution

$\overline{X} < M < Z$

Shakehandwithlife.in

# 2D: Correlational Analysis

# Why we need Visualization

- Visualization is the most powerful tool for EDA
- human brain process image >>> numbers/text



$$\gamma_1 = \mathrm{E}\left[\left(\frac{X-\mu}{\sigma}\right)^3\right] = \frac{\mu_3}{\sigma^3} = \frac{\mathrm{E}\left[(X-\mu)^3\right]}{(\mathrm{E}[(X-\mu)^2])^{3/2}} = \frac{\kappa_3}{\kappa_2^{3/2}}$$

# I still don't believe!

7563950684 73
6586630375 76
8603726586 02
8465891078 30

# I still don't believe!

7563**3**9505068847**3**

658663**3**0**3**7576

860**3**72658602

84658910783**3**0

# So what I can leverage?

# Chart Suggestions—A Thought-Starter

**Variable Width Column Chart**
Two Variables per Item

**Table or Table with Embedded Charts**
Many Categories

**Bar Chart**
Many Items

**Column Chart**
Few Items

Few Categories

One Variable per Item

Among Items

**Circular Area Chart**
Cyclical Data

**Line Chart**
Non-Cyclical Data

Many Periods

**Column Chart**
Single or Few Categories

**Line Chart**
Many Categories

Few Periods

Over Time

**Comparison**

**Column Histogram**
Few Data Points

**Line Histogram**
Many Data Points

Single Variable

**Scatter Chart**
Two Variables

**Bubble Chart**
Three Variables

**Relationship**

## What would you like to show?

**Distribution**

**Scatter Chart**
Two Variables

**3D Area Chart**
Three Variables

**Composition**

Changing Over Time

Static

Few Periods

Many Periods

Only Relative Differences Matter
**Stacked 100% Column Chart**

Relative and Absolute Differences Matter
**Stacked Column Chart**

Only Relative Differences Matter
**Stacked 100% Area Chart**

Relative and Absolute Differences Matter
**Stacked Area Chart**

Simple Share of Total
**Pie Chart**

Accumulation or Subtraction to Total
**Waterfall Chart**

Components of Components
**Stacked 100% Column Chart with Subcomponents**

# Let's look at
# some misleading visualizations...

# Amount of Environmental Waste Created and Recycled



3,500,000
3,000,000
2,500,000
2,000,000
1,500,000

2000  2005  2007  2009  2010  2011  2012  2013

■ Total waste generated   ■ Trashed   ■ Recycled

3
2
1

A  B

Item A
Item B
Item C
Item D

Internet Explorer vs Murder Rate

— Murders in US    Internet Explorer Market Share



Reconstructed Temperature

Medieval Warm Period

Little Ice Age

2004 ✳



"MASSIVE INCREASE IN HOUSE PRICES THIS YEAR!"

Average house price (£)

# Visualization Tools

# Google Data Studio

- Beta since 2016
- went out of beta on 25th Sept
- Hundreds of data connectors
- Run on cloud
- Even high school students could create dynamic interactive data visualization



Google Data Studio

# Let's do some exercise...

Google Data Studio:
https://datastudio.google.com

Data Source:
https://www.kaggle.com/vikalpdongre/us-flights-data-2008

Importing Data to Google Cloud BigQuery:
https://cloud.google.com/bigquery/docs/loading-data-cloud-storage-csv

# Exercise 1:
# Find out the relationship of month and flight delay

# Exercise 2:

## So is it because there are more flights in Dec?

# Exercise 3:
# Then what is the source of flight delay?

# Exercise 4:

**We now know that weather is the reason of delay…Then is it related to where you are flying?**

# Exercise 5:

# But what if I really want a Christmas holiday... When should I fly?

# Exercise 6:

## If you still interested - which airline we should go for ...

# Conclusions – some tips...

- Always start with message you want to deliver
- use only one color tone in one graph
- highlight the thing you want to highlight
- only show data that are relevant
- only show data that can be process by human brain
- Tell the story you want vs tell the right story?

## -Don't risk deceiving yourself-