

Visualisasi Data dan Informasi

Pertemuan #2 - Mengenal Tools Visualisasi

Tujuan Pembelajaran

- Mengela library yang digunakan pada Python
- Mengenal visualisasi variabel dan statistik

Outline

- Visualisasi Variabel

1. Pie Chart
2. Bar Chart
3. Line Graphs
4. Scatter Plot
5. Heatmap

Visualisasi Statistik

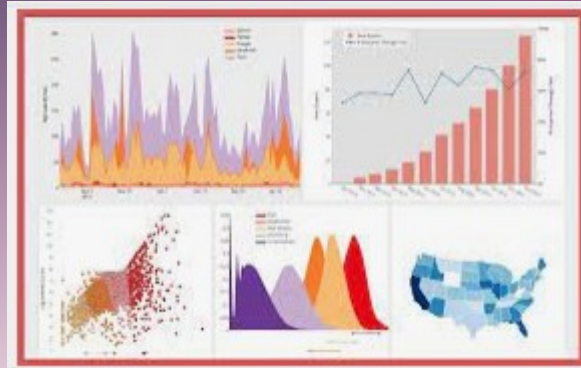
1. Histogram
2. Correlation
3. Descriptive Statistik
4. Grouping (Pivot)
5. Anova

Library

- Pustaka Python adalah potongan kode yang dapat digunakan kembali yang mungkin ingin Anda sertakan dalam program.
- Library yang digunakan
 1. Matplotlib
 2. Pandas
 3. Scatter
 4. Numpy

Matplotlib

- Matplotlib adalah lintas platform, visualisasi data, dan pustaka plot grafis untuk Python dan ekstensi numeriknya NumPy.



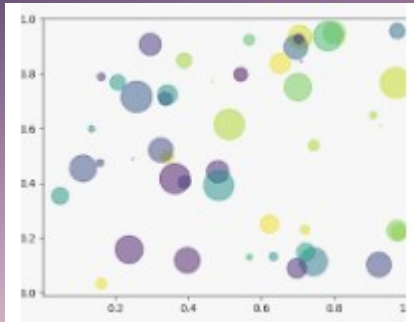
Pandas

- Pandas adalah paket Python open source yang paling banyak digunakan untuk ilmu data/analisis data dan tugas pembelajaran mesin. Itu dibangun di atas paket lain bernama Numpy, yang menyediakan dukungan untuk array multi-dimensi.



Scatter

- Plot sebar adalah jenis bagan yang biasanya digunakan untuk mengamati dan menampilkan hubungan antar variabel secara visual. Nilai-nilai variabel diwakili oleh titik-titik.



Numpy

- NumPy adalah pustaka Python yang digunakan untuk bekerja dengan array. Ini juga memiliki fungsi untuk bekerja dalam domain aljabar linier, transformasi fourier, dan matriks.
- NumPy adalah singkatan dari Numerik Python.



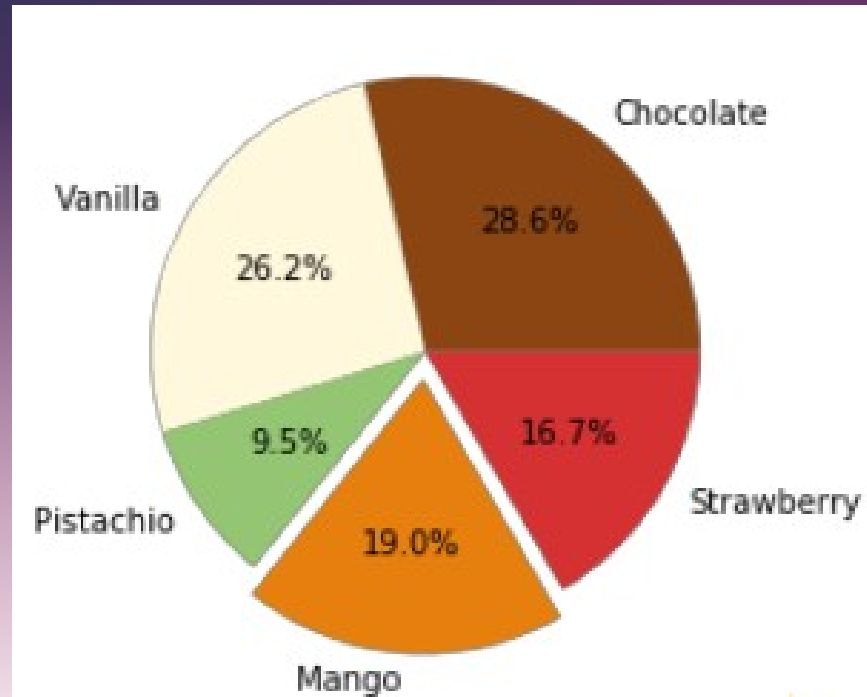
Visualisasi

- Visualisasi berperan peran penting dalam bidang machine learning dan data science. Seringkali kita perlu menyaring informasi kunci yang ditemukan dalam sejumlah data data menjadi bentuk yang bermakna dan mudah dicerna.
- Visualisasi yang baik dapat menceritakan sebuah cerita tentang data Anda dengan cara yang tidak dapat dilakukan oleh sebuah kalimat.
- Di lab ini kita akan mengeksplorasi beberapa teknik visualisasi yang umum. Lab ini akan menggunakan toolkit seperti Matplotlib's Pyplot dan Seaborn untuk membuat gambar informatif yang memberikan informasi dan pengetahuan mengenai dataset.

Visualisasi Variabel

Pie Chart

- Pie chart digunakan untuk menunjukkan seberapa banyak dari setiap jenis kategori dalam dataset berbanding dengan keseluruhan.
 - Variabel label berisi tupel rasa es krim
 - Variabel voting berisi tupel voting.
 - Data tersebut mewakili jumlah voting rasa es krim favorit.

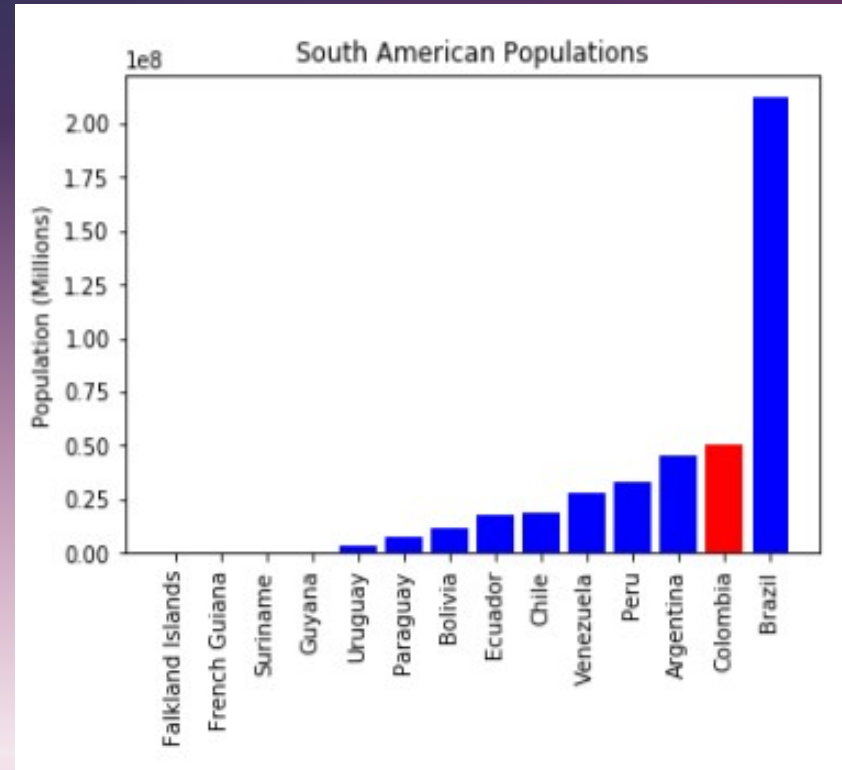


Visualisasi Variabel

Bar Chart

- Bar Chart adalah merupakan tools visualisassi yang dapat digunakan untuk membandingkan data kategorikal.
- Mirip dengan diagram lingkaran, diagram ini dapat digunakan untuk membandingkan kategori data satu sama lain.
- Diagram batang dapat menampilkan lebih banyak kategori data daripada diagram lingkaran.

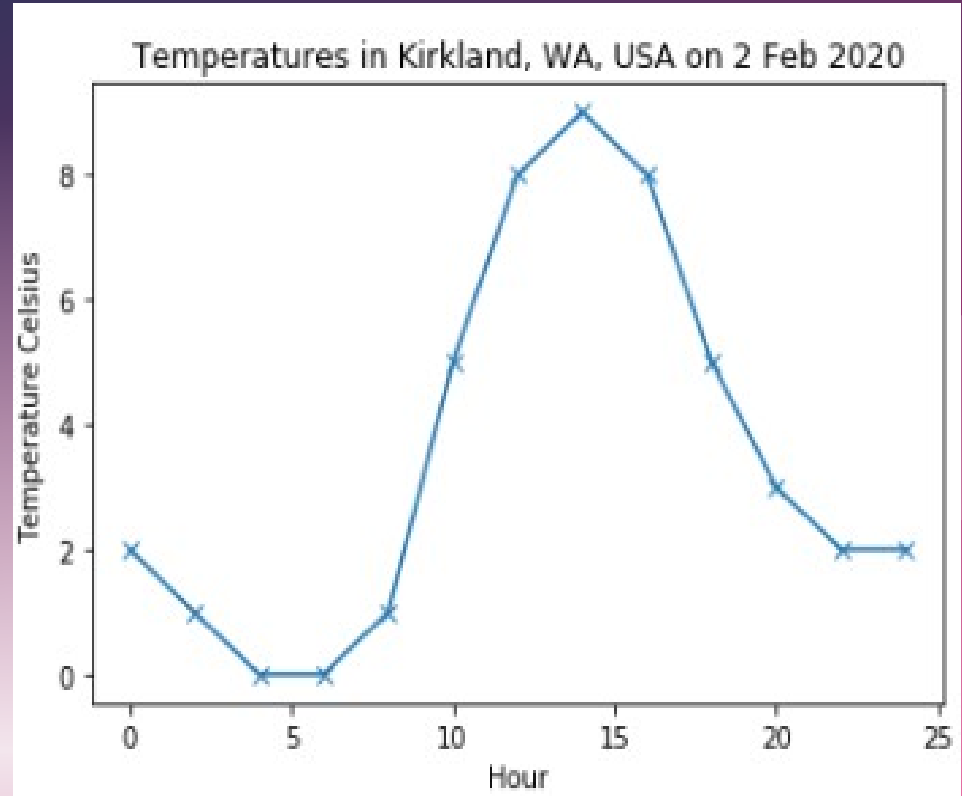
- Mari kita mulai dengan melihat diagram batang yang menunjukkan populasi setiap negara di Amerika Selatan.
- Visualisasi ditunjukkan dengan cara mengurutkan dari negara yang memiliki populasi terbesar ke populasi terendah.
- Highilght ditunjukkan untuk negara Colombia



Visualisasi Variabel

Line Graph

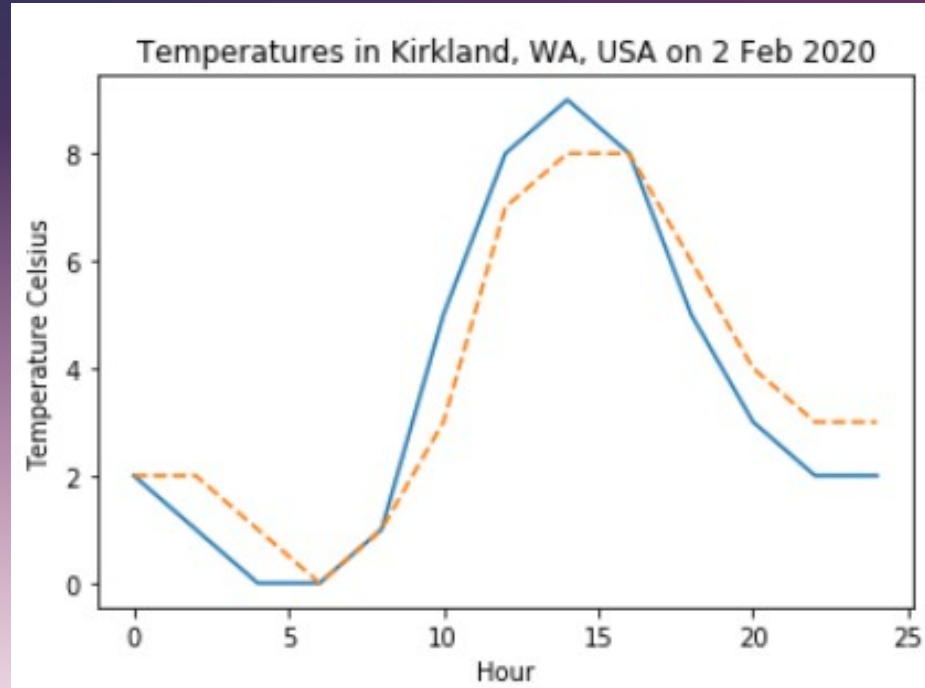
- Line Graph adalah bentuk visualisasi lainnya selain diagram lingkaran dan diagram batang.
- Diagram garis lebih berguna untuk menunjukkan bagaimana kemajuan data selama beberapa periode.
- Misalnya, grafik garis dapat berguna dalam membuat grafik temperatur dari waktu ke waktu, harga saham dari waktu ke waktu, berat menurut hari, atau metrik berkelanjutan lainnya.



Visualisasi Variabel

Line Graph

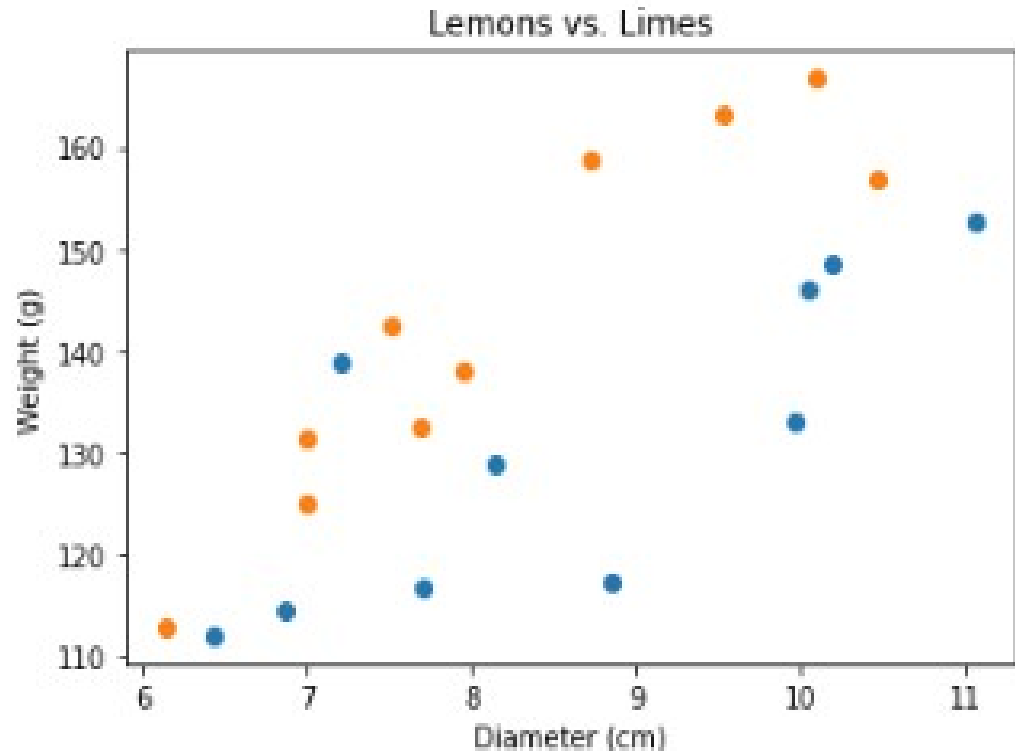
- Kita bahkan dapat memiliki beberapa garis pada grafik yang sama didalam satu gambar
- Biasanya kita mengilustrasikan dua line graph untuk menggambarkan dua data yaitu data aktual dan data prediksi.



Visualisasi Variabel

Scatter Plot

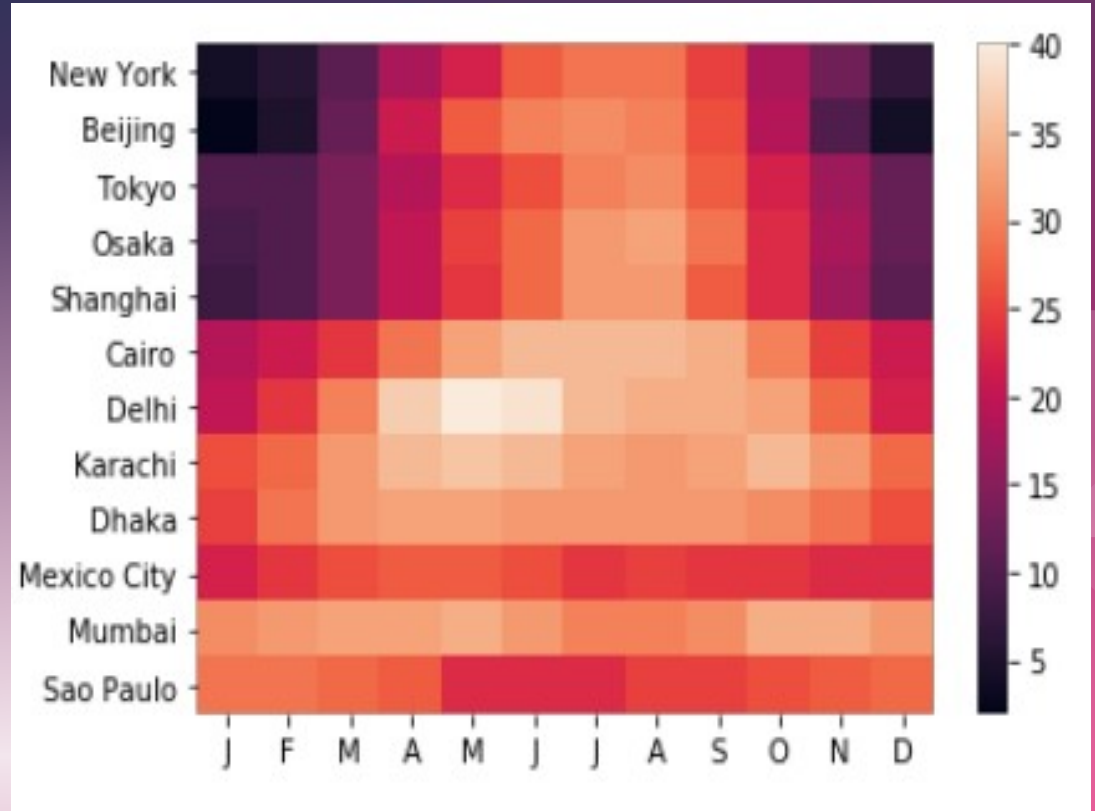
- Scatter plot berfungsi baik untuk data dengan dua komponen numerik.
- Scatter plot dapat memberikan informasi yang berguna terutama mengenai pola atau pencilan.
- Pada contoh di bawah ini, kita memiliki data yang terkait dengan perbedaan lemon dan lime berdasarkan karakteristik fisiologis.
 - Berat (g)
 - Diameter (cm)



Visualisasi Variabel

Heatmap

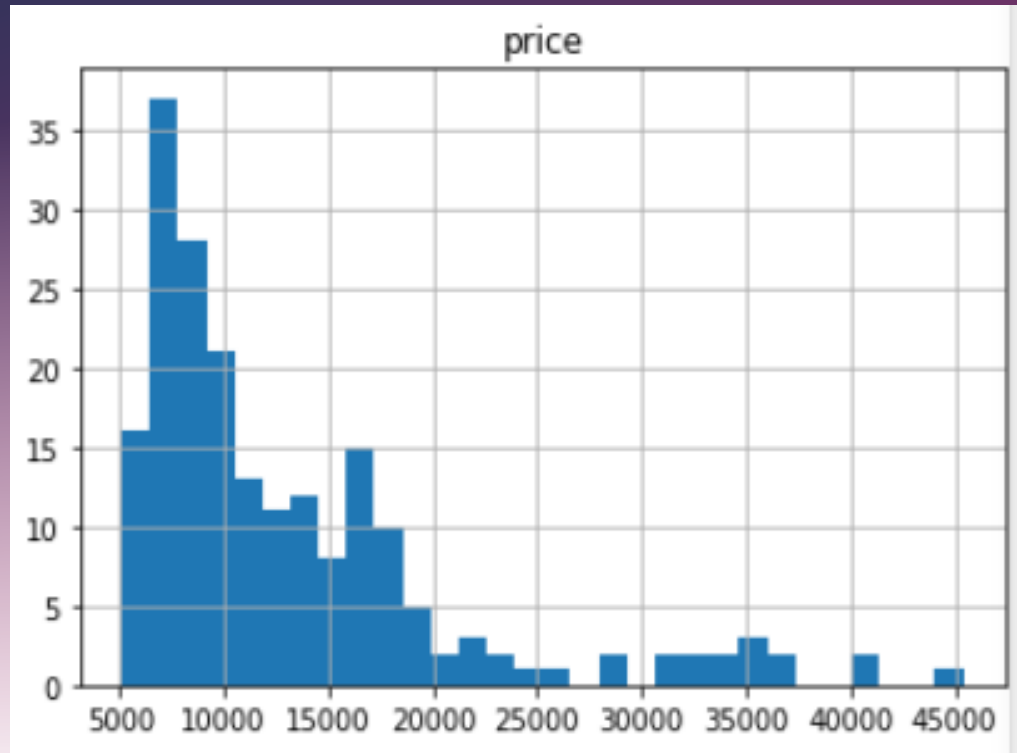
- Heatmap adalah jenis visualisasi yang menggunakan kode warna untuk mewakili nilai / kepadatan relatif data di seluruh permukaan.
- Warna-warna ini kemudian dapat digunakan untuk memeriksa data secara visual guna menemukan kelompok dengan nilai serupa dan mendeteksi tren dalam data.



Visualisasi Statistik

Histogram

- Histogram adalah salah satu visualisasi yang cukup penting dalam memahami distribusi pada data kita. Pandas Histogram menyediakan method yang memudahkan kita untuk membuat histogram.
- Plot histogram secara tradisional hanya membutuhkan satu dimensi data.
- Ini dimaksudkan untuk menunjukkan jumlah nilai atau kumpulan nilai secara serial.



Visualisasi Statistik

Histogram

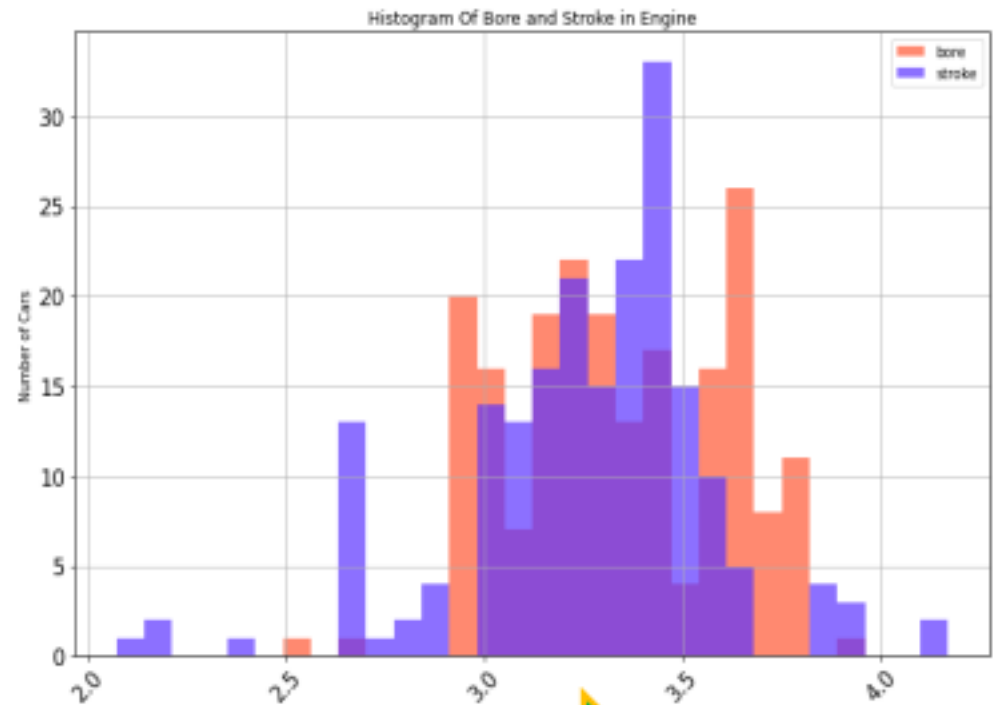
- Data yang digunakan adalah data spesifikasi mobil dari berbagai merk

symboling	normalized-losses	make	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	height	curb-weight	engine-type	num-of-cylinders	engine-size	fuel-system	bore	stroke	compre
0	3	122	alfa-romero	std	two	convertible	rwd	front	88.6	0.811148	0.890278	48.8	2548	dohc	four	130	mpfi	3.47	2.68
1	3	122	alfa-romero	std	two	convertible	rwd	front	88.6	0.811148	0.890278	48.8	2548	dohc	four	130	mpfi	3.47	2.68
2	1	122	alfa-romero	std	two	hatchback	rwd	front	94.5	0.822681	0.909722	52.4	2823	ohcv	six	152	mpfi	2.68	3.47
3	2	164	audi	std	four	sedan	fwd	front	99.8	0.848630	0.919444	54.3	2337	ohc	four	109	mpfi	3.19	3.40
4	2	164	audi	std	four	sedan	4wd	front	99.4	0.848630	0.922222	54.3	2824	ohc	five	136	mpfi	3.19	3.40

Visualisasi Statistik

Histogram

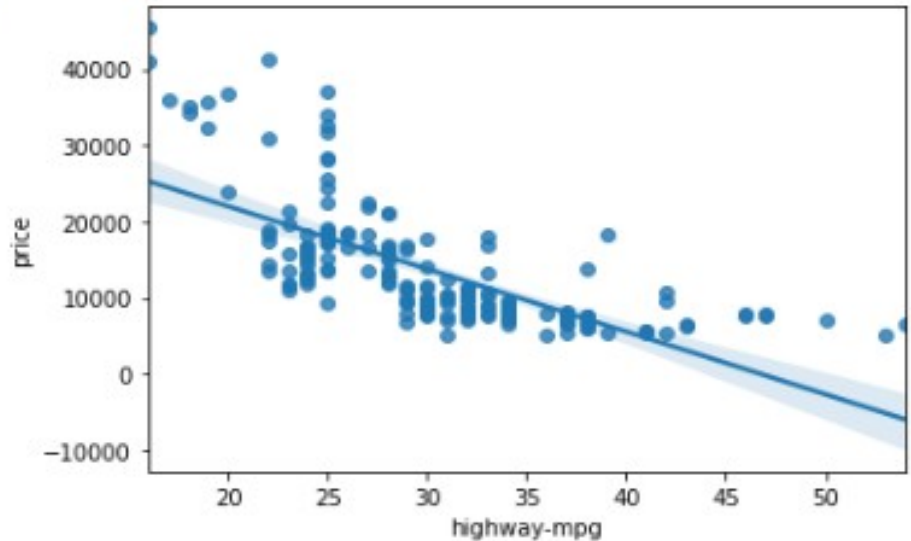
- Kita juga dapat memplot beberapa grup secara berdampingan. Di sini saya ingin melihat dua histogram, histogram price akan dikelompokkan berdasarkan roda penggerak dari kendaraan (fwd – berpengerak roda depan, 4wd – berpengerak 4 roda, atau rwd – pengerak belakang).



Visualisasi Statistik

Correlation & Causation

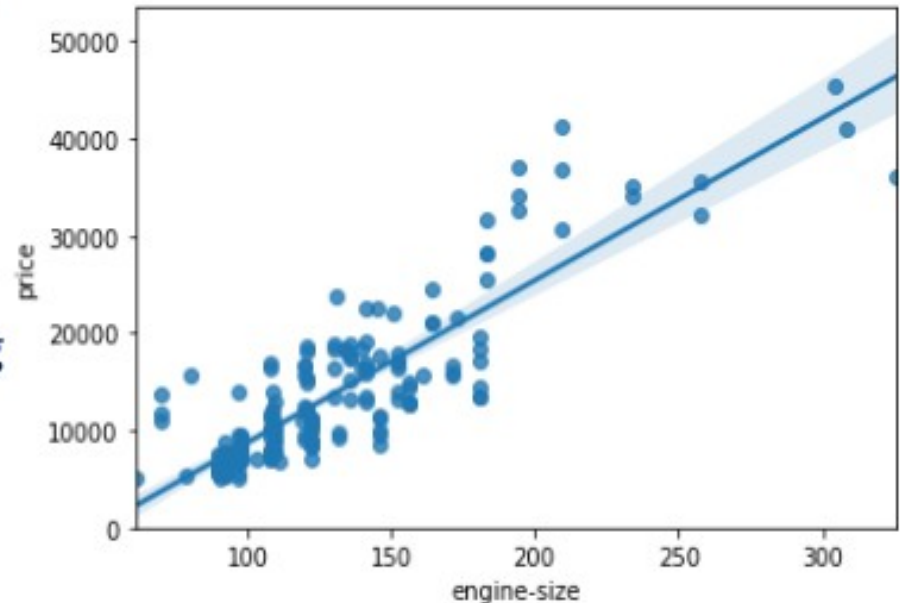
- Korelasi merupakan suatu pengukuran sejauh mana nilai saling ketergantungan antar variabel.
- Causation merupakan hubungan antara sebab dan akibat antara dua variabel
- Penting untuk mengetahui perbedaan antara keduanya dan bahwa korelasi tidak mendeskripsikan sebab-akibat.
- Menentukan korelasi jauh lebih sederhana menentukan sebab memerlukan analisis lebih lanjut



Visualisasi Statistik

Correlation & Causation

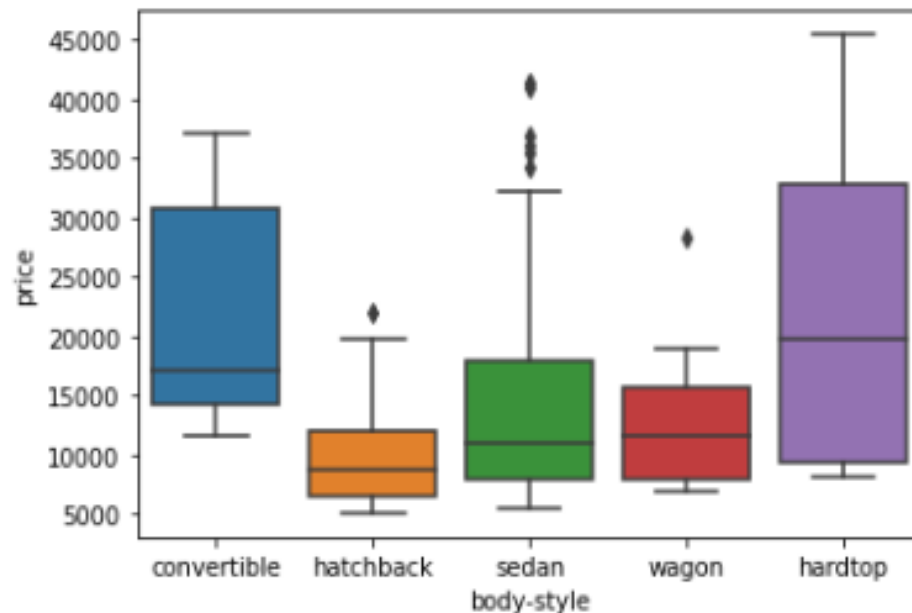
- Korelasi Pearson
- Pearson Correlation adalah metode default dari fungsi "corr". Kita dapat menghitung Korelasi Pearson dari variabel 'int64' atau 'float64'. Terkadang kita ingin mengetahui signifikansi dari estimasi korelasi, kita dapat menggunakan p-value.
- Korelasi Pearson mengukur ketergantungan linier antara dua variabel X dan Y.



Visualisasi Statistik

Variabel Kategori Statistik

- Ini adalah variabel yang menggambarkan 'karakteristik' dari unit data, dan dipilih dari sekelompok kategori. Variabel kategori dapat memiliki tipe "objek" atau "int64". Cara yang baik untuk memvisualisasikan variabel kategori adalah dengan menggunakan boxplot.
- Boxplot menggambarkan variable variable statistic seperti quartil 1, median / quartil 2, quartil 3, nilai maksimum, nilai minimum, dan outlier.

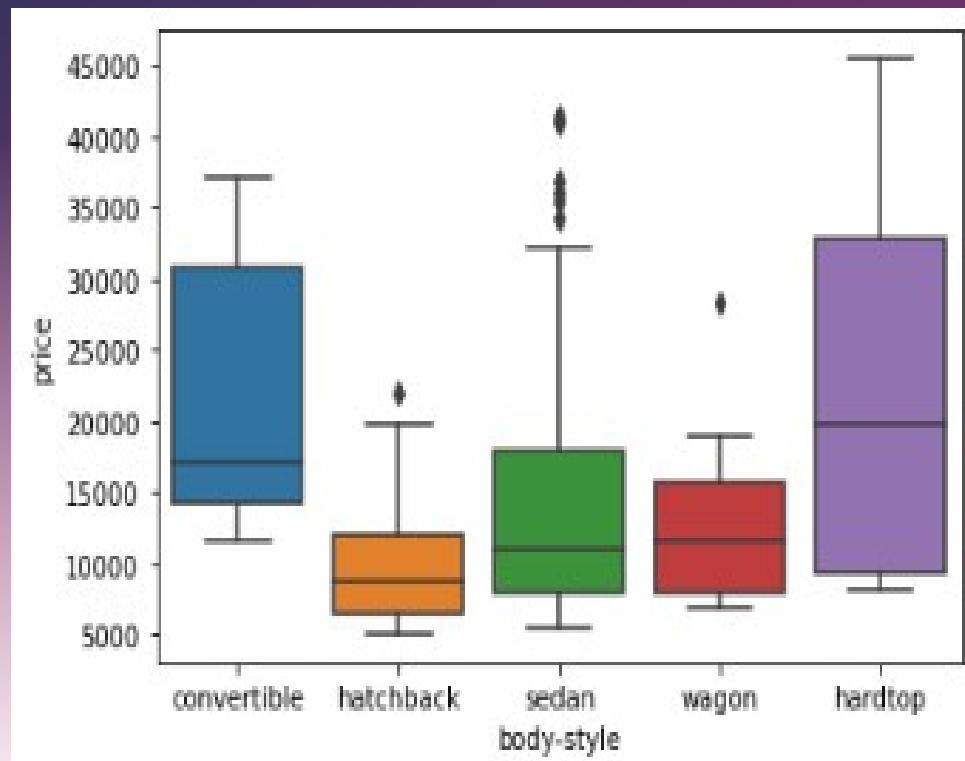


Visualisasi Statistik

Descriptive Statistic

- Fungsi deskripsikan secara otomatis menghitung statistik dasar untuk semua variabel kontinu.
- Analisis yang bisa kita dapatkan dari deskriptif statistik adalah
 - Jumlah variabel
 - Rata-rata
 - Standard deviasi
 - Nilai minimal
 - IQR (Interquartile Range: 25%, 50% and 75%)
 - Nilai Maximal

yearling normalized-losses wheel-base length width height



Visualisasi Statistik

Grouping

- Method "groupby" digunakan untuk mengelompokkan data menurut kategori yang berbeda. Data dikelompokkan berdasarkan satu atau beberapa variabel dan analisis dilakukan pada kelompok individu.
- Sebagai contoh, mari kita kelompokkan berdasarkan variabel "roda penggerak". Kita melihat bahwa ada 3 kategori roda penggerak yang berbeda.

- ```
df['drive-wheels'].unique()
array(['rwd', 'fwd', '4wd'], dtype=object)
```

# Visualisasi Statistik

## Grouping

- Anda juga dapat mengelompokkan dengan beberapa variabel. Misalnya, mari kita kelompokkan berdasarkan 'roda penggerak' dan 'body-style'.
- Ini mengelompokkan dataframe dengan kombinasi unik 'drive-wheels' dan 'body-style'. Kita dapat menyimpan hasilnya dalam variabel 'grouped\_test1'.

|    | drive-wheels | body-style  | price        |
|----|--------------|-------------|--------------|
| 0  | 4wd          | hatchback   | 7603.000000  |
| 1  | 4wd          | sedan       | 12647.333333 |
| 2  | 4wd          | wagon       | 9095.750000  |
| 3  | fwd          | convertible | 11595.000000 |
| 4  | fwd          | hardtop     | 8249.000000  |
| 5  | fwd          | hatchback   | 8396.387755  |
| 6  | fwd          | sedan       | 9811.800000  |
| 7  | fwd          | wagon       | 9997.333333  |
| 8  | rwd          | convertible | 23949.600000 |
| 9  | rwd          | hardtop     | 24202.714286 |
| 10 | rwd          | hatchback   | 14337.777778 |
| 11 | rwd          | sedan       | 21711.833333 |
| 12 | rwd          | wagon       | 16994.222222 |

# Visualisasi Statistik

## ANOVA

- Analysis of Varians (ANOVA) adalah metode statistik yang digunakan untuk menguji apakah ada perbedaan yang signifikan antara rata-rata dua kelompok atau lebih.
- ANOVA mengembalikan dua parameter
  - F-Score:
  - P-Value
- F-Score: ANOVA mengasumsikan rata-rata semua kelompok adalah sama, anova akan menghitung seberapa jauh rata-rata yang sebenarnya menyimpang dari asumsi, dan melaporkannya sebagai F-Score.
- Skor yang lebih besar berarti ada perbedaan yang lebih besar antara rata-rata.
- P-Value: Nilai-P menunjukkan seberapa signifikan secara statistik nilai skor yang dihitung.

# Visualisasi Statistik

## ANOVA

- Jika variabel harga pada dataset mobil sangat berkorelasi dengan variabel lainya, ANOVA akan mengembalikan skor F-Score yang cukup besar dan nilai-p yang kecil.
- ANOVA menganalisis perbedaan antara kelompok yang berbeda dari variabel yang sama, fungsi groupby akan berguna dalam kasus ANOVA.
- Mari kita lihat apakah jenis 'roda penggerak' mempengaruhi 'harga',

**#JADIJAGOANDIGITAL**  
**TERIMA KASIH**