

Visualisasi

Ali Akbar Septiandri

Universitas Al-Azhar Indonesia

aliakbars@live.com

October 10, 2019

Selayang Pandang

- ① Visualisasi Data
- ② Tipe-tipe Grafik
- ③ Visualisasi Efektif

Bahan Bacaan

- 1 Few, S. (2012). Show me the numbers: Designing tables and graphs to enlighten. Analytics Press.
- 2 `https://nces.ed.gov/forum/pdf/NCES_table_design.pdf`

“There are three kinds of lies:
lies, damned lies, and **statistics**.”

Visualisasi Data

Sebagai analis data yang baik,
penting sekali untuk memvisualisasikan data.

Mengapa?

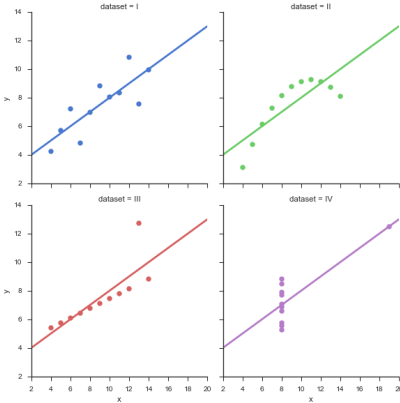
Figure 1 displays bar charts showing the minimum inhibitory concentration (MIC) of Penicillin (P), Streptomycin (S), and Neomycin (N) for various bacterial strains. The y-axis represents the MIC on a logarithmic scale from 0.001 to 1000. The x-axis lists the bacterial strains. The legend indicates: Penicillin (P) in blue, Streptomycin (S) in red, and Neomycin (N) in black. The chart is divided into sections for different bacterial groups: Mycobacterium tuberculosis, Pseudomonas aeruginosa, Aerobacter aerogenes, Klebsiella pneumoniae, Escherichia coli, Salmonella schottmuelleri, Brucella abortus, Salmonella (Eberthella) typhosa, Streptococcus fecalis, Streptococcus viridans, Streptococcus hemolyticus, Diplococcus pneumoniae, Proteus vulgaris, Bacillus anthracis, Staphylococcus aureus, and Staphylococcus albus. The MIC values are indicated by the height of the bars. A red circle highlights the MIC values for Streptococcus fecalis, Streptococcus hemolyticus, and Diplococcus pneumoniae.

Bacterial Strain	Penicillin (P)	Streptomycin (S)	Neomycin (N)
<i>Mycobacterium tuberculosis</i>	1000	1	1
<i>Pseudomonas aeruginosa</i>	1000	1	0.5
<i>Aerobacter aerogenes</i>	1000	1	1
<i>Klebsiella pneumoniae</i>	1000	1	1
<i>Escherichia coli</i>	100	0.1	0.1
<i>Salmonella schottmuelleri</i>	10	0.1	0.1
<i>Brucella abortus</i>	0.1	1	0.1
<i>Salmonella (Eberthella) typhosa</i>	0.1	0.1	0.1
<i>Streptococcus fecalis</i>	0.1	0.1	0.1
<i>Streptococcus viridans</i>	0.1	10	10
<i>Streptococcus hemolyticus</i>	0.1	10	10
<i>Diplococcus pneumoniae</i>	0.1	10	10
<i>Proteus vulgaris</i>	0.1	0.1	0.1
<i>Bacillus anthracis</i>	0.1	0.1	0.1
<i>Staphylococcus aureus</i>	0.1	0.1	0.1
<i>Staphylococcus albus</i>	0.1	0.1	0.1

Really a streptococcus!
(realized ~20 years later)

Adapted from Brian Schmotzer

Visualisasi vs Summary Statistics

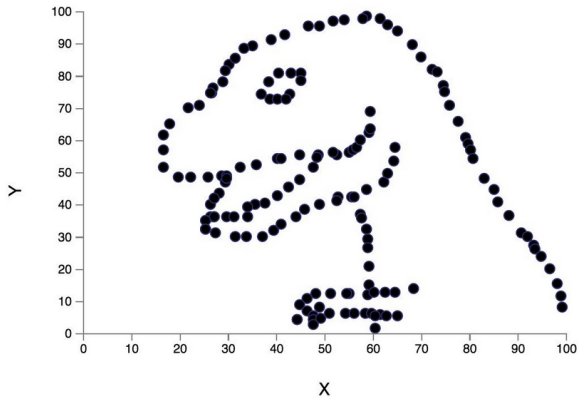


Gambar: *The infamous Anscombe's quartet*

Keempat data tersebut punya
 $\mu_x, \mu_y, \sigma_x, \sigma_y, N, \rho$ yang sama!

*Don't trust summary statistics.
Always **visualize** your data!*

N = 157 ; X mean = 50.7333 ; X SD = 19.5661 ; Y mean = 46.495 ; Y SD = 27.2828 ;
Pearson correlation = -0.1772



Alberto Cairo @albertocairo · 15 Aug 2016

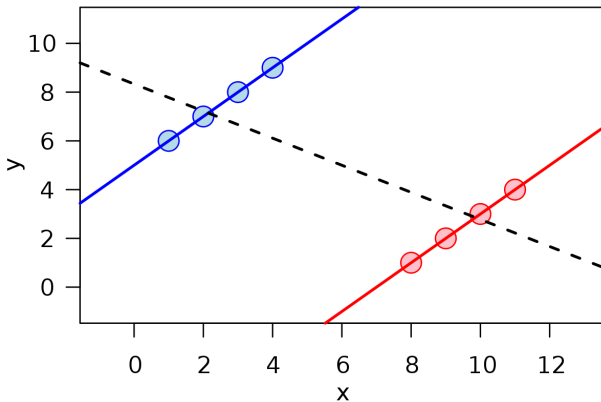
Don't trust summary statistics. Always visualize your data first robertgrantstats.co.uk/drawmydata.html pic.twitter.com/5j94Dw9UAf

13

895

984

Simpson's Paradox



Gambar: Kesalahan dalam interpretasi pola karena tidak dibagi per grup

Dengan visualisasi, bukan berarti bebas masalah.

Kita tetap bisa berbohong!

WEDDING BUDGET BREAKDOWN

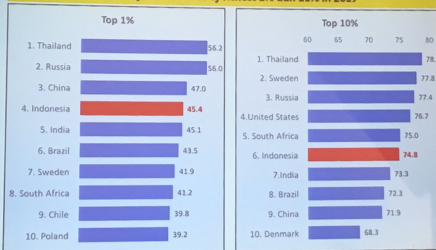


THE BIGGEST
CHUNK OF YOUR
WEDDING WILL
TYPICALLY BE
SPENT ON
CATERING!



The world's most unequal countries

Share of total wealth of richest 1% dan 10% in 2017

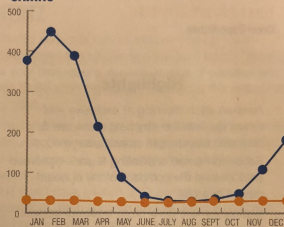


Source: Credit Suisse Global Wealth Databook 2017.

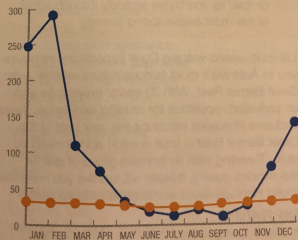


TEMPERATURE & RAIN CHART

CAIRNS

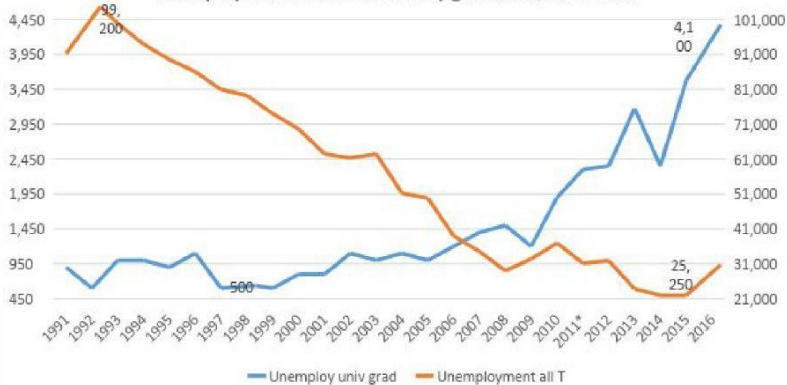


TOWNSVILLE



● RAINFALL (mm) ● TEMPERATURE (max. °C)

Figure 1: Trends in the national unemployment level and the unemployment level of university graduates, 1991-2016





Ainun Najib

@ainunnajib

Following



A good example of “How to Lie with Statistics”. Maaf mas dokter [@Gamal_Albinsaid](#) , we’re good friends, namun kali ini anda keliru dan mengelirukan masyarakat.



dr. Gamal Albinsaid [@Gamal_Albinsaid](#)

3. Ini bukan soal angka, ini soal luka, Luka Indonesiaku.

Saya katakan ini karena saya menyadari betapa beratnya tugas yang menunggu kita semua di depan sana....

Show this thread

3:02 PM - 7 Jan 2019 from [Central Region, Singapore](#)

5,253 Retweets 3,948 Likes

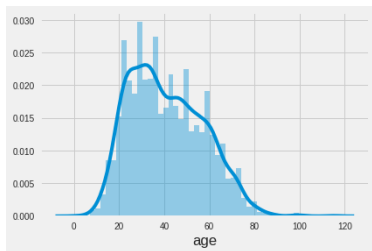


253 5.3K 3.9K

Gambar: Dapat digunakan untuk keuntungan politik. Sumber: [@ainunnajib](#)

Tipe-tipe Grafik

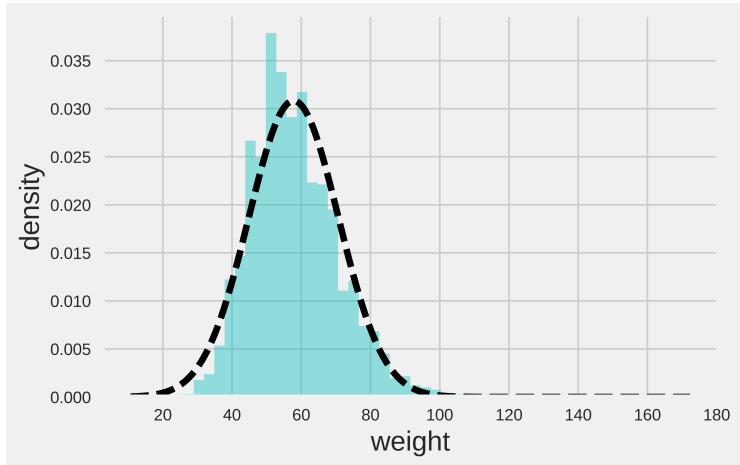
Histogram



Gambar: Contoh histogram dengan *kernel density estimation*

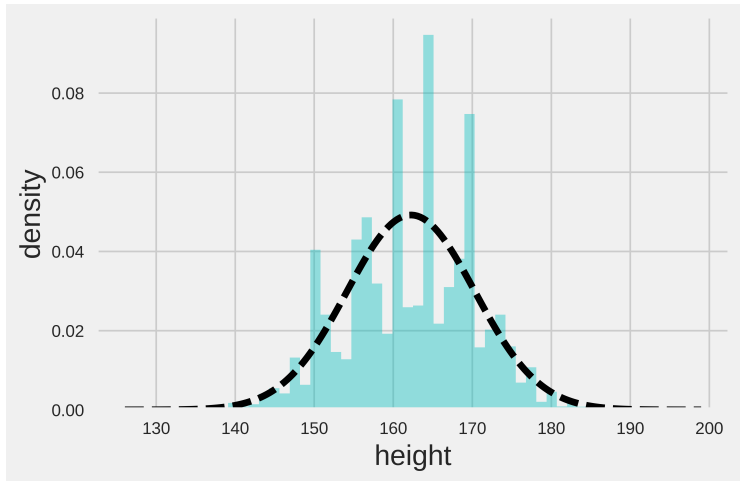
- ① Digunakan untuk melihat distribusi dari variabel
- ② Dibagi berdasarkan *bins*
- ③ Sangat bergantung pada jumlah *bins* yang digunakan!

Histogram



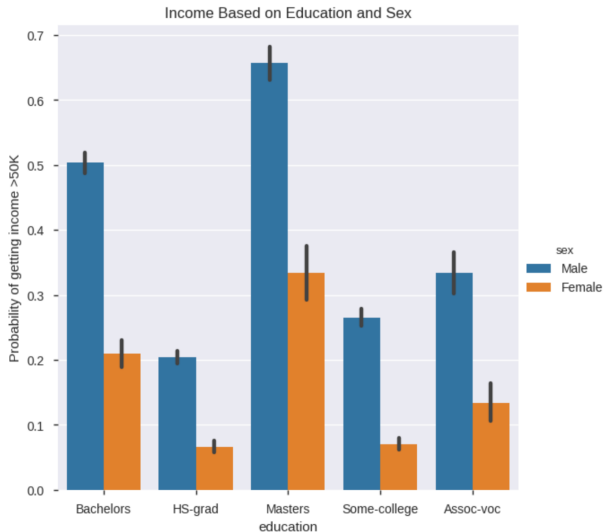
Gambar: Apa yang aneh dari histogram ini?

Histogram

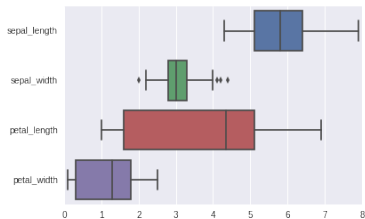


Gambar: Apa yang aneh dari histogram ini?

Bar Plot



Box Plot



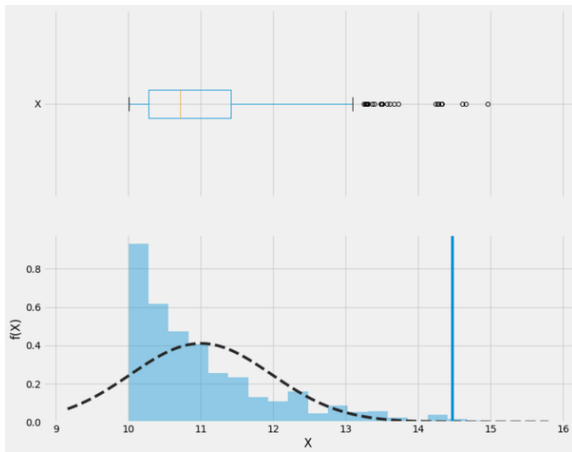
Gambar: Box plot untuk membandingkan atribut

- 1 Menggambarkan jangkauan dan persentil
- 2 Dapat digunakan untuk membandingkan atribut
- 3 Membantu menemukan pencilan

Potensi Masalah pada Box Plot

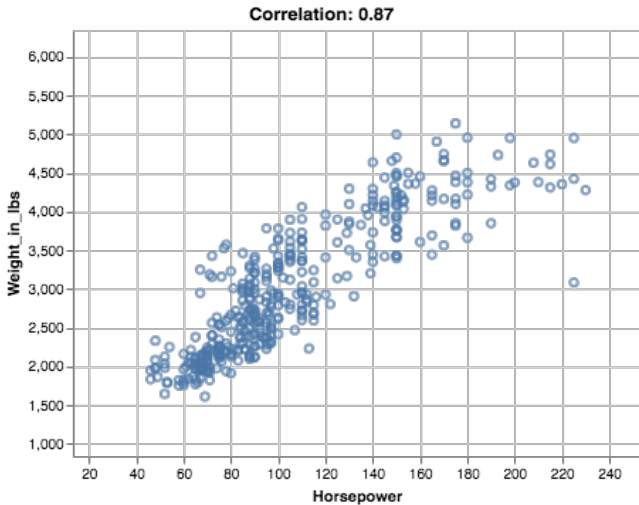
https://twitter.com/randal_olson/status/895678675814764544

Box Plot vs. Histogram

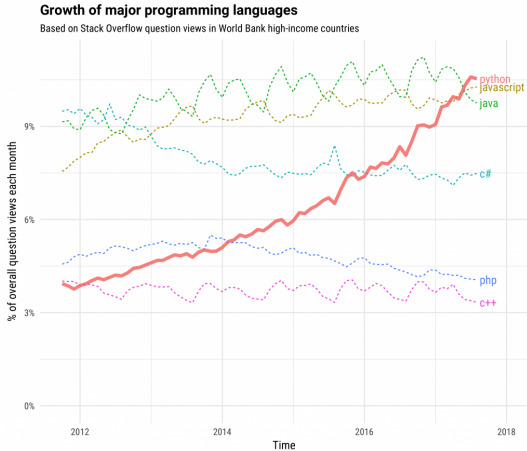


Gambar: Perbandingan box plot dan histogram untuk data yang sama

Scatter Plot

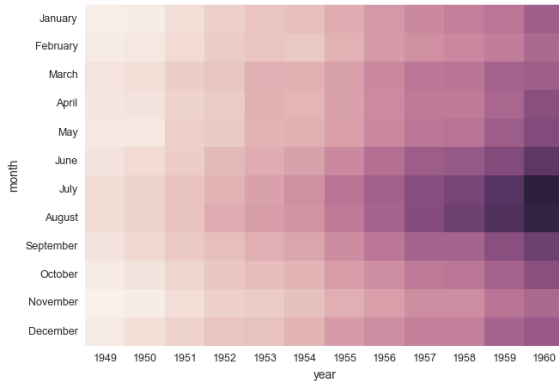


Timeseries



Gambar: Perubahan persentase pengunjung pertanyaan berdasarkan waktu [Robinson, 2017]

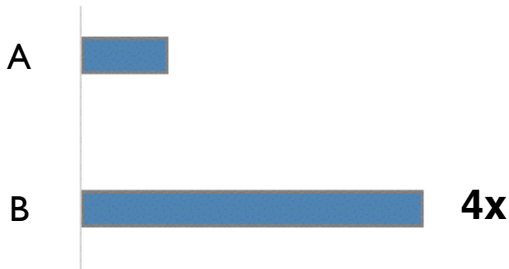
Heatmap



Gambar: Aktivitas penerbangan berdasarkan dua dimensi waktu

Visualisasi Efektif

How much longer?



How much steeper slope?

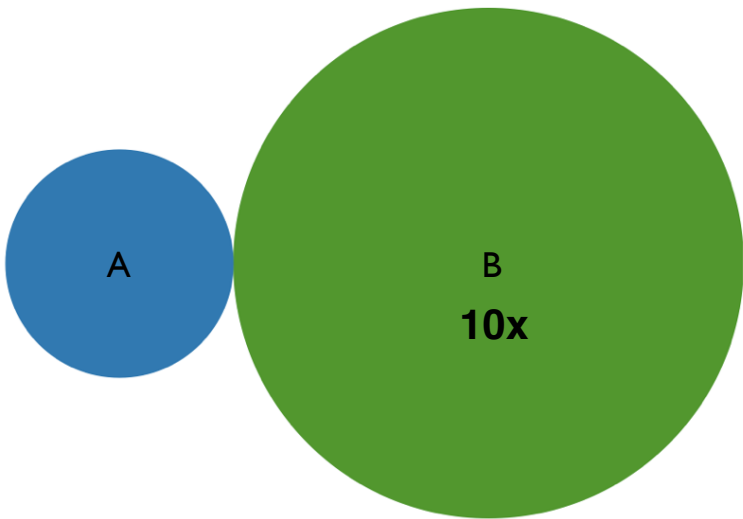


A
4x



B

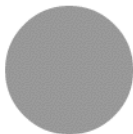
How much larger area?



How much darker?



A



B

2x

How much bigger value?



A



B

4x



Most
Efficient



Least
Efficient

Position



Length



Slope



Angle



Area



Intensity



Color



Shape

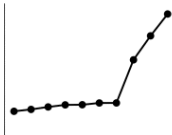
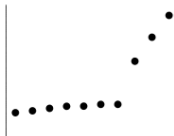


Quantitative

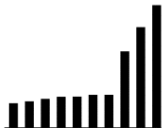
Ordered

Categories

Most Effective



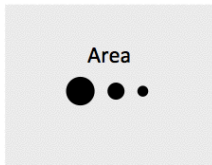
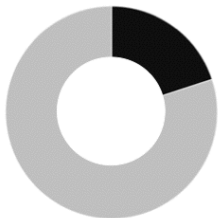
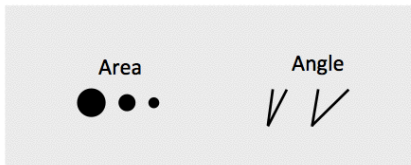
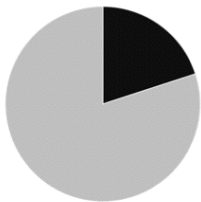
Position



Length

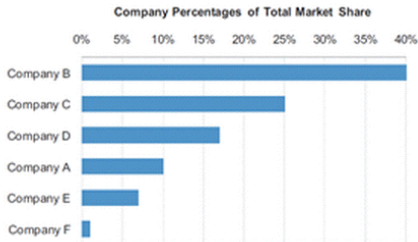
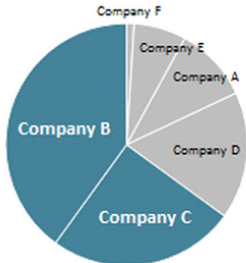


Less Effective



Pie vs. Bar Charts

65% of the market is controlled by companies B and C



Least Effective

SANFORD AND SELNICK

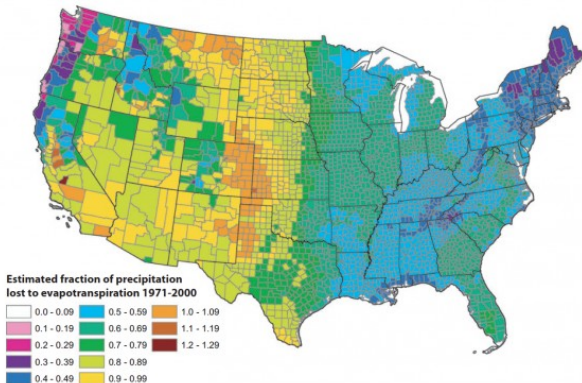


FIGURE 13. Estimated Mean Annual Ratio of Actual Evapotranspiration (ET) to Precipitation (P) for the Conterminous U.S. for the Period 1971-2000. Estimates are based on the regression equation in Table 1 that includes land cover. Calculations of ET/P were made first at the 800-m resolution of the PRISM climate data. The mean values for the counties (shown) were then calculated by averaging the 800-m values within each county. Areas with fractions >1 are agricultural counties that either import surface water or mine deep groundwater.

CS109 Data Science

Lecture 3: Exploratory Data Analysis

Referensi



Michael Waskom (2015)

Anscombe's quartet

http://seaborn.pydata.org/examples/anscombes_quartet.html



David Robinson (6 September 2017)

The Incredible Growth of Python

<https://stackoverflow.blog/2017/09/06/incredible-growth-python/>

Terima kasih