# Osamelci
http://goo.gl/forms/VjRKEX0eEY

# Osamelci



FIGURE 7-32 Estimated weather data from Weather Underground

[Yau: Visualize this.]

# Iskanje osamelcev

[Yau: Visualize this.]

# Iskanje osamelcev - primer



FIGURE 7-34  FlowingData subscriber counts over time

[Yau: Visualize this.]

# Iskanje osamelcev - primer



FIGURE 7-35 Histogram showing distribution of subscriber counts
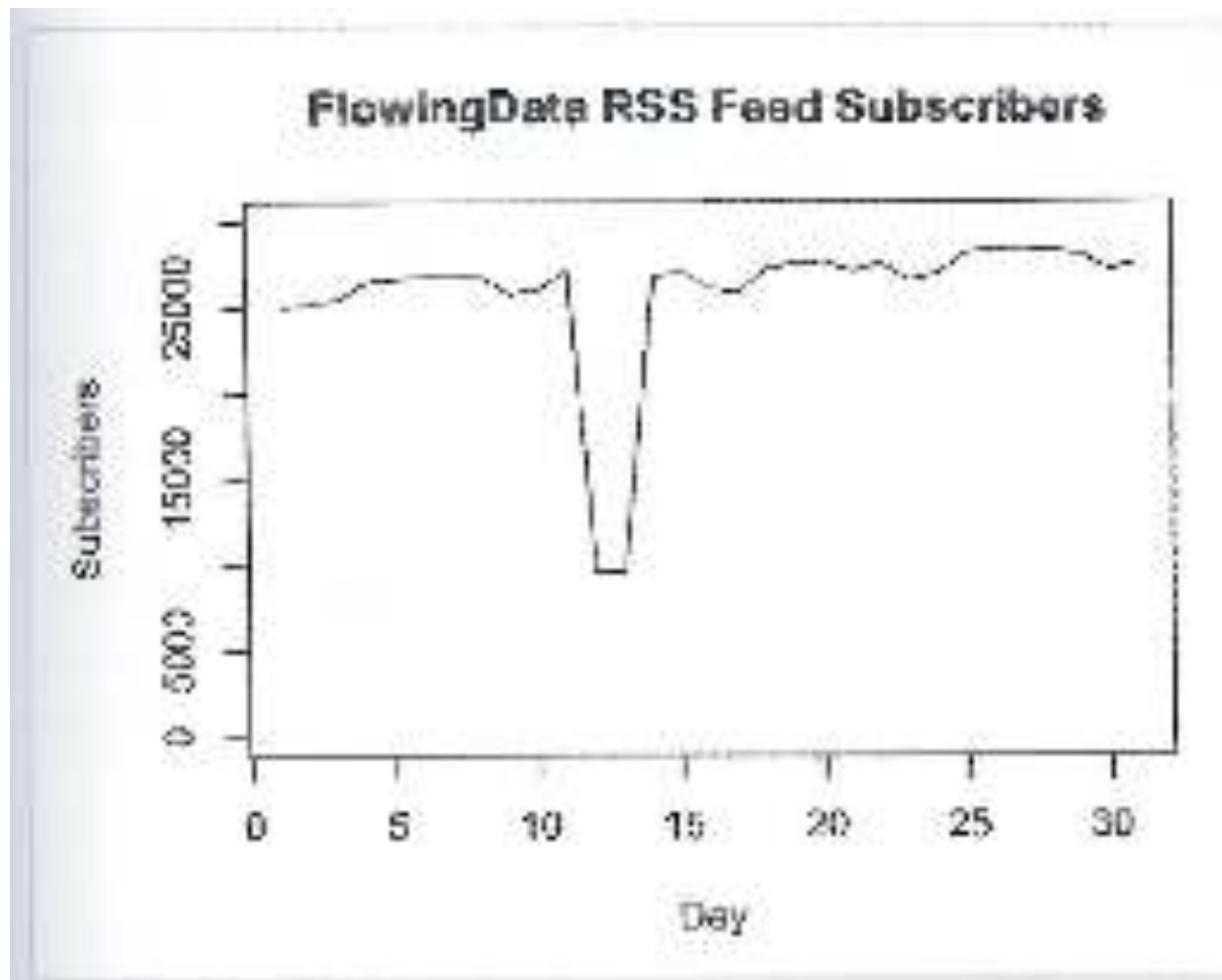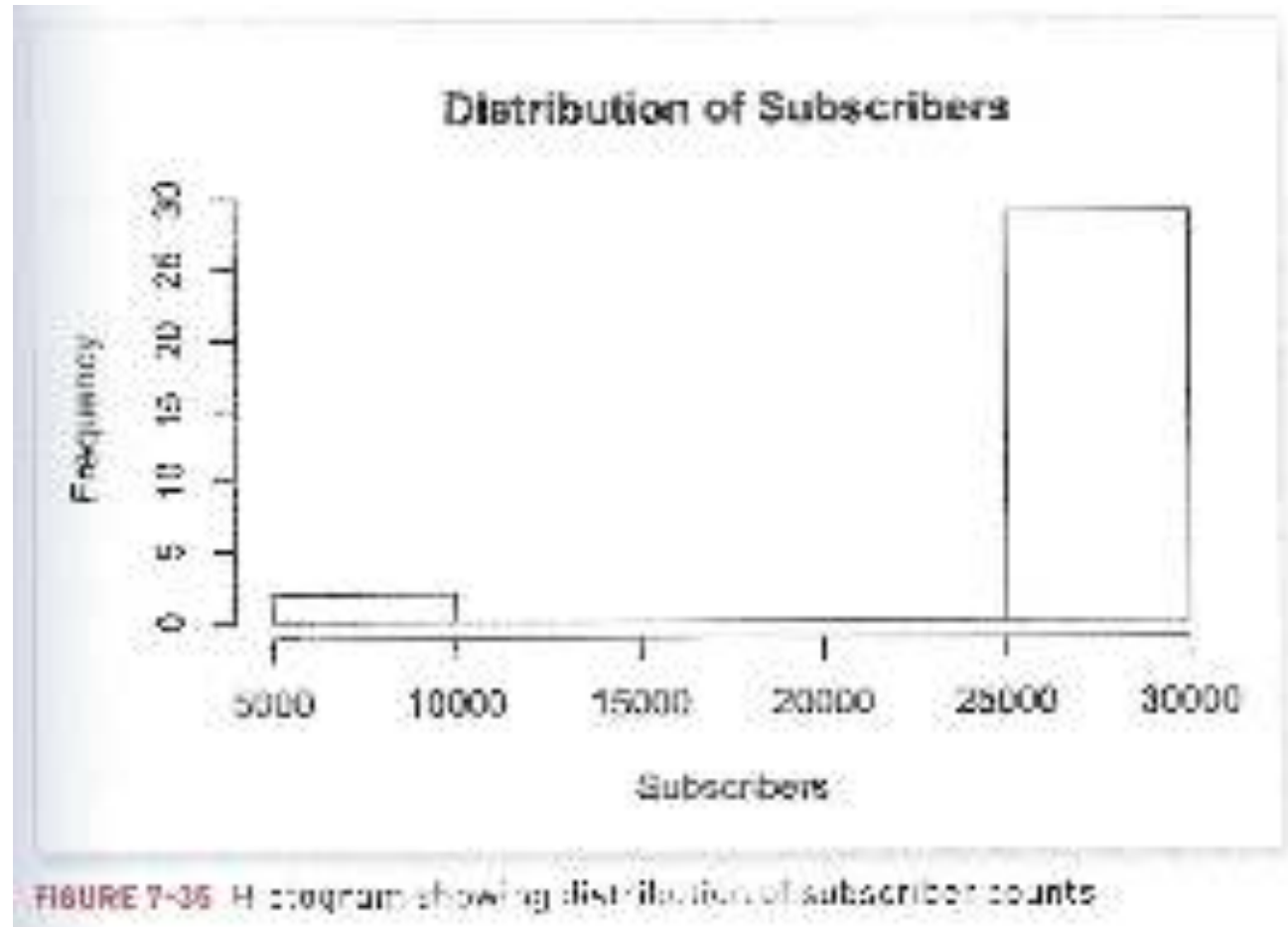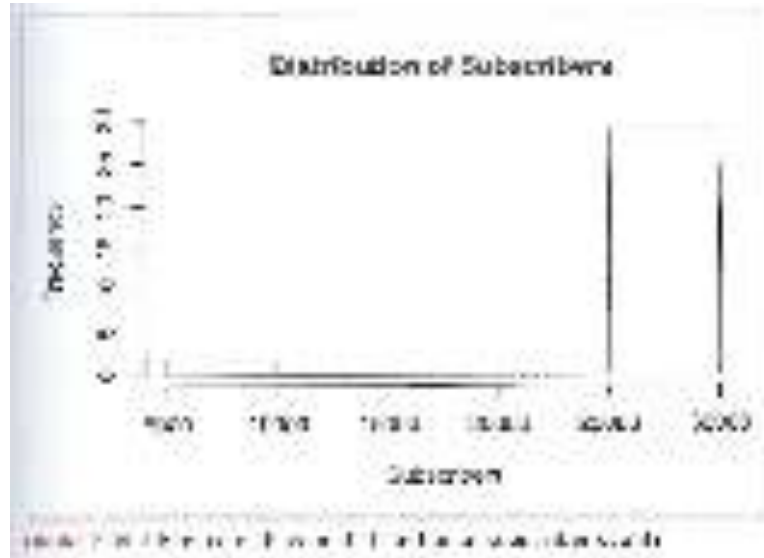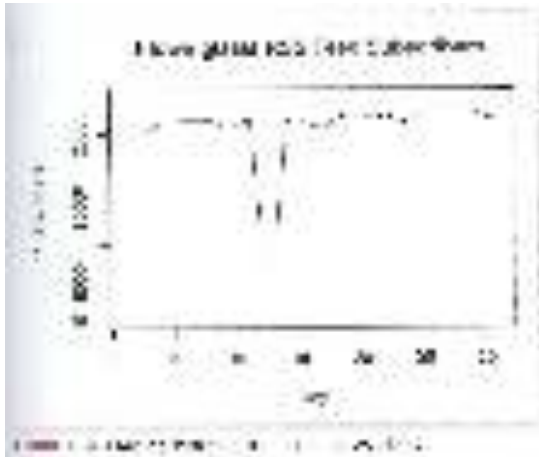
[Yau: Visualize this.]

# Iskanje osamelcev - primer

[Yau: Visualize this.]
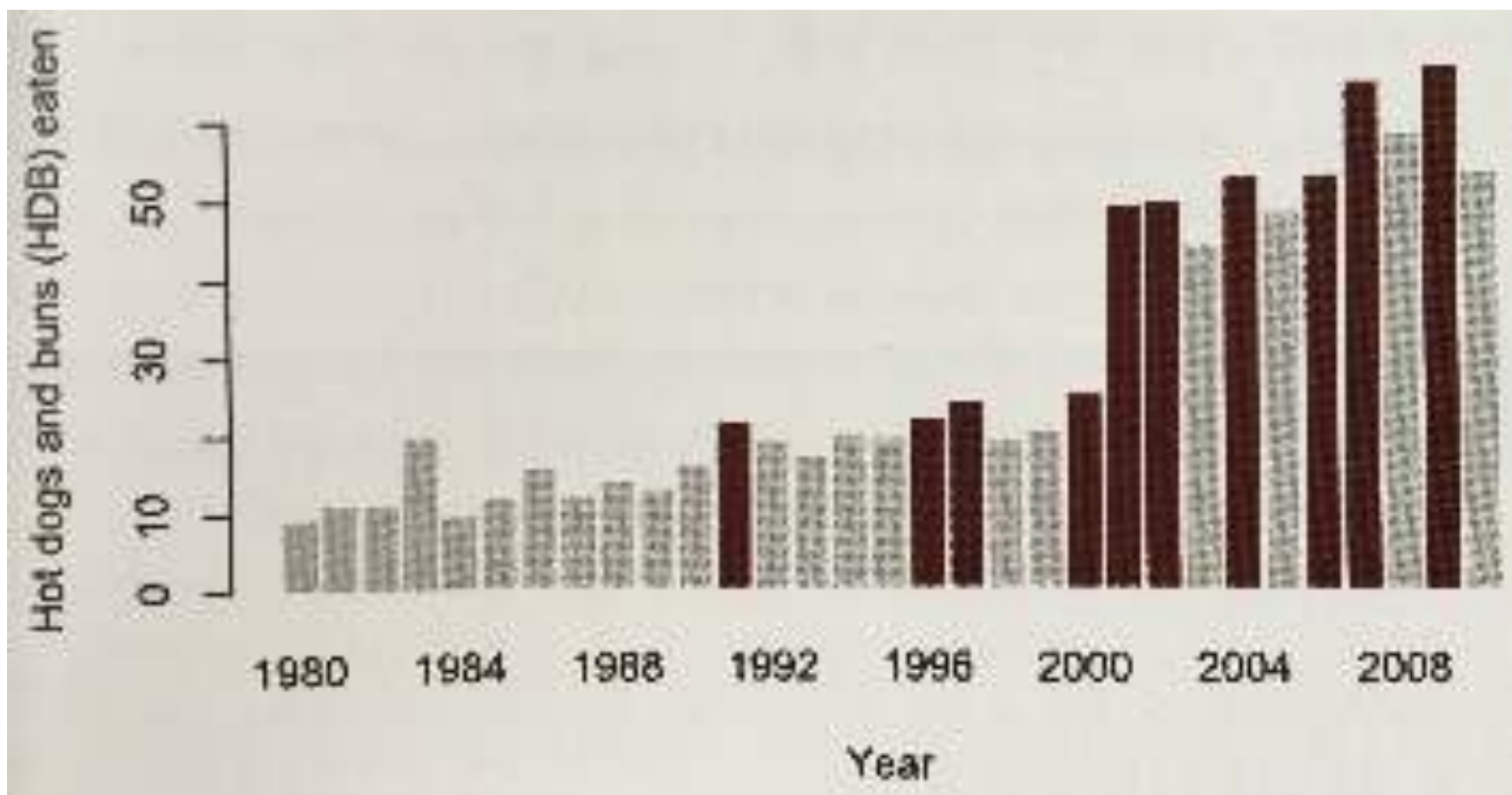
# Iskanje osamelcev - primer

[Yau: Visualize this.]

# Iskanje osamelcev – primer II

Tekmovanje (10-12 min) "Nathan's Hot Dog Eating Contest" v Coney Island, New York, ob prazniku "Independence day", vsako leto od 1972.

[Yau: Visualize this.]

# Iskanje osamelcev – primer II

Tekmovanje (10-12 min) "Nathan's Hot Dog Eating Contest" v Coney Island, New York, ob prazniku "Independence day", vsako leto od 1972.

[Yau: Visualize this.]

# Iskanje osamelcev – primer II

Tekmovanje (10-12 min) "Nathan's Hot Dog Eating Contest" v Coney Island, New York, ob prazniku "Independence day", vsako leto od 1972.

Leta 2001 prvič tekmuje Takeru Kobayashi (l.r. 1978), ki se zelo resno loti priprav in uspe podvojiti rekord Iz 25.25 na 50 hot dog-ov.

[Yau: Visualize this.]
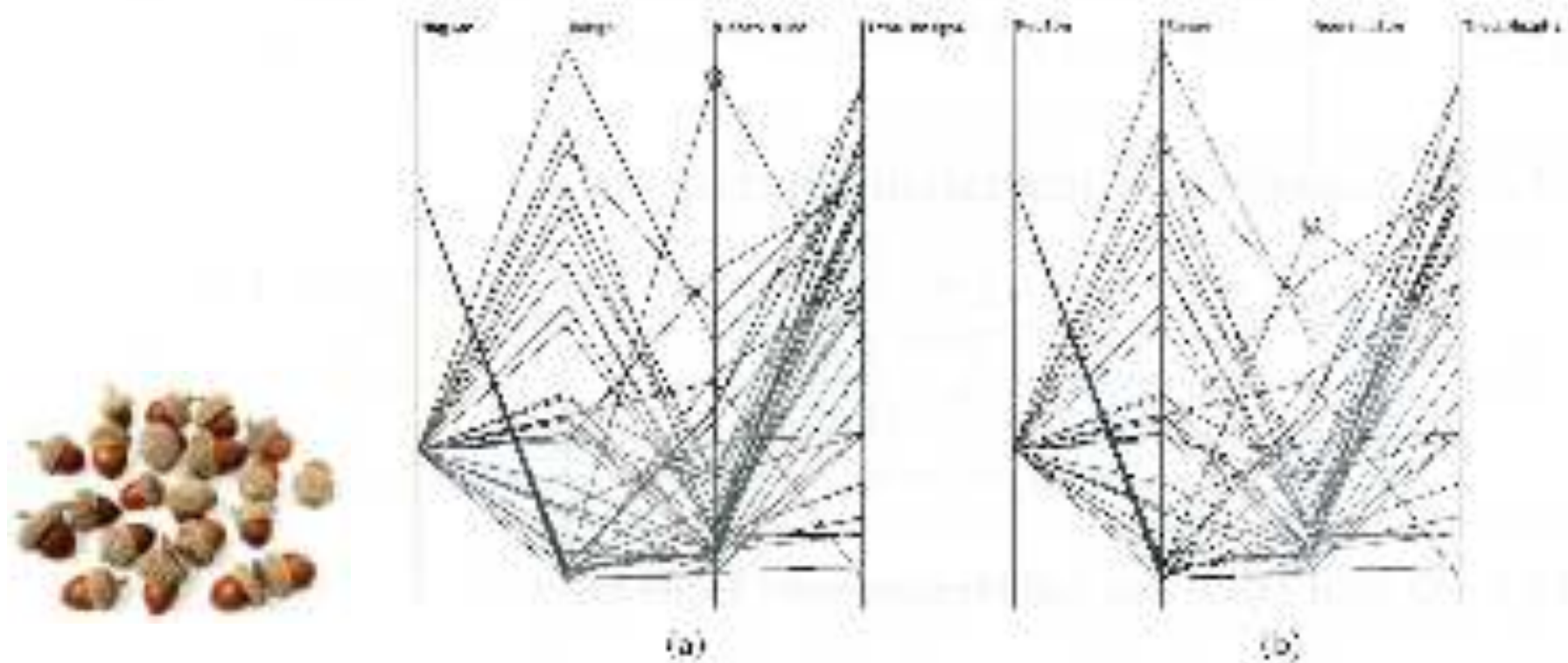
# Iskanje osamelcev – primer III



Figure 13.1. Parallel coordinates view of data set describing acorn attributes, with a single outlier (circled) in the acorn size dimension (a) in its original position and (b) with the distance artificially shortened. [378]. (Image © 1997 IEEE.)

[Ward et al.: Interactive Data Visualization.]
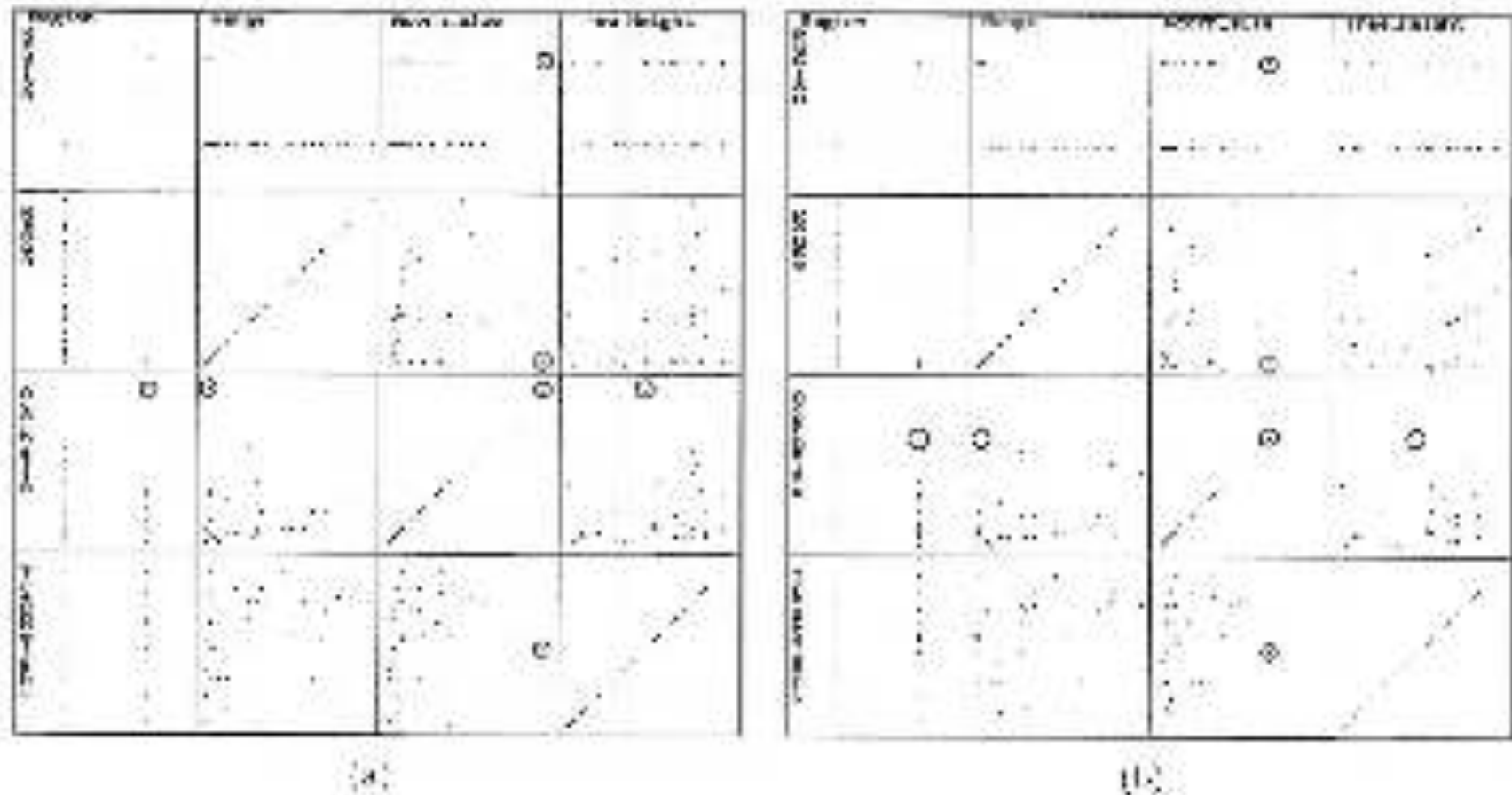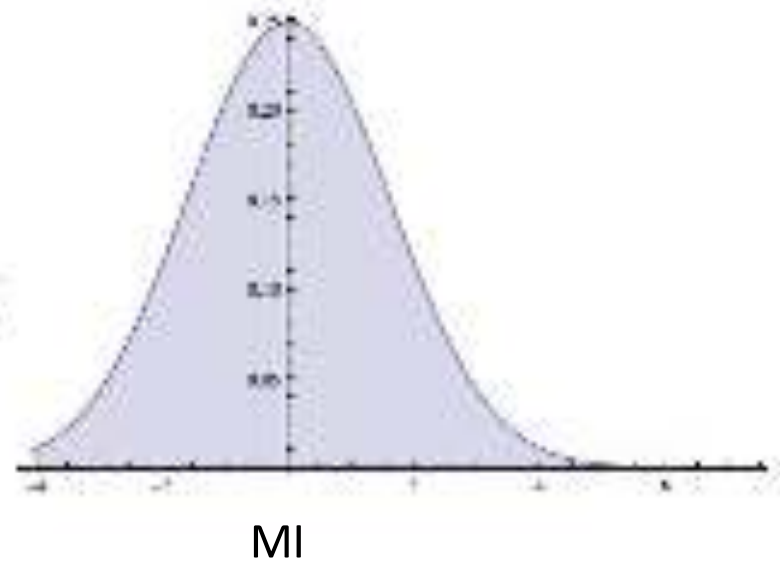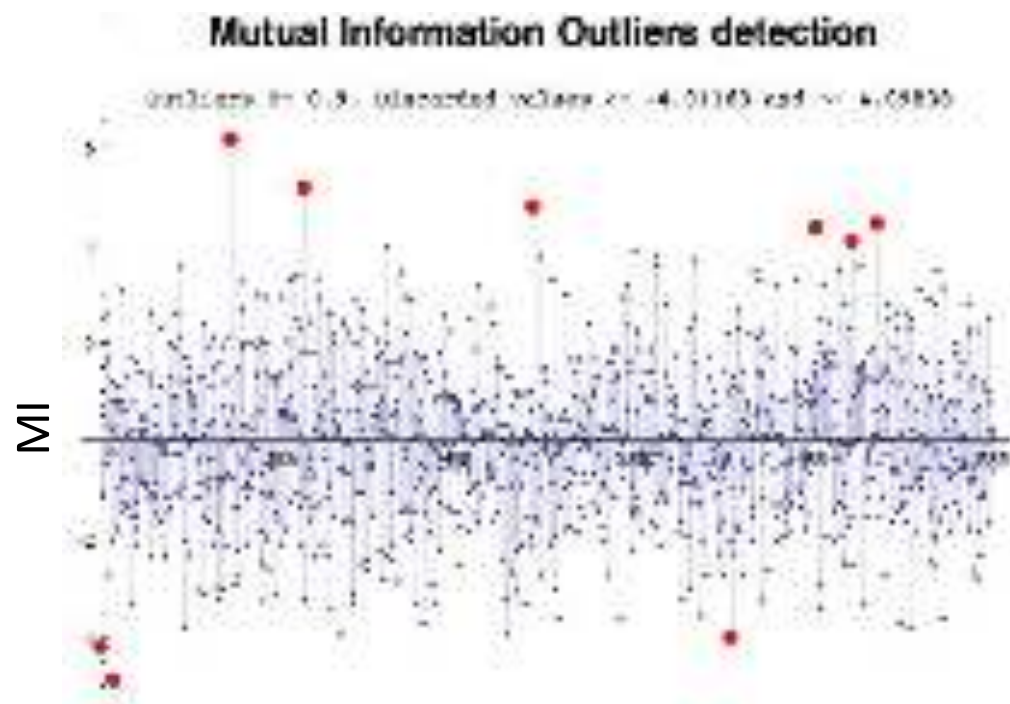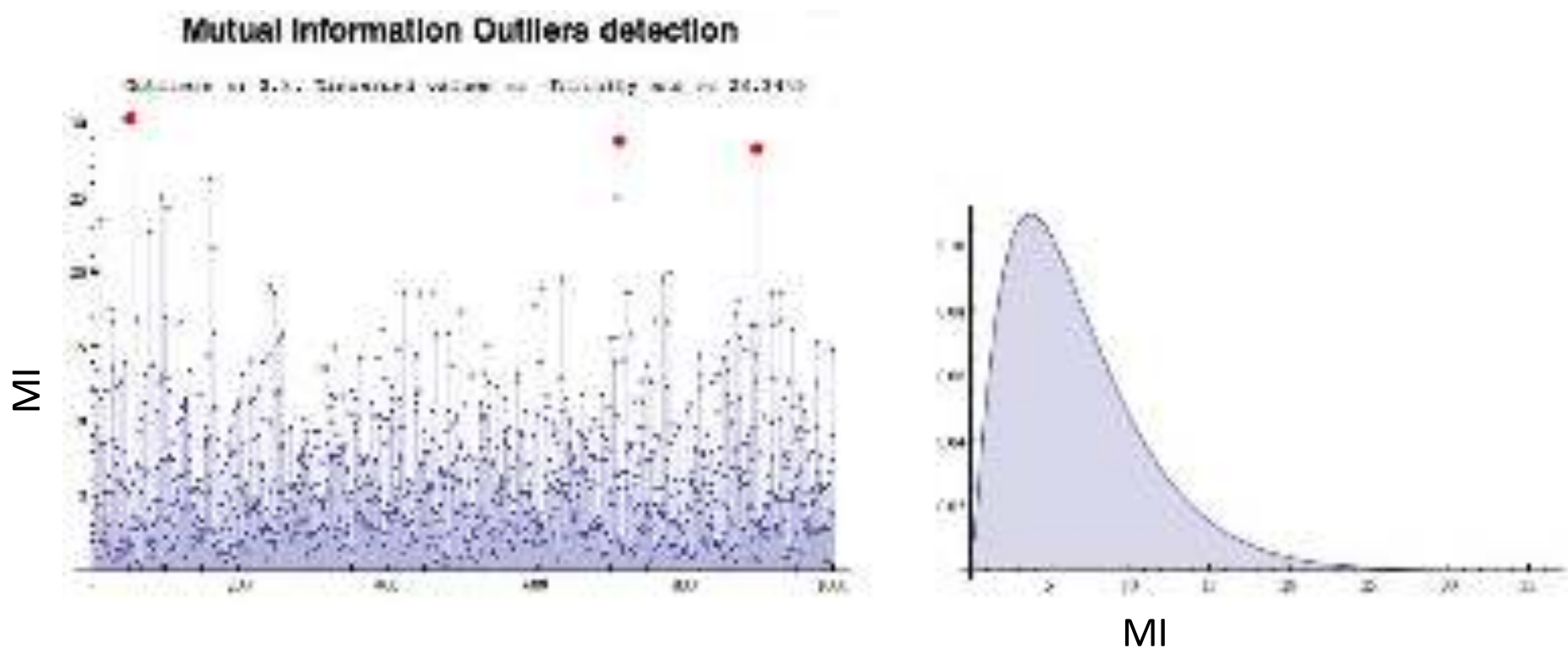
# Iskanje osamelcev – primer III



Figure 13.2  Identifying outliers with scatterplot matrices: same data as previous figure, using scatterplot matrices [575]. (Image © 1987 IEEE.)

[Ward et al.: Interactive Data Visualization.]

# Iskanje osamelcev – različne distribucije



Mutual Information Outliers detection

MI

MI

[http://www.analyticbridge.com]

# Iskanje osamelcev – različne distribucije



MI

[http://www.analyticbridge.com]

# Iskanje osamelcev: Q-Q plot
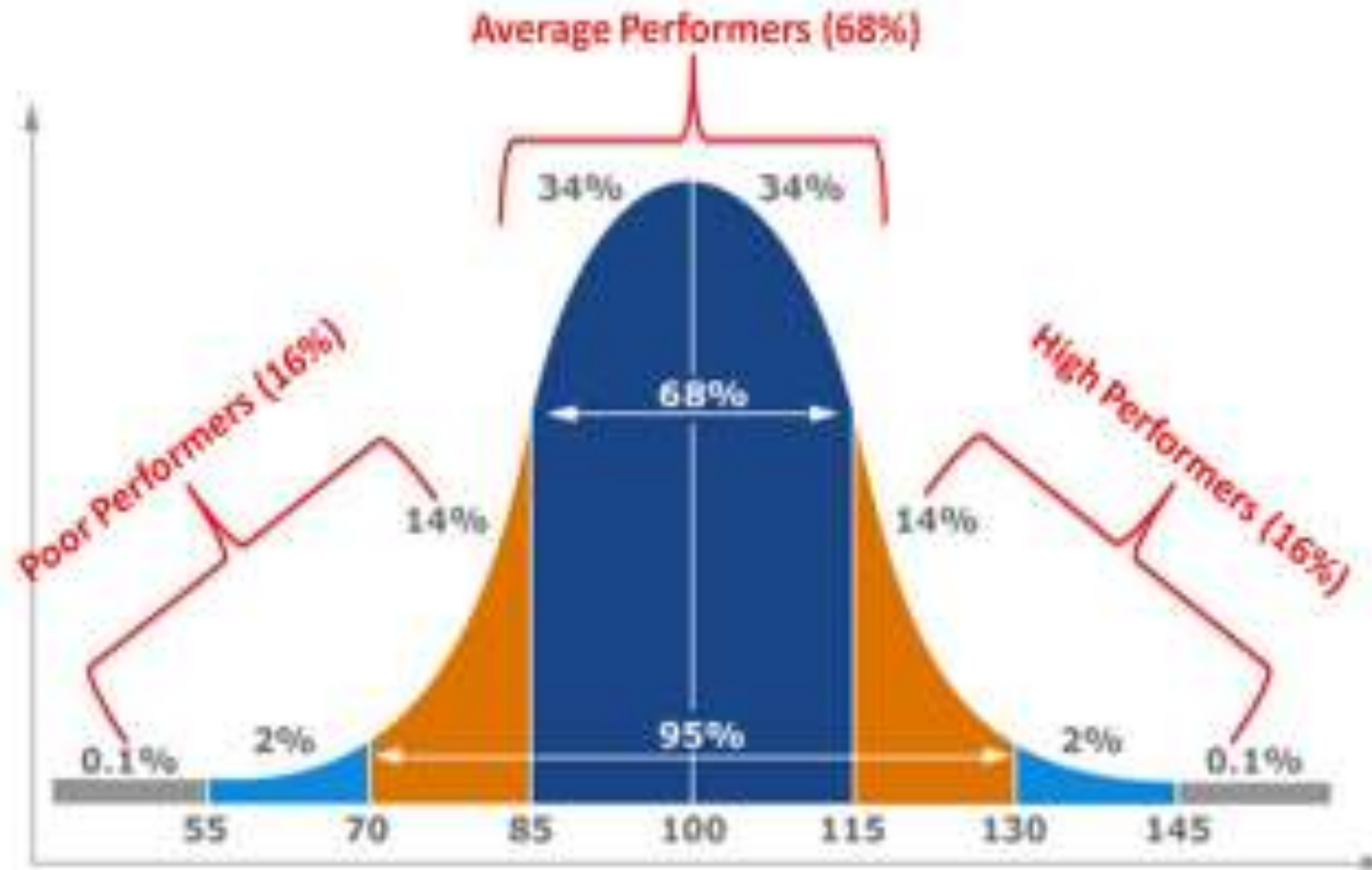


Sample Plot

The points on this plot form a nearly linear pattern, which indicates that the normal distribution is a good model for this data set.

[http://www.itl.nist.gov/div898/handbook/eda/section3/normprpl.htm]

# Iskanje osamelcev: Z-score

[http://www.compensationcafe.com/2014/04/ding-dong-the-wicked-bell-curve-is-dead.html]

# Iskanje osamelcev: Z-score

Z-Scores and Modified Z-Scores

The Z-score of an observation is defined as

$$Z_i = \frac{Y_i - \bar{Y}}{s}$$

with $\bar{Y}$ and $s$ denoting the sample mean and sample standard deviation, respectively. In other words, data is given in units of how many standard deviations it is from the mean.

Although it is common practice to use Z-scores to identify possible outliers, this can be misleading (partiucarly for small sample sizes) due to the fact that the maximum Z-score is at most $(n-1)/\sqrt{n}$

$$s = \sqrt{\frac{1}{N-1} \sum_{i=1}^{N} (Y_i - \bar{Y})^2}$$



probability

Z-score [s]

[http://www.itl.nist.gov/div898/handbook/eda/section3/eda35h.htm]

# Iskanje osamelcev: Z-score

Iglewicz and Hoaglin recommend using the modified Z-score

$$M_i = \frac{0.6745(x - \tilde{x})}{MAD}$$

with MAD denoting the median absolute deviation and $\tilde{x}$ denoting the median.

These authors recommend that modified Z-scores with an absolute value of greater than 3.5 be labeled as potential outliers.

5. median absolute deviation - the median absolute deviation (MAD) is defined as

$$MAD = median(|Y_i - \tilde{Y}|)$$

where $\tilde{Y}$ is the median of the data and |Y| is the absolute value of Y. This is a variation of the average absolute deviation that is even less affected by extremes in the tail because the data in the tails have less influence on the calculation of the median than they do on the mean.

http://www.itl.nist.gov/div898/handbook/eda/section3/eda35h.htm

# Iskanje osamelcev: Z-score
# Koliko osamelcev?

*Formal Outlier Tests*

A number of formal outlier tests have proposed in the literature. These can be grouped by the following characteristics:

- What is the distributional model for the data? We restrict our discussion to tests that assume the data follow an approximately normal distribution.

- Is the test designed for a single outlier or is it designed for multiple outliers?

- If the test is designed for multiple outliers, does the number of outliers need to be specified exactly or can we specify an upper bound for the number of outliers?

http://www.itl.nist.gov/div898/handbook/eda/section3/eda35h.htm

# Iskanje osamelcev: Z-score
# Koliko osamelcev?

The text in the scanned images is too faded and degraded to read reliably.

http://www.itl.nist.gov/div898/handbook/eda/section3/eda35h.htm

# Iskanje osamelcev – več skupin

[Mira et la. (2012) RODHA: Robust Outlier Detection using Hybrid Approach]

# Iskanje osamelcev: Z-score
# Samo en osamelec?

http://www.itl.nist.gov/div898/handbook/eda/section3/eda35h1.htm

# Iskanje osamelcev: Z-score
# Točno k osamelcev?

http://www.itl.nist.gov/div898/handbook/eda/section3/eda35h1.htm

# Iskanje osamelcev: model-based

http://scikit-learn.org/stable/modules/outlier_detection.html

# Iskanje osamelcev: model-based

http://scikit-learn.org/stable/modules/outlier_detection.html

# Feedback – predavanje 3
http://goo.gl/forms/NC0tQUEA5b