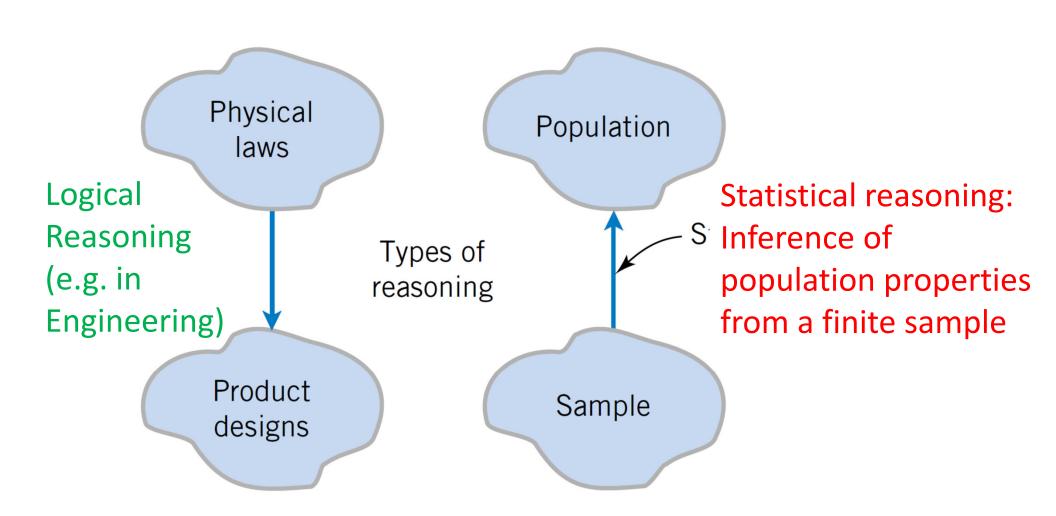
Descriptive statistics:
Populations, Samples
Histograms, Quartiles
Sample mean and
variance

# Two types of reasoning



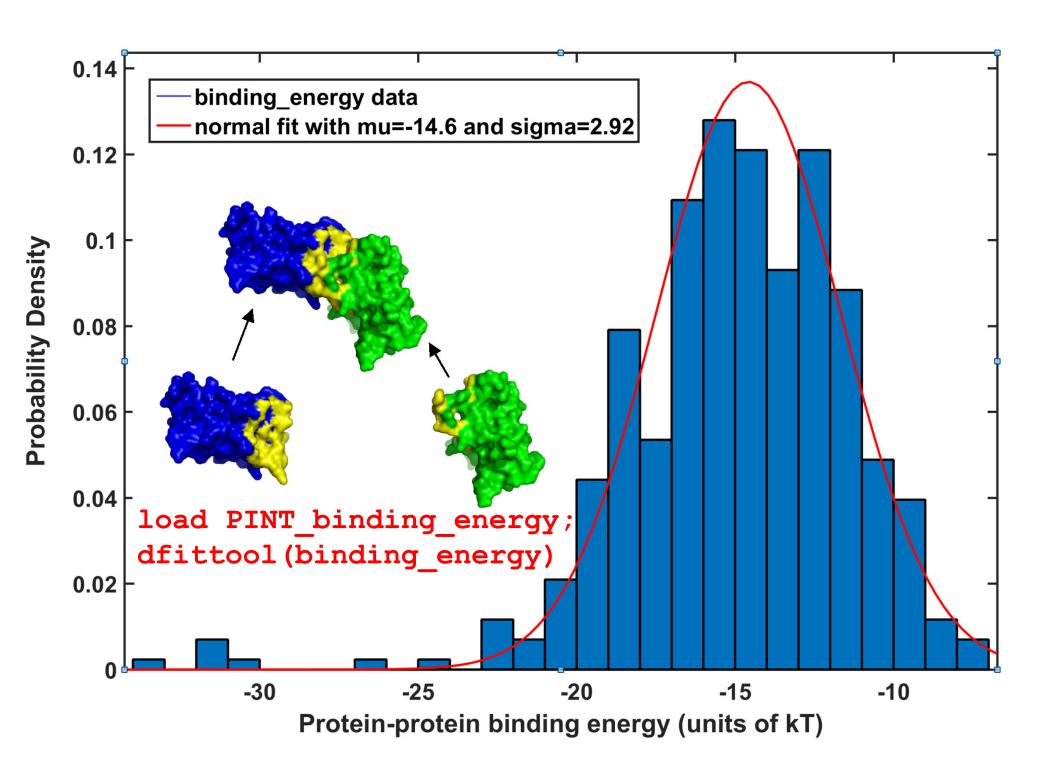
#### **Numerical Summaries of Data**

- Data are the numerical observations of a phenomenon of interest.
- The totality of all observations is a population.
  - Population can be infinite
     (e.g. abstract random variables)
  - It can be very large (e.g. 7 billion humans or all patients who have cancer of a given type)
- A (usually small) portion of the population collected for analysis is a random sample.
- We want to use sample to infer facts about populations
- The inference is not perfect but gets better and better as sample size increases.

#### Some Definitions

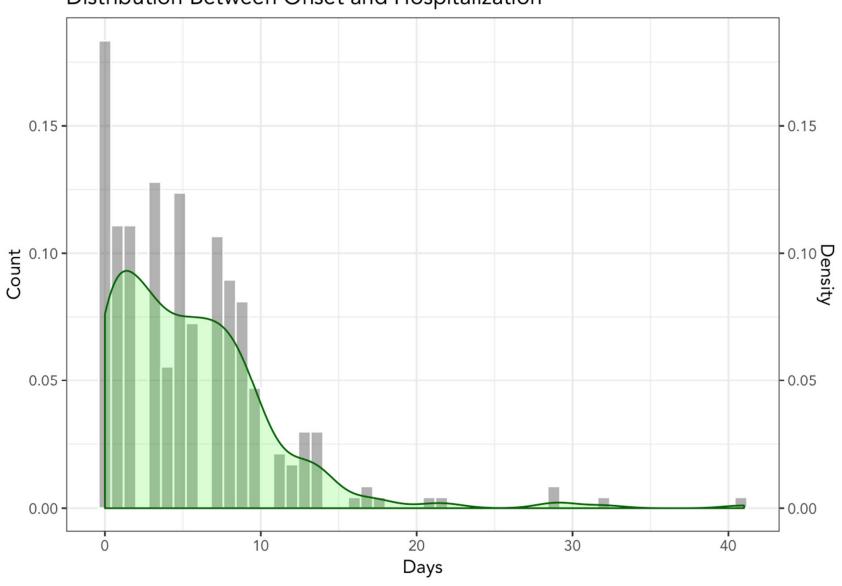
- The random variables  $X_1, X_2,...,X_n$  are a random sample of size n if:
  - a) The  $X_i$  are independent random variables.
  - b) Every  $X_i$  has the same probability distribution.
- Such  $X_1, X_2,...,X_n$  are also called independent and identically distributed (or i. i. d.) random variables

# Ways to describe a sample: Histogram approximates PDF (or PMF)



# PDF of time between COVID-19 symptoms onset and hospitalization in IL, April 2020



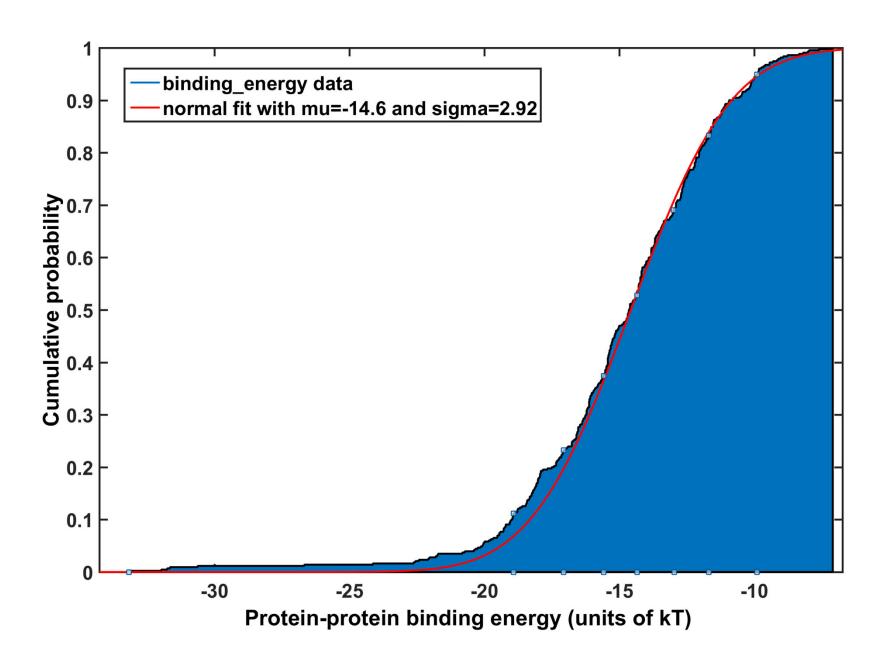


# Histograms with Unequal Bin Widths

- If the data is tightly clustered in some regions and scattered in others, it is visually helpful to use narrow bin widths in the clustered region and wide bin widths in the scattered areas.
- To <u>approximate the PDF</u>, the rectangle area, not the height, must be proportional to the bin relative frequency.

Rectangle height = 
$$\frac{\text{bin relative frequency}}{\text{bin width}}$$

# **Cumulative Frequency Plot**



# Median, Quartiles, Percentiles

- The median  $q_2$  divides the sample into two equal parts: 50% (n/2) of sample points below  $q_2$  and 50% (n/2) points above  $q_2$
- The three quartiles partition the data into four equally sized counts or segments.
  - -25% of the data is less than  $q_1$ .
  - -50% of the data is less than  $q_2$ , the median.
  - -75% of the data is less than  $q_3$ .
- There are 100 percentiles. n-th percentile  $p_n$  is defined so that n% of the data is less than  $p_n$

#### Box-and-Whisker Plot

- A box plot is a graphical display showing Spread,
   Outliers, Center, and Shape (SOCS).
- It displays the 5-number summary: min,  $q_1$ , median,  $q_3$ , and max.

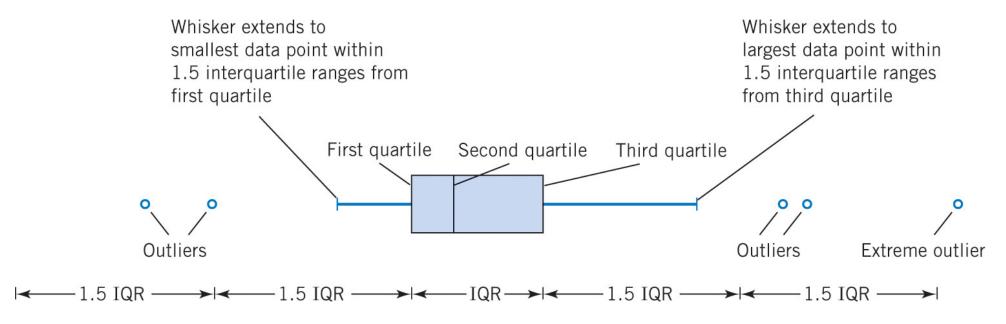
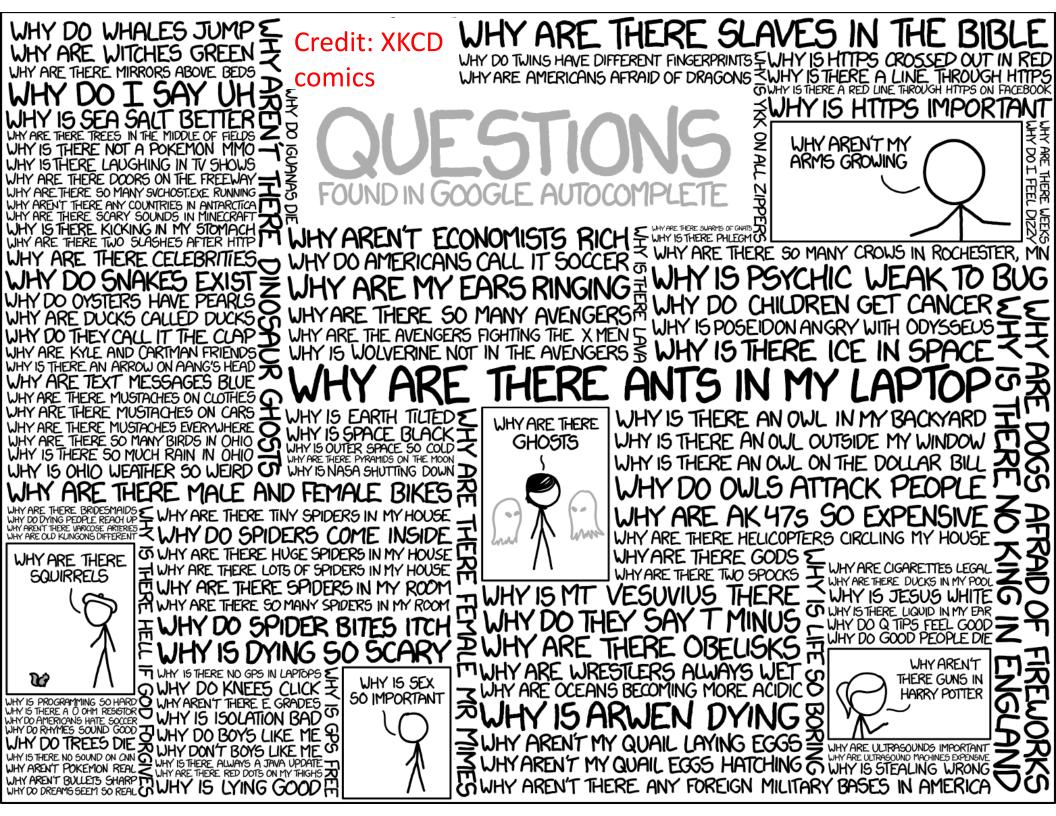


Figure 6-13 Description of a box plot.

Sec 6-4 Box Plots



#### Matlab exercise #1

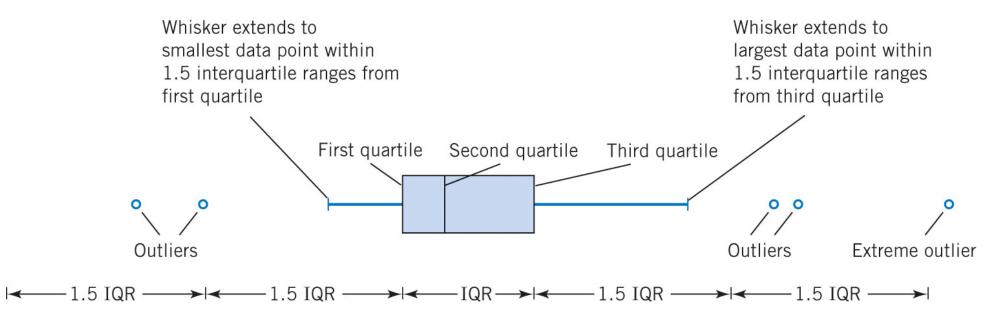
- Find the median and lower & upper quartiles of a n=1000 sample drawn from a standard normal distribution
- <u>Do not use</u> built-in Matlab functions for this exercise!
- Hint: use [a,b]=sort(r1); to rank order your sample. The variable a returns r1 sorted in the increasing order.
- How to find the median and both quartiles from a?

# How to find the median & quartiles

- % Example: find median and lower quartile of
- % a sample with n=100 drawn from uniform
- r1=randn(1000,1);
- [a,b]=sort(r1);
- med=(a(500)+a(501))./2
- sum(r1<med) % verify</li>
- q1=(a(250)+a(251))./2
- sum(r1<q1) % verify</li>
- q3=(a(750)+a(751))./2
- sum(r1<q3) % verify</li>

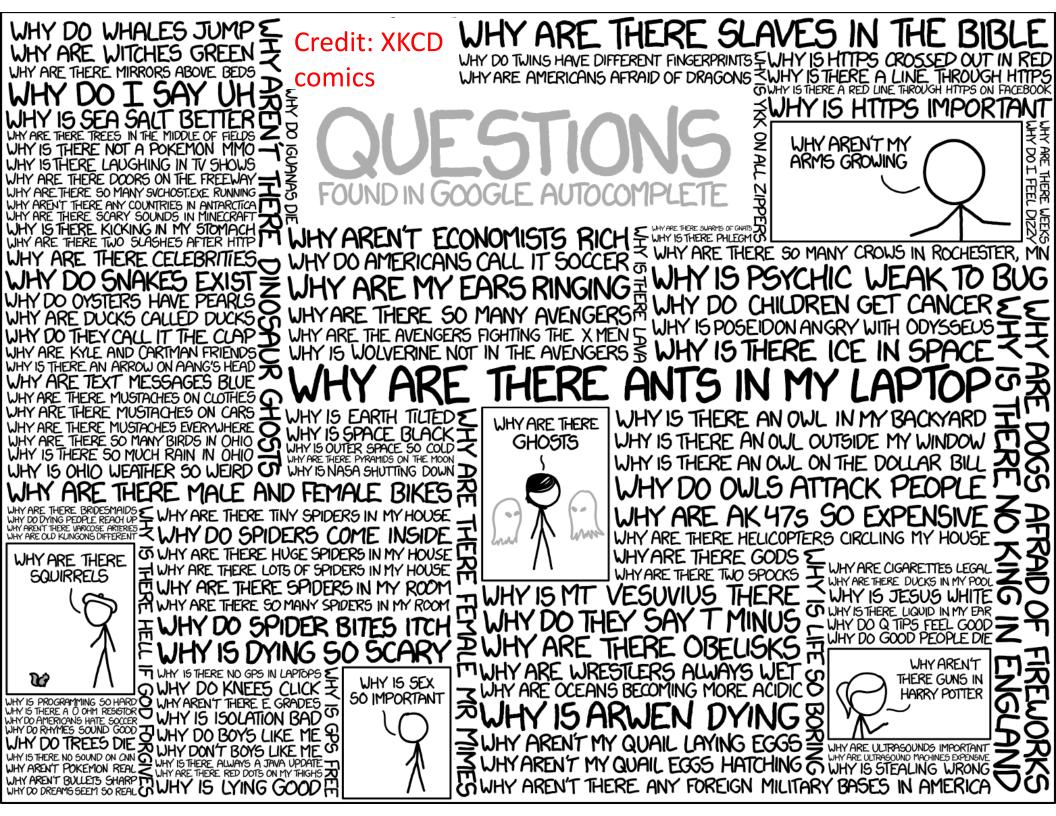
#### Matlab exercise #2:

- Generate a sample with n= 1000 following standard normal distribution
- Calculate median, first, and third quartiles
- Calculate IQR and find ranges shown below
- Find and count left and right outliers
- Do not use built-in Matlab functions for this!
- Make box and whisker plot: use boxplot



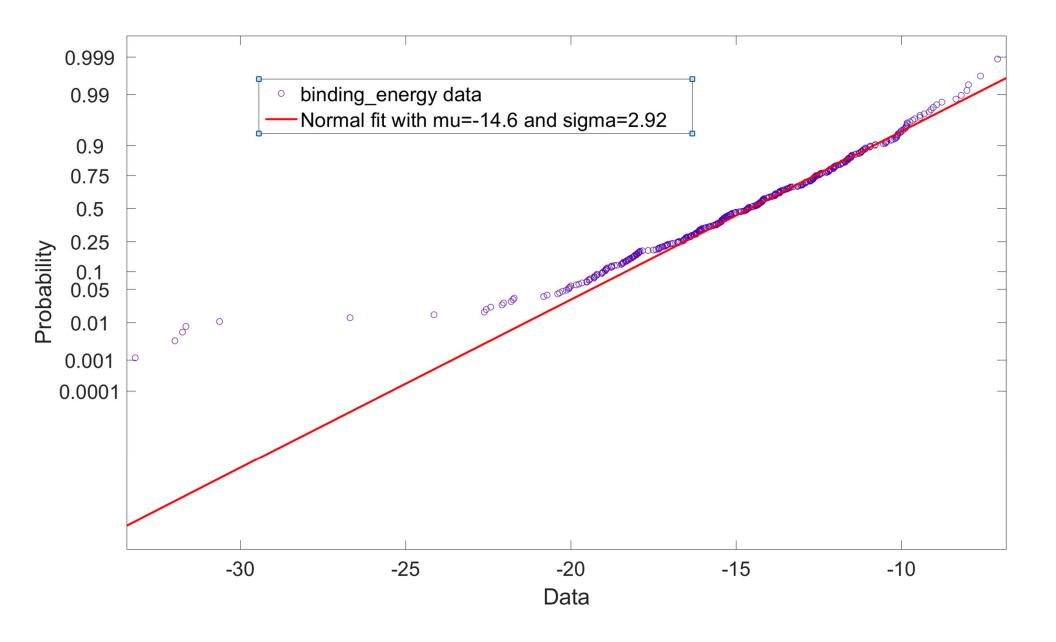
# How many right outliers one expects in a sample of n=1000 following normal distribution?

- % find the third quartile of a standard distribution
- norminv(0.75) %ans = 0.6745
- % Calculate IQR Interquartile Range
- IQR=2.\*norminv(0.75) % 1.3490
- % Calculate 0.5\*IQR+1.5\*IQR the right whisker position
- whisker=0.5.\*IQR+1.5\*IQR %ans = 2.6980
- % Find the probability to be above the right whisker
- 1-normcdf(whisker) %ans = 0.00349
- % Find number of right outliers in a sample of 1000 points
- 1000.\*(1-normcdf(whisker)) %ans = 3.49



# **Probability Plots**

- How do we know if a particular probability distribution is a reasonable model for a data set?
- A histogram of a large data set reveals the shape of a distribution. The histogram of a small data set does not provide a clear picture.
- A probability plot is helpful for all data set size.
   How good is the model based on a particular probability distribution can be verified using a subjective visual examination.



# How To Build a Probability Plot

- Sort the data observations in ascending order:  $X_{(1)}, X_{(2)}, ..., X_{(n)}$ .
- Empirically determined cumulative frequency  $Prob(x \le x_{(j)}) = j/n$ . To correct for discreteness of  $x_{(j)}$  better use  $Prob(x \le x_{(j)}) = (j-0.5)/n$
- If you believe that CDF(x) describes your random variable (j-0.5)/n should be close to  $CDF(x_{(j)})$
- Probability plot is  $x_{(j)} \cdot [(j-0.5)/n]/CDF(x_{(j)})$  plotted versus the observed value  $x_{(j)}$ .
- If the fit is good one gets a straight line
- Deviations can be seen especially at tails.

# **Probability Plot Variations**

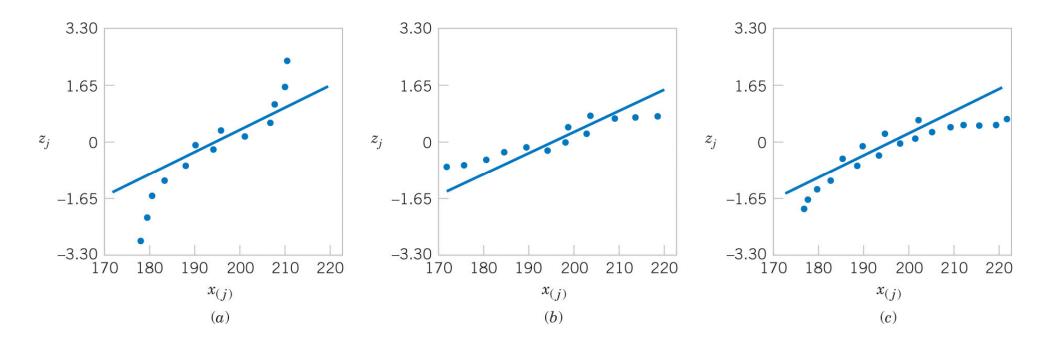


Figure 6-21 Normal probability plots indicating a non-normal distribution.

- (a) Light tailed distribution (squeezed together)
- (b) Heavy tailed distribution (stretched out)
- (c) Right skewed distribution (left end squeezed, right end stretched)

