# Correlation Analysis (for Categorical Data)

❑ **$X^2$ (chi-square) test:**

$$\overset{\text{observed}}{\underset{\text{expected}}{\chi^2 = \sum_{i}^{n} \frac{(O_i - E_i)^2}{E_i}}}$$

❑ Null hypothesis: The two distributions are independent

❑ The cells that contribute the most to the $X^2$ value are those whose actual count is very different from the expected count

   ❑ The larger the $X^2$ value, the more likely the variables are related

❑ Note:  Correlation does not imply causality

   ❑ # of hospitals and # of car-theft in a city are correlated

   ❑ Both are causally linked to the third variable: population

# Chi-Square Calculation: An Example

|                         | Play chess | Not play chess | Sum (row) |
|-------------------------|------------|----------------|-----------|
| Like science fiction    | 250 (X1)   | 200 (X2)       | 450       |
| Not like science fiction| 50 (X3)    | 1000 (X4)      | 1050      |
| Sum(col.)               | 300        | 1200           | 1500      |

❑ Null hypothesis: The two distributions are independent

❑ What does that mean?

❑ The ratio between people who play chess vs not play chess is the same for both groups of like science fiction and not like science fiction

❑ X1:X2=X3:X4=300:1200

❑ X1:X3=X2:X4=450:1050

❑ X1+X2=450       X3+X4=1050

❑ X1+X3=300       X2+X4=1200

2

# Chi-Square Calculation: An Example

|  | Play chess | Not play chess | Sum (row) |
|---|---|---|---|
| Like science fiction | 250 (90) | 200 (360) | 450 |
| Not like science fiction | 50 (210) | 1000 (840) | 1050 |
| Sum(col.) | 300 | 1200 | 1500 |

How to derive 90?
450/1500 * 300 = 90

We can reject the null hypothesis of independence at a confidence level of 0.001

❑ $X^2$ (chi-square) calculation (numbers in parenthesis are expected counts calculated based on the data distribution in the two categories)

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

❑ It shows that like_science_fiction and play_chess are correlated in the group

# Chi-Square Calculation: An Example

|  | A | B | C | D | Sum (row) |
|---|---|---|---|---|---|
| 1 |  |  |  |  | 200 |
| 0 |  |  |  |  | 1000 |
| Sum(col.) | 300 | 300 | 300 | 300 | 1200 |

❑ Degree of freedom

  ❑ (#categories_in_variable_A -1)(#categories_in_variable_B -1)

  ❑ number of values that are free to vary

# Chi-Square Calculation: An Example

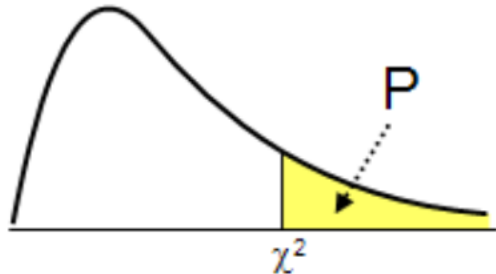|  | Play chess | Not play chess | Sum (row) |
|---|---|---|---|
| Like science fiction | 250 (90) | 200 (360) | 450 |
| Not like science fiction | 50 (210) | 1000 (840) | 1050 |
| Sum(col.) | 300 | 1200 | 1500 |

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

We can reject the null hypothesis of independence at a confidence level of 0.001

❑ Degree of freedom =?

**Values of the Chi-squared distribution**



| DF | P | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 0.995 | 0.975 | 0.20 | 0.10 | 0.05 | 0.025 | 0.02 | 0.01 | 0.005 | 0.002 | 0.001 |
| 1 | 0.0000393 | 0.000982 | 1.642 | 2.706 | 3.841 | 5.024 | 5.412 | 6.635 | 7.879 | 9.550 | 10.828 |
| 2 | 0.0100 | 0.0506 | 3.219 | 4.605 | 5.991 | 7.378 | 7.824 | 9.210 | 10.597 | 12.429 | 13.816 |
| 3 | 0.0717 | 0.216 | 4.642 | 6.251 | 7.815 | 9.348 | 9.837 | 11.345 | 12.838 | 14.796 | 16.266 |
| 4 | 0.207 | 0.484 | 5.989 | 7.779 | 9.488 | 11.143 | 11.668 | 13.277 | 14.860 | 16.924 | 18.467 |
| 5 | 0.412 | 0.831 | 7.289 | 9.236 | 11.070 | 12.833 | 13.388 | 15.086 | 16.750 | 18.907 | 20.515 |
| 6 | 0.676 | 1.237 | 8.558 | 10.645 | 12.592 | 14.449 | 15.033 | 16.812 | 18.548 | 20.791 | 22.458 |

# Variance for Single Variable (Numerical Data)

❑ The variance of a random variable $X$ provides a measure of how much the value of $X$ deviates from the mean or expected value of $X$:

$$\sigma^2 = \text{var}(X) = E[(X-\mu)^2] = \begin{cases} \sum_x (x-\mu)^2 f(x) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} (x-\mu)^2 f(x)dx & \text{if } X \text{ is continuous} \end{cases}$$

   ❑ where $\sigma^2$ is the variance of X, $\sigma$ is called *standard deviation*

     $\mu$ is the mean, and $\mu$ = E[X] is the expected value of X

   ❑ That is, variance is the expected value of the square deviation from the mean

   ❑ It can also be written as: $\sigma^2 = \text{var}(X) = E[(X-\mu)^2] = E[X^2] - \mu^2 = E[X^2] - [E(x)]^2$

❑ Sample variance

$$s^2 = \frac{1}{N}\sum_i^n (x_i - \hat{\mu})^2 \qquad\qquad s^2 = \frac{1}{n-1}\sum_i^n (x_i - \hat{\mu})^2$$

# Covariance for Two Variables

❏ Covariance between two variables $X_1$ and $X_2$

$$\sigma_{12} = E[(X_1 - \mu_1)(X_2 - \mu_2)] = E[X_1 X_2] - \mu_1 \mu_2 = E[X_1 X_2] - E[X_1]E[X_2]$$

where $\mu_1 = E[X_1]$ is the respective mean or **expected value** of $X_1$; similarly for $\mu_2$

❏ Sample covariance between $X_1$ and $X_2$:

$$\hat{\sigma}_{12} = \frac{1}{n}\sum_{i=1}^{n}(x_{i1} - \widehat{\mu_1})(x_{i2} - \widehat{\mu_2})$$

❏ **Positive covariance:** If $\sigma_{12} > 0$

❏ **Negative covariance:** If $\sigma_{12} < 0$

# Covariance for Two Variables

❑ **Independence**: If $X_1$ and $X_2$ are independent, $\sigma_{12} = 0$ but the reverse is not true

    ❑ Some pairs of random variables may have a covariance 0 but are not independent

    ❑ Only under some additional assumptions (e.g., the data follow multivariate normal distributions) does a covariance of 0 imply independence

# Example:  Calculation of Covariance

❑   Suppose two stocks $X_1$ and $X_2$ have the following values in one week:

   ❑   (2, 5), (3, 8), (5, 10), (4, 11), (6, 14)

❑   Question:  If the stocks are affected by the same industry trends, will their prices rise or fall together?

❑   Covariance formula

$$\sigma_{12} = E[(X_1 - \mu_1)(X_2 - \mu_2)] = E[X_1 X_2] - \mu_1 \mu_2 = E[X_1 X_2] - E[X_1]E[X_2]$$

❑   Its computation can be simplified as:   $\sigma_{12} = E[X_1 X_2] - E[X_1]E[X_2]$

   ❑   $E(X_1)$ = (2 + 3 + 5 + 4 + 6)/ 5 = 20/5 = 4

   ❑   $E(X_2)$ = (5 + 8 + 10 + 11 + 14) /5 = 48/5 = 9.6

   ❑   $\sigma_{12}$ = (2×5 + 3×8 + 5×10 + 4×11 + 6×14)/5 − 4 × 9.6 = 4

❑   Thus, $X_1$ and $X_2$ rise together since $\sigma_{12} > 0$

# Correlation between Two Numerical Variables

❑ **Correlation** between two variables $X_1$ and $X_2$ is the standard covariance, obtained by normalizing the covariance with the standard deviation of each variable
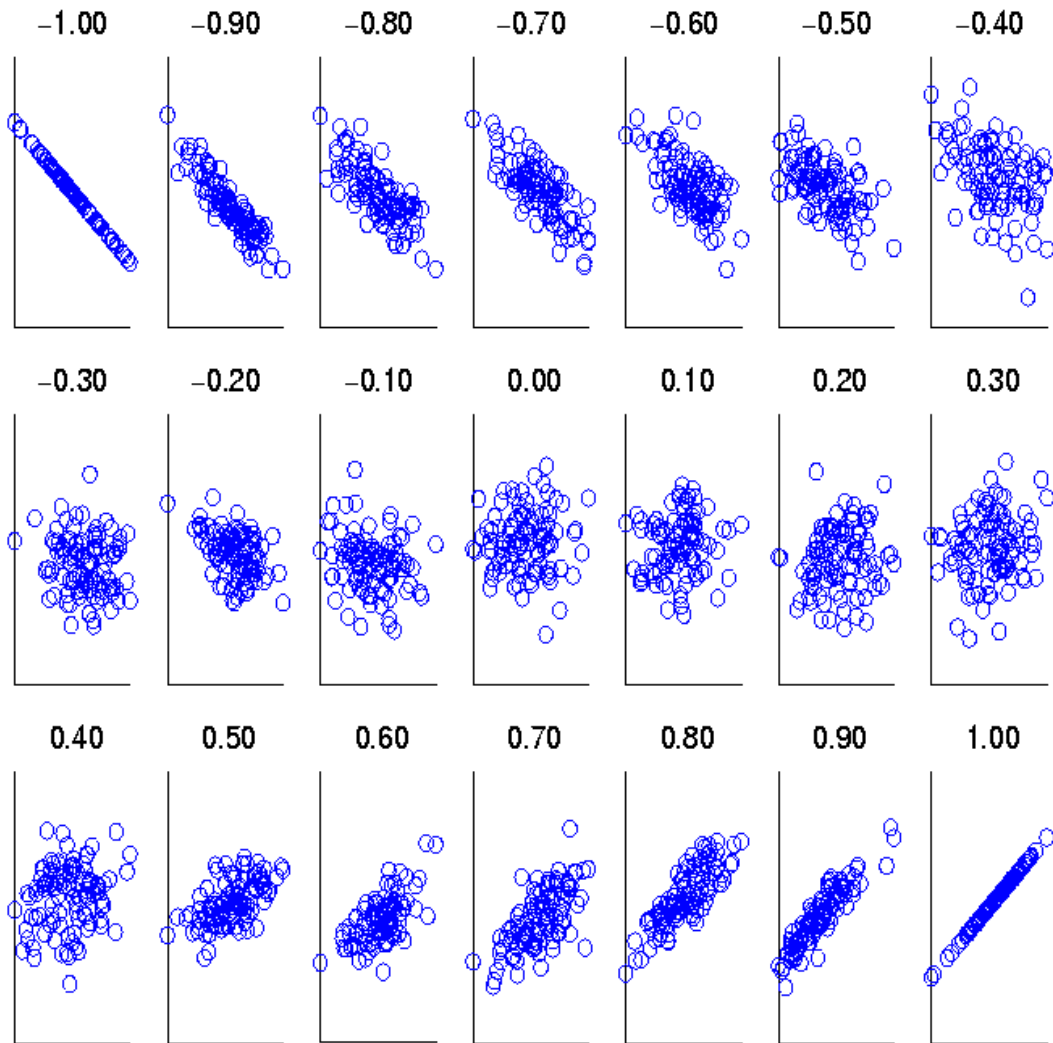
$$\rho_{12} = \frac{\sigma_{12}}{\sigma_1 \sigma_2} = \frac{\sigma_{12}}{\sqrt{\sigma_1^2 \sigma_2^2}}$$

❑ **Sample correlation** for two attributes $X_1$ and $X_2$:

$$\hat{\rho}_{12} = \frac{\hat{\sigma}_{12}}{\hat{\sigma}_1 \hat{\sigma}_2} = \frac{\sum_{i=1}^{n}(x_{i1} - \hat{\mu}_1)(x_{i2} - \hat{\mu}_2)}{\sqrt{\sum_{i=1}^{n}(x_{i1} - \hat{\mu}_1)^2 \sum_{i=1}^{n}(x_{i2} - \hat{\mu}_2)^2}}$$

where n is the number of tuples, $\mu_1$ and $\mu_2$ are the respective means of $X_1$ and $X_2$, $\sigma_1$ and $\sigma_2$ are the respective standard deviation of $X_1$ and $X_2$

❑ If $\rho_{12} > 0$: A and B are positively correlated ($X_1$'s values increase as $X_2$'s)

    ❑   The higher, the stronger correlation

❑ If $\rho_{12} = 0$: independent (under the same assumption as discussed in co-variance)

❑ If $\rho_{12} < 0$: negatively correlated

# Visualizing Changes of Correlation Coefficient



- ❑ Correlation coefficient value range: [−1, 1]

- ❑ A set of scatter plots shows sets of points and their correlation coefficients changing from −1 to 1

# Covariance Matrix

❑ The variance and covariance information for the two variables $X_1$ and $X_2$ can be summarized as 2 X 2 covariance matrix as

$$\Sigma = E[(\mathbf{X} - \mu)(\mathbf{X} - \mu)^T] = E[(\begin{matrix} X_1 - \mu_1 \\ X_2 - \mu_2 \end{matrix})(X_1 - \mu_1 \quad X_2 - \mu_2)]$$

$$= \begin{pmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] \end{pmatrix}$$

$$= \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}$$

❑ Generalizing it to *d* dimensions, we have,

$$D = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{d1} & x_{d2} & \cdots & x_{dd} \end{pmatrix} \quad \mathbf{\Sigma} = E[(\mathbf{X} - \mu)(\mathbf{X} - \mu)^T] = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{pmatrix}$$
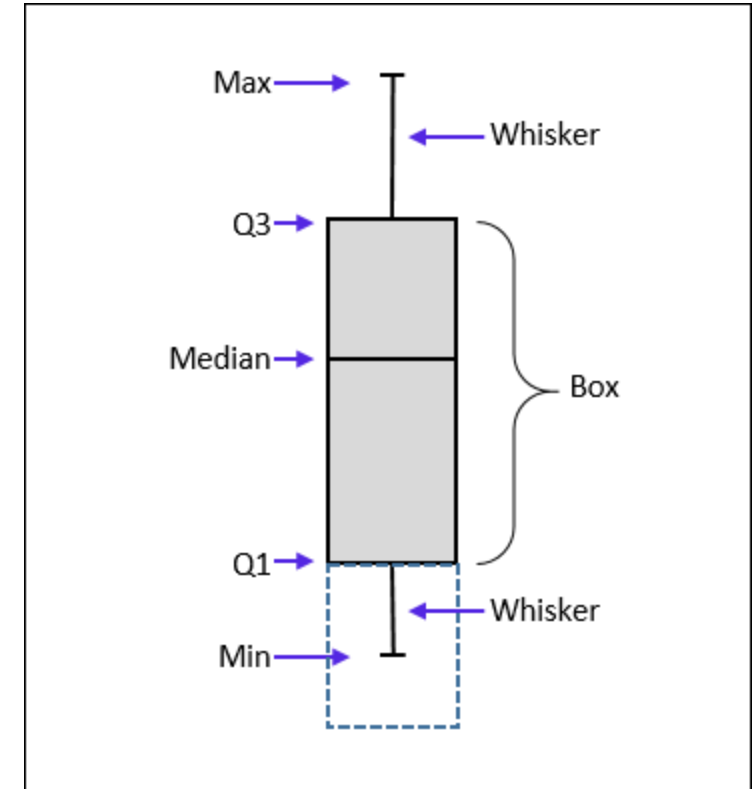
# Graphic Displays of Basic Statistical Descriptions

❑ **Boxplot**: graphic display of five-number summary

❑ **Histogram**: x-axis are values, y-axis repres. frequencies

❑ **Quantile plot**: each value $x_i$ is paired with $f_i$ indicating that approximately 100 $f_i$ % of data are $\leq x_i$

❑ **Quantile-quantile (q-q) plot**: graphs the quantiles of one univariant distribution against the corresponding quantiles of another

❑ **Scatter plot**: each pair of values is a pair of coordinates and plotted as points in the plane

# Measuring the Dispersion of Data: Quartiles & Boxplots

- **Quartiles**: $Q_1$ (25th percentile), $Q_3$ (75th percentile)

- **Inter-quartile range**: IQR = $Q_3 - Q_1$

- **Five number summary**: min, $Q_1$, median, $Q_3$, max

- **Boxplot**: Data is represented with a box

  - $Q_1$, $Q_3$, IQR:  The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR

  - Median ($Q_2$) is marked by a line within the box

  - Whiskers: two lines outside the box extended to Minimum and Maximum

  - Outliers: points beyond a specified outlier threshold, plotted individually

    - **Outlier**: usually, a value higher/lower than 1.5 x IQR

# Histogram Analysis

- ❑ Histogram: Graph display of tabulated frequencies, shown as bars

- ❑ Differences between histograms and bar charts

  - ❑ Histograms are used to show distributions of variables while bar charts are used to compare variables

  - ❑ Histograms plot binned quantitative data while bar charts plot categorical data

  - ❑ Bars can be reordered in bar charts but not in histograms

  - ❑ Differs from a bar chart in that it is the area of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width



Olympic Medals of all Times (till 2012 Olympics)

Bar chart

15

# Histograms Often Tell More than Boxplots



- ❑ The two histograms shown in the left may have the same boxplot representation
  - ❑ The same values for: min, Q1, median, Q3, max
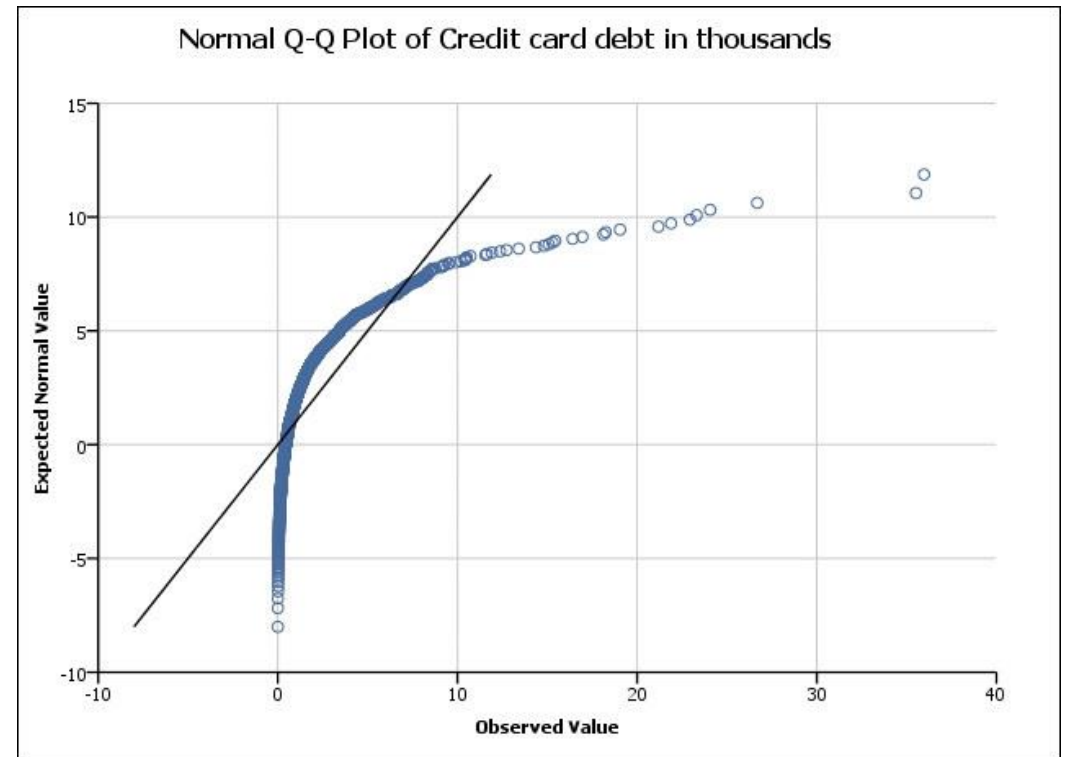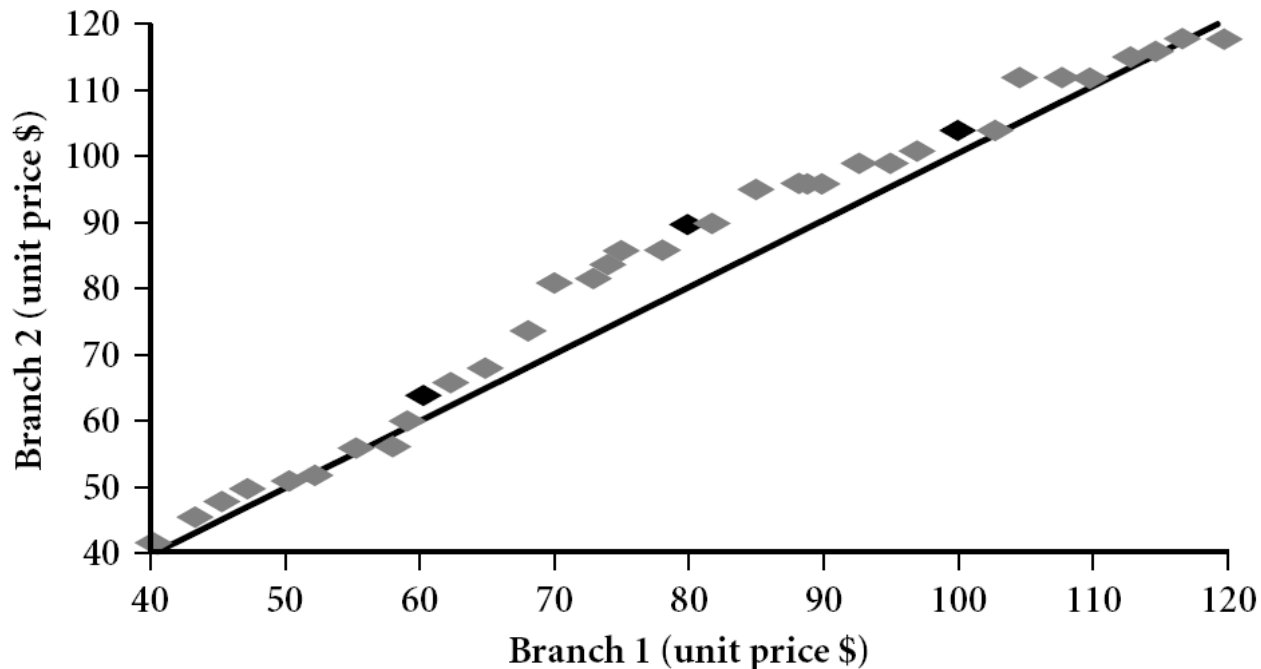- ❑ But they have rather different data distributions

# Quantile Plot

❑ Displays all of the data (allowing the user to assess both the overall behavior and unusual occurrences)

❑ Plots **quantile** information

   ❑ For a data $x_i$ data sorted in increasing order, $f_i$ indicates that approximately 100 $f_i$% of the data are below or equal to the value $x_i$
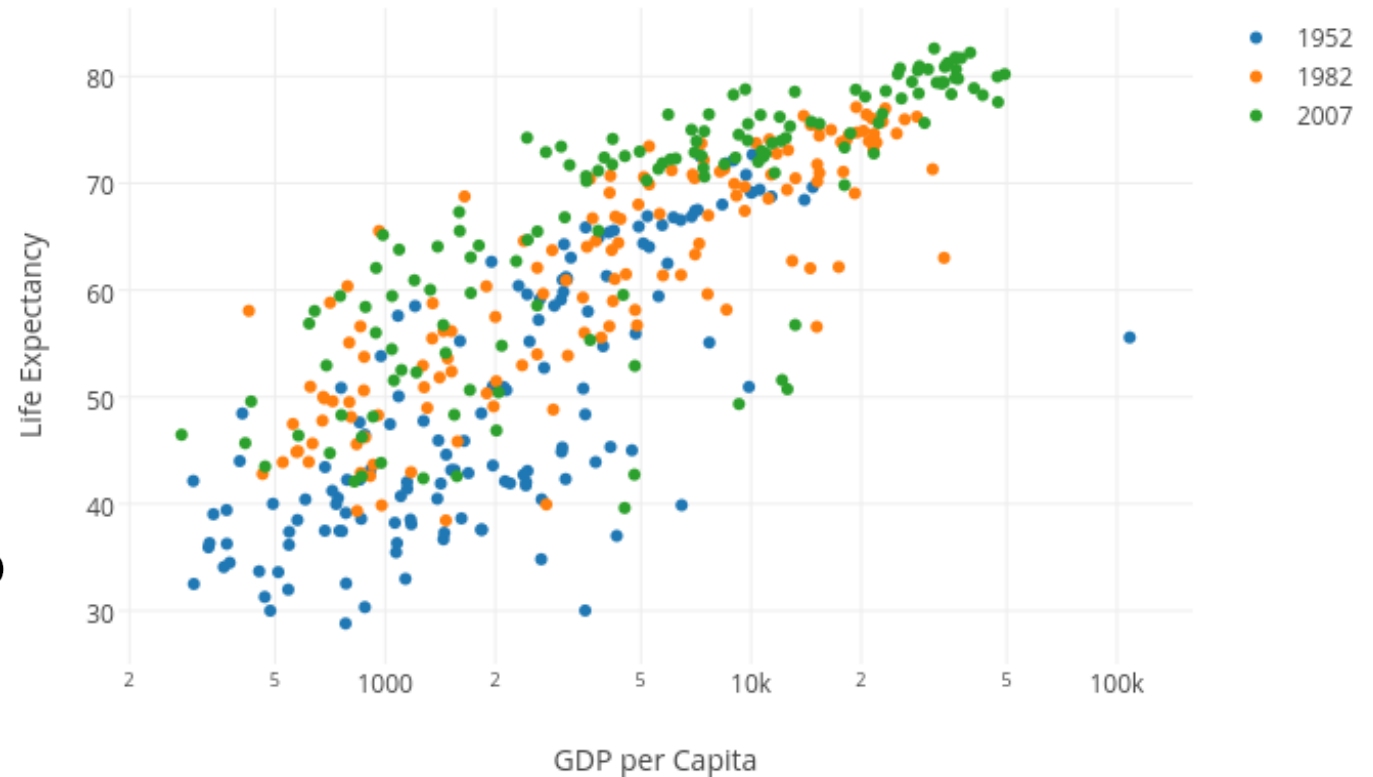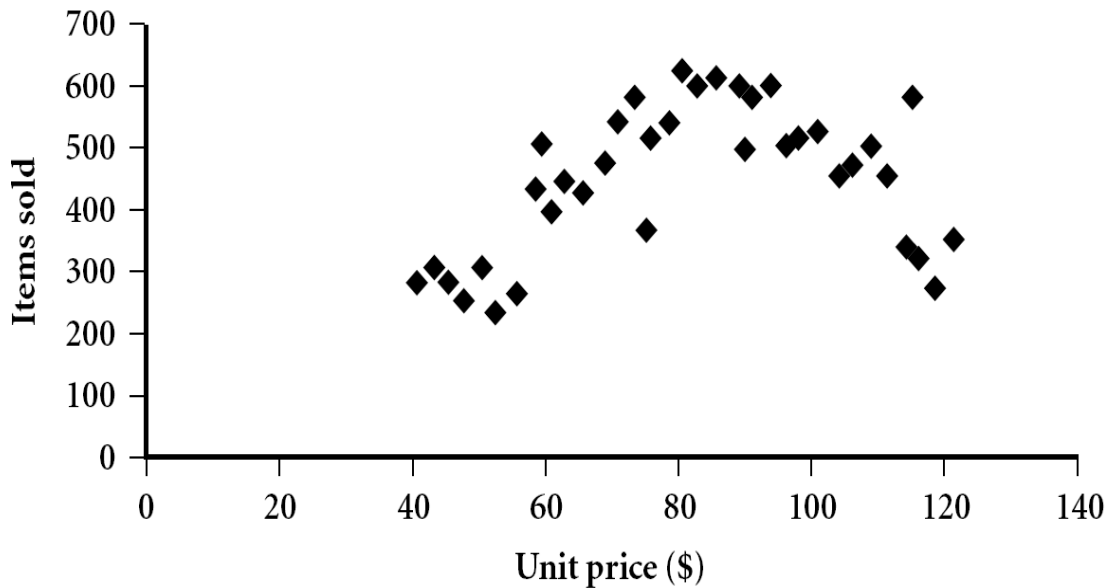
# Quantile-Quantile (Q-Q) Plot

❑ Graphs the quantiles of one univariate distribution against the corresponding quantiles of another

❑ View: Is there is a shift in going from one distribution to another?

❑ Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile. Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2
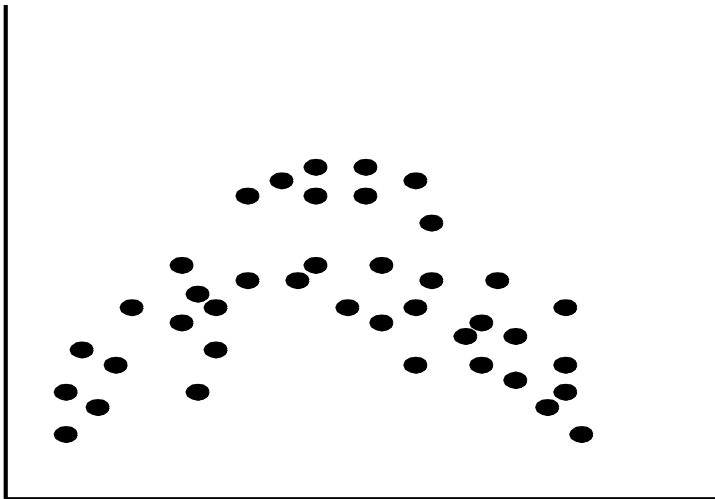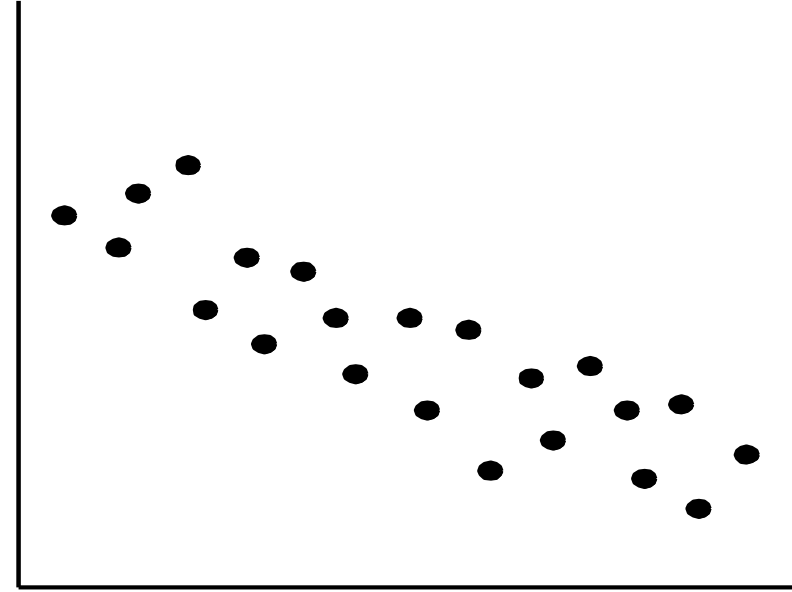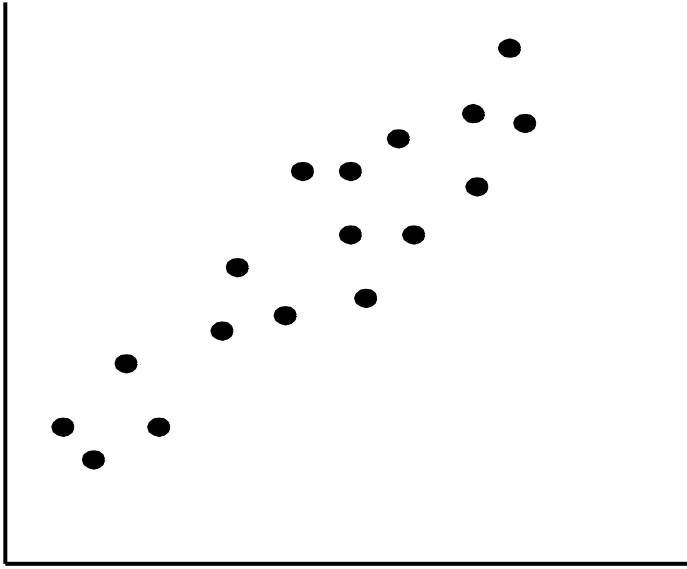
# Scatter plot

❑ Provides a first look at bivariate data to see clusters of points, outliers, etc.

❑ Each pair of values is treated as a pair of coordinates and plotted as points in the plane

# Positively and Negatively Correlated Data



- ❑ The left half fragment is positively correlated
- ❑ The right half is negative correlated

# Uncorrelated Data