# MultiVariate Normal Distribution

Jing Qin

15/03/2022

# How table 4.2 is built?

| Sample size $n$ | Significance levels $\alpha$ | | |
|---|---|---|---|
| | .01 | .05 | .10 |
| 5 | .8299 | .8788 | .9032 |
| 10 | .8801 | .9198 | .9351 |
| 15 | .9126 | .9389 | .9503 |
| 20 | .9269 | .9508 | .9604 |
| 25 | .9410 | .9591 | .9665 |
| 30 | .9479 | .9652 | .9715 |
| 35 | .9538 | .9682 | .9740 |
| 40 | .9599 | .9726 | .9771 |
| 45 | .9632 | .9749 | .9792 |
| 50 | .9671 | .9768 | .9809 |
| 55 | .9695 | .9787 | .9822 |
| 60 | .9720 | .9801 | .9836 |
| 75 | .9771 | .9838 | .9866 |
| 100 | .9822 | .9873 | .9895 |
| 150 | .9879 | .9913 | .9928 |
| 200 | .9905 | .9931 | .9942 |
| 300 | .9935 | .9953 | .9960 |

**Table 4.2** Critical Points for the Q–Q Plot Correlation Coefficient Test for Normality

1. Generate multiple datasets with N(0,1) for n=100

2. Make Q-Q plots for each of the dataset and derive $r_Q$ respectively

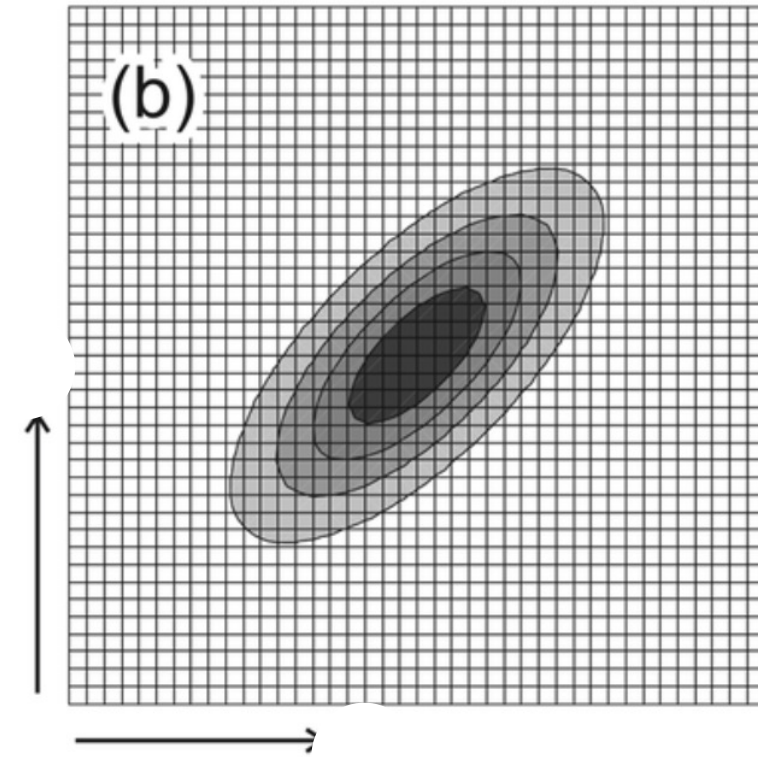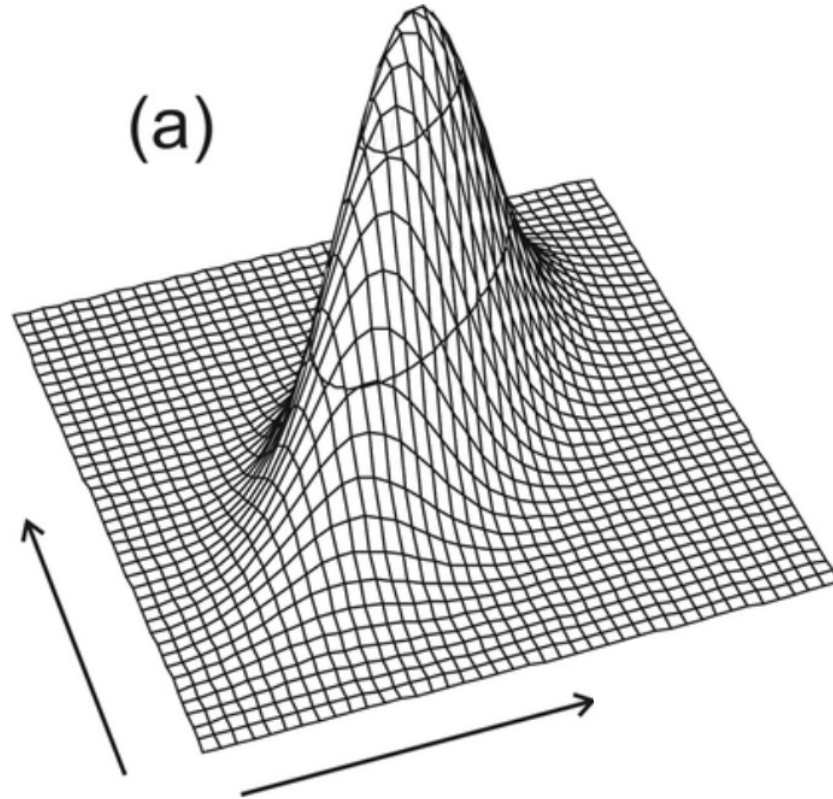3. Collect all the $r_Q$ and find the critical value for given significant level 0.05.

# Bivariate normal distribution $p = 2$

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$$

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$$

(a)

(b)

$$\frac{1}{2\pi\sqrt{\sigma_1^2\sigma_2^2(1-\rho^2)}} \times$$

$$\exp\left\{ -\frac{1}{2(1-\rho^2)} \left[ \frac{(x_1-\mu_1)^2}{\sigma_1^2} + \frac{(x_2-\mu_2)^2}{\sigma_2^2} - 2\rho\left(\frac{x_1-\mu_1}{\sigma_1}\right)\left(\frac{x_2-\mu_2}{\sigma_2}\right) \right] \right\}$$

# Towards general: vector and matrix form

$$\frac{1}{2\pi\sqrt{\sigma_1^2\sigma_2^2(1-\rho^2)}} \times$$

$$\exp\left\{-\frac{1}{2(1-\rho^2)}\left[\frac{(x_1-\mu_1)^2}{\sigma_1^2} + \frac{(x_2-\mu_2)^2}{\sigma_2^2} - 2\rho\left(\frac{x_1-\mu_1}{\sigma_1}\right)\left(\frac{x_2-\mu_2}{\sigma_2}\right)\right]\right\}$$
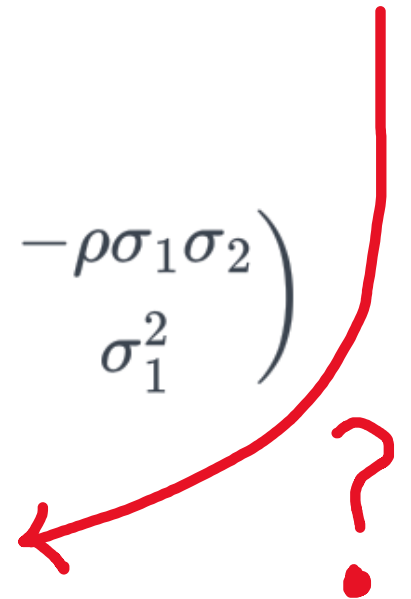
$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \qquad |\Sigma| = \sigma_1^2\sigma_2^2(1-\rho^2)$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \qquad \Sigma^{-1} = \frac{1}{\sigma_1^2\sigma_2^2(1-\rho^2)}\begin{pmatrix} \sigma_2^2 & -\rho\sigma_1\sigma_2 \\ -\rho\sigma_1\sigma_2 & \sigma_1^2 \end{pmatrix}$$

$$\frac{1}{(2\pi)^{2/2}|\Sigma|^{1/2}} \exp\left\{-(\boldsymbol{x}-\boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{x}-\boldsymbol{\mu})/2\right\}$$

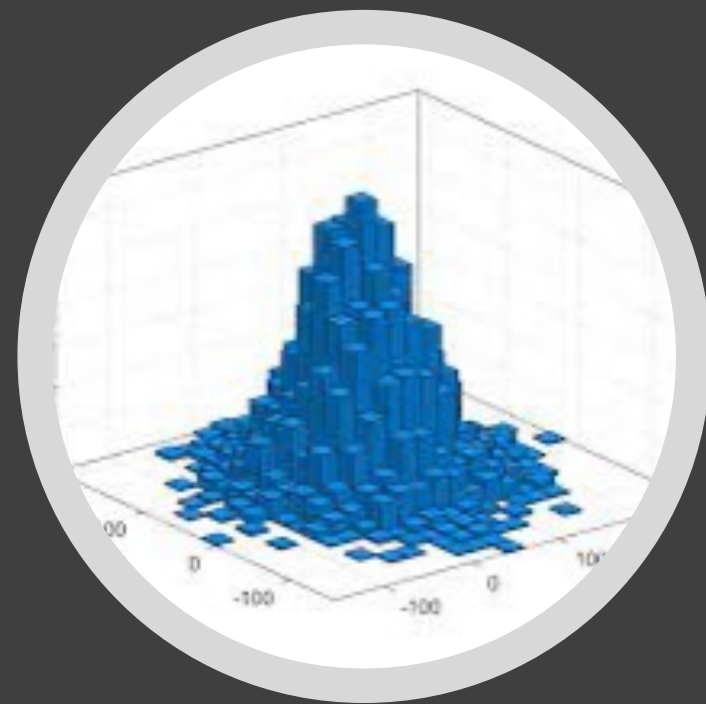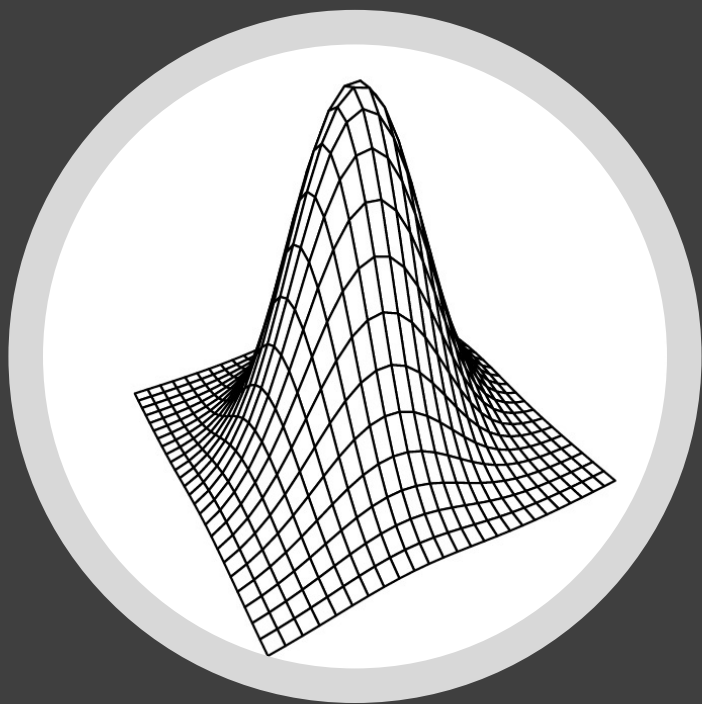For general $p \geq 2$, joint PDF (4-4)     $X \sim N_p(\boldsymbol{\mu}, \Sigma)$

$$\frac{1}{(2\pi)^{2/2}|\Sigma|^{1/2}} \exp\left\{-(\boldsymbol{x} - \boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{x} - \boldsymbol{\mu})/2\right\}$$

This is why we need the vectors!

$$\frac{1}{(2\pi)^{p/2}|\Sigma|^{1/2}} \exp\left\{-(\boldsymbol{x} - \boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{x} - \boldsymbol{\mu})/2\right\}$$

Consider $\boldsymbol{X} = (X_1, X_2, \ldots, X_p)'$ and $\boldsymbol{x} = (x_1, x_2, \ldots, x_p)'$

$E(\boldsymbol{X}) = \boldsymbol{\mu}$ and $Cov(\boldsymbol{X}) = \Sigma$.

In practice: *Is my data normally distributed*?

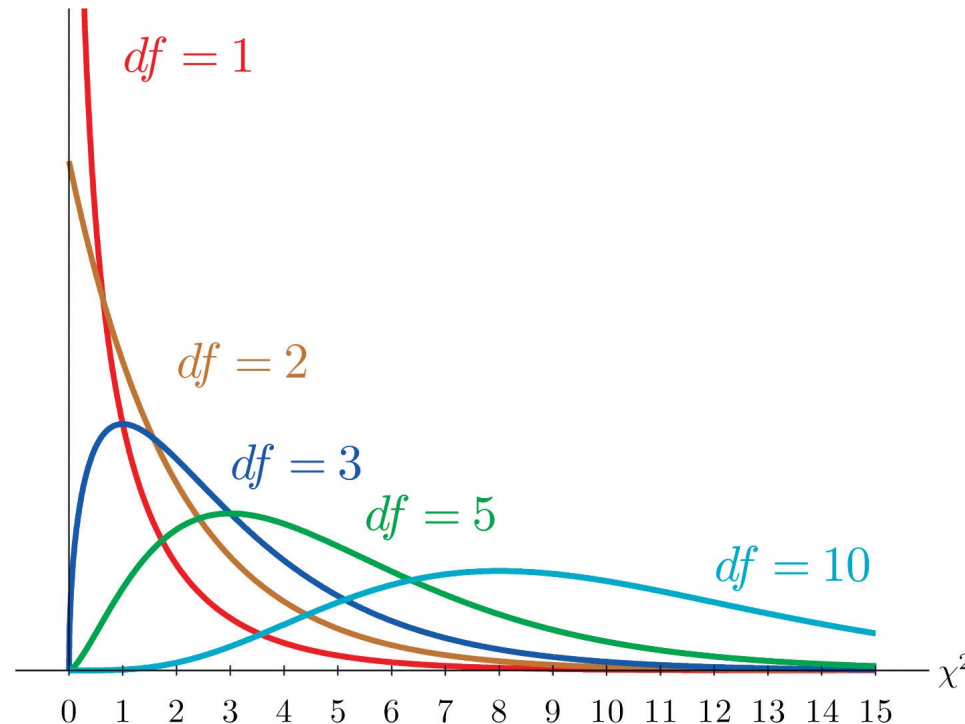Example: radiation data with door open+closed (t4-1.dat; t4-5.dat)

# Quadratic form (4-8)

$$\boldsymbol{X} \sim N_p(\boldsymbol{\mu}, \Sigma) \implies \frac{1}{(2\pi)^{p/2}|\Sigma|^{1/2}} \exp\left\{-\boxed{(\boldsymbol{x} - \boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{x} - \boldsymbol{\mu})}/2\right\}$$
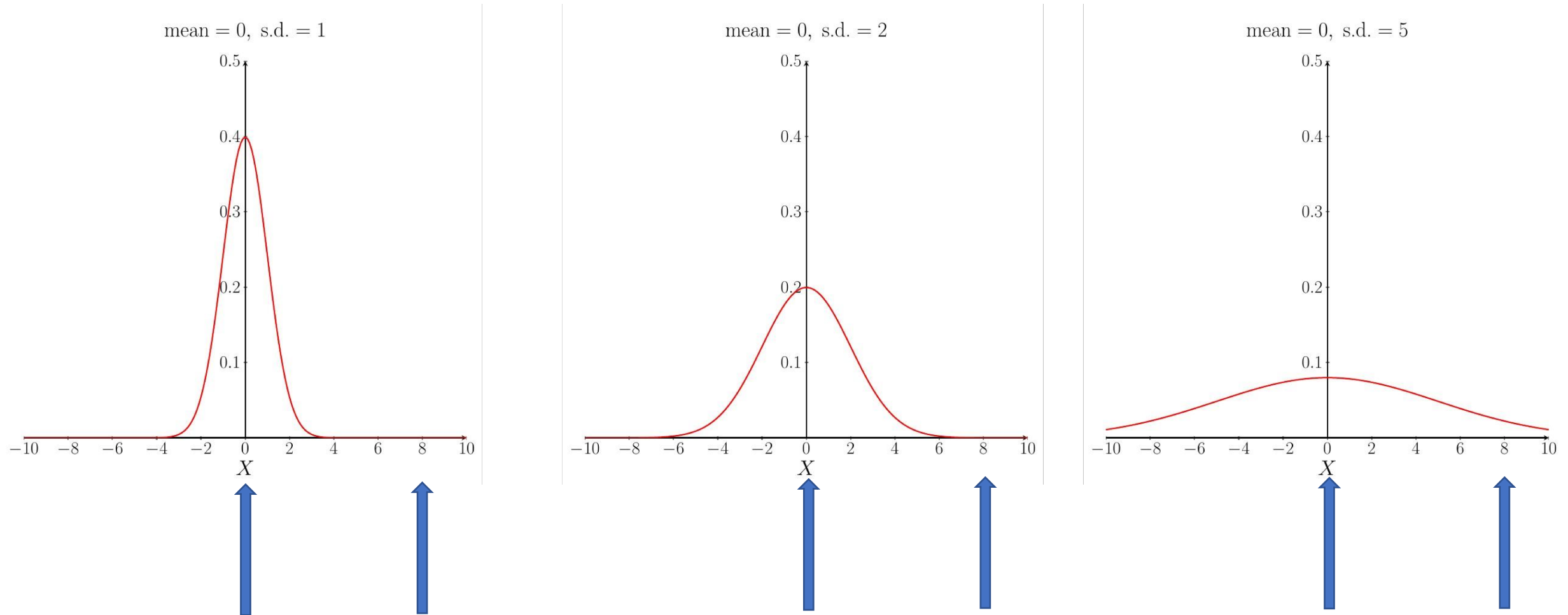
$p$-dimension                    1-dimension

- (4-8) Assume $\boldsymbol{X} \sim N_p(\boldsymbol{\mu}, \Sigma)$, we have $(\boldsymbol{X} - \boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{X} - \boldsymbol{\mu}) \sim \chi_p^2$
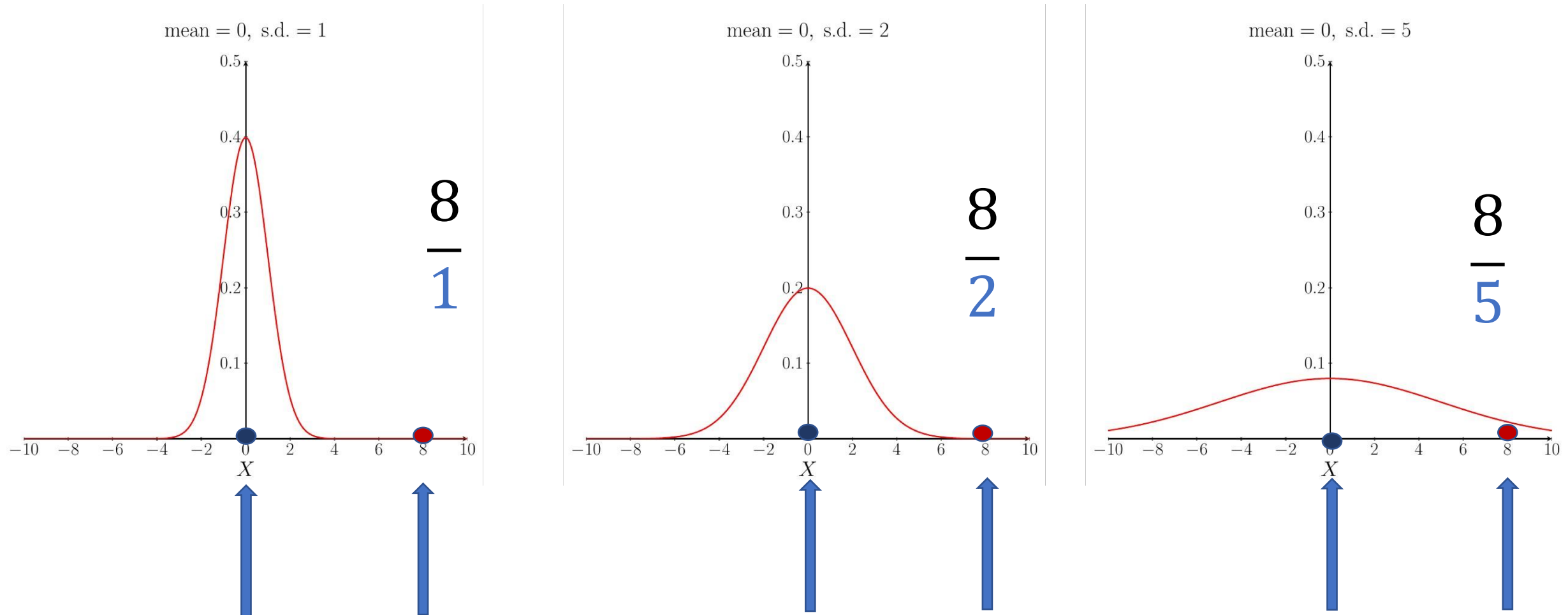
# Statistical/Mahalanobis distance

# Statistical/Mahalanobis distance

Bigger variability, smaller difference



mean = 0,  s.d. = 1

$$\frac{8}{1}$$

mean = 0,  s.d. = 2

$$\frac{8}{2}$$

mean = 0,  s.d. = 5

$$\frac{8}{5}$$

# Statistical/Mahalanobis distance

- (2-17, 4-3) The quadratic form $(x - \mu)'\Sigma^{-1}(x - \mu)$ is referred to as squared statistical/Mahalanobis distance. R cmd `mahalanobis()`

From $x$ to $\mu$

$$d_E(x, y) = \sqrt{(x - y)^T \cdot (x - y)}$$

$$d_M(x, y) = \sqrt{(x - y)^T \cdot S^{-1} \cdot (x - y)}$$

$$= \sqrt{[x_1 - y_1 \quad x_2 - y_2] \begin{bmatrix} \dfrac{1}{\sigma_1^2} & 0 \\ 0 & \dfrac{1}{\sigma_2^2} \end{bmatrix} \begin{bmatrix} x_1 - y_1 \\ x_2 - y_2 \end{bmatrix}}$$

$$= \sqrt{\left[\dfrac{x_1 - y_1}{\sigma_1^2} \quad \dfrac{x_2 - y_2}{\sigma_2^2}\right] \begin{bmatrix} x_1 - y_1 \\ x_2 - y_2 \end{bmatrix}}$$

$$= \sqrt{\dfrac{(x_1 - y_1)^2}{\sigma_1^2} + \dfrac{(x_2 - y_2)^2}{\sigma_2^2}}$$



Mahalanobis

Euclidean

A

B

X

$$d^2 = constant$$

- **Result (4.7)** The solid ellipsoid of $x$ values satisfying

$$(x - \mu)' \Sigma^{-1} (x - \mu) \leq c^2 = \chi_p^2(\alpha)$$

has probability $1 - \alpha$.

$qchisq(\alpha, p)$

| $1 - \alpha$ | Observed count | Expected count |
|:---:|:---:|:---:|
| 0.25 | 17 | 10.5 |
| 0.50 | 29 | 21 |
| 0.75 | 33 | 31.5 |

Expected number of observations versus data. Note

n = 42

# Q-Q plot, again

**Dimension?**

- Result (4.7) Assume $\boldsymbol{X} \sim N_p(\boldsymbol{\mu}, \Sigma)$, we have $\boxed{(\boldsymbol{X} - \boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{X} - \boldsymbol{\mu})} \sim \chi_p^2$

# Q-Q plot, again

**Dimension?**

- Result (4.7)

Assume $\boldsymbol{X} \sim N_p(\boldsymbol{\mu}, \Sigma)$, we have $\boxed{(\boldsymbol{X} - \boldsymbol{\mu})' \Sigma^{-1} (\boldsymbol{X} - \boldsymbol{\mu})} \sim \chi_p^2$
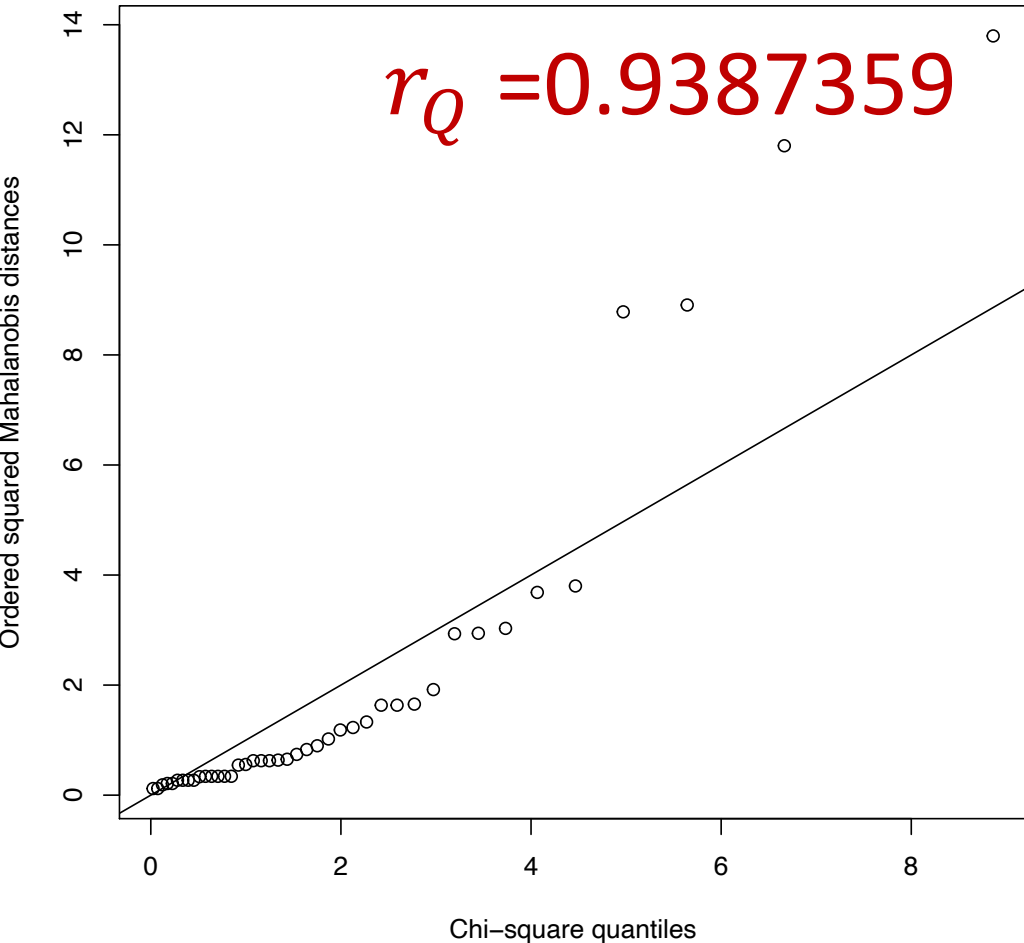
$r_Q$ =?

Can we use Table 4.2 again?



**Table 4.2** Critical Points for the Q–Q Plot Correlation Coefficient Test for Normality

| Sample size $n$ | Significance levels $\alpha$ | | |
|---|---|---|---|
| | .01 | .05 | .10 |
| 5 | .8299 | .8788 | .9032 |
| 10 | .8801 | .9198 | .9351 |
| 15 | .9126 | .9389 | .9503 |
| 20 | .9269 | .9508 | .9604 |
| 25 | .9410 | .9591 | .9665 |
| 30 | .9479 | .9652 | .9715 |
| 35 | .9538 | .9682 | .9740 |
| 40 | .9599 | .9726 | .9771 |
| 45 | .9632 | .9749 | .9792 |
| 50 | .9671 | .9768 | .9809 |
| 55 | .9695 | .9787 | .9822 |
| 60 | .9720 | .9801 | .9836 |
| 75 | .9771 | .9838 | .9866 |
| 100 | .9822 | .9873 | .9895 |
| 150 | .9879 | .9913 | .9928 |
| 200 | .9905 | .9931 | .9942 |
| 300 | .9935 | .9953 | .9960 |

# Q-Q plot, again

- Result (4.7)

Assume $X \sim N_p(\mu, \Sigma)$, we have $\boxed{(X - \mu)'\Sigma^{-1}(X - \mu)} \sim \chi_p^2$



$r_Q$ =0.9387359

1. Generate multiple datasets with $\chi_p^2$ for n=42

2. Make Q-Q plots for each of the dataset

   and derive $r_Q$ respectively

3. Collect all the $r_Q$ and find the critical value

   for some given significant level.

FindCrikChi2update.R

```
> source("FindCrikChi2.R")
> result1[[2]]
[1] 0.9948543
```
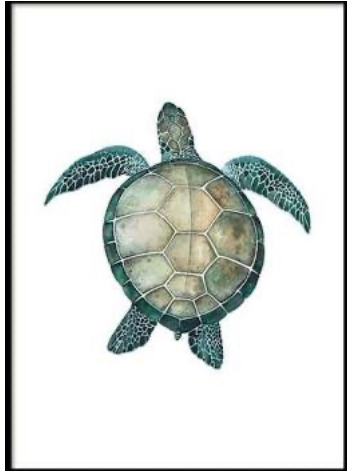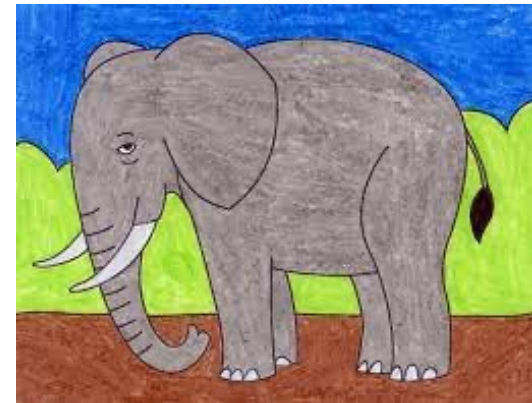
Conclusion ?

# Is *quadratic form* enough for assessing normality?

- Result (4.7) Assume $\boldsymbol{X} \sim N_p(\boldsymbol{\mu}, \Sigma)$, we have $(\boldsymbol{X} - \boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{X} - \boldsymbol{\mu}) \sim \chi_p^2$

**This animal is a turtle.**    **This is an animal has 4 legs.**

# Check MVN, continued

- (Result 4.4) Any subset of a MVN distributed random vector is normally distributed.

**Example 4.5  (The distribution of a subset of a normal random vector)**

If $\mathbf{X}$ is distributed as $N_5(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, find the distribution of $\begin{bmatrix} X_2 \\ X_4 \end{bmatrix}$. We set

$$\mathbf{X}_1 = \begin{bmatrix} X_2 \\ X_4 \end{bmatrix}, \quad \boldsymbol{\mu}_1 = \begin{bmatrix} \mu_2 \\ \mu_4 \end{bmatrix}, \quad \boldsymbol{\Sigma}_{11} = \begin{bmatrix} \sigma_{22} & \sigma_{24} \\ \sigma_{24} & \sigma_{44} \end{bmatrix}$$

# In summary, basic track of checking MVN

- Test univariate normality for each marginal distribution with QQ-plot.

- Test bivariate normality for each pair of attributes. For example, a matrix of scatterplots and QQ-plot based on $(\boldsymbol{X} - \boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{X} - \boldsymbol{\mu}) \sim \chi_2^2$

- Test over all MVN using QQ-plot based on $(\boldsymbol{X} - \boldsymbol{\mu})'\Sigma^{-1}(\boldsymbol{X} - \boldsymbol{\mu}) \sim \chi_p^2$

- Linear pattern in QQ-plot can be evaluated through hypothesis test.

**Advanced track: well, it is by taking care of all the possible subsets of the attributes...or PCA**

# Other useful properties of MVN

- (Result 4.3) Let $A$ be a $(q \times p)$ numeric matrix, then

$$AX \sim N_q(A\mu, A\Sigma A')$$

- Exercise 2 Find the mean vector and the total variance of $AX$.
  Given that $A = \begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \end{pmatrix}$ and $X = (X_1, X_2, X_3)'$.

Further we know $X \sim N_3(\mu, \Sigma)$, where $\mu = (1, 2, 1)'$ and

$$\Sigma = \begin{pmatrix} 3 & -1 & 1 \\ -1 & 1 & 0 \\ 1 & 0 & 2 \end{pmatrix}$$

# Other useful properties of MVN

**Result 4.5.**

(a) If $\mathbf{X}_1$ and $\mathbf{X}_2$ are independent, then $\text{Cov}(\mathbf{X}_1, \mathbf{X}_2) = \mathbf{0}$, a $q_1 \times q_2$ matrix of
$(q_1 \times 1)$ $(q_2 \times 1)$
zeros.

(b) If $\begin{bmatrix} \mathbf{X}_1 \\ \hline \mathbf{X}_2 \end{bmatrix}$ is $N_{q_1+q_2}\left( \begin{bmatrix} \boldsymbol{\mu}_1 \\ \hline \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \vdots & \Sigma_{12} \\ \hline \Sigma_{21} & \vdots & \Sigma_{22} \end{bmatrix} \right)$, then $\mathbf{X}_1$ and $\mathbf{X}_2$ are independent if and only if $\Sigma_{12} = \mathbf{0}$.

(c) If $\mathbf{X}_1$ and $\mathbf{X}_2$ are independent and are distributed as $N_{q_1}(\boldsymbol{\mu}_1, \Sigma_{11})$ and $N_{q_2}(\boldsymbol{\mu}_2, \Sigma_{22})$, respectively, then $\begin{bmatrix} \mathbf{X}_1 \\ \hline \mathbf{X}_2 \end{bmatrix}$ has the multivariate normal distribution

$$N_{q_1+q_2}\left( \begin{bmatrix} \boldsymbol{\mu}_1 \\ \hline \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \vdots & \mathbf{0} \\ \hline \mathbf{0}' & \vdots & \Sigma_{22} \end{bmatrix} \right)$$