# Lecture 33 One-way ANOVA Examples

BIO210 Biostatistics

Xi Chen

Fall, 2024

School of Life Sciences
Southern University of Science and Technology

南方科技大学生命科学学院
SUSTech · SCHOOL OF
**LIFE SCIENCES**

## The Iris Flower Dataset

- Introduced by Ronald Fisher in his 1936 paper: **The use of multiple measurements in taxonomic problems**.
- Extensively used in the machine learning community for testing classification methods. https://en.wikipedia.org/wiki/Iris_flower_data_set
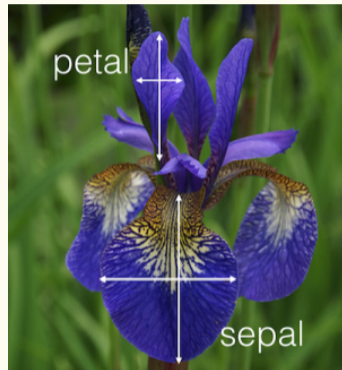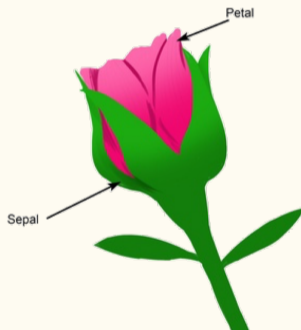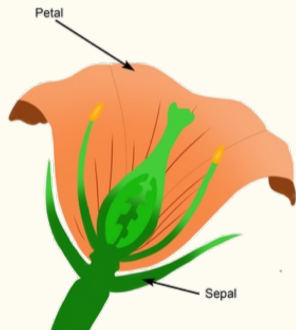


Setosa



Versicolor



Virginica

# The Iris Flower Dataset

THE USE OF MULTIPLE MEASUREMENTS IN TAXONOMIC PROBLEMS

By R. A. FISHER, Sc.D., F.R.S.

I. DISCRIMINANT FUNCTIONS

II. ARITHMETICAL PROCEDURE

Table I

## The Iris Flower Dataset - Formatting

Typical data input format ($m \times n$ matrix):

$n$ features

$m$ observations

$$\begin{array}{cccccccc}
1 & 3 & d & k & 7 & a & \cdots \\
8 & 2 & c & 8 & 1 & c & \cdots \\
7 & 4 & e & x & 1 & d & \cdots \\
9 & 6 & z & y & 5 & e & \cdots \\
5 & 8 & x & z & 8 & f & \cdots \\
\vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots
\end{array}$$

observations: subjects of interest, ~~samples~~ of interest;

features: characteristics describing the observation and they vary among observations.

| sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) | species |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 7.0 | 3.2 | 4.7 | 1.4 | versicolor |
| 6.4 | 3.2 | 4.5 | 1.5 | versicolor |
| 6.9 | 3.1 | 4.9 | 1.5 | versicolor |
| 6.3 | 3.3 | 6.0 | 2.5 | virginica |
| 5.8 | 2.7 | 5.1 | 1.9 | virginica |
| 7.1 | 3.0 | 5.9 | 2.1 | virginica |

## The Iris Flower Dataset - Plotting

- Software choices
- **R**
- **Python**
- SAS
- Stata
- SPSS
- Minitab

# Fisher's Least Significant Difference (LSD)

When doing $post\ hoc$ pairwise $\boldsymbol{t}$-tests, use the following **test statistic (equal variance) for all comparisons:**

$$\boldsymbol{t} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\text{MSW}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}, \text{where } \nu = n - k$$

Note the difference between $s_p^2$ and $\text{MSW}$

## One-way/factor ANOVA

- **One-way/factor ANOVA**: samples can be distinguished by one facotr:
- Brands of tyres
- Species
- *etc.*

- **Two-way/factor ANOVA**: samples can be distinguished by two facotrs:
- Brands of tyres + colours
- Species + location
- *etc.*