

SUMMARIZING DATA & REVISITING PROBABILITY (NOS 2101)PROBABILITY:

Probability is the chance of occurrence anything

$$P(A) = s/p$$

Where  $s$  is sample size or no of positive outcomes and  $p$  is the population size or total no of outcomes.

Probability distribution:

Which describes how the values of a random variable are distributed.

Binominal distribution: The collection of all possible outcomes of a sequence of coin tossing.

Normal distribution:

The means of sufficiently large samples of a data population.

Example :

Probability of ace of Diamond in a pack of 52 cards when 1 card is pulled out at random.

Other mathematical variables.

A random variable is a real-valued function defined on the points of a sample space.

RANDOM VARIABLES ARE TWO BROAD CATEGORIES:

- Random variable with discrete values
- Bivariate Random variable

DISCRETE:

Specified finite or countable list of values, endowed with a probability mass function, characteristic of a probability distribution;

CONTINUOUS:

Any numerical value in an interval or collection of intervals, via a probability density function that is characteristic of a probability distribution;

MIXTURE OF BOTH TYPES:

The realizations of a random variable, that is, the results of randomly choosing values according to the variable's probability distribution function, are called random variates.

2. Do ANOVA test of 3 different datasets which are subset of:

sepal length

iris.1 ← iris [1:10, 1:1] row  
 iris.2 ← iris [11:20, 1:1] column  
 iris.3 ← iris [21:30, 1:1]  
 names(iris.1)

1:10

iris.1	iris.2	iris.3
5.1	5.4	5.4
4.9	4.8	5.1
4.7	4.8	4.6
4.6	4.3	5.1
5.0	5.8	4.8
5.4	5.7	5.0
4.6	5.4	5.0
5.0	5.1	5.2
4.4	5.7	5.2
4.9	5.1	4.7
<u>4.86</u>	<u>5.21</u>	<u>5.01</u>

mean =  $\bar{x}$

sd =

0.291357

0.4817791

↓

iris.2

0.2469818

Total sum of squares

SST =

$$\sum_{i=1}^n \sum_{j=1}^k (x_{ij} - \bar{x})^2$$

$$\bar{x} = 5.1 + 4.9 + 4.7 + 4.6 + 5.0 + 5.4 + 4 + 5.0 + 4.4 + 4.9 + 5.4$$

Stacking:

```
iris.new <- data.frame(cbind(iris.1,  
                             iris.2, iris.3))
```

```
iris.stacked <- stack(iris.new)
```

```
names(iris.stacked)
```

```
[1] "values" "ind"
```

```
iris.anova <- aov(values ~ ind, data =  
                  iris.stacked)
```

```
sum(iris.stacked$values - sst) ^ 2
```

```
sst = 4.018667
```

```
sst <- mean(iris.stacked$values)
```

```
sst = 5.026667
```

$$\bar{x} \quad SST = 4.026667$$

mean (i.e. stacked  
& values)

$$= \sum (x_{ij} - \bar{x})^2$$

$$SST = 4.018667$$

Find (SSTR) :- Treatment sum of squares

$$10 \times (4.86 - 5.026667)^2 + [10 \times (5.21 - 5.026667)^2] \\ + [10 \times (5.01 - 5.026667)^2]$$

$$10 \times (-0.166667)^2 + 10 \times (0.183333)^2 +$$

$$10 \times (-0.016667)^2$$

$$10 \times 0.02777789 + 10 \times 0.03361099 +$$

$$0.0002777889$$

$$= 0.002777889 + 0.3361099 +$$

$$0.2777789 = \underline{0.616666689}$$

$$SSTR = 0.617$$

Error sum of squares  
(SS<sub>E</sub>)

$$\sum \sum (x_{ij} - \bar{x}_j)^2$$

iris.1

5.1 4.9 4.7 4.6 5.0 5.4 4.6 5.0  
4.4 4.9

next mean(iris.1)

m1

4.86

iris.1 - m1

0.24 0.04 -0.16 -0.26 0.14 0.54 -0  
0.14 -0.46 0.04

0.0576 0.0016 0.0256 0.0676 0.0196

0.2916 0.0676 0.0196 0.2116 0.00

↑  
(iris.1 - m1)<sup>2</sup> =

$$\text{sum}((\text{iris.1} - m_1)^2) = 0.764$$

$$\text{sum}((\text{iris.2} - m_2)^2) = 2.089$$

$$\text{sum}((\text{iris.3} - m_3)^2) = 0.549$$

$$m_3 = 5.0$$

$$m_2 = 5.2$$

$$0.764 + 2.089 + 0.549$$

$$SSE = \underline{\underline{3.402}}$$

$$(iris.2 - m2) =$$

$$0.19 \quad -0.41 \quad -0.41 \quad -0.91 \quad 0.59 \quad 0.49$$

$$0.19 \quad -0.11 \quad 0.49 \quad -0.11$$

$$(iris.2 - m2)^2 =$$

$$0.0361 \quad 0.1681 \quad 0.1681 \quad 0.8281$$

$$0.3481 \quad 0.2401 \quad 0.0361 \quad 0.0121$$

$$0.2401 \quad 0.0121$$

$$\text{sum}((iris.2 - m2)^2)$$

$$= 2.089$$

$$(iris.3 - m3)$$

$$0.39 \quad 0.09 \quad -0.41 \quad 0.09 \quad -0.21 \quad -0.01$$

$$-0.01 \quad 0.19 \quad 0.19 \quad -0.31$$

$$(iris.3 - m3)^2 =$$

$$0.1521 \quad 0.0081 \quad 0.1681 \quad 0.0081$$

$$0.0441 \quad 0.0001 \quad 0.0001 \quad 0.0361$$

$$0.0361 \quad 0.0961$$

Ar  
1.  
2.  
3.  
4.  
5.  
6.  
7.  
8.

$$SST = SSTR + SSE$$

$$4.018667 = 0.617 + 3.402$$
$$4.019$$

$$\text{round}(4.018667, 3)$$
$$= 4.019$$

$$\text{Total mean square} = \frac{4.019}{N-1} = \frac{4.019}{30-1}$$
$$= 0.1385862$$

$$\text{Mean square treatment: (MSTR)}$$
$$= \frac{SSTR}{c-1} = \frac{0.617}{(3-1)} = \boxed{0.3085}$$
$$= \underline{0.3083}$$

$$p\text{-val} = P_f(3.354131, 2, 27)$$
$$p\text{-val}$$
$$= 0.95$$



Mean square error  
(MSE)

$$= \frac{SSE}{N-c} = \frac{3.402}{30-3} = \underline{0.126}$$

F. distribution:-

$$\frac{MSTR}{MSE} = \frac{0.3085}{0.126} = 2.448413 \approx 2.447.$$

critical F. value:-

$$\begin{aligned} df_1 &= c-1 = 2 \\ df_2 &= N-c = 27 \\ &CV \\ F_{2,27} &= \end{aligned}$$

For p-value:-

$$= P(F, df=n-1)$$

$$af(.95, df=2, df2=27)$$

$$p = 3.354131$$

The 95% percentile of the F-distribution with (2, 27) degrees of freedom is 3.354131

$$2.447 - 3.354 + 3.1^*$$

P-value look up:

Given an observed value of a test statistic  $t_{obs}$  that is supposed (under the null hypothesis) be 'F' distributed with  $df_1$  numerator degrees of freedom and  $df_2$  denominator degrees of freedom the p-value for a lower-tailed test

a) a lower-tailed test:

is the lower tail area given by

$$P_f(t_{obs}, df_1, df_2)$$

$$P_f(2.447, df_1=2, df_2=27)$$

$$= 0.8944789$$

$$1 - P_f = \underline{0.1055211}$$

$$= 0.105$$

b) an upper-tailed test is the upper tail area given by

$$1 - pf(t_{obs}, df1, df2)$$

c) a two-tailed test is twice whichever one-tailed p-value is smaller

$$p_{low} \leftarrow pf(t_{obs}, df1, df2)$$

$$2 \times \min(p_{low}, 1 - p_{low})$$

Example:- in DeGroot & Schervish, the (obviously made up) "observed" value of the test statistic they call it  $v$ ) is 3.0. the numerator degrees of freedom is 5 and denominator degrees of freedom is 20.