

## Module 3 summary: Probability: The Language of Uncertainty

## Contingency Tables

- Consider the relationship between the cut of a diamond and its clarity.

```
> tally(~ cut + clarity, margins=TRUE, data=Diamonds)
```

cut	clarity								Total
	I1	SI1	SI2	VS1	VS2	VS1	VS2	IF	
Fair	210	408	466	170	261	17	69	9	1610
Good	96	1560	1081	648	978	186	286	71	4906
Very Good	84	3240	2100	1775	2591	789	1235	268	12082
Premium	205	3575	2949	1989	3357	616	870	230	13791
Ideal	146	4282	2598	3589	5071	2047	2606	1212	21551
Total	741	13065	9194	8171	12258	3655	5066	1790	53940

3

## Contingency Tables

- To explore the relationship between two categorical variables for the object, we create contingency tables, also known as cross-tabulations.
- Contingency tables present one categorical variable as the rows and the other categorical variable as the columns.
- These tables are used to calculate the conditional probabilities or percentages, which makes it easier for us to explore potential associations between variables.

2

## Contingency Tables

- Can calculate the conditional column percentages using the following code:
- Use round() to reduce the decimal points, otherwise R uses 6 decimal points.**

```
> round(tally(~ cut | clarity, margins=TRUE, format = "proportion", data=Diamonds), 3)
```

cut	clarity							
	I1	SI1	SI2	VS1	VS2	VS1	VS2	IF
Fair	0.283	0.031	0.051	0.021	0.021	0.005	0.014	0.005
Good	0.130	0.119	0.118	0.079	0.080	0.051	0.056	0.040
Very Good	0.113	0.248	0.228	0.217	0.211	0.216	0.244	0.150
Premium	0.277	0.274	0.321	0.243	0.274	0.169	0.172	0.128
Ideal	0.197	0.328	0.283	0.439	0.414	0.560	0.514	0.677
Total	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

4

## Interpretation of Contingency table

- These are conditional column probabilities.
- If we add all the probabilities for a column together, they will equal 1, e.g. sum column 1 = (0.283, 0.130, 0.113, 0.227, 0.197) = 1.00.
- Comparing the worst clarity I1 with the best, IF (Flawless). Probability of IF diamonds having Ideal cut is 0.667 vs I1 of 0.197.

5

## Clustered Bar Charts using contingency table

```
> barplot(table, main = "Diamond Cut Quality by
Clarity", ylab="Proportion within
Clarity", ylim=c(0,8), legend=rownames(table), beside=TRUE,
args.legend=c(x = "top", horiz=TRUE, title="Cut"), xlab="Clarity")
> grid()
```

- Notice how grid() was added after the plot was produced in R. Grid lines can help the viewer quickly read off and compare values on the plot axes.

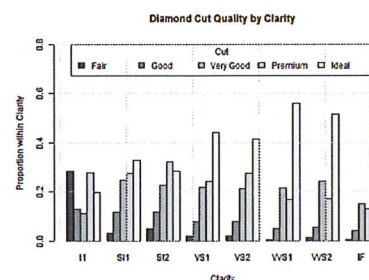
7

## Clustered Bar Charts using contingency table

- Clustered bar charts are a great way to visualise two qualitative variables.
- First, must create contingency table in an object called table.
- This will make it easy for us to create the clustered bar chart.
- `> table<-tally( ~ cut | clarity, format = "proportion", data=Diamonds)`

6

When the clarity of a diamond increases, the quality of the cut also tends to increase.



8

## Probability:

- **Probability** is defined as the **proportion** of times a random **event** occurs in a very large number of **trials**. Its **value is between 0 and 1**.
- The probability of an event =  $f/n$ , where  $f$  is the frequency or number of times an event occurs and  $n$  is the total sample size.
- As the sample size  $n$  increases, the sample will begin to approximate the true population probability.

9

## Rules

$$Pr(2 \text{ serves}) = 4984/17042 = .292$$

$$Pr(< 1 \text{ serve}) + Pr(1 \text{ serve}) + Pr(2 \text{ serves}) + Pr(3 \text{ serves or more}) = 1$$

$$\frac{3368 + 5445 + 4984 + 3244}{17042} = \frac{17042}{17042} = 1$$

11

Table 10 of the Australian Bureau of Statistics  
2011 -2012 National Health Survey

FRUIT INTAKE  
Fruit Intake - Table  
Table 1. Self-reported Fruit Intake-2011-12, Persons '000

	Age group (years)						
	18-24	25-34	35-44	45-54	55-64	65-74	75 +
<b>Males</b>							
Usual daily intake of fruit							
< 1 serve	202	403	463	306	202	129	73
1 serve	297	503	626	479	345	243	142
2 serves	254	379	349	303	309	207	174
3 serves or +	158	202	215	266	263	195	164
Total	1122	1606	1661	1433	1269	824	651
<b>Females</b>							
Usual daily intake of fruit							
< 1 serve	159	264	240	281	191	83	61
1 serve	417	602	659	486	307	209	165
2 serves	318	603	448	603	483	323	248
3 serves or +	138	231	281	324	343	242	214
Total	1232	1655	1587	1514	1324	857	688
<b>Persons</b>							
Usual daily intake of fruit							
< 1 serve	361	667	703	587	393	212	134
1 serve	714	1105	1285	965	652	452	307
2 serves	572	982	797	906	792	530	422
3 serves or +	279	430	478	590	603	437	386
Total	1926	3209	3263	3048	2437	1631	1249

10

## Rules

- Two events are **mutually exclusive** if, when one event occurs, the other cannot and vice versa.
- Mutually exclusive sets have **no intersection**:  $Pr(A \cap B) = 0$ . We use  $\cap$  to denote an intersection.
- Example: The levels of fruit consumption are mutually exclusive. A person cannot occupy more than one category at a particular time.

12

## Rules Intersection, Union and Complement

$$Pr(1 \text{ serve} \cap 2 \text{ serves}) = 0$$

$$Pr(< 1 \text{ serve} \cap \text{Male}) = \frac{2471}{17042} = .145$$

$$Pr(1 \text{ serve} \cup < 1 \text{ serve}) = \frac{3368 + 5445}{17042} = .517$$

$$Pr(\leq 1 \text{ serve}) = Pr(1 \text{ serve}) + Pr(2 \text{ serves}) + Pr(3 \text{ serves or more})$$

$$= \frac{5445 + 4954 + 3241}{17042} = .802$$

or...

$$Pr(\leq 1 \text{ serve}) = 1 - Pr(< 1 \text{ serve}) = 1 - \frac{3368}{17042} = .802$$

13

## Conditional Probability

The probability that an event, B, will occur given that another event, A has already occurred.

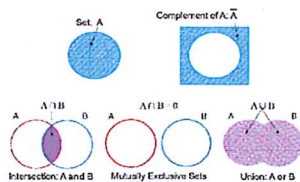
$$Pr(B|A) = \frac{Pr(A \cap B)}{Pr(A)}$$

Using an example...

$$Pr(< 1 \text{ serve} | \text{Male}) = \frac{Pr(\text{Male} \cap < 1 \text{ serve})}{Pr(\text{Male})} = \frac{2471/17042}{8842/17042} = .280$$

15

## Rules



14

## Conditional probability

Using conditional probability to check independence.

The two events A and B are independent if and only if  $Pr(A|B) = Pr(A)$  or  $Pr(B|A) = Pr(B)$ .

Use this rule to reconfirm that gender and fruit consumption are dependent:

$$Pr(< 1 \text{ serve} | \text{Male}) = .280$$

$$Pr(< 1 \text{ serve}) = \frac{3368}{17042} = .198$$

16

## Permutations

- There are six candidates. You need to vote for the top three (order matters).
- How many possible ways can you assign your votes, 1<sup>st</sup>, 2<sup>nd</sup> and 3<sup>rd</sup> preference? This is an example of a permutation problem, 3 possible ways:

Veronica Paskett	Milagros Depocio	Loraine Muntz	Thuy Silverberg	Myriam Hakes	Maude Dimery
1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	-	-	-
-	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	-	-
-	2 <sup>nd</sup>	1 <sup>st</sup>	-	3 <sup>rd</sup>	-

17

## Combinations

- In combination, we need to know how many possible combinations of selecting four out of ten friends exist (with combinations the order does not matter), 3 possibilities are:

Leah	Rosale	Marlena	Tarra	Graham	Gilberto	Marcos	Gladis	Otha	Jeremiah
Ticket	-	Ticket	-	Ticket	-	-	-	-	Ticket
-	-	-	-	-	-	Ticket	Ticket	Ticket	Ticket
-	Ticket	-	Ticket	-	Ticket	-	-	-	Ticket

19

## Permutations

Calculation of number of possible permutation of k observations out of n observations;

$$P(n, k) = \frac{n!}{(n-k)!}$$

The n! is known as the factorial of a number. For example  $6! = 6 \times 5 \times 4 \times 3 \times 2 \times 1 = 720$ .

Solving the voting problem:

$$P(6, 3) = \frac{6!}{(6-3)!} = \frac{6!}{3!} = \frac{720}{6} = 120$$

Using R:

```
> factorial(6)/factorial(6-3)
[1] 120
```

18

## Combinations

$$C(n, k) = \frac{n!}{(n-k)!k!}$$

This is known as the choose formula or the binomial coefficient (we will revisit this in Module 4). Solving, we find:

$$C(n, k) = \frac{n!}{(n-k)!k!} = \frac{10!}{(10-4)!4!} = \frac{10!}{6!4!} = \frac{3628800}{17280} = 210$$

Using R:

```
> choose(10,4)
[1] 210
```

20