

The background is a light gray gradient. It is decorated with numerous realistic water droplets of various sizes, some with highlights and shadows, scattered across the frame. In the upper center, there is a faint, circular, textured pattern that resembles a lens flare or a subtle watermark.

BUSINESS ANALYTICS

RELATIONSHIPS/ASSOCIATIONS

RELATIONSHIPS AMONG VARIABLES

- THE PRIMARY INTEREST IN DATA ANALYSIS IS USUALLY IN *RELATIONSHIPS* BETWEEN VARIABLES.
 - THE MOST USEFUL NUMERICAL SUMMARY MEASURE IS CORRELATION.
 - THE MOST USEFUL GRAPH IS A SCATTERPLOT.
 - TO BREAK DOWN A NUMERICAL VARIABLE BY A CATEGORICAL VARIABLE, IT IS USEFUL TO CREATE SIDE-BY-SIDE BOX PLOTS.
 - EXCEL'S® PIVOT TABLE BREAKS DOWN ONE VARIABLE BY OTHERS SO THAT ALL SORTS OF RELATIONSHIPS CAN BE UNCOVERED VERY QUICKLY.

RELATIONSHIPS AMONG CATEGORICAL VARIABLES

- THE MOST MEANINGFUL WAY TO EXAMINE RELATIONSHIPS BETWEEN TWO CATEGORICAL VARIABLES IS WITH COUNTS AND CORRESPONDING CHARTS OF THE COUNTS.
 - YOU CAN FIND COUNTS OF THE CATEGORIES OF EITHER VARIABLE SEPARATELY, AS WELL AS COUNTS OF THE *JOINT* CATEGORIES OF THE TWO VARIABLES.
 - CORRESPONDING PERCENTAGES OF TOTALS AND CHARTS HELP TELL THE STORY.
- IT IS CUSTOMARY TO DISPLAY ALL SUCH COUNTS IN A TABLE CALLED A **CROSSTABS** (FOR CROSSTABULATIONS). THIS IS ALSO SOMETIMES CALLED A **CONTINGENCY TABLE**.



SMOKING DRINKING.XLSX

- **OBJECTIVE:** TO USE A CROSSTABS TO EXPLORE THE RELATIONSHIP BETWEEN SMOKING AND DRINKING.
- **SOLUTION:** DATA SET LISTS THE SMOKING AND DRINKING HABITS OF 8761 ADULTS.
- CATEGORIES HAVE BEEN CODED “N,” “O,” “H,” “S,” AND “D” FOR “NON,” “OCCASIONAL,” “HEAVY,” “SMOKER,” AND “DRINKER.”

	A	B	C
1	Person	Smoking	Drinking
2	1	NS	OD
3	2	NS	HD
4	3	OS	HD
5	4	HS	ND
6	5	NS	OD
7	6	NS	ND
8	7	NS	OD
9	8	NS	ND
10	9	OS	HD
11	10	HS	HD



EXAMPLE 3.1: SMOKING DRINKING.XLSX (SLIDE 2 OF 2)

- TO CREATE THE CROSSTABS, ENTER THE CATEGORY HEADINGS IN EXCEL AND USE THE COUNTIFS FUNCTION TO FILL THE TABLE WITH COUNTS OF JOINT CATEGORIES.
- NEXT, SUM ACROSS ROWS AND DOWN COLUMNS TO GET TOTALS.
- THEN EXPRESS THE COUNTS AS PERCENTAGES OF ROW AND PERCENTAGES OF COLUMN.

	E	F	G	H	I
1	Crosstabs from COUNTIFS formulas				
2					
3		NS	OS	HS	Total
4	ND	2118	435	163	2716
5	OD	2061	1067	552	3680
6	HD	733	899	733	2365
7	Total	4912	2401	1448	8761
8					
9	Shown as percentages of row				
10		NS	OS	HS	Total
11	ND	78.0%	16.0%	6.0%	100.0%
12	OD	56.0%	29.0%	15.0%	100.0%
13	HD	31.0%	38.0%	31.0%	100.0%
14					
15	Shown as percentages of column				
16		NS	OS	HS	
17	ND	43.1%	18.1%	11.3%	
18	OD	42.0%	44.4%	38.1%	
19	HD	14.9%	37.4%	50.6%	
20	Total	100.0%	100.0%	100.0%	

RELATIONSHIPS AMONG CATEGORICAL VARIABLES AND A NUMERICAL VARIABLE

- THE **COMPARISON PROBLEM** IS AN IMPORTANT PROBLEMS IN DATA ANALYSIS. IT OCCURS WHENEVER YOU WANT TO COMPARE A NUMERICAL MEASURE ACROSS TWO OR MORE SUBPOPULATIONS.
 - EXAMPLES:
 - THE SUBPOPULATIONS ARE MALES AND FEMALES, AND THE NUMERICAL MEASURE IS SALARY.
 - THE SUBPOPULATIONS ARE DIFFERENT REGIONS OF THE COUNTRY, AND THE NUMERICAL MEASURE IS THE COST OF LIVING.
 - THE SUBPOPULATIONS ARE DIFFERENT DAYS OF THE WEEK, AND THE NUMERICAL MEASURE IS THE NUMBER OF CUSTOMERS GOING TO A PARTICULAR FAST-FOOD CHAIN.

RELATIONSHIPS AMONG NUMERICAL VARIABLES

- TO STUDY RELATIONSHIPS AMONG NUMERICAL VARIABLES, A NEW TYPE OF CHART, CALLED A SCATTERPLOT, AND TWO NEW SUMMARY MEASURES, CORRELATION AND COVARIANCE, ARE USED.
- THESE MEASURES CAN BE APPLIED TO ANY VARIABLES THAT ARE DISPLAYED NUMERICALLY.
- HOWEVER, THEY ARE APPROPRIATE ONLY FOR TRULY NUMERICAL VARIABLES, NOT FOR CATEGORICAL VARIABLES THAT HAVE BEEN CODED NUMERICALLY.

SCATTERPLOTS

- A **SCATTERPLOT** IS A SCATTER OF POINTS, WHERE EACH POINT DENOTES THE VALUES OF AN OBSERVATION FOR TWO SELECTED VARIABLES.
 - IT IS A GRAPHICAL METHOD FOR DETECTING RELATIONSHIPS BETWEEN TWO NUMERICAL VARIABLES.
 - THE TWO VARIABLES ARE OFTEN LABELED GENERICALLY AS X AND Y, SO A SCATTERPLOT IS SOMETIMES CALLED AN **X-Y CHART**.
 - THE PURPOSE OF A SCATTERPLOT IS TO MAKE A RELATIONSHIP (OR THE LACK OF IT) APPARENT.



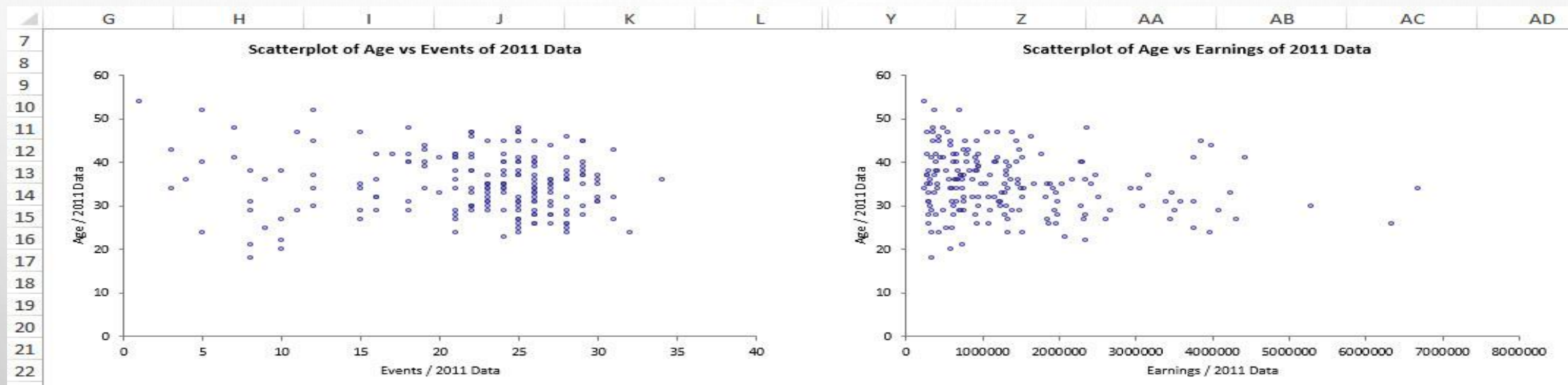
GOLFSTATS.XLSX

- **OBJECTIVE:** TO USE SCATTERPLOTS TO SEARCH FOR RELATIONSHIPS IN THE GOLF DATA.
- **SOLUTION:** DATA SET INCLUDES AN OBSERVATION (STATS) FOR EACH OF THE TOP 200 EARNERS ON THE PGA TOUR.
- USING EXCEL YOU CAN CREATE A SCATTERPLOT FOR TWO VARIABLES SUCH AS AGE AND EVENTS, OR AGE AND EARNINGS.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Rank	Player	Age	Events	Rounds	Cuts Made	Top 10s	Wins	Earnings	Yards/Drive	Driving Accuracy	Greens in Regulation	Putting Average	Sand Save Pct
2	1	Luke Donald	34	19	67	17	14	2	6,683,215	284.1	64.3	67.3	1.7	59.1
3	2	Webb Simpson	26	26	98	23	12	2	6,347,354	296.2	61.9	69.8	1.731	52
4	3	Nick Watney	30	22	77	19	10	2	5,290,674	301.9	58.2	66.9	1.738	48.1
5	4	K.J. Choi	41	22	75	18	8	1	4,434,691	285.6	62	65.9	1.787	55.6
6	5	Dustin Johnson	27	21	71	17	6	1	4,309,962	314.2	57.2	68.4	1.759	41.5
7	6	Matt Kuchar	33	24	88	22	9	0	4,233,920	286.2	64.7	67	1.735	58.9
8	7	Bill Haas	29	26	92	22	7	1	4,088,637	296.6	63.6	69.4	1.775	43.9
9	8	Steve Stricker	44	19	69	18	5	2	3,992,785	288.8	62.5	66	1.71	52.1
10	9	Jason Day	24	21	73	18	10	0	3,962,647	302.6	54.7	64.9	1.737	61
11	10	David Toms	45	23	79	16	7	1	3,858,090	279.1	71.8	66.6	1.749	55.9



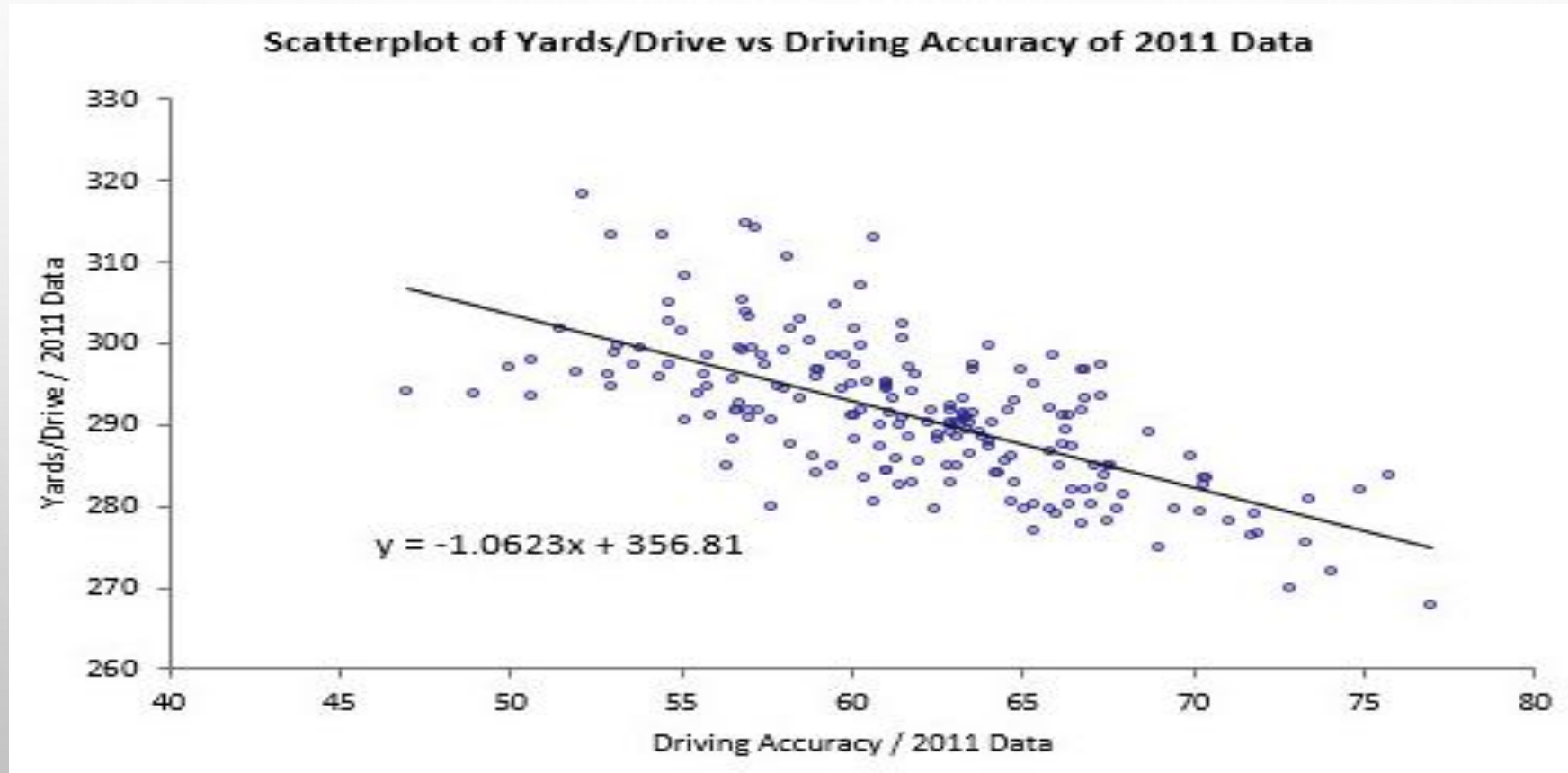
GOLFSTATS.XLSX



TREND LINES IN SCATTERPLOTS

- ONCE YOU HAVE A SCATTERPLOT, EXCEL ENABLES YOU TO SUPERIMPOSE ONE OF SEVERAL TREND LINES ON THE SCATTERPLOT.
 - A **TREND LINE** IS A LINE OR CURVE THAT “FITS” THE SCATTER AS WELL AS POSSIBLE.
 - THIS COULD BE A STRAIGHT LINE, OR IT COULD BE ONE OF SEVERAL TYPES OF CURVES.
- ON THE LAYOUT TAB FOR THE SCATTERPLOT CLICK ON TRENDLINE AND CHOOSE THE APPROPRIATE ONE. (IN EXCEL 2013 ON THE DESIGN TAB CHOOSE ADD CHART ELEMENT).

SCATTERPLOT WITH TREND LINE AND EQUATION SUPERIMPOSED



CORRELATION AND COVARIANCE

(SLIDE 1 OF 4)

- CORRELATION AND COVARIANCE MEASURE THE STRENGTH AND DIRECTION OF A *LINEAR* RELATIONSHIP BETWEEN TWO NUMERICAL VARIABLES.
 - THE RELATIONSHIP IS “STRONG” IF THE POINTS IN A SCATTERPLOT CLUSTER TIGHTLY AROUND SOME STRAIGHT LINE.
 - IF THIS STRAIGHT LINE RISES FROM LEFT TO RIGHT, THE RELATIONSHIP IS *POSITIVE* AND THE MEASURES WILL BE POSITIVE NUMBERS.
 - IF IT FALLS FROM LEFT TO RIGHT, THE RELATIONSHIP IS *NEGATIVE* AND THE MEASURES WILL BE NEGATIVE NUMBERS.
 - THE TWO NUMERICAL VARIABLES MUST BE “PAIRED” VARIABLES.
 - THEY MUST HAVE THE SAME NUMBER OF OBSERVATIONS, AND THE VALUES FOR ANY OBSERVATION SHOULD BE NATURALLY PAIRED.

CORRELATION AND COVARIANCE

(SLIDE 2 OF 4)

- **COVARIANCE** IS ESSENTIALLY AN AVERAGE OF PRODUCTS OF DEVIATIONS FROM MEANS.

$$\text{Covar}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$

- EXCEL HAS A BUILT-IN COVAR FUNCTION
- COVARIANCE HAS A SERIOUS LIMITATION AS A DESCRIPTIVE MEASURE BECAUSE IT IS VERY SENSITIVE TO THE *UNITS* IN WHICH X AND Y ARE MEASURED.

CORRELATION AND COVARIANCE

(SLIDE 3 OF 4)

- **CORRELATION** IS A UNITLESS QUANTITY THAT IS UNAFFECTED BY THE MEASUREMENT SCALE.

$$\text{Correl}(X, Y) = \frac{\text{Covar}(X, Y)}{\text{Stdev}(X) \times \text{Stdev}(Y)}$$

- THE CORRELATION IS ALWAYS BETWEEN -1 AND +1.
 - THE CLOSER IT IS TO EITHER OF THESE TWO EXTREMES, THE CLOSER THE POINTS IN A SCATTERPLOT ARE TO A STRAIGHT LINE.
- EXCEL HAS A BUILT-IN *CORREL* FUNCTION AND THE BUILT IN ADD-IN DATA ANALYSIS CAN CALCULATE CORRELATION ON MULTIPLE VARIABLES.

CORRELATION AND COVARIANCE

(SLIDE 4 OF 4)

- THREE IMPORTANT POINTS ABOUT SCATTERPLOTS, CORRELATIONS, AND COVARIANCES:
 - A CORRELATION IS A SINGLE-NUMBER SUMMARY OF A SCATTERPLOT. IT NEVER CONVEYS AS MUCH INFORMATION AS THE FULL SCATTERPLOT.
 - YOU ARE USUALLY ON THE LOOKOUT FOR LARGE CORRELATIONS, THOSE NEAR -1 OR $+1$.
 - DO NOT EVEN TRY TO INTERPRET COVARIANCES NUMERICALLY EXCEPT POSSIBLY TO CHECK WHETHER THEY ARE POSITIVE OR NEGATIVE. FOR INTERPRETIVE PURPOSES, CONCENTRATE ON CORRELATIONS.