# BUSINESS ANALYTICS

DR. BRENDA MULLALLY

# OUTLIERS

- AN **OUTLIER** IS A VALUE OR AN ENTIRE OBSERVATION (ROW) THAT LIES WELL OUTSIDE OF THE NORM.
  - SOME STATISTICIANS DEFINE AN OUTLIER AS ANY VALUE MORE THAN THREE STANDARD DEVIATIONS FROM THE MEAN, BUT THIS IS ONLY A RULE OF THUMB.
- EVEN IF VALUES ARE NOT UNUSUAL BY THEMSELVES, THERE STILL MIGHT BE UNUSUAL *COMBINATIONS* OF VALUES.
- WHEN DEALING WITH OUTLIERS, IT IS BEST TO RUN THE ANALYSES TWO WAYS: WITH THE OUTLIERS AND WITHOUT THEM.

# MISSING VALUES

- MOST REAL DATA SETS HAVE GAPS IN THE DATA.

- THERE ARE TWO ISSUES: HOW TO DETECT THESE **MISSING VALUES** AND WHAT TO DO ABOUT THEM.

- THE MORE IMPORTANT ISSUE IS WHAT TO DO ABOUT THEM:

  - ONE OPTION IS TO SIMPLY IGNORE THEM. THEN YOU WILL HAVE TO BE AWARE OF HOW THE SOFTWARE DEALS WITH MISSING VALUES.

  - ANOTHER OPTION IS TO FILL IN MISSING VALUES WITH THE AVERAGE OF NON MISSING VALUES, BUT THIS ISN'T USUALLY A VERY GOOD OPTION.

  - A THIRD OPTION IS TO EXAMINE THE NONMISSING VALUES IN THE *ROW* OF A MISSING VALUE; THESE VALUES MIGHT PROVIDE CLUES ON WHAT THE MISSING VALUE SHOULD BE.

# EXCEL TABLES FOR FILTERING, SORTING, AND SUMMARIZING

- TABLES ARE A TOOL INTRODUCED IN EXCEL 2007.

- YOU NOW HAVE THE ABILITY TO DESIGNATE A RECTANGULAR DATA SET AS A TABLE AND THEN EMPLOY A NUMBER OF POWERFUL TOOLS FOR ANALYZING TABLES.

- THESE TOOLS INCLUDE:
  - FILTERING
  - SORTING
  - SUMMARIZING

# EXAMPLE 2.7:CATALOG MARKETING.XLSX

- **OBJECTIVE**: TO ILLUSTRATE EXCEL TABLES FOR ANALYZING THE HYTEX DATA.

- **SOLUTION**: DATA SET CONTAINS DATA ON 1000 CUSTOMERS OF HYTEX, A FICTIONAL DIRECT MARKETING COMPANY.

- DESIGNATE THE DATA SET AS A TABLE BY SELECTING ANY CELL IN THE DATA SET AND CLICKING THE TABLE BUTTON ON THE INSERT RIBBON.

- USE THE DROPDOWN ARROWS NEXT TO THE VARIABLE NAMES TO FILTER IN MANY DIFFERENT WAYS.

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Person | Age | Gender | Own Home | Married | Close | Salary | Children | History | Catalogs | Region | State | City | First Purchase | Amount Spent |
| 2 | 1 | 1 | 0 | 0 | 0 | 1 | $16,400 | 1 | 1 | 12 | South | Florida | Orlando | 10/23/2008 | $218 |
| 3 | 2 | 2 | 0 | 1 | 1 | 0 | $108,100 | 3 | 3 | 18 | Midwest | Illinois | Chicago | 5/25/2006 | $2,632 |
| 4 | 3 | 2 | 1 | 1 | 1 | 1 | $97,300 | 1 | NA | 12 | South | Florida | Orlando | 8/18/2012 | $3,048 |
| 5 | 4 | 3 | 1 | 1 | 1 | 1 | $26,800 | 0 | 1 | 12 | East | Ohio | Cleveland | 12/26/2009 | $435 |
| 6 | 5 | 1 | 1 | 0 | 0 | 1 | $11,200 | 0 | NA | 6 | Midwest | Illinois | Chicago | 8/4/2012 | $106 |
| 7 | 6 | 2 | 0 | 0 | 0 | 1 | $42,800 | 0 | 2 | 12 | West | Arizona | Phoenix | 3/4/2010 | $759 |
| 8 | 7 | 2 | 0 | 0 | 0 | 1 | $34,700 | 0 | NA | 18 | Midwest | Kansas | Kansas City | 6/11/2012 | $1,615 |
| 9 | 8 | 3 | 0 | 1 | 1 | 0 | $80,000 | 0 | 3 | 6 | West | California | San Francisco | 8/17/2006 | $1,985 |
| 10 | 9 | 2 | 1 | 1 | 0 | 1 | $60,300 | 0 | NA | 24 | Midwest | Illinois | Chicago | 5/29/2012 | $2,091 |
| 11 | 10 | 3 | 1 | 1 | 1 | 0 | $62,300 | 0 | 3 | 24 | South | Florida | Orlando | 6/9/2008 | $2,644 |

# CATALOG MARKETING.XLSX



| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Person | Age | Gender | Own Home | Married | Close | Salary | Children | History | Catalogs | Region | State | City | First Purchase | Amount Spent |
| 2 | 1 | 1 | 0 | 0 | 0 | 1 | $16,400 | 1 | 1 | 12 | South | Florida | Orlando | 10/23/2008 | $218 |
| 3 | 2 | 2 | 0 | 1 | 1 | 0 | $108,100 | 3 | 3 | 18 | Midwest | Illinois | Chicago | 5/25/2006 | $2,632 |
| 4 | 3 | 2 | 1 | 1 | 1 | 1 | $97,300 | 1 | NA | 12 | South | Florida | Orlando | 8/18/2012 | $3,048 |
| 5 | 4 | 3 | 1 | 1 | 1 | 1 | $26,800 | 0 | 1 | 12 | East | Ohio | Cleveland | 12/26/2009 | $435 |
| 6 | 5 | 1 | 1 | 0 | 0 | 1 | $11,200 | 0 | NA | 6 | Midwest | Illinois | Chicago | 8/4/2012 | $106 |
| 7 | 6 | 2 | 0 | 0 | 0 | 1 | $42,800 | 0 | 2 | 12 | West | Arizona | Phoenix | 3/4/2010 | $759 |
| 8 | 7 | 2 | 0 | 0 | 0 | 1 | $34,700 | 0 | NA | 18 | Midwest | Kansas | Kansas City | 6/11/2012 | $1,615 |
| 9 | 8 | 3 | 0 | 1 | 1 | 0 | $80,000 | 0 | 3 | 6 | West | California | San Francisco | 8/17/2006 | $1,985 |
| 10 | 9 | 2 | 1 | 1 | 0 | 1 | $60,300 | 0 | NA | 24 | Midwest | Illinois | Chicago | 5/29/2012 | $2,091 |

# FILTERING

- FINDING RECORDS THAT MATCH PARTICULAR CRITERIA IS CALLED *FILTERING*.

- ONE WAY TO FILTER IS TO CREATE AN EXCEL TABLE, WHICH AUTOMATICALLY PROVIDES DROPDOWN ARROWS NEXT TO THE FIELD NAMES THAT ALLOW YOU TO FILTER.

- THERE ARE ALSO THREE WAYS TO FILTER ON ANY RECTANGULAR DATA SET WITH VARIABLE NAMES:
    1. USE THE FILTER BUTTON FROM THE SORT & FILTER DROPDOWN LIST ON THE HOME RIBBON.
    2. USE THE FILTER BUTTON FROM THE SORT & FILTER GROUP ON THE DATA RIBBON.
    3. RIGHT-CLICK ANY CELL IN THE DATA SET AND SELECT FILTER. YOU GET SEVERAL OPTIONS, THE MOST POPULAR OF WHICH IS FILTER BY SELECTED CELL'S VALUE.

# CATALOG MARKETING.XLSX

- **OBJECTIVE**: TO INVESTIGATE THE TYPES OF FILTERS THAT CAN BE APPLIED TO THE HYTEX DATA.

- **SOLUTION**: THERE IS ALMOST NO LIMIT TO THE FILTERS YOU CAN APPLY, BUT HERE ARE A FEW POSSIBILITIES:
  - FILTER ON ONE OR MORE VALUES IN A FIELD.
  - FILTER ON MORE THAN ONE FIELD.
  - FILTER ON A CONTINUOUS NUMERICAL FIELD.
  - *TOP 10* AND *ABOVE/BELOW AVERAGE* FILTERS.
  - FILTER ON A TEXT FIELD.
  - FILTER ON A DATE FIELD.
  - FILTER ON COLOR OR ICON.
  - USE A CUSTOM FILTER.

# EXAMPLE 2.7
# CATALOG MARKETING.XLSX

RESULTS FROM A TYPICAL FILTER

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Person | Age | Gender | Own Home | Married | Close | Salary | Children | History | Catalogs | Region | State | City | First Purchase | Amount Spent |
| 155 | 154 | 2 | 0 | 1 | 1 | 0 | $96,800 | 3 | NA | 24 | Midwest | Kentucky | Louisville | 4/28/2012 | $3,082 |
| 163 | 162 | 2 | 0 | 1 | 1 | 1 | $62,200 | 3 | NA | 24 | Midwest | Indiana | Indianapolis | 6/7/2008 | $2,119 |
| 245 | 244 | 2 | 1 | 1 | 1 | 0 | $82,400 | 2 | 3 | 24 | Midwest | Indiana | Indianapolis | 3/25/2011 | $2,035 |
| 370 | 369 | 2 | 1 | 1 | 1 | 0 | $113,400 | 3 | 3 | 18 | Midwest | Kentucky | Louisville | 11/25/2011 | $1,790 |
| 430 | 429 | 2 | 1 | 1 | 1 | 1 | $113,000 | 2 | 2 | 18 | Midwest | Kentucky | Louisville | 6/15/2011 | $1,554 |
| 570 | 569 | 2 | 1 | 1 | 1 | 1 | $70,400 | 2 | NA | 12 | Midwest | Indiana | Indianapolis | 4/12/2007 | $1,127 |
| 764 | 763 | 2 | 0 | 1 | 1 | 1 | $85,500 | 2 | 2 | 18 | Midwest | Kentucky | Louisville | 7/3/2011 | $895 |
| 790 | 789 | 2 | 1 | 1 | 1 | 1 | $74,500 | 2 | 2 | 12 | Midwest | Indiana | Indianapolis | 3/7/2012 | $824 |
| 804 | 803 | 2 | 0 | 1 | 1 | 1 | $72,200 | 2 | 2 | 18 | Midwest | Kentucky | Louisville | 5/29/2011 | $715 |
| 851 | 850 | 2 | 1 | 1 | 1 | 1 | $77,100 | 2 | 2 | 6 | Midwest | Indiana | Indianapolis | 6/17/2012 | $568 |
| 1002 | Total | | | | | | $84,750 | | | | | | | | $14,709 |

# 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 |

- 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.
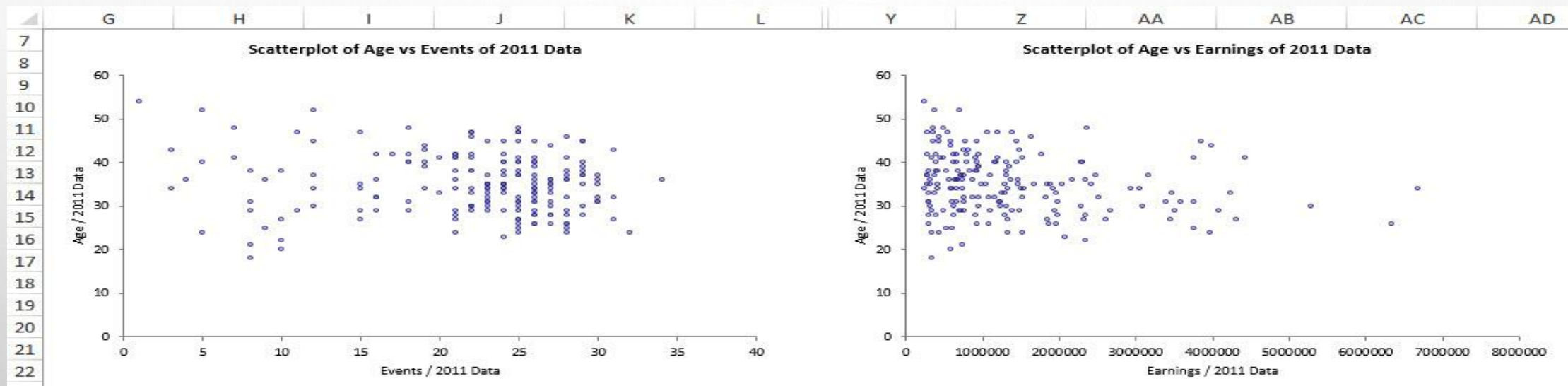
# GOLFSTATS.XLSX

- **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.

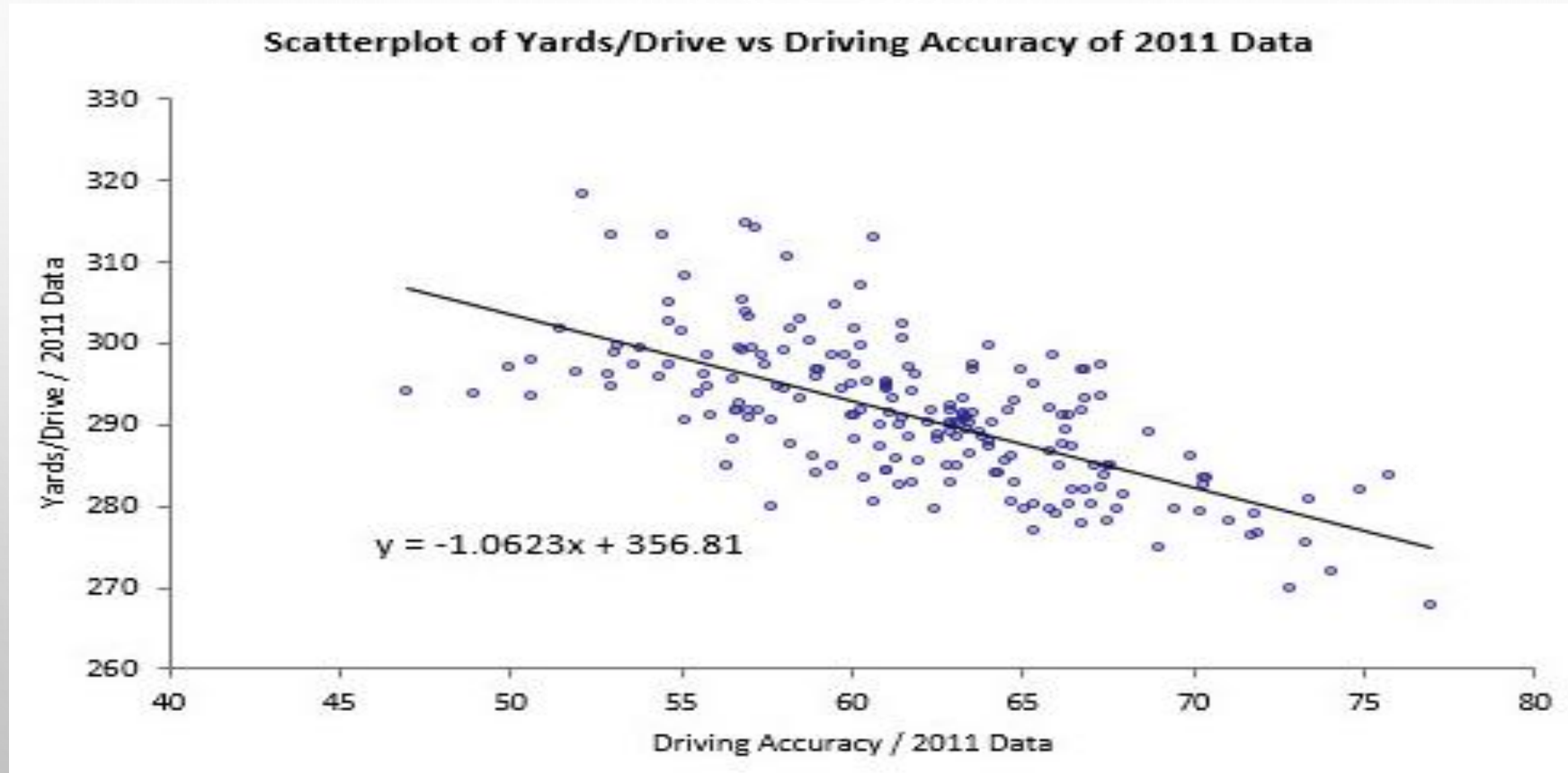| | A | B | C | D | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Rank | Player | Age | Events | Rounds | Cuts Made | Top 10s | Wins | Earnings | Yards/Drive | Driving Accuracy | Greens in Regulation | Putting Average | Sand Save Pct |
| 2 | 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



Scatterplot of Yards/Drive vs Driving Accuracy of 2011 Data

$y = -1.0623x + 356.81$

# 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.

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

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

- 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.