# Session 6: Categorical and Linear Models

## 6.1 – Introduction to Linear Models

| | |
|---|---|
| **Bringing Models to Data:**<br>• **Categorical**<br><br>• **Linear Models**<br>$\quad (x) = ax + b),$<br><br>• **Fitting Lines to Data**  | • **Non Linear Models** <br><br>• **Big Coefficient**<br>$y = a_1 X_1 + a_2 X_2 \quad$ Does $a_1$ or $a_2$ dominate?<br>Only works where we have the data. Often we don't have.<br><br>• **New Reality**, where we don't have the data but things are a lot better than we have today. |

## 6.2 – Categorical Models

| | |
|---|---|
| **Bin Reality into different categories:**<br>Example: Amazon early on – Info company or Delivery company.<br>"Lump to Live" people categorize to simplify the world, making faster decision.<br>Example: What to eat? Green or Not Green.<br>Broccoli, Grasshopper, Banana, Candy bar, Orange , Asparagus, Pear, Strawberry.<br>Simple rule, not always selective enough. |  Point here is that to reduce the large deviations from the mean, we should apply categorization to the grouping to reduce total variation. *Total variation* is the sum of square of the differences. $(\sum variation\ above = 53{,}200)$ |

**Quiz:** 5 students take an exam: Amy scores 100; Ben scores 90; Carla scores 70; Demitri scores 65; Emile scores 85. What is the Total Variation of these scores?  (1) 82, (2) 324, (3) 830, (4) 940

**Analysis:** $\mu = \frac{1}{5}(100 + 90 + 70 + 65 + 85) = 82$ and

$$\sigma^2 = \sum \{(100 - 82)^2 + 90 - 82)^2 + 70 - 82)^2 + 65 - 82)^2 + 85 - 82)^2\} = 830$$

**Ans:** (3) 830

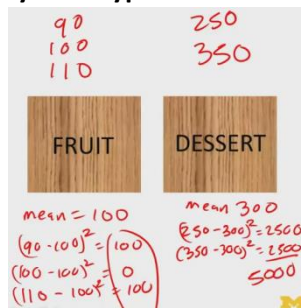| | | |
|---|---|---|
| **Breaking into categories by food type:**<br>Fruit (Pear, Apple, Banana) and Dessert (Cake, Pie)<br><br>In the chart:<br>Fruit (mean=100, variation = 200)<br>Dessert (mean=300, variation=5,000) |  | **How much did I explain?**<br>Combined fruit & dessert variation = 53,200. Fruit variation = 200, Dessert variation = 5,000<br>Explanation:<br>$\frac{53{,}200 - 5{,}000 + 200)}{53{,}200} =$<br>$\frac{48{,}000}{53{,}200} = 0.90226$ |

| | |
|---|---|
| **R-squared:** % variation explained<br>$\quad 1 - \dfrac{total\ subgroup\ variation}{total\ grouped\ variation} = 1 - \dfrac{5200}{53{,}200} = 90.2\%$<br>• R-squared near 1 ➜ model explains a lot<br>• R-squared near 0 ➜ model explains little.<br>• Data quality impacts a good model's R-squared | **Many subcategories:** (fruit, dessert, vegetables, grains).<br>Experts have a lot of 'right' categories. |

**Finally:** Even if you have a lot of explanation, it doesn't always follow you have a good model. ==*"Correlation is not causation".*==
Example: Schools with equestrian teams happen to do better. But it may be due to other factors such as parental support in these districts as opposed to the criteria that was used (money, class sizes, etc.)

**Summary:** Simplest method of explanation is by categorization and the R-squared metric. More categories means there are more possible explanations. The next section on linear models basically creates a box for each value of x.
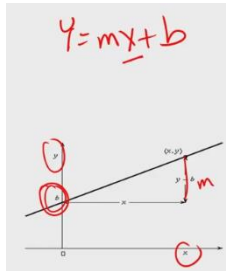
## 6.3 – Linear Models

| | | |
|---|---|---|
| **Linear Model:** Variable X (independent) and variable Y (dependent) that is a function of x, $y = F(x)$ Used for: <br> • Prediction <br> • Understand data | $Y = mx + b$ | **Example:** Purchase TV where x = length of diagonal, and y = cost of TV. <br> Linear Model: $Cost = 15 * (length) + 100$ <br> **Concerns:** <br> *Sign* – (does y increase or decrease in x?), <br> *Magnitude* – (how much does y increase for each one unit increase in x?) <br> How much would a 30-in TV cost? From model <br> $C = 15(30) + 100 = \$550$ |

**Quiz:** According to a psychiatrist, the probability that one developers haa mental disorder is linearly related to the number of stressful life events one experiences. This psychiatrist uses an equation to predict mental disorders: P(mental disorder)=0.1(# stressful life events)+.005. If Jeff has had 5 stressful life events, what is the probability he develops a mental disorder, according to this equation?  (1) 50.5%, (2) 60%, (3) 49.5%, (4) 5.05%
**Analysis:** $P = 0.1(5) + 0.005 = 0.505 = 50.5\%$
**Ans:** (1) 50.5%
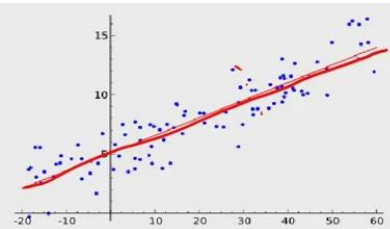
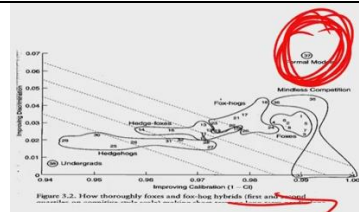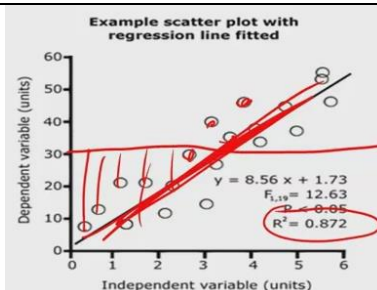| | |
|---|---|
| **Linear Model:** Best linear fit to a data set. Offset from line is a measure of variation. | Robyn Dawes 1979: *"The Robust Beauty of Improper Linear Models in Decision Making"* <br> **Example:** 43 bank loan officers predict which 30 of 60 firms would go bankrupt. They see the financial statements. Bankers 75% accurate, Linear Model – ratio of assets to liabilities 80% accurate. |
| Mehl  (1954) 20 studies of clinicians, Sawyer (1966) 45 studies of predictions in the social world. <br><br> ***Experts NEVER did significantly better than the linear models.*** | ***Recall Tetlock:*** Formal models do better than experts. |

## 6.4 – Fitting Lines to Data

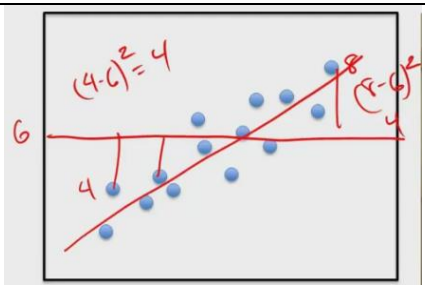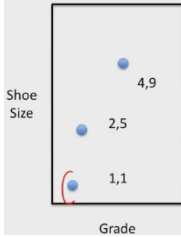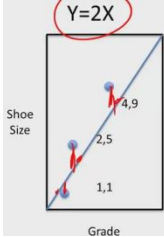| | |
|---|---|
| **Recall R-squared** Note mean of this data and the variation compared to variation of best fit line $R^2 = 0.872$ Explains 87.2% of variation. <br><br> y = 8.56 x + 1.73 <br> $F_{1,19}$= 12.63 <br> $R^2$ = 0.872 | Data: <br> mean = 6 <br> Find variation from mean. <br><br> Compare with variation from line. <br><br> $(4-6)^2 = 4$ <br> $(r-6)^2$ |

| | | | |
|---|---|---|---|
| **Example:** (grade, shoe size)<br>Note variation of dependent<br>variable (shoe size, mean=5)<br>$(1,\mathbf{1})$ - $(1-5)^2 = 16$<br>$(2,\mathbf{5})$ - $(5-5)^2 = 0$<br>$(4,\mathbf{9})$ - $(9-5)^2 = 16$<br>Variation = 32 |  | **Find best linear fit:** Assume y=2x.<br>(x,y,2x)<br>$(1,\mathbf{1},2)$ - $(2-1)^2 = 1$<br>$(2,\mathbf{5},4)$ - $(4-5)^2 = 1$<br>$(4,\mathbf{9},8)$ - $(8-9)^2 = 1$<br>Variation = 3 so model explains<br>$1 - \dfrac{Var_l}{Var_m} = 1 - \dfrac{3}{32} = 90.625\ \%$ |  |

| | |
|---|---|
| **Best Line calculation:**<br>$y = mx + b$<br><br>$x = 1$: $(m+b-1)^2 = m^2 + 2mb + b^2 - 2m - 2b + 1$<br><br>$x = 2$: $(2m+b-5)^2 = 4m^2 + 4mb + b^2 - 20m - 10b + 25$<br><br>$x = 4$: $(4m+b-9)^2 = 16m^2 + 8mb + b^2 - 72m - 18b + 81$<br><br>$sum = 21m^2 + 14mb + 3b^2 - 94m - 30b + 107$ | **Find b, m:**<br>$b = -1$, & $m = {}^8\!/_3$<br>(x, y, model) $\qquad y = \frac{8}{3}x - 1$<br>$\left(1, 1, {}^5\!/_3\right)$: $\left(\frac{5}{3} - 1\right)^2 = \left({}^2\!/_3\right)^2$<br>$\left(2, 5, {}^{13}\!/_3\right)$: $\left(\frac{13}{3} - \frac{15}{3}\right)^2 = \left({}^2\!/_3\right)^2$<br>$\left(4, 9, {}^{29}\!/_3\right)$: $\left(\frac{29}{3} - \frac{27}{3}\right)^2 = \left({}^2\!/_3\right)^2$<br>$sum = 3\left(\frac{4}{9}\right) = {}^4\!/_3$ |

| | |
|---|---|
| **R-squared:**<br>$1 - \dfrac{\left({}^4\!/_3\right)}{32} = 95.8\%$ | **Multiple variables:** $y = ax_1 + bx_2 + c$<br>Recall: Sign and magnitude of x coefficients.<br>Example:<br>Y = Test Score, Q = IQ, T = Teacher Quality, Z = Class Size,<br>$Y = a + bQ + cT + dZ$<br>Note: We expect: $b > 0$, $c > 0$, $d < 0$, but may be wrong<br>**Studies:**<br>Class Size (78 studies, d>0 in 4, d<0 in 13, d~0 in 61)<br>Kindergarten Teacher (teacher quality matters a lot – Chetty (2010)) |

| |
|---|
| **Summary:** Linear models help us explain variation in data. Sign and Magnitude are indicators of the influence of corresponding independent variable. |

## 6.5 – Reading Regression Output

| | | |
|---|---|---|
| **Understanding the Data:**<br>Multiple variables:<br>$y = m_1 x_1 + m_2 x_2 + b$<br><br>Chart data implies:<br>$y = 20x_1 + 10x_2 + 25$<br><br>SE=standard error of coefficient, P-value indicates probability of sign error (1.1% in chart for $m_2$). |  | **Example:**<br>Y = Test Score, T = Teacher Quality, Z = Class Size,<br>$Y = cT + dZ + b$ and expect $c > 0$, $d < 0$<br><br>SE – standard error in coefficients<br>'Intercept', SE =2 ➔ b∈ [23,27] @ 68% confidence<br>$m_1$, SE=1 ➔ $m_1 \in [19,21]$ @ 68%<br>$m_2$, SE=4 ➔ $m_2 \in [6,16]$ @ 68%, note >0 but large error |

| |
|---|
| **Quiz:** A researcher attempts to use a linear model to understand tennis. In her model, points the dependent variable is "points won" and there are several independent variables, such as: serve speed, backhand speed, accuracy, player height and dominant hand. The researcher finds that left-handedness has a positive coefficient and a p-value of 0.20. What does this mean? (1) The coefficient indicates that left-handedness is a disadvantage to tennis players, and the p-value confirms this. (2) The coefficient indicates that left-handedness is an advantage to tennis players, but the p-value indicates that there is some chance that left-handedness is a disadvantage. (3) The |

coefficient indicates that left-handedness is always an advantage to tennis players, and the p-value confirms this. (4) The coefficient indicates that left-handedness has the greatest influence on points won of any variable, but the p-value indicates that there is some chance that this is not the case.

**Ans**: (2)

*Explanation: The positive coefficient means that left-handedness increases the points won. In other words, being left-handed helps the tennis player. The p-value means there is a 20% chance that this is wrong, or in other words, that being left-handed actually has a negative effect.*

| | |
|---|---|
| **Summary:**<br>R-Squared – how much of the data is explained,<br>Observations – how many data points,<br>Standard Error – How much variation in the data to begin with<br>Linear Model – Intercept (sign and magnitude, SE and probability of sign error)<br>Linear Model – Coefficients of independent variables (sign and magnitude, SE & probability of sign error)<br>**Recall:** *Models are generally better than we are!* |  |

## 6.6 – From Linear to Nonlinear

| | | |
|---|---|---|
| **Neuman:** *Study of non-linear functions is like the study of non elephants.*<br>**How to handle nonlinear functions?**<br>*1. Piecewise linear approximations.*<br> |  | **Example:**<br>break clustered data into groups and piecewise linear fit in each section (spline fit). |
| *2. Add non-linear terms:*<br>Can substitute to create a linear form of a non-linear relationship.<br><br>e.g., in $y = a\sqrt{x} + b$ let $z = \sqrt{x}$<br><br>Then we have a linear form $y = az + b$ | |  |

## 6.7 – The Big Coefficient vs The New Reality

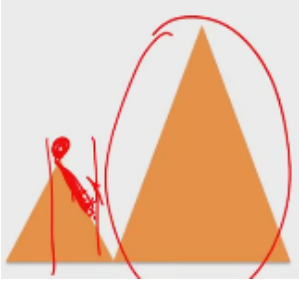| | |
|---|---|
| **Big Coefficient:** does one coefficient dominate in $y = a_1x_1 + a_2x_2 + b$ ? if so, then that has the major influence in explaining y. | **Evidence Based:** (1) medicine, (2) philanthropy, (3) education, (4) management<br>In each, look for dominant influence. |
| **Process**:<br>(1) construct model, (2) gather data, (3) identify important variables, (4) change those variables to test and/or optimize. | **Big Data approach:** (1) gather data, (2) find pattern(s), (3) identify important variables, (4) change those variables<br>**Observations:** (1) Appears to ignore model step, BUT big data does not obviate the uses of models. Identification of a pattern does not explain where it came from, the why. Correlation is not Causation, (2) Linear models tell sign and magnitude of changes in independent variables within data range. |

| | | | |
|---|---|---|---|
| **Examples:**<br>*Feedbacks*: (1) antilock brakes on cars. People adapt and may shorten safe distance between cars. (2) Class size-Performance: data in the range of [20,30] students may not extrapolate to classes of smaller sizes. |  | *Multiple Peaks:* Local maximum (or minimum) can mislead an optimization process. |  |
| **The New Reality:** Need to consider that some big change is necessary to avoid local maximum trap.<br><br>*Big Coefficient:*<br>Tax Cigarettes, Increase HOV lanes, Oat Bran Pretzels<br>*New Reality:*<br>Universal Health Care, Huge Rail System, Fitness Regime | | **Big Coefficient Logic: $447 B**<br> | |
| **New Reality Logic:** 1956 - US Highway system with 41,000 miles of highways for $25B (CPI: $207B today or at $10 million/mile equivalent in today's dollar of $410B). | | **Evidence-based** methods really useful. Put your money on the big coefficients but remember these can blind you to new reality opportunities. ***Models can help optimize in the box as well as think outside the box.*** | |

**Quiz:** Big coefficients may be limiting if they prevent us from achieving a new, better reality. However, programs intended to create new realities are not always successful. Which of the following is a possible explanation? (1) New realities often occur when the coefficients of linear regressions are miscalculated. (2) New realities often don't have a lot of evidentiary support since they are a new way of doing things. (3) New realities focus on existing empirical data rather than trying something completely different. (4) None of the above, new realities are always better than using big coefficients.
**Ans:** (2)
*Explanation: The correct response is that we don't have much empirical support for new realities, by definition. The third option, that new realities fail to try something completely different, is a limitation of big coefficients rather than of new realities. The other two options are incorrect as well: there is no evidence that miscalculated coefficients somehow often lead to new realities; nor can we say that new realities are always better than big coefficients.*