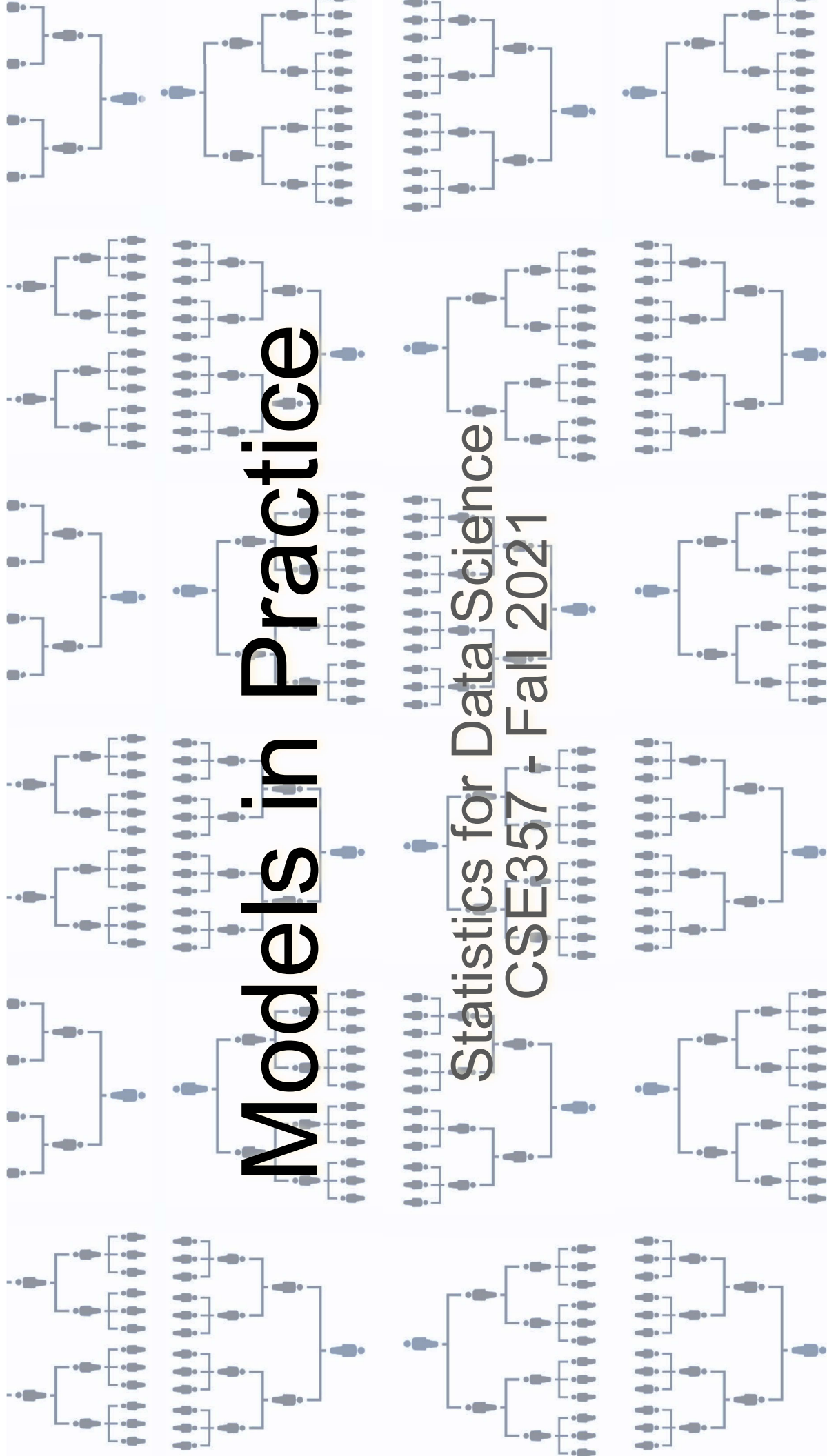


# Models in Practice

Statistics for Data Science  
CSE357 - Fall 2021



# Statistical Models in Practice

Accuracy Metrics

Moderation Analysis

Mediation Analysis

Multi-Level Analysis (Fixed Effect Model)

# Accuracy Metrics

|                  |              | Predicted condition  |   |
|------------------|--------------|--|---|
|                  |              | Positive (PP)  | Negative (PN)   |
| Actual condition | Positive (P) | True positive (TP),<br>hit   | False negative (FN),<br>type II error, miss,<br>underestimation |
|                  | Negative (N) | False positive (FP),<br>type I error, false alarm,<br>overestimation | True negative (TN),<br>correct rejection                        |

(Wikimedia, 2021)

Total population  
= P + N

# Accuracy Metrics

(Wikimedia, 2021)

Sources: [3][4][5][6][7][8][9][10] view · talk · edit

|                  |              | Predicted condition  |   |
|------------------|--------------|--|---|
|                  |              | Positive (PP)  | Negative (PN)   |
| Actual condition | Positive (P) | True positive (TP),<br>hit   | False negative (FN),<br>type II error, miss,<br>underestimation |
|                  | Negative (N) | False positive (FP),<br>type I error, false alarm,<br>overestimation             | True negative (TN),<br>correct rejection                        |
|                  |              | Positive predictive value (PPV),<br>precision<br>= $\frac{TP}{PP} = 1 - FDR$     |   |
|                  |              | False discovery rate (FDR)<br>= $\frac{FP}{PP} = 1 - PPV$                        |   |
|                  |              | $F_1$ score<br>= $\frac{2PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$ |   |

# Accuracy Metrics

Sources: [3][4][5][6][7][8][9][10] view · talk · edit

|                  |                             | Predicted condition   |   |
|------------------|-----------------------------|---|---|
|                  |                             | Positive (PP)   | Negative (PN)   |
| Actual condition | Total population<br>= P + N |   |   |
|                  | Positive (P)                | <b>True positive (TP),</b><br>hit   | <b>False negative (FN),</b><br>type II error, miss,<br>underestimation  |
|                  | Negative (N)                | <b>False positive (FP),</b><br>type I error, false alarm,<br>overestimation         | <b>True negative (TN),</b><br>correct rejection   |
|                  |                             | Positive predictive value (PPV),<br>precision<br>$= \frac{TP}{PP} = 1 - FDR$        | True positive rate (TPR), recall,<br>sensitivity (SEN), probability of detection,<br>hit rate, power<br>$= \frac{TP}{P} = 1 - FNR$  |
|                  |                             | False discovery rate (FDR)<br>$= \frac{FP}{PP} = 1 - PPV$                           | False negative rate (FNR),<br>miss rate<br>$= \frac{FN}{P} = 1 - TPR$   |
|                  |                             | $F_1 \text{ score} = \frac{2PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$ | $F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$ <p>(harmonic mean of precision and recall)</p> |

# Accuracy Metrics

Sources: [3][4][5][6][7][8][9][10] view · talk · edit

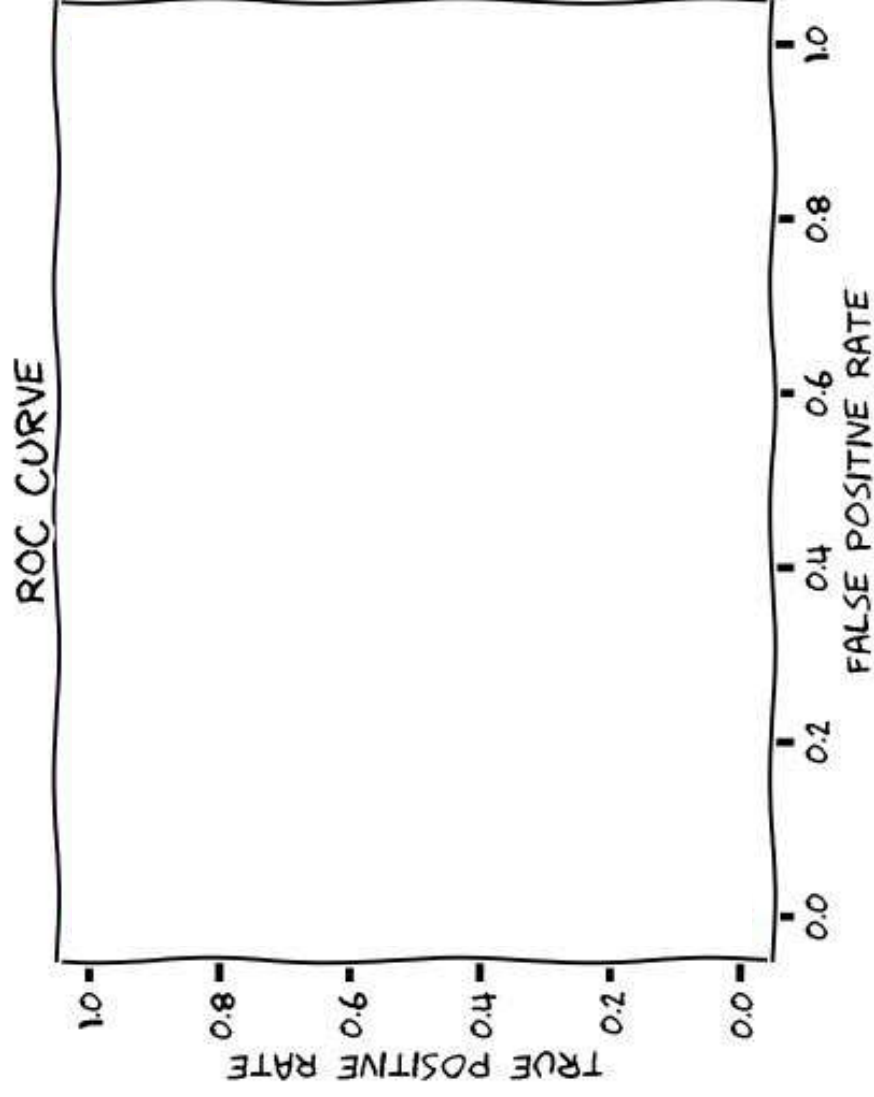
|   |                             | Predicted condition  |  |
|---|-----------------------------|--|--|
|   |                             | Positive (PP)  | Negative (PN)  |
| Actual condition                                  | Total population<br>= P + N |  |  |
|   | Positive (P)                | True positive (TP),<br>hit   | False negative (FN),<br>type II error, miss,<br>underestimation  |
|   | Negative (N)                | False positive (FP),<br>type I error, false alarm,<br>overestimation             | True negative (TN),<br>correct rejection   |
| Prevalence<br>$= \frac{P}{P + N}$                 |                             | Positive predictive value (PPV),<br>precision<br>$= \frac{TP}{PP} = 1 - FDR$     | True positive rate (TPR), recall,<br>sensitivity (SEN), probability of detection,<br>hit rate, power<br>$= \frac{TP}{P} = 1 - FNR$ |
| Accuracy (ACC)<br>$= \frac{TP + TN}{P + N}$       |                             | False discovery rate (FDR)<br>$= \frac{FP}{PP} = 1 - PPV$                        | False negative rate (FNR),<br>miss rate<br>$= \frac{FN}{P} = 1 - TPR$  |
| Balanced accuracy<br>(BA) = $\frac{TPR + TNR}{2}$ |                             | $F_1$ score<br>$= \frac{2PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$ |  |

# Accuracy Metrics

Sources: [3][4][5][6][7][8][9][10] view · talk · edit

| (Wikimedia, 2021)   |  | Predicted condition   |   | Informedness, bookmaker informedness (BM)   | Prevalence threshold (PT)  |
|---|--|---|---|---|--|
|   |  | Positive (PP)   | Negative (PN)   | $= \text{TPR} + \text{TNR} - 1$   | $= \frac{\sqrt{\text{TPR} \times \text{FPR}} - \text{FPR}}{\text{TPR} - \text{FPR}}$                           |
| Actual condition  | Positive (P)   | True positive (TP),<br>hit  | False negative (FN),<br>type II error, miss,<br>underestimation   | True positive rate (TPR), recall,<br>sensitivity (SEN), probability of detection,<br>hit rate, power<br>$= \frac{\text{TP}}{\text{P}} = 1 - \text{FNR}$ | False negative rate (FNR),<br>miss rate<br>$= \frac{\text{FN}}{\text{P}} = 1 - \text{TPR}$                     |
|   | Negative (N)   | False positive (FP),<br>type I error, false alarm,<br>overestimation                              | True negative (TN),<br>correct rejection  | False positive rate (FPR),<br>probability of false alarm, fall-out<br>$= \frac{\text{FP}}{\text{N}} = 1 - \text{TNR}$                                   | True negative rate (TNR),<br>specificity (SPC), selectivity<br>$= \frac{\text{TN}}{\text{N}} = 1 - \text{FPR}$ |
|   | Prevalence<br>$= \frac{\text{P}}{\text{P} + \text{N}}$   | Positive predictive value (PPV),<br>precision<br>$= \frac{\text{TP}}{\text{PP}} = 1 - \text{FDR}$ | False omission rate<br>(FOR)<br>$= \frac{\text{FN}}{\text{PN}} = 1 - \text{NPV}$  | Positive likelihood ratio (LR+)<br>$= \frac{\text{TPR}}{\text{FPR}}$  | Negative likelihood ratio (LR-)<br>$= \frac{\text{FNR}}{\text{TNR}}$   |
| Accuracy (ACC)<br>$= \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}}$ | False discovery rate (FDR)<br>$= \frac{\text{FP}}{\text{PP}} = 1 - \text{PPV}$   | Negative predictive value (NPV) $= \frac{\text{TN}}{\text{PN}} = 1 - \text{FOR}$                  | Markedness (MK), deltaP ( $\Delta p$ )<br>$= \text{PPV} + \text{NPV} - 1$   | Diagnostic odds ratio (DOR) $= \frac{\text{LR}+}{\text{LR}-}$   |  |
| Balanced accuracy<br>(BA) $= \frac{\text{TPR} + \text{TNR}}{2}$         | $F_1$ score<br>$= \frac{2\text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$ | Fowlkes–Mallows index<br>(FM) $= \sqrt{\text{PPV} \times \text{TPR}}$                             | Matthews correlation coefficient (MCC)<br>$= \frac{\sqrt{\text{TPR} \times \text{TNR} \times \text{PPV} \times \text{NPV}} - \sqrt{\text{FNR} \times \text{FPR} \times \text{FOR} \times \text{FDR}}}{\text{TP} + \text{FN} + \text{FP}}$ | Threat score (TS), critical success index (CSI), Jaccard index<br>$= \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$                               |  |

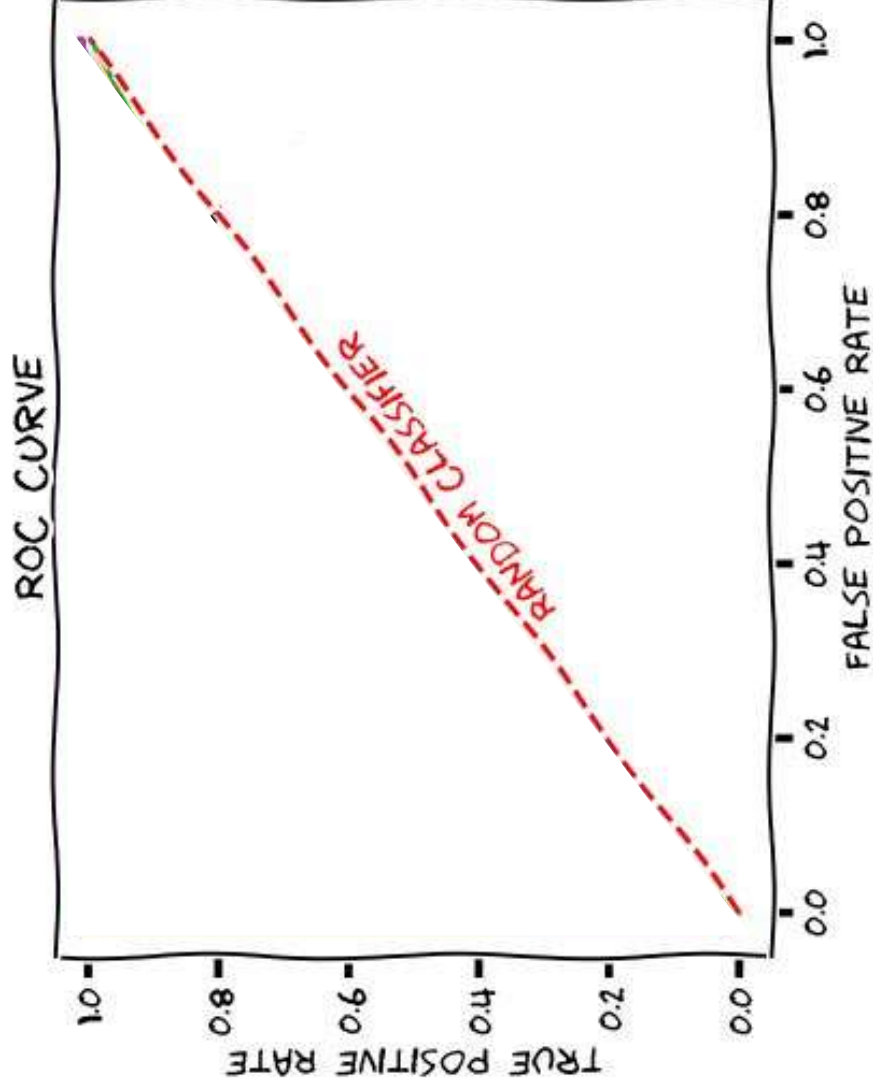
# Accuracy Metrics



(Glass Box: Measuring Performance, Retrieved 2021)



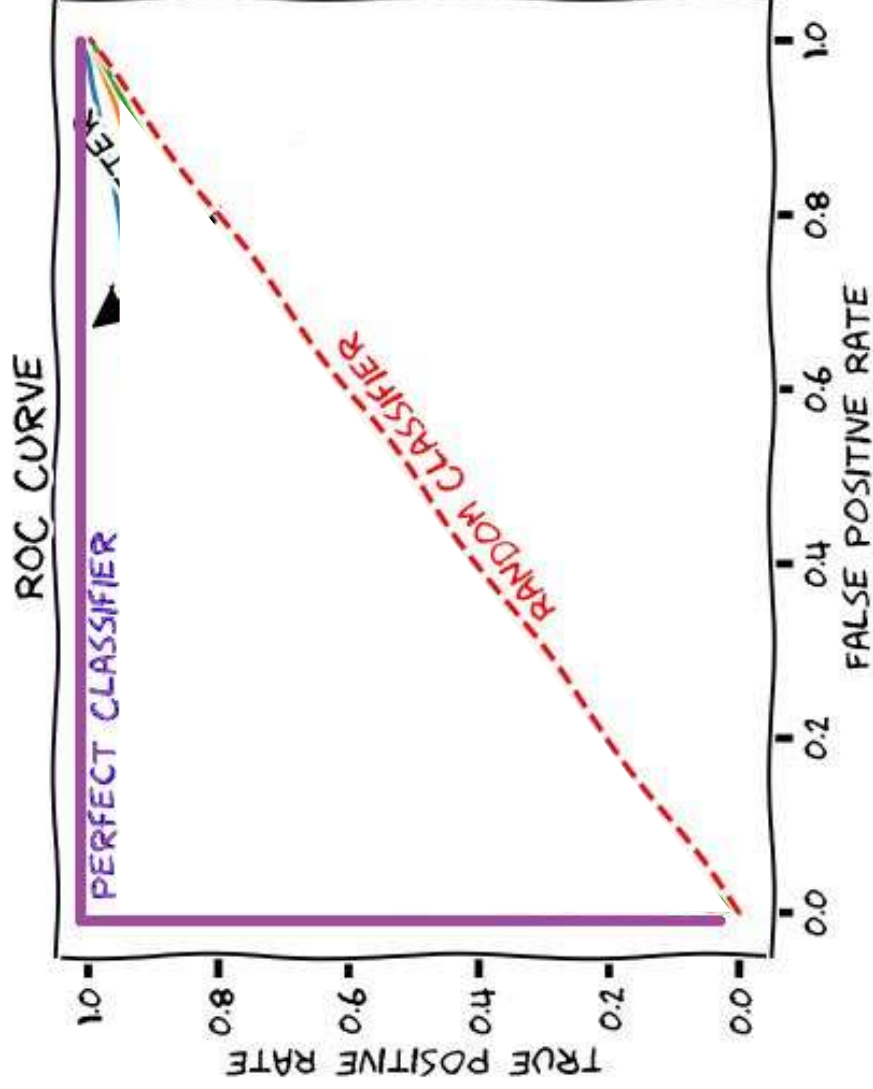
# Accuracy Metrics



Logistic regression gives us a probability which we threshold (often at 0.5) to get predicted true(1) versus predicted false(0).

Change FP and TP rate by varying the threshold for predicting true or false.

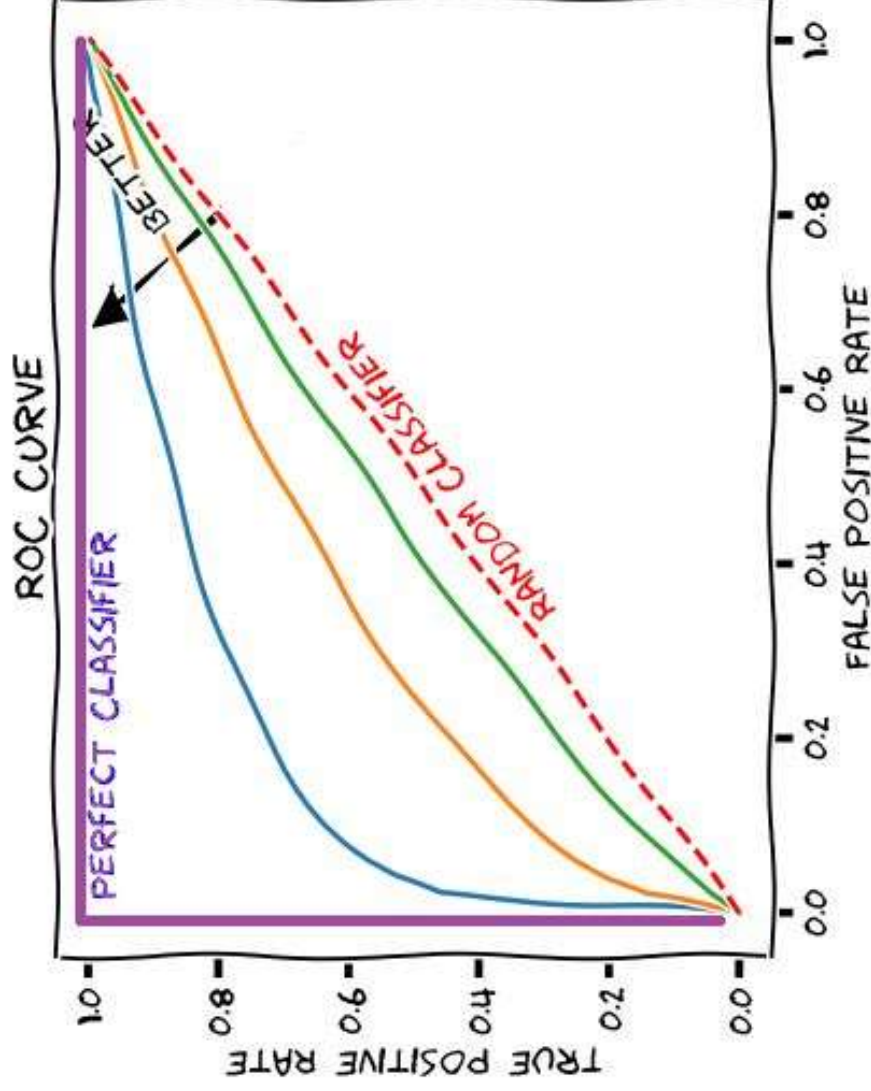
# Accuracy Metrics



Logistic regression gives us a probability which we threshold (often at 0.5) to get predicted true(1) versus predicted false(0).

Change FP and TP rate by varying the threshold for predicting true or false.

# Accuracy Metrics



Logistic regression gives us a probability which we threshold (often at 0.5) to get predicted true(1) versus predicted false(0).

Change FP and TP rate by varying the threshold for predicting true or false.

# Mediation Analysis

Path Analyses (a type of “structured equation modeling”)

How much does **M** mediate the effect of **X** on **Y**?

$$Y = \beta_0 c' b + c' X + b M + \epsilon c' b$$

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$

# Mediation Analysis

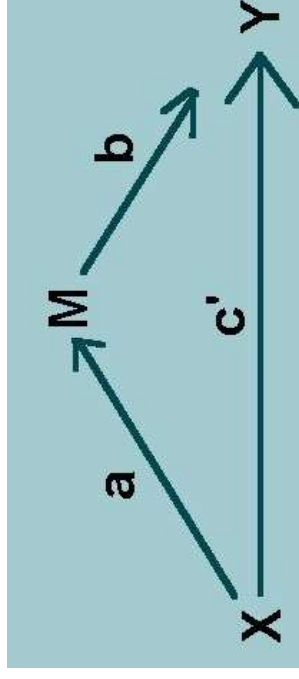
Path Analyses (a type of “structured equation modeling”)

How much does **M** mediate the effect of **X** on **Y**?

$$Y = \beta_{0c} + cX + \epsilon_c$$

$$X = \beta_{0a} + aM + \epsilon_a$$

$$Y = \beta_{0c'b} + c'X + bM + \epsilon_{c'b}$$



(Kenney, 2015)

<http://davidakenny.net/cm/mediate.htm>

# Mediation Analysis

Path Analyses (a type of “structured equation modeling”)

How much does **M** mediate the effect of **X** on **Y**?

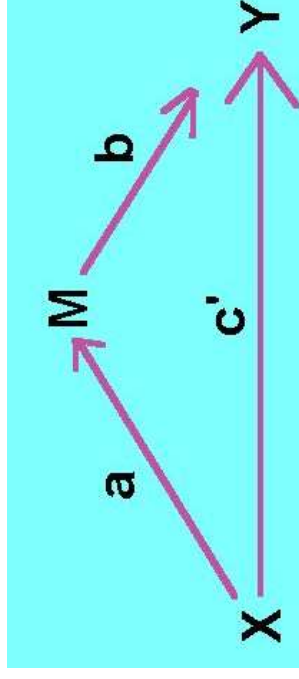
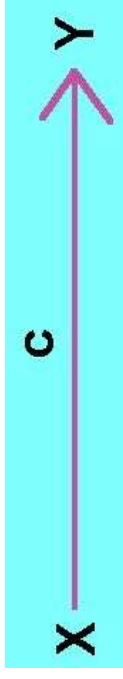
$$Y = \beta_{0c} + cX + \epsilon_c$$

$$M = \beta_{0a} + aX + \epsilon_a$$

$$Y = \beta_{0c'b} + c'X + bM + \epsilon_{c'b}$$

Effect size: often reported as  $c - c'$ .

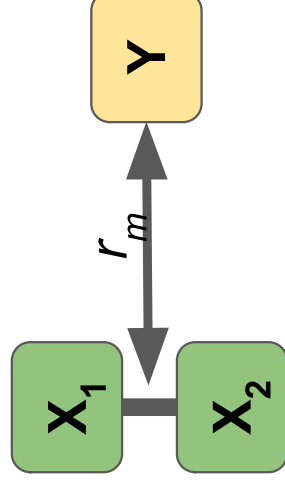
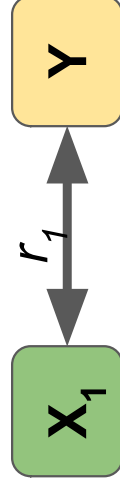
Used for \*basic\* causal inference.



(Kenney, 2015)

<http://davidakenny.net/cm/mediate.htm>

# Moderation (interaction)

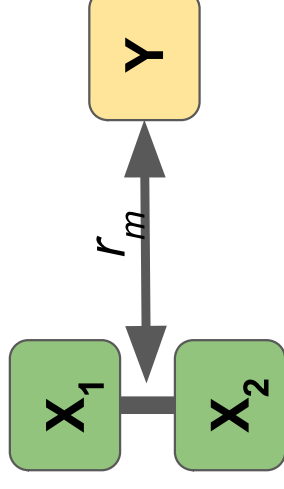
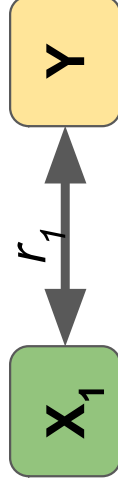


When  $r_1 \neq r_m$ ,  $X_2$  moderates the relationship between  $X_1$  and  $Y$ .

Examples:

?

# Moderation (interaction)



When  $r_1 \neq r_m$ ,  $X_2$  moderates the relationship between  $X_1$  and  $Y$ .

Examples:

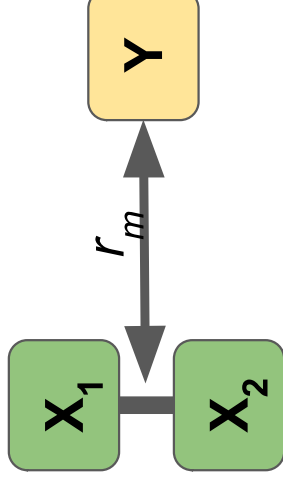
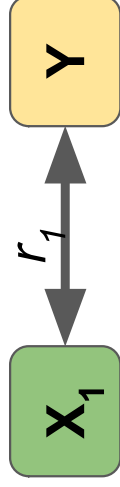
Y: Attend church?     $X_1$ : Agreeableness,     $X_2$ : From US?

Movie Reviews:

Y: Rated Depressing,     $X_1$ : "death" in review,     $X_2$ : Silly Horror Movie?



# Moderation (interaction)



When  $r_1 \neq r_m$ ,  $X_2$  moderates the relationship between  $X_1$  and  $Y$ .

More precisely moderation analyses fit the model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_M X_1 X_2 + \beta_2 X_2 + \epsilon$$

(Element-wise multiplication)

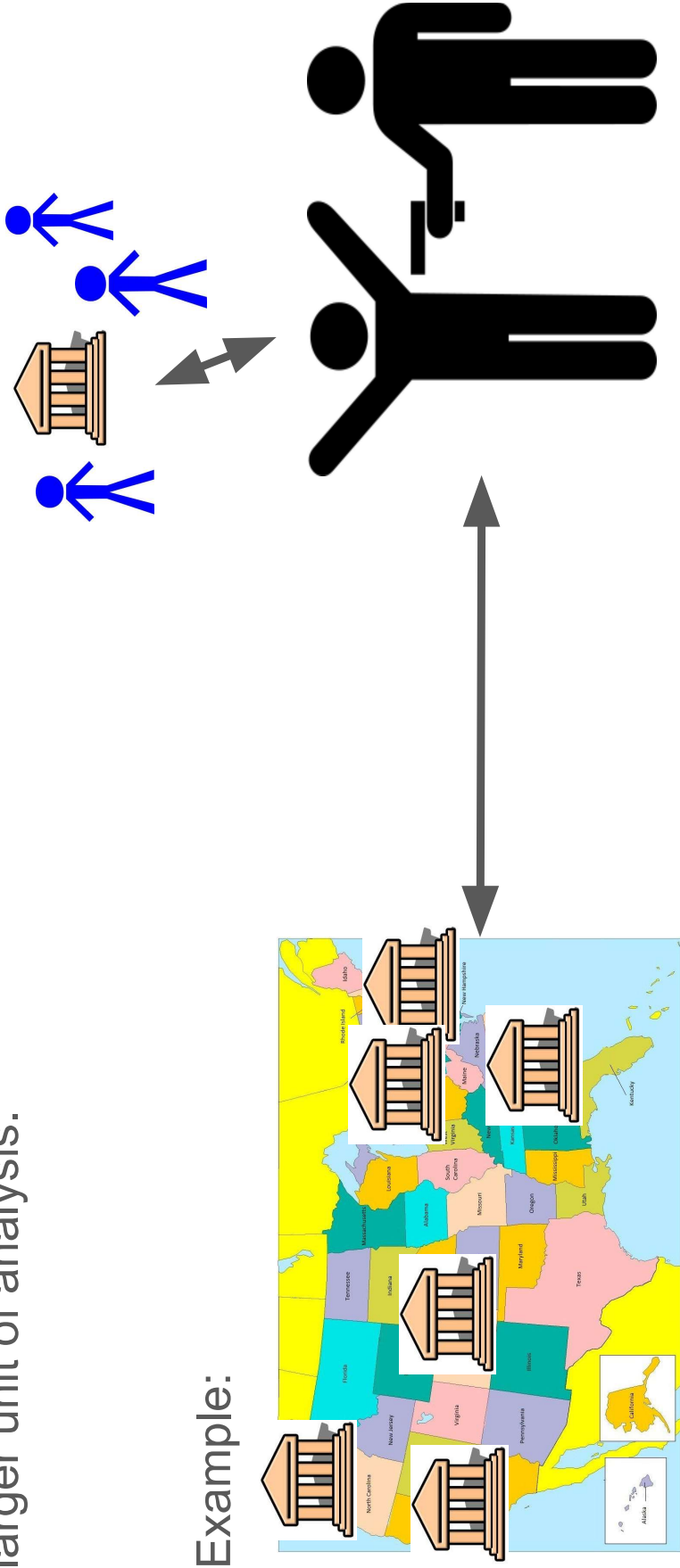
$X_1 X_2$ : The interaction term.

$\beta_M$  can then be tested for significance using the same t-test we use for any individual coefficient in multiple linear regression

# Ecological Fallacy

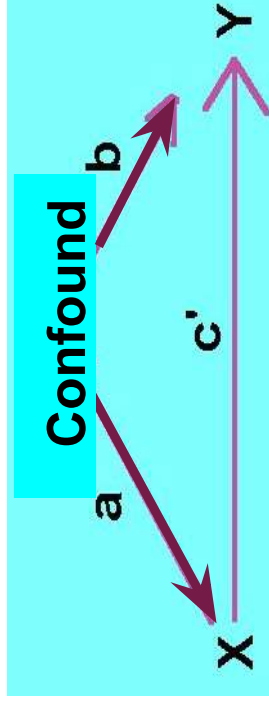
The assumption that an effect at one unit of analysis will hold for a smaller or larger unit of analysis.

Example:



# Multi-Level Models

Problem: Sometimes variables at one unit of analysis are *confounded* by a variable at another level.



# Multi-Level Models

Problem: Sometimes variables at one unit of analysis are *confounded* by a variable at another level.

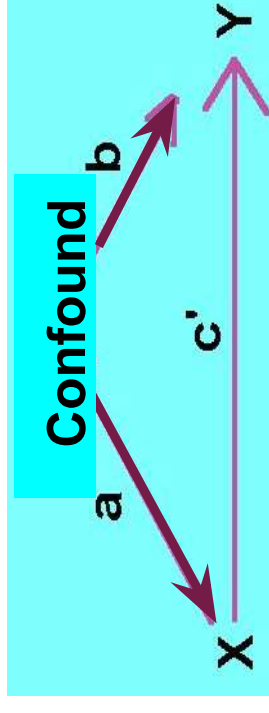
Examples:

Pot heads are more likely to say “hella”  
but really californians are more like to say “hella” and be potheads.

X = use of “hella”

Y = pot-head or not

Confound = from california?



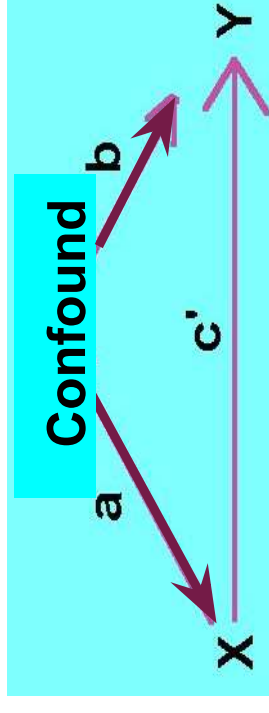
# Multi-Level Models

Problem: Sometimes variables at one unit of analysis are *confounded* by a variable at another level.

Examples:

Pot heads are more likely to say “hella”  
but really californians are more like to say “hella” and be potheads.

Females are more likely to post pictures of food  
but really both food posts and females are more common on Pinterest.



# Multi-Level Models

Problem: Sometimes variables at one unit of analysis are *confounded* by a variable at another level.

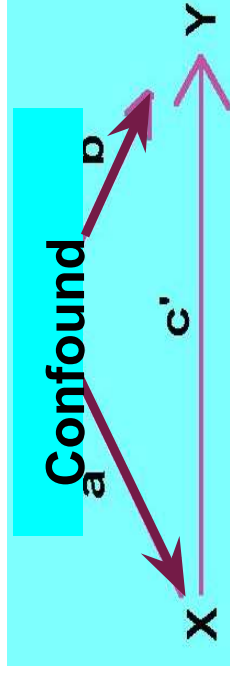
Examples:

Pot heads are more likely to say “hella”  
but really californians are more like to say “hella” and be potheads.

Females are more likely to post pictures of food  
but really both food posts and females are more common on Pinterest.

Solution: include aggregate confounding variable as a covariate in multiple linear regression. (also useful for prediction)

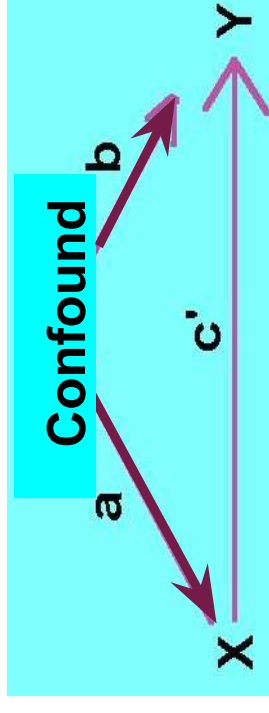
$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{m1} + \epsilon_i$$



# Multi-Level Models

Problem: Sometimes variables at one unit of analysis are *confounded* by a variable at another level.

$$Y = \beta_0 + \beta_1 X_1 + \beta_A A + \epsilon$$



A : aggregate indicator variable (is in region or not? Pinterest usage).

Solution: include aggregate confounding variable as a covariate in multiple linear regression. (also useful for prediction)

# Fixed-Effects Model (a multi-level model useful for time data)

$$y_{it} - \bar{y}_i = (X_{it} - \bar{X}_i) \beta + (\alpha_i - \bar{\alpha}_i) + (u_{it} - \bar{u}_i) \implies \ddot{y}_{it} = \ddot{X}_{it} \beta + \ddot{u}_{it}$$

$$\bar{y}_i = \frac{1}{T} \sum_{t=1}^T y_{it} \quad \bar{X}_i = \frac{1}{T} \sum_{t=1}^T X_{it} \quad \bar{u}_i = \frac{1}{T} \sum_{t=1}^T u_{it}$$

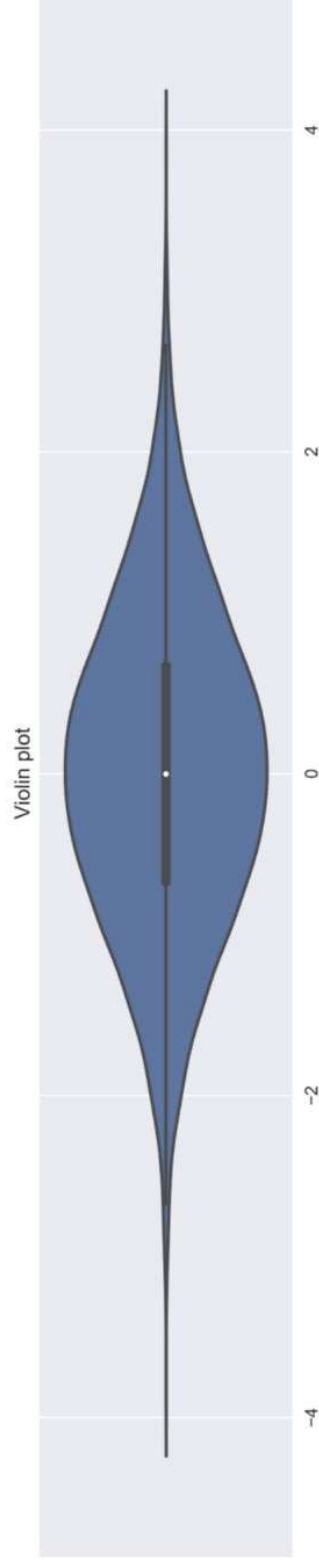
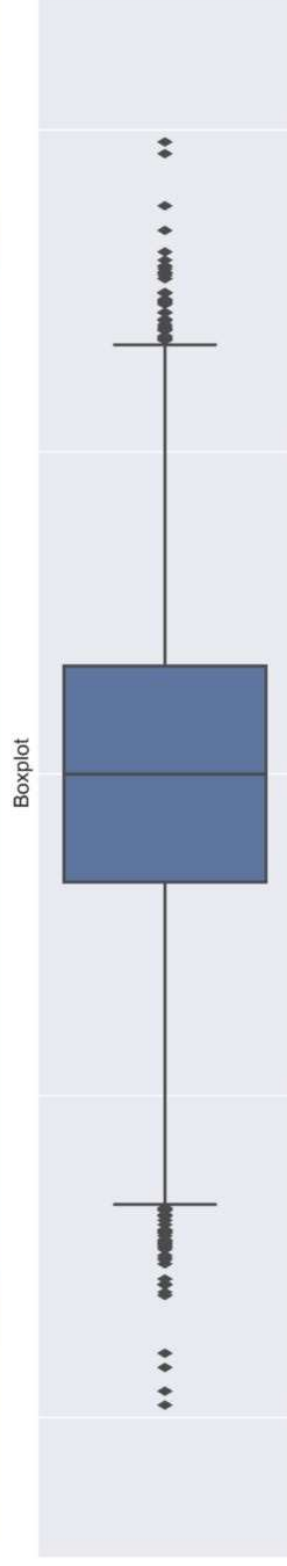
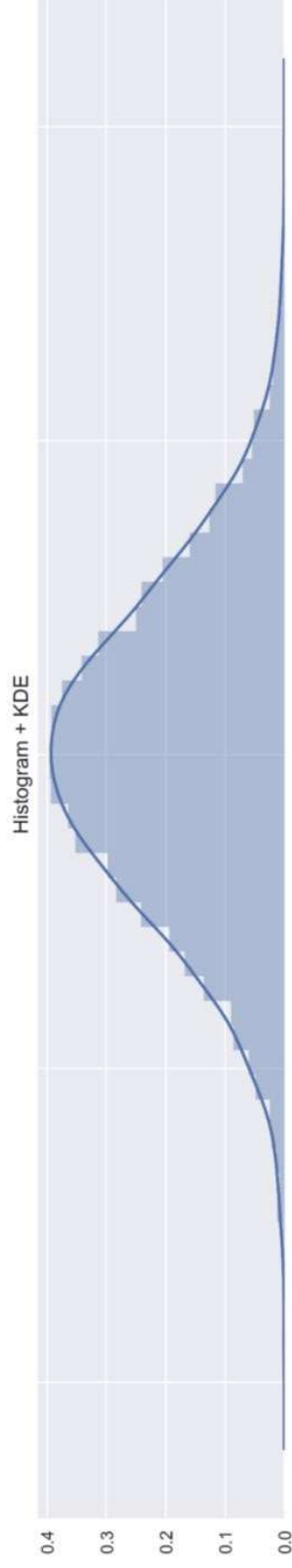
$$y_{it} = X_{it} \beta + \alpha_i + u_{it}$$

- $y_{it}$  is the dependent variable observed for individual  $i$  at time  $t$ .
- $X_{it}$  is the time-variant  $1 \times k$  (the number of independent variables) regressor vector.
- $\beta$  is the  $k \times 1$  matrix of parameters.
- $\alpha_i$  is the unobserved time-invariant individual effect. For example, the innate ability for individuals or historical and institutional factors for countries.
- $u_{it}$  is the error term.

Unlike  $X_{it}$ ,  $\alpha_i$  cannot be directly observed.



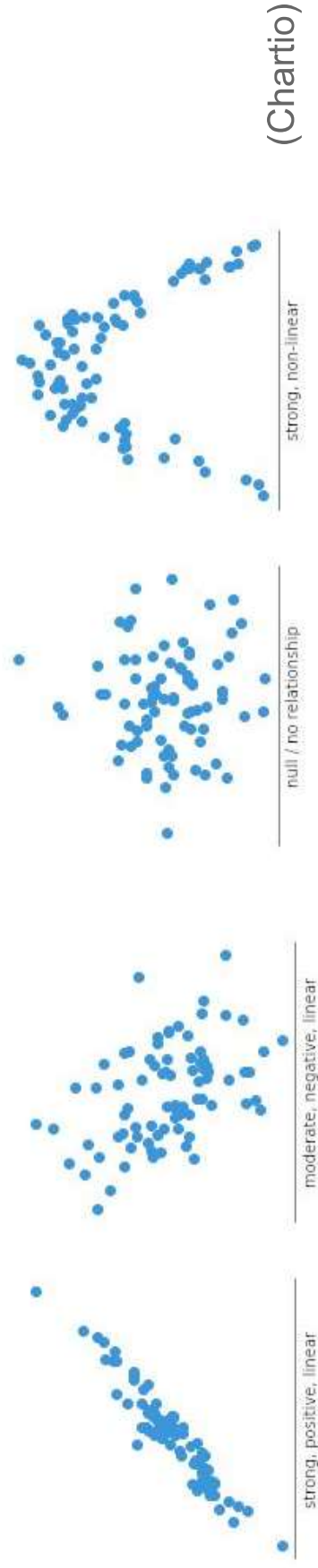
# Useful Plots: For distributions



([Lewinson, 2019](#))

# Useful Plots: Correlation

**Scatter Plot:** for two variables expected to be associated (with optional regression line)

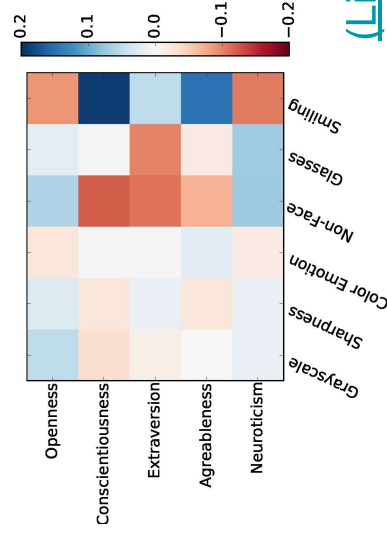


**Correlation Matrix:** for comparing associations between many variables (use Bonferroni correction if hyp testing)

|    | FriendSize | Intelligence Quotient | Income | Sat W/ Life | Depression |
|----|------------|-----------------------|--------|-------------|------------|
| F1 | 0.03       | 0.04                  | 0.12   | 0.02        | -0.1       |
| F2 | 0.04       | -0.26                 | -0.19  | -0.09       | 0.11       |
| F3 | -0.07      | -0.13                 | 0.02   | -0.02       | -0.02      |
| F4 | -0.03      | 0.27                  | -0.08  | -0.12       | 0.11       |
| F5 | -0.01      | 0.23                  | 0.29   | 0.07        | -0.21      |

**Fig 3. Individual factor correlations with outcomes.** Note how F4 which captures the use of swear words negatively correlates with Satisfaction with Life (SWL).

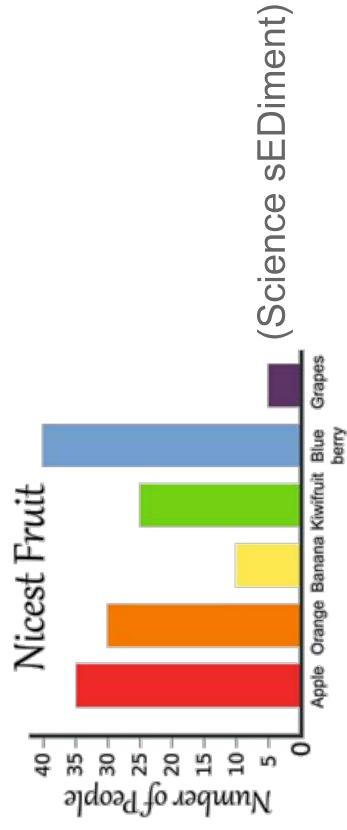
<https://doi.org/10.1371/journal.pone.0201703.g003>



(Liu et al., 2016)

# Useful Plots: Any Values

**Bar Plot:** To visually compare values under multiple conditions.



**Line Plot:** When one variable has a natural ordering (e.g. time)



(plot source: NYT U.S. Coronavirus Data)

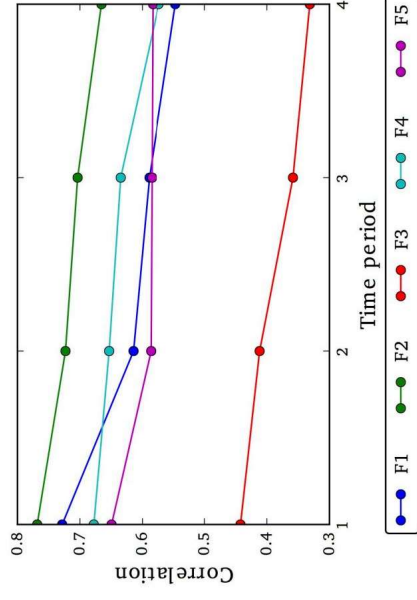
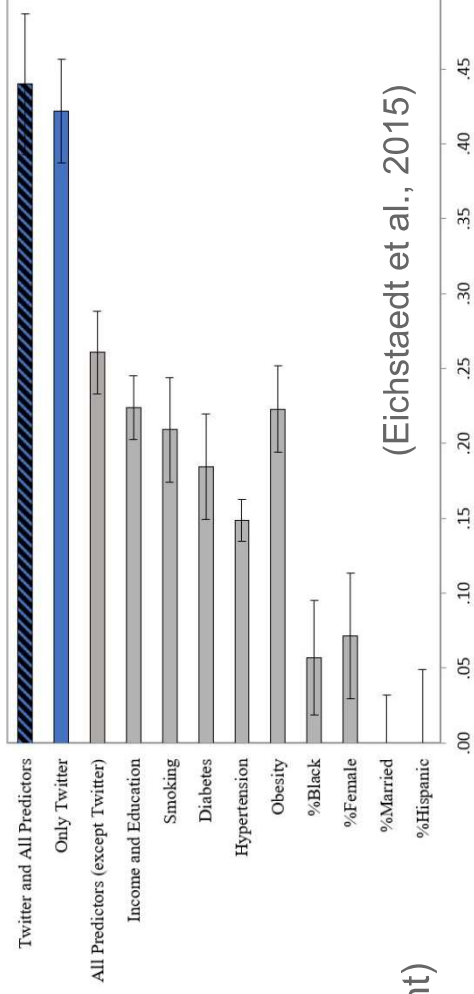
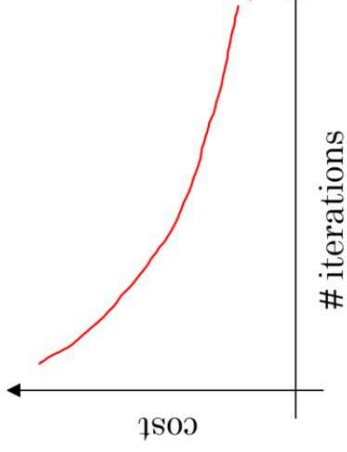


Fig 6. Test re-test validity of our learned factors.  
<https://doi.org/10.1371/journal.pone.0201703.g006>

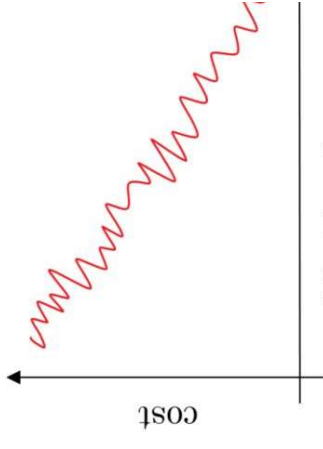
# Useful Plots: Prediction

**Learning Curve:** for plotting error from gradient descent.

for a model with  
convex optimization  
(i.e. linear regression)

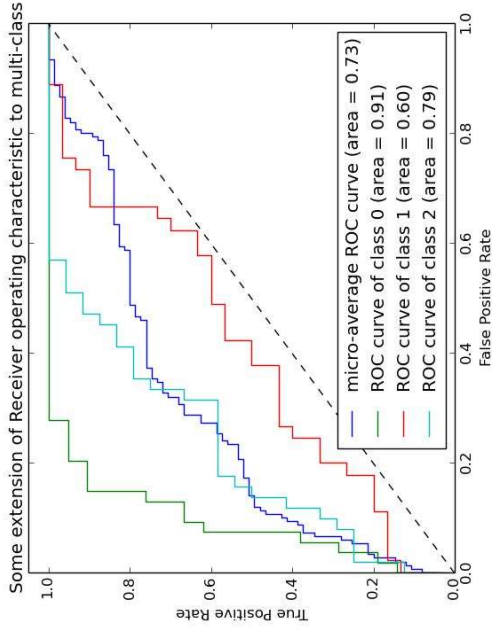


for a model with  
non-convex  
optimization (i.e.  
most deep learning)

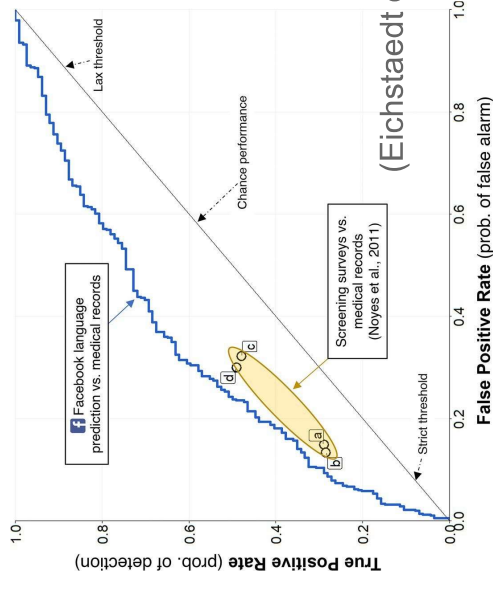


([Dabura, 2017](#))

**ROC Plot:** for visualizing true-positive to false-positive rates (used for AUC metric)



([PLOT ROC](#))



([Eichstaedt et al., 2018](#))