# GARBAGE IN, GARBAGE OUT

Steven Hyde, Eric Bachura, & Joseph Harrison

#### ML Basics

- Unsupervised Exploratory, Inductive
- Supervised Confirmatory, Deductive





#### Traditional Statistics vs Machine Learning

- <u>Statistical Inference</u>
- Small sample size
  - More assumption
- Lots of human effort
- Abstract constructs
- Messy real world data
- Theory driven

- Predictive Accuracy
- Large sample sizes
  - Very few assumption
- No human effort
- Simple constructs
- Tightly controlled data
- Data driven



### For Example:

 Examines the impact of CEO narcissism on firm strategy and performance (Chatterjee & Hambrick, 2007). Measured through the prominence of CEO photograph; CEO prominence in company press releases; proportion of first-person singular pronouns by CEO; measure of relative pay. N=111

#### Vs

• Identify hand-written numbers, with a training set in the tens of thousands (Qiao et al., 2018)



#### Black Box Problem

 "Theories without methodological implications are likely to be little more than idle speculation with minimal empirical import. And methods without theoretical substance can be sterile, representing technical sophistication in isolation." Van Maanen and colleagues (2007: 1146) state,

#### LATEST TRICKS

Rotating objects in an image confuses DNNs, probably because they are too different from the types of image used to train the network.



Even natural images can fool a DNN, because it might focus on the picture's colour, texture or background rather than picking out the salient features a human would recognize.

#### Manhole cover



Pretzel



(Heaven 2019) nature

### Phase 1: Pre Implementation Data Handling

- Is ML needed?
- Establish ground truth
  - Data source aligns with construct at right level
  - Training data validated
  - Sufficient variance
  - Data is large enough
- Theoretical pre-pruning



### Theoretical pre-pruning

- Traditional approach- kitchen sink first, prune based on data
- Put theory first
- Don't give the algorithm the "background color" if you don't want it to use it.
  - Be as inclusive as possible, but only with features for which prior theory or logic would suggest there is a conceptual link





# Phase 2: ML Implementation

- Select appropriate ML technique for construct
  - Binary, categorical, or continuous
  - Algorithm works with construct
    - Horse race if unknown

# Phase 3: Post-Implementation

- Report everything
  - Hyperparameters too
- Refine your model



# Demonstration: Need for Affiliation

- McClelland's need are subconscious trait-like motive
- Affiliation is the need to establishing and maintaining relationships (McClelland, 1987; McClelland and Boyatzis, 1982)
- Thematic Apperception Test (TAT)

### Training data

- Quarterly earnings calls of the SP 500, 2008-2016
  - 700 randomly selected transcripts and 364 from the manual
  - 2 independent coders used Winter's (1994) scoring manual
  - Training set (80%) and hold out (20%)



#### Feature selection – Pre pruning

- LIWC, over 90 categories, representing on average 86%
- Pruned model:

#### Motivation

(drives, affiliation, achievement, power, reward, risk, authentic, comparisons, number, and quantifiers)

#### Attainment

Goal

(tone, affect words, positive emotions, negative emotions, anxiety, anger, sadness, swear words, discrepancy, tentative, certainty, and clout)

#### Relations to others

(personal pronouns, social, family, friend, work, leisure, home, sexual, prepositions, assent, and total words count)

### ML Techniques

- Scikit-learn, regression task
- Horse race: NN\*, RF, SVM, & LIWC

Туре	All Features	Pruned Features
NN	Structure: Multi-layer perceptron, 2 hidden layers with 36 nodes each Solver: Stochastic gradient descent Activation: Hyperbolic tangent function Maximum iterations: 10000	<ul> <li>Structure: Multi-layer perceptron,</li> <li>2 hidden layers with 18 nodes</li> <li>each</li> <li>Solver: Optimized gradient</li> <li>descent</li> <li>Activation: Rectified linear unit</li> <li>Maximum iterations: 10000</li> </ul>
RF	Structure: 1000 trees Decision criterion: MAE Features considered on splitting: All available	Structure: 1000 trees Decision Criterion: MSE Features considered on splitting: Pre-pruned features
SVM	Kernel: Polynomial Regularization value: 1	Kernel: Polynomial Regularization value: 10

#### Results

Model	R <sup>2</sup>	RMSE	MSE	MAE	
Average Performance Across ML Models					
Pre-pruned LIWC categories (Pruned)	0.879	0.590	0.435	0.251	
All LIWC categories (All)	0.838	0.671	0.551	0.268	
Affiliation dictionary (LIWC)	0.005	5.230	27.350	4.315	
Performance of each ML Model					
NN-Pruned	0.955	0.194	0.038	0.124	
NN-All	0.915	0.266	0.071	0.147	
RF-Pruned	0.901	0.288	0.083	0.151	
RF-All	0.886	0.309	0.095	0.176	
SVM-Pruned	0.724	0.480	0.230	0.084	
SVM-All	0.550	0.613	0.376	0.143	

#### Findings

- Pre-pruned models consistently outperformed unpruned models: mean +9.7% in R<sup>2</sup>, -19% RMSE, -35.2% MSE, & -23% MAE
  - LWIC was a poor measure (R<sup>2</sup> = .005, RMSE = 5.230, MSE = 27.350, MAE = 4.315)
- Order of accuracy: NN, RF, and SMV.
- Pre-Pruning more beneficial for SVM and NN, least for RF.







#### Post Implementation: What's next?

• Continued refinement of the model

#### Phase I: Pre-implementation Data Handling

Step	Considerations
1: Establish the relevance of ML for the	<ul> <li>Will the application of ML techniques improve</li> </ul>
focal construct	measurement relative to simpler techniques?
2: Establish the "ground truth" of the	Does the data source align with the construct?
training data	Has the training data been validated?
	Is there sufficient variance in the training data?
	<ul> <li>Is the training sample sufficiently large to designate a</li> </ul>
	holdout sample?
3: Apply theoretical pre-pruning	<ul> <li>What available features are conceptually linked to the</li> </ul>
	construct?

Phase II: ML Implementation				
Step	Considerations			
4: Select appropriate ML techniques for the focal construct	<ul> <li>Is the construct binary, categorical, or continuous?</li> <li>Which ML algorithm most accurately predicts the construct?</li> <li>What hyperparameters optimize the accuracy of the model?</li> </ul>			

Phase III: Post-implementation Reporting and Refinement



#### Questions

