

Methods Bites

Machine Learning Workshop

Iasmin Goes & Daniel Weitzel

Department of Political Science

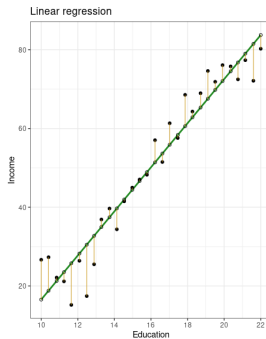
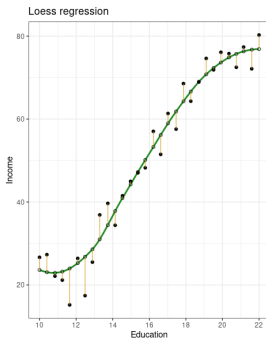
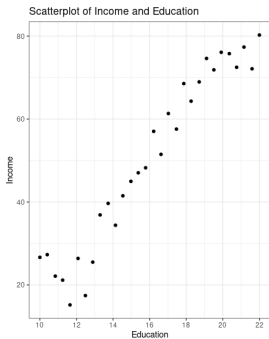
Colorado State University

25 April 2023

Why we analyze data

- ▶ **Goal:** understand the *systematic* relationship between an outcome (dependent variable) and predictors (independent variables)
- ▶ We want to know the function f in $Y = f(X) + \epsilon$
 - ▶ f is the unknown form with which X provides systematic information about Y

Why we analyze data



Difference between “traditional” statistics and ML

- ▶ Goal
 - ▶ **TS:** test hypotheses, make inferences about population parameters
 - ▶ **ML:** build predictive models for new data
- ▶ Approach
 - ▶ **TS:** pre-defined assumptions about the population distribution; these assumptions, in turn, inform model selection
 - ▶ **ML:** no assumptions about distribution or functional form; models learn patterns and relationships from the data
- ▶ Data
 - ▶ **TS:** good for smaller datasets that are a representative sample of a larger population
 - ▶ **ML:** good for larger and more complex data sets, also unstructured data (images, text, audio)

Difference between “traditional” statistics and ML

- ▶ Interpretability
 - ▶ **TS:** a set of pre-defined assumptions allows for higher interpretability; results are easier to understand
 - ▶ **ML:** can be “black boxes”, models can be difficult to interpret
- ▶ Training
 - ▶ **TS:** usually no split between training and test data: the analysis is done on one full dataset
 - ▶ **ML:** the dataset is split into training, (cross-)validation, and test set to prevent overfitting of models to idiosyncratic features of the data

When ML makes sense

- ▶ Large data sets (number of observations)
- ▶ Large number of predictors and/or no theory about f
- ▶ Accurate predictions are more valuable than causal inference
- ▶ Complex non-linearity in the data
- ▶ Unstructured data
- ▶ Goal is feature generation (making new variables)

Types of ML

▶ **Classification (categorical outcome)**

- ▶ Naive Bayes
- ▶ Clustering algorithms (kNN)
- ▶ Logistic regression
- ▶ Random forest
- ▶ Gradient boosting machine
- ▶ Support vector machine
- ▶ Neural networks

▶ **Regression (continuous outcome)**

- ▶ Lasso, ridge, and linear regression
- ▶ Random forest
- ▶ Gradient boosting machine
- ▶ Support vector machine
- ▶ Neural networks

Types of ML

- ▶ **Supervised learning:** the computer is trained on a labeled dataset and learns to make predictions or classifications based on new data
- ▶ **Unsupervised learning:** the computer is given an unlabeled dataset and is tasked with finding patterns or relationships within the data
- ▶ **Reinforcement learning:** the computer learns to make decisions based on a reward signal, which is provided when it performs an action that leads to a positive outcome

The Fingerprints of Fraud: Evidence from Mexico's 1988 Presidential Election

FRANCISCO CANTÚ *University of Houston*

This paper investigates the opportunities for non-democratic regimes to rely on fraud by documenting the alteration of vote tallies during the 1988 presidential election in Mexico. In particular, I study how the alteration of vote returns came after an electoral reform that centralized the vote-counting process. Using an original image database of the vote-tally sheets for that election and applying Convolutional Neural Networks (CNN) to analyze the sheets, I find evidence of blatant alterations in about a third of the tallies in the country. This empirical analysis shows that altered tallies were more prevalent in polling stations where the opposition was not present and in states controlled by governors with grassroots experience of managing the electoral operation. This research has implications for understanding the ways in which autocrats control elections as well as for introducing a new methodology to audit the integrity of vote tallies.

Supervised learning

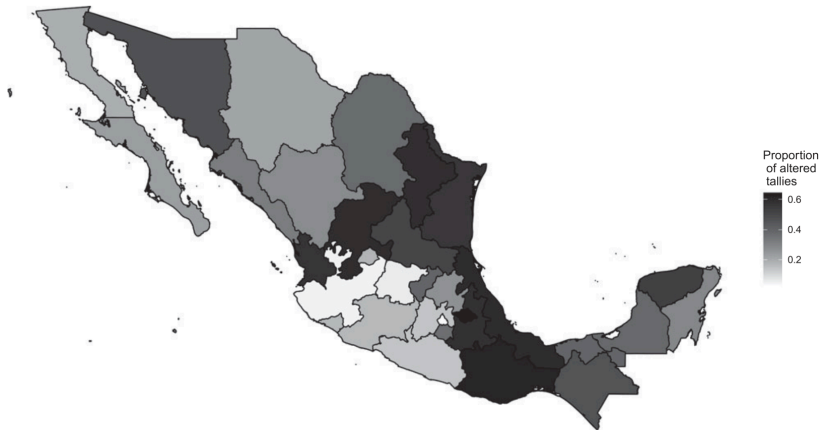
FIGURE 1. Examples of Vote Tallies with Alteration in Their Numbers. Mexico, 1988

A		
VOTACION RECIBIDA EN LA URNA (con número)	VOTOS ENCUENTRADOS EN OTRAS URNAS (con número)	(con número)
131	131	
97	7	
128	138	
00		
128	138	

B		
VOTACION RECIBIDA EN LA URNA (con número)	VOTOS ENCUENTRADOS EN OTRAS URNAS (con número)	(con número)
29		
120		
131		
1		
10		
37		
1		
22		
2		
273		
14		
287		

Supervised learning

FIGURE 3. Rates of Tallies Classified as Altered by State



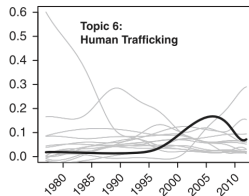
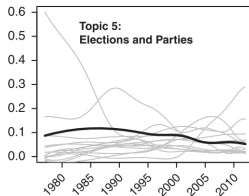
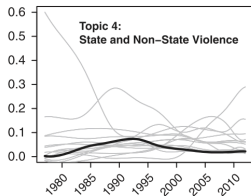
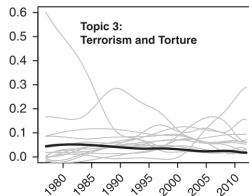
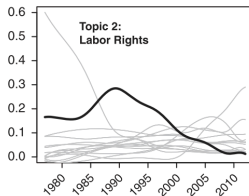
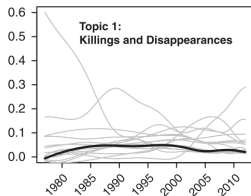
Notes: This figure shows the proportion of tallies in every state classified by the CNN as altered.

The Politics of Scrutiny in Human Rights Monitoring: Evidence from Structural Topic Models of US State Department Human Rights Reports*

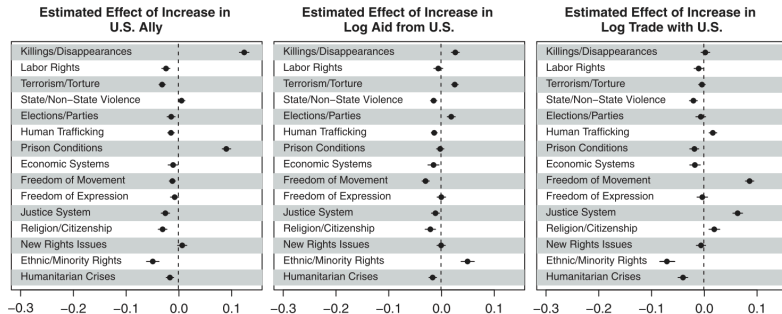
BENJAMIN E. BAGOZZI AND DANIEL BERLINER

Human rights monitoring reports play important roles both in the international human rights regime and in productions of human rights data. However, human rights reports are produced by organizations subject to formal and informal pressures that may influence the topics considered salient for attention and scrutiny. We study this potential using structural topic models (STMs), a method used for identifying the latent topical dimensions of texts and assessing the effects of covariates on these dimensions. We apply STMs to a corpus of 6298 State Department Country Reports on Human Rights Practices (1977–2012), identifying a plausible set of topics including killings and disappearances, freedoms of expression and movement, and labor rights, among others. We find that these topics vary markedly both over time and space. We also find that while US domestic politics play no systematic role in shaping topic prevalence, US allies tend to receive more attention to violations of physical integrity rights. These results challenge extant research, and illustrate the usefulness of STM methods for future study of foreign policy documents. Our findings also highlight the importance of topical attention shifts in documents that monitor and evaluate countries.

Unsupervised learning



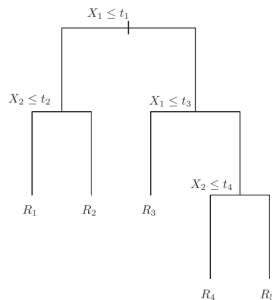
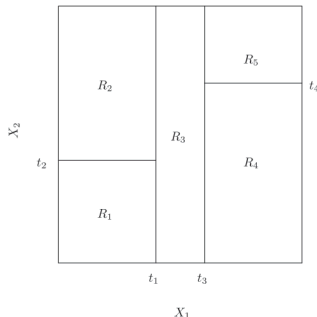
Unsupervised learning



In more detail: tree-based ML

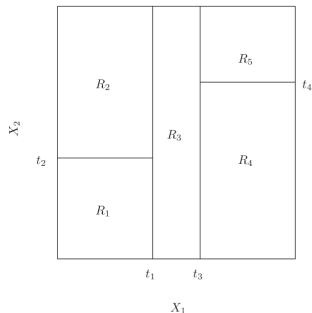
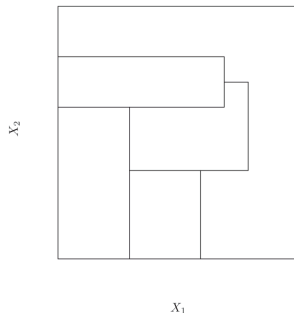
A decision tree

- **Goal:** split covariate space into regions, with each region corresponding to a unique covariate combination
 - The model then makes one prediction for all observations within this region



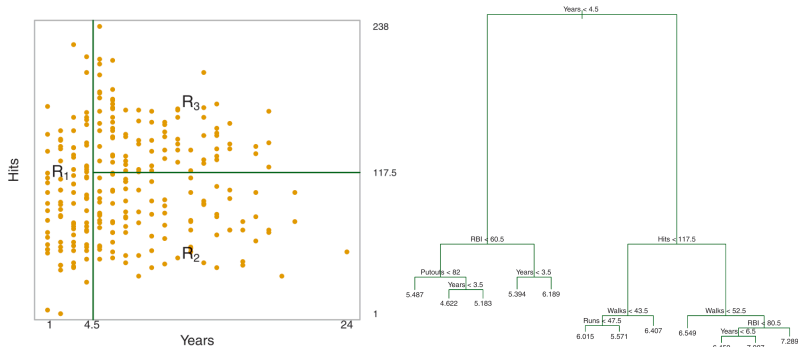
In theory, regions could have any shape

But it is computationally infeasible to consider every possible partition (left), which is why the model works with high-dimensional rectangles (right). This is called *recursive binary splitting*



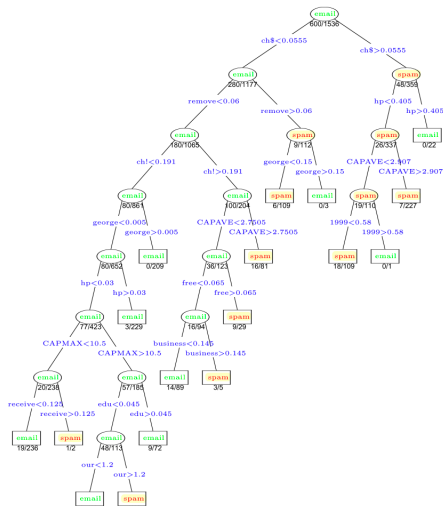
For example, the salary of baseball players

This is a *regression* tree: the outcome is quantitative



Another example: legit emails vs. spam

This is a *classification* tree: the outcome is categorical/qualitative



Random forests

- ▶ **Why forests?** Because one tree can be sensitive to data changes
- ▶ **Why random?** Because each binary split only considers a random sample of predictors
 - ▶ If there is a very strong predictor in the dataset, we don't want *all* trees to use this predictor in the first split
- ▶ The model then aggregates the results based on the predictions of most trees

Things you can explain with tree-based models

- ▶ Civil war onset (Muchlinski et al 2016)
- ▶ Supreme Court rulings (Kaufman, Kraft and Sen 2019)
- ▶ Women's legislative representation and the allocation of government expenditures (Funk, Paul and Philips 2022)
- ▶ Negative campaigning and voting (Montgomery and Olivella 2016)
- ▶ Democracy (Weitzel et al 2023a) and democratic backsliding (Weitzel et al 2023b)
- ▶ Variation in GDP data reported across different sources (Goes 2023)

Advantages of tree-based models

- ▶ No need to develop theoretical expectations
- ▶ No need to make assumptions about...
 - ▶ predictor variables
 - ▶ functional form of predictors (linear, log, squared)
- ▶ More honesty about the lack of causality
 - ▶ Regressions are not causal either, but people often interpret them causally

ML steps

1. Data cleaning and preparation

- ▶ Missing values, feature engineering (one-hot encoding, ranging, standardizing)

2. Choose a model and train it

- ▶ Select appropriate ML algorithm
- ▶ Split the data into training, (cross-)validation, and testing sets to evaluate the model's performance
- ▶ Train the algorithm on the prepared training set
- ▶ Be aware of *data leakage*

3. Evaluate the model

- ▶ Evaluate the model's performance on the (cross-)validation set
- ▶ Fine-tune it as necessary to improve accuracy

An applied example

Predicting Democracy scores

- ▶ Are democracy scores subjective or objective?
- ▶ Idea: Use a democracy score (liberal democracy) as the outcome variable and train it on a model with purely objective indicators.
- ▶ Predict democracy scores based on this model and compare to observed democracy scores.

The code

- ▶ Entire code available on Github.
- ▶ We use the V-Dem package to load V-Dem 13.
- ▶ Pre-processing steps involve cleaning the data, transforming variables, visualizations of the data.

Libraries

```
devtools::install_github("vdeminstitute/vdemdata",  
                          force = TRUE)
```

```
library("tidyverse") # For data processing  
library("vdemdata")  # The data set we will use  
library("h2o")        # the machine learning package  
library("randomizr")  # for grouped fold assignment  
library("naniar")     # Missing data visualization
```

Train set

- ▶ We reduce the data set to the id, outcome, and predictor variables.
- ▶ We also remove all rows in which the outcome (Liberal Democracy) is missing.
- ▶ We reduce the training set to years before 2012.

```
df_vdem_train <-  
  df_vdem |>  
  ungroup() |>  
  dplyr::select(all_of(ids),  
                all_of(preds),  
                all_of(outcome)) |>  
  drop_na(all_of(outcome)) |>  
  dplyr::filter(year <= 2011)
```

Cross-validation set

- ▶ We add cross-validation folds to the training data.
- ▶ There are six equally sized folds.
- ▶ Stratification is based on country.
- ▶ This allows us to train the model and validate it on data that it has never learned about.

```
df_vdem_train$folds <- cluster_ra(  
  clusters = df_vdem_train$country_id,  
  conditions = c("Fold_1", "Fold_2",  
                 "Fold_3", "Fold_4",  
                 "Fold_5", "Fold_6"))
```

Test set

- ▶ We generate a test set that we use to assess the quality of our model.
- ▶ This is data that never was part of the training process. The model does not know it.
- ▶ Our test set are all years after 2011.

```
df_vdem_predict <-  
  df_vdem |>  
  ungroup() |>  
  filter(year > 2011) |>  
  dplyr::select(all_of(ids),  
                all_of(preds),  
                all_of(outcome))
```

H2O Config

- ▶ We use H2O for the random forest
- ▶ Configuration uses n-1 cores and 20GB of RAM
- ▶ H2O works in R, Python, and Flow.
- ▶ Other packages also great!

Initialization

```
h2o.no_progress()  
h2o.init(nthreads=-1, max_mem_size = "20g")  
h2o.removeAll()
```

H2O data objects

```
train_h2o  <- as.h2o(df_vdem_train)  
test_h2o   <- as.h2o(df_vdem_predict)
```

Estimating the model

```
model_rf <-  
  h2o.randomForest(  
    model_id = "ld_1",  
    x = predictors,  
    y = outcome,  
    fold_column = "folds",  
    training_frame = train_h2o,  
    ntrees = 400,  
    mtries = 4,  
    col_sample_rate_per_tree = 0.8,  
    seed = 1904,  
    keep_cross_validation_predictions = TRUE)
```

Goodness of fit

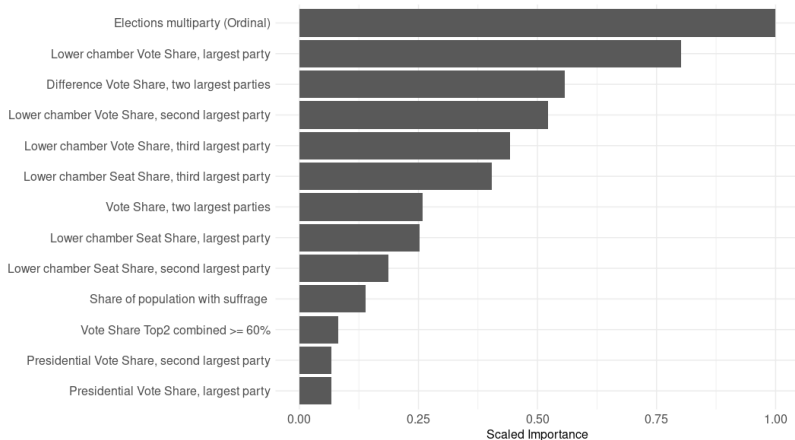
- ▶ After estimation we can look at the performance metrics in the training and cross-validation data set.
- ▶ Common are R2 and Mean Squared Error.

```
> h2o.r2(model_rf_libdem, train = TRUE, xval = TRUE)
      train      xval
0.9539659 0.7451879
> h2o.performance(model_rf_libdem, train = TRUE)
H2ORegressionMetrics: drf
** Reported on training data. **
** Metrics reported on Out-Of-Bag training samples **

MSE:  0.002712985
RMSE: 0.05208633
MAE:  0.03132654
RMSLE: 0.04023765
Mean Residual Deviance : 0.002712985
```


What matters

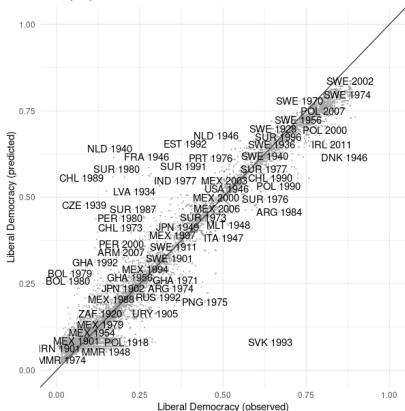
- ▶ Variable Importance Plots show us the scaled importance of individual predictors for the random forest.



The predictions

A Prediction on Training Dataset

N = 16,911, 91% Subset of available V-Dem Data



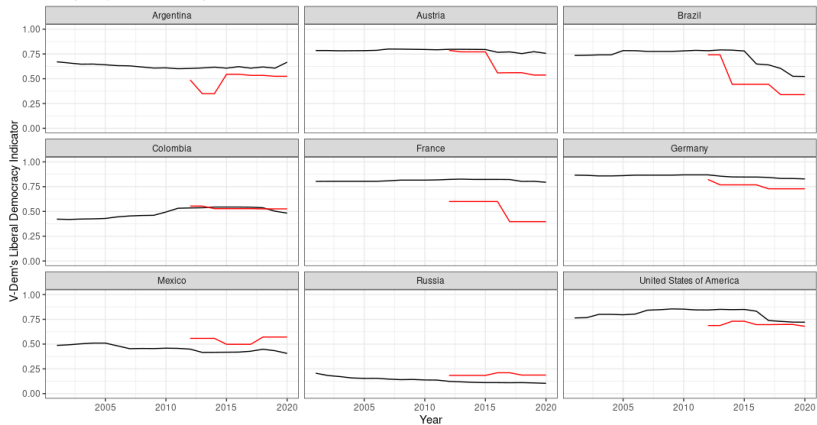
B Prediction on Test Dataset

N = 1,611, 9% Subset of available V-Dem Data



The predictions

Comparing observed and predicted values



Why is the performance so poor?

- ▶ Statistical reasons:
 - ▶ More feature engineering necessary!
 - ▶ More model tuning is necessary, we have not adjusted any of the hyperparameters of the random forest.
- ▶ Theoretical reasons:
 - ▶ Our outcome is liberal democracy and all our predictors are indicators of an electoral democracy.

Additional Resources

Books, packages, and articles

- ▶ [Introduction to Statistical Learning](#), great introduction
- ▶ [Elements of Statistical Learning](#), more advanced
- ▶ [List of R packages for Machine Learning](#), CRAN Task View List
- ▶ [Machine Learning Methods That Economists Should Know About](#), academic article with an overview