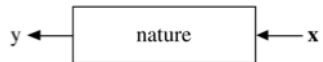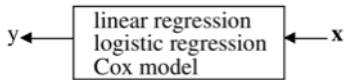# Class slides for Thursday, Sept 24: Armed conflict, part 2
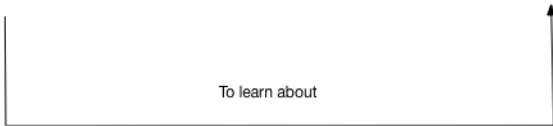
Matthew J. Salganik

COS 597E/SOC 555 Limits to prediction
Fall 2020, Princeton University

Social science (Data modeling)

## Social science (Data modeling)



## Computer science (Algorithmic modeling)

Social science (Data modeling)

y ← [ linear regression / logistic regression / Cox model ] ← x        y ← [ nature ] ← x

To learn about

Computer science (Algorithmic modeling)

y ← [ unknown ] ← x        Compare ŷ and y to learn about algorithm

decision trees
neural nets

Third way (Prediction for understanding)

y ← [ unknown ] ← x        y ← [ nature ] ← x

decision trees
neural nets

Compare ŷ and y

Study structure of model

**The long and the short of the problem**
Civil wars and internal armed conflicts, 1946–2012

Legend:
- East Asia, South-East Asia & Oceania
- South & Central Asia
- Sub-Saharan Africa
- Middle East & North Africa
- Europe
- Americas

Total deaths*‡, '000

100 DEADLIEST, RANKED BY NUMBER OF COMBATANT DEATHS*, '000

Highest → Lowest

Selected conflicts (with combatant deaths, '000):
CHINA, 800.0; VIETNAM, 198.9; CAMBODIA, 186.8; AFGHANISTAN, 334.4; ALGERIA, 91.3; SOUTH VIETNAM, 82.5; CAMEROON‡, 77.5; GREECE, 77.0; ETHIOPIA (Eritrea), 122.7; ANGOLA‡, 71.1; UGANDA, 67.4; LEBANON, 86.3; MOZAMBIQUE, 62.9; SRI LANKA, 62.2; SUDAN, 57.3; PAKISTAN (East), 59.0; ANGOLA, 39.5; ETHIOPIA, 48.4; INDONESIA (East Timor), 37.7; NIGERIA, 37.5; IRAN (Kurdistan), 34.0; SOMALIA, 31.6; MYANMAR, 25.2; EL SALVADOR, 28.7; TURKEY (Kurdistan), 27.0; CHAD, 22.1; GUATEMALA, 21.0; NICARAGUA, 20.3; INDIA (Kashmir), 19.3; PHILIPPINES (Mindanao), 22.1; PHILIPPINES, 19.5; COLOMBIA, 19.1; IRAQ, 19.4; PAKISTAN, 19.2; MYANMAR (Karen), 17.4; MYANMAR (Kachin), 18.0; RUSSIA (Chechnya), 17.8; ALGERIA, 18.3; INDONESIA, 17.3; YEMEN (North)‡, 15.0; MYANMAR (Shan), 15.0; CONGO‡, 15.5; SYRIA, 15.9; LAOS, 14.0; SUDAN (South), 13.7; ZIMBABWE, 13.0; CONGO-BRAZZAVILLE, 14.2; MOZAMBIQUE, 11.7; PERU, 11.6; CHINA (Tibet), 10.0; ETHIOPIA (Ogaden), 10.8; SIERRA LEONE‡, 10.0; NEPAL, 10.0; RWANDA‡, 9.2; ISRAEL (Palestinian Territories), 9.1; SOVIET UNION (Ukraine), 8.0; INDIA (Punjab/Khalistan), 7.6; BURUNDI, 8.6; KENYA, 6.5; GUINEA-BISSAU, 6.7; SOUTH AFRICA (Natal), 6.7; TAJIKISTAN, 4.8; INDIA, 6.3; MOROCCO (Western Sahara), 6.4; MALAYSIA, 5.4; CAMEROON, 5.0; YEMEN (South), 6.7; INDIA, 4.8; PHILIPPINES, 4.6; INDONESIA (West Papua), 4.8; AZERBAIJAN‡, 4.8; MADAGASCAR, 4.6; SOVIET UNION (Lithuania), 4.3; PAKISTAN (Baluchistan), 4.3; SYRIA, 3.8; SERBIA (Croatia), 3.9; YEMEN‡, 3.7; MYANMAR (Arakan), 3.4; CONGO, 3.4; IRAN, 3.7; BOSNIA (Croatia)‡, 3.6; IRAN (Kurdistan), 3.3; ETHIOPIA (Oromo), 2.8; HYDERABAD, 3.2; SERBIA (Kosovo)‡, 2.6; LIBERIA, 2.6; RUSSIA (Caucasus Emirate), 2.4; PARAGUAY, 2.0; THAILAND, 2.5; ARGENTINA, 2.3; SERBIA (Kosovo)‡, 2.6; INDONESIA (Aceh), 2.3; LIBYA, 2.3; INDIA (Nagaland), 2.2; SOUTH AFRICA, 2.2; GEORGIA, 2.2; YEMEN (North), 2.0; DOMINICAN REPUBLIC, 1.5; INDONESIA, 1.8; UNITED KINGDOM (Northern Ireland), 1.5

Sources: PRIO; Uppsala University

*Based on over 250 conflicts, 1946–2012   †Deaths in battle of government troops and troops of politically organised rebels; conflicts restarted within 10 years counted as continuous.   ‡Including foreign intervention

https://www.economist.com/content/inner-turmoil

# Comparing Random Forest with Logistic Regression for Predicting Class-Imbalanced Civil War Onset Data

**David Muchlinski**

*School of Social and Political Science, University of Glasgow, Glasgow, UK*
*e-mail: david.muchlinski@glasgow.ac.uk (corresponding author)*

**David Siroky**

*Department of Political Science, Arizona State University, Tempe, AZ*
*e-mail: david.siroky@asu.edu*

**Jingrui He**

*Department of Computer Science and Engineering, Arizona State University, Tempe, AZ*
*e-mail: jingrui.he@asu.edu*

**Matthew Kocher**

*Department of Political Science, Yale University, New Haven, CT*
*e-mail: mathew.kocher@yale.edu*

▶ Two goals: (1) compare random forest to logistic regression for predicting civil war onset (2) learn from random forest about civil war onset.

▶ Two goals: (1) compare random forest to logistic regression for predicting civil war onset (2) learn from random forest about civil war onset.

▶ Regarding goal (1):

   ▶ This feels like a weird hybrid. If you are going to argue for prediction, why not go all in? (Schelling wind tunnel story)

- ▶ Two goals: (1) compare random forest to logistic regression for predicting civil war onset (2) learn from random forest about civil war onset.
- ▶ Regarding goal (1):
    - ▶ This feels like a weird hybrid. If you are going to argue for prediction, why not go all in? (Schelling wind tunnel story)
    - ▶ I find the comparison between random forest and logistics regression misleading because two things are varying: number of predictors and learning algorithm.

▶ Two goals: (1) compare random forest to logistic regression for predicting civil war onset (2) learn from random forest about civil war onset.

▶ Regarding goal (1):

   ▶ This feels like a weird hybrid. If you are going to argue for prediction, why not go all in? (Schelling wind tunnel story)

   ▶ I find the comparison between random forest and logistics regression misleading because two things are varying: number of predictors and learning algorithm. Better way might be to keep the predictor set fixed.

- ▶ Two goals: (1) compare random forest to logistic regression for predicting civil war onset (2) learn from random forest about civil war onset.
- ▶ Regarding goal (1):
  - ▶ This feels like a weird hybrid. If you are going to argue for prediction, why not go all in? (Schelling wind tunnel story)
  - ▶ I find the comparison between random forest and logistics regression misleading because two things are varying: number of predictors and learning algorithm. Better way might be to keep the predictor set fixed.
  - ▶ A lot of work has already happened in both feature selection and feature engineering (e.g., GDP growth and GDP per capita) before random forest is ever used

- ▶ Two goals: (1) compare random forest to logistic regression for predicting civil war onset (2) learn from random forest about civil war onset.
- ▶ Regarding goal (1):
    - ▶ This feels like a weird hybrid. If you are going to argue for prediction, why not go all in? (Schelling wind tunnel story)
    - ▶ I find the comparison between random forest and logistics regression misleading because two things are varying: number of predictors and learning algorithm. Better way might be to keep the predictor set fixed.
    - ▶ A lot of work has already happened in both feature selection and feature engineering (e.g., GDP growth and GDP per capita) before random forest is ever used
    - ▶ They did not model the panel data structure (e.g., no country fixed effects). Why? If the goal is to help policy makers, then we should allow country fixed effects.

- ► Two goals: (1) compare random forest to logistic regression for predicting civil war onset (2) learn from random forest about civil war onset.
- ► Regarding goal (1):
    - ► This feels like a weird hybrid. If you are going to argue for prediction, why not go all in? (Schelling wind tunnel story)
    - ► I find the comparison between random forest and logistics regression misleading because two things are varying: number of predictors and learning algorithm. Better way might be to keep the predictor set fixed.
    - ► A lot of work has already happened in both feature selection and feature engineering (e.g., GDP growth and GDP per capita) before random forest is ever used
    - ► They did not model the panel data structure (e.g., no country fixed effects). Why? If the goal is to help policy makers, then we should allow country fixed effects.
    - ► Starts like a paper about onset of civil wars, ends like a paper about random forest

- ▶ Two goals: (1) compare random forest to logistic regression for predicting civil war onset (2) learn from random forest about civil war onset.
- ▶ Regarding goal (1):
  - ▶ This feels like a weird hybrid. If you are going to argue for prediction, why not go all in? (Schelling wind tunnel story)
  - ▶ I find the comparison between random forest and logistics regression misleading because two things are varying: number of predictors and learning algorithm. Better way might be to keep the predictor set fixed.
  - ▶ A lot of work has already happened in both feature selection and feature engineering (e.g., GDP growth and GDP per capita) before random forest is ever used
  - ▶ They did not model the panel data structure (e.g., no country fixed effects). Why? If the goal is to help policy makers, then we should allow country fixed effects.
  - ▶ Starts like a paper about onset of civil wars, ends like a paper about random forest
  - ▶ Notice that models did not maximize scoring function (AUC)

**Fig. 4** Plot of variable importance by mean decrease in Gini Score.



**Fig. 5** Partial dependence plots.

Fig. 4 Plot of variable importance by mean decrease in Gini Score.



Fig. 5 Partial dependence plots.

▶ I'm not sure I believe this. What are the units?

**Variable Importance for Random Forests**

GDP Growth
GDP per Capita
Life Expectancy
Western Europe and US Dummy
Infant Mortality
Trade as Percent of GDP
Mountainous Terrain
Illiteracy Rate
Population (logged)
Linguistic Hetrogeneity
Anocracy
Median Regional Polity Score
Primary Commodity Exports (Squared)
Democracy
Military Power
Population Density
Political Instability
Ethnic Fractionalization
Secondary Education
Primary Commodity Exports

Mean Decrease in Gini Score (OOB Estimates)

**Fig. 4** Plot of variable importance by mean decrease in Gini Score.



**Fig. 5** Partial dependence plots.

▶ I'm not sure I believe this. What are the units?
▶ Recall Arvind's concerns about correlated predictors

**Fig. 4** Plot of variable importance by mean decrease in Gini Score.

**Fig. 5** Partial dependence plots.

- ▶ I'm not sure I believe this. What are the units?
- ▶ Recall Arvind's concerns about correlated predictors
- ▶ If we did this with a slightly different implementation of random forest would we get the same results? What about another ML method?

Assignment 2

- It is very hard to tell what is happening in the comments and replies without digging in, hence your assignment. But imagine what we happen if the code was not available!

- It is very hard to tell what is happening in the comments and replies without digging in, hence your assignment. But imagine what we happen if the code was not available!
- Can we use predictive models to learn about the data generating process?

- ▶ It is very hard to tell what is happening in the comments and replies without digging in, hence your assignment. But imagine what we happen if the code was not available!
- ▶ Can we use predictive models to learn about the data generating process?
- ▶ What is a good way to measure generalization performance when there is measurement uncertainty?

# What Is Civil War?

## CONCEPTUAL AND EMPIRICAL COMPLEXITIES
## OF AN OPERATIONAL DEFINITION

NICHOLAS SAMBANIS

*Department of Political Science*
*Yale University*

The empirical literature on civil war has seen tremendous growth because of the compilation of quantitative data sets, but there is no consensus on the measurement of civil war. This increases the risk of making inferences from unstable empirical results. Without ad hoc rules to code its start and end and differentiate it from other violence, it is difficult, if not impossible, to define and measure civil war. A wide range of variation in parameter estimates makes accurate predictions of war onset difficult, and differences in empirical results are greater with respect to war continuation.

*Keywords:*   civil war; Correlates of War; data sets; coding rules

Significant differences across civil war lists are mainly due to disagreement on three questions: What threshold of violence distinguishes civil war from other forms of internal armed conflict? How do we know when a civil war starts and ends? How can we distinguish between intrastate, interstate, and extrastate wars? Answers to these questions are not only relevant for the purposes of accurate coding, but they also reveal the degree to which we share a common understanding of the concept of civil war.

12 different measures of civil wars. Does this matter? Should we standardize on one definition? Think back to ImageNet.

▶ What is a good way to measure generalization performance when the data has structure (focus on time)?

- ▶ What is a good way to measure generalization performance when the data has structure (focus on time)?
- ▶ Should we use the future to predict the past? Should we use the present to predict the present?

Go to file/function · Addins

**data ×**

Filter — Cols: `<` `<` 1 - 50 `>` `>>`

| cowcode | year | warstds | ptime | yrint | autonomy | rf | popdense | auto98 | dem98 | pol98 | army85 | milgnp92 | milper | milex | trade | logpop | ienergy | ienercap | inienc | gdpmkt | gdpcap | irgdp | decade | regy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 700 | 1945 | 0 | 12 | 0 | 0.005150819 | 35.48197 | 118.554791 | 3.963861 | 3.975114 | 0.011253747 | 129472.9 | 3.364459 | 121.0874 | 2026920 | 72.88137 | 15.39613 | 65.53742 | 5.40e-06 | -13.58918 | 69136410776 | 3576.21532 | 4546.011 | 3 | 4582.636 |
| 700 | 1946 | 0 | 24 | 1 | 0.005150819 | 31.41992 | 117.756342 | 3.964134 | 3.911079 | -0.013254829 | 129413.0 | 3.395091 | 121.8854 | 2023556 | 72.90009 | 15.41061 | 64.19439 | 5.39e-06 | -13.58218 | 67276776506 | 3543.94032 | 4521.105 | 3 | 4564.097 |
| 700 | 1947 | 0 | 36 | 2 | 0.005150819 | 31.40798 | 118.280616 | 3.973155 | 3.940862 | -0.032293758 | 130411.0 | 3.396670 | 122.7806 | 2021188 | 72.96288 | 15.41321 | 64.45685 | 5.39e-06 | -13.58279 | 67529305184 | 3545.35165 | 4517.728 | 3 | 4558.854 |
| 700 | 1948 | 0 | 48 | 3 | 0.005150819 | 31.42972 | 118.325869 | 3.982209 | 3.929385 | -0.052823053 | 126781.7 | 3.398607 | 118.2564 | 1937612 | 73.10245 | 15.40248 | 62.79688 | 5.40e-06 | -13.58111 | 65957355084 | 3538.99228 | 4505.616 | 3 | 4552.190 |
| 700 | 1949 | 0 | 60 | 4 | 0.005150819 | 34.40010 | 118.312296 | 3.973679 | 3.949860 | -0.023818451 | 130979.2 | 3.383485 | 122.2451 | 2022832 | 72.83019 | 15.41403 | 64.85889 | 5.38e-06 | -13.58636 | 68168435790 | 3554.02928 | 4528.104 | 3 | 4570.362 |
| 700 | 1950 | 0 | 72 | 5 | 0.005150819 | 34.47221 | 118.752253 | 3.963514 | 3.971818 | 0.006504856 | 130616.5 | 3.350695 | 121.5477 | 2026570 | 72.76833 | 15.40876 | 65.95401 | 5.42e-06 | -13.58743 | 69416252793 | 3583.41291 | 4546.038 | 3 | 4589.689 |
| 700 | 1951 | 0 | 84 | 6 | 0.005150819 | 34.46567 | 118.900216 | 3.973817 | 3.964128 | -0.009688565 | 129142.7 | 3.353969 | 119.6265 | 2001987 | 72.86190 | 15.40223 | 65.25755 | 5.42e-06 | -13.58528 | 68807529646 | 3580.55359 | 4549.637 | 3 | 4586.529 |
| 700 | 1952 | 0 | 96 | 7 | 0.005150819 | 31.45453 | 118.660591 | 3.972832 | 3.965914 | -0.006917568 | 129067.3 | 3.352165 | 119.6122 | 1974899 | 72.82532 | 15.40364 | 65.09238 | 5.40e-06 | -13.58923 | 68688755336 | 3574.92453 | 4547.422 | 3 | 4586.157 |
| 700 | 1953 | 0 | 108 | 8 | 0.005150819 | 33.43161 | 118.917614 | 3.967038 | 3.972325 | 0.005286929 | 130133.3 | 3.352917 | 120.3301 | 1988674 | 72.84402 | 15.40568 | 65.85418 | 5.42e-06 | -13.58931 | 69545362107 | 3589.93183 | 4556.311 | 3 | 4593.182 |
| 700 | 1954 | 0 | 120 | 9 | 0.005150819 | 33.48062 | 118.635862 | 3.962076 | 3.978004 | 0.015927615 | 131839.4 | 3.347502 | 122.4269 | 2035846 | 72.69640 | 15.41166 | 66.34346 | 5.41e-06 | -13.58497 | 70037114788 | 3587.75744 | 4559.297 | 3 | 4592.352 |
| 700 | 1955 | 0 | 132 | 10 | 0.005150819 | 33.49078 | 118.639812 | 3.968333 | 3.971692 | 0.003358665 | 129915.1 | 3.347625 | 119.9212 | 1979018 | 72.77357 | 15.40546 | 65.77326 | 5.41e-06 | -13.58784 | 69452466977 | 3585.27307 | 4553.980 | 3 | 4588.174 |
| 700 | 1956 | 0 | 144 | 11 | 0.005150819 | 35.38244 | 118.115543 | 3.970045 | 3.943897 | -0.026148613 | 129301.5 | 3.384583 | 120.0053 | 1966458 | 72.92703 | 15.41439 | 63.96239 | 5.39e-06 | -13.58264 | 67285007988 | 3549.27050 | 4520.337 | 3 | 4564.339 |
| 700 | 1957 | 0 | 156 | 12 | 0.005150819 | 35.49392 | 119.079787 | 3.964757 | 3.973885 | 0.009127687 | 130448.2 | 3.346094 | 120.8477 | 1995574 | 72.79605 | 15.40779 | 66.07430 | 5.42e-06 | -13.58623 | 69782069345 | 3587.01581 | 4554.107 | 3 | 4589.405 |
| 700 | 1958 | 0 | 168 | 13 | 0.005150819 | 35.48778 | 118.650265 | 3.964497 | 3.974992 | 0.010495344 | 131197.9 | 3.347627 | 122.6380 | 2040801 | 72.70017 | 15.41062 | 66.80191 | 5.43e-06 | -13.58330 | 70335295752 | 3590.12954 | 4558.496 | 3 | 4592.079 |
| 700 | 1959 | 0 | 180 | 14 | 0.005150819 | 35.47105 | 118.824075 | 3.968870 | 3.971809 | 0.002938786 | 131245.8 | 3.353200 | 121.8285 | 2028679 | 72.77011 | 15.40810 | 66.44923 | 5.42e-06 | -13.58810 | 70118706834 | 3590.63452 | 4557.006 | 3 | 4591.583 |
| 700 | 1960 | 0 | 192 | 15 | 0.00000000 | 1.00000 | 117.804310 | 10.000000 | 0.000000 | -10.00000000 | 129057.9 | 3.347722 | 118.8071 | 1939181 | 72.69955 | 16.11969 | 0.23300 | 0.00e+00 | -17.57641 | 121000000 | 120.80670 | 3800.417 | 1 | 2317.327 |
| 700 | 1961 | 0 | 204 | 16 | 0.00000000 | 1.00000 | 15.661860 | 10.000000 | 0.000000 | -10.00000000 | 47000.0 | 3.545582 | 60.0000 | 14240 | 11.15703 | 16.13917 | 0.26600 | 0.00e+00 | -17.46342 | 1240000000 | 120.92500 | 3800.593 | 1 | 2317.327 |
| 700 | 1962 | 0 | 216 | 17 | 0.00000000 | 1.00000 | 15.982750 | 10.000000 | 0.000000 | -10.00000000 | 47000.0 | 3.546275 | 77.0000 | 14184 | 12.55081 | 16.15945 | 0.14100 | 0.00e+00 | -17.23132 | 1230000000 | 118.01740 | 3800.673 | 1 | 2317.327 |
| 700 | 1963 | 0 | 228 | 18 | 0.00000000 | 1.00000 | 16.321039 | 10.000000 | 0.000000 | -10.00000000 | 47000.0 | 3.548764 | 90.0000 | 13125 | 14.22764 | 16.18039 | 0.34400 | 0.00e+00 | -17.24751 | 961000000 | 90.27886 | 3799.289 | 1 | 2317.327 |
| 700 | 1964 | 0 | 240 | 19 | 0.00000000 | 1.00000 | 16.675520 | 10.000000 | 0.000000 | -10.00000000 | 47000.0 | 3.339405 | 100.0000 | 14544 | 26.03551 | 16.20188 | 0.37500 | 0.00e+00 | -17.18271 | 800000000 | 73.57040 | 3798.504 | 1 | 2317.327 |
| 700 | 1965 | 0 | 252 | 20 | 0.00000000 | 1.00000 | 17.045191 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.340246 | 100.0000 | 14544 | 26.94445 | 16.22381 | 0.43000 | 0.00e+00 | -17.06778 | 1010000000 | 90.56831 | 3799.648 | 1 | 2496.514 |
| 700 | 1966 | 0 | 264 | 21 | 0.00000000 | 1.00000 | 17.429300 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.337526 | 90.0000 | 13865 | 32.67108 | 16.24609 | 0.46500 | 0.00e+00 | -17.01181 | 1400000000 | 123.18010 | 3801.726 | 1 | 2496.514 |
| 700 | 1967 | 0 | 276 | 22 | 0.00000000 | 1.00000 | 17.827311 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.343529 | 90.0000 | 16819 | 27.14286 | 16.26867 | 0.66200 | 1.00e-07 | -16.68116 | 1670000000 | 143.94250 | 3803.416 | 1 | 2496.514 |
| 700 | 1968 | 0 | 288 | 23 | 0.00000000 | 1.00000 | 18.238970 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.346049 | 75.0000 | 17932 | 20.98274 | 16.29150 | 0.70900 | 1.00e-07 | -16.63540 | 1370000000 | 115.46970 | 3801.935 | 1 | 2496.514 |
| 700 | 1969 | 0 | 300 | 24 | 0.00000000 | 1.00000 | 18.664190 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.342436 | 70.0000 | 20067 | 24.11003 | 16.31454 | 0.51300 | 0.00e+00 | -16.98202 | 1410000000 | 115.76040 | 3801.840 | 1 | 2496.514 |
| 700 | 1970 | 0 | 312 | 25 | 0.00000000 | 1.00000 | 19.103189 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.346659 | 70.0000 | 20915 | 25.07887 | 16.33779 | 0.58000 | 0.00e+00 | -16.88252 | 1750000000 | 140.39391 | 3803.689 | 2 | 2764.380 |
| 700 | 1971 | 0 | 324 | 26 | 0.00000000 | 1.00000 | 19.549879 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.351174 | 70.0000 | 19401 | 21.71281 | 16.36091 | 0.85300 | 1.00e-07 | -16.51990 | 1830000000 | 143.63580 | 3804.501 | 2 | 2764.380 |
| 700 | 1972 | 0 | 336 | 27 | 0.00000000 | 1.00000 | 20.014311 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.344355 | 83.0000 | 29844 | 27.06314 | 16.38439 | 0.64100 | 0.00e+00 | -16.82911 | 1600000000 | 122.25420 | 3802.985 | 2 | 2764.380 |
| 700 | 1973 | 0 | 348 | 28 | 0.00000000 | 1.00000 | 20.498091 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.348400 | 84.0000 | 32289 | 32.86908 | 16.40827 | 0.74200 | 1.00e-07 | -16.70668 | 1730000000 | 129.67650 | 3803.839 | 2 | 2764.380 |
| 700 | 1974 | 0 | 360 | 29 | 0.00000000 | 1.00000 | 21.002331 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.350272 | 91.0000 | 32400 | 27.69231 | 16.43257 | 0.83400 | 1.00e-07 | -16.61409 | 2160000000 | 157.39259 | 3806.147 | 2 | 2764.380 |
| 700 | 1975 | 0 | 372 | 30 | 0.00000000 | 1.00000 | 21.527700 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.349988 | 80.0000 | 34733 | 28.86598 | 16.45728 | 0.82500 | 1.00e-07 | -16.64965 | 2369999872 | 168.59000 | 3807.223 | 2 | 3053.327 |
| 700 | 1976 | 0 | 384 | 31 | 0.00000000 | 1.00000 | 22.075081 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.348045 | 130.0000 | 40555 | 26.94836 | 16.48239 | 1.06200 | 1.00e-07 | -16.42223 | 2560000000 | 177.53149 | 3808.534 | 2 | 3053.327 |
| 700 | 1977 | 0 | 396 | 32 | 0.00000000 | 1.00000 | 22.645620 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 3.344108 | 142.0000 | 42422 | 28.08696 | 16.50790 | 0.98600 | 1.00e-07 | -16.52200 | 2950000128 | 199.99580 | 3810.478 | 2 | 3053.327 |
| 700 | 1978 | 1 | 399 | 33 | 0.00000000 | 1.00000 | 23.236290 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 4.261217 | 143.0000 | 52289 | 26.48608 | 16.53365 | 0.99100 | 0.00e+00 | -16.54269 | 3300000000 | 217.79089 | 3812.273 | 2 | 3053.327 |
| 700 | 1979 | 0 | 399 | 34 | 0.00000000 | 1.00000 | 23.842110 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 4.000952 | 110.0000 | 59400 | 24.71180 | 16.55939 | 0.95900 | 1.00e-07 | -16.60126 | 3700000000 | 237.85249 | 3814.280 | 2 | 3053.327 |
| 700 | 1980 | 0 | 399 | 35 | 0.00000000 | 1.00000 | 24.459810 | 7.000000 | 0.000000 | -7.00000000 | 47000.0 | 4.004853 | 89.0000 | 58000 | 47.57342 | 16.58497 | 0.76200 | 0.00e+00 | -16.85678 | 3640000000 | 228.32120 | 3813.710 | 1 | 3427.958 |

| | warstds | ager | agexp | anoc | army85 | autch98 | auto4 | autonomy | avgnabo | centpol3 | coldwar | decade1 | decade2 | decade3 | decade4 | dem | dem4 | demch98 | dlan |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | peace | 34.46177 | 8.510845 | 0 | 129472.9 | 0 | 3.925812 | 0.005150819 | 0.4329550 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 3.995929 | 0 | 70.0 |
| 2 | peace | 34.34635 | 8.478997 | 0 | 129413.0 | 0 | 10.000000 | 0.000000000 | 0.0450517 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 3 | peace | 77.00000 | 8.481015 | 0 | 130431.0 | 0 | 10.000000 | 0.000000000 | 0.0300345 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 4 | peace | 78.00000 | 8.451628 | 0 | 126781.7 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 5 | peace | 79.00000 | 8.500172 | 0 | 130979.2 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 6 | peace | 80.00000 | 8.528873 | 0 | 130616.5 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 7 | peace | 81.00000 | 8.546965 | 0 | 129142.7 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 8 | peace | 82.00000 | 8.550921 | 0 | 129067.3 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 9 | peace | 83.00000 | 8.530715 | 0 | 130133.3 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 10 | peace | 84.00000 | 8.544636 | 0 | 131839.4 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 11 | peace | 85.00000 | 8.532668 | 0 | 129951.5 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 12 | peace | 86.00000 | 8.488276 | 0 | 129301.5 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 13 | peace | 87.00000 | 8.534172 | 0 | 130448.2 | 0 | 10.000000 | 0.000000000 | 0.2500000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 14 | peace | 88.00000 | 8.529870 | 0 | 131597.9 | 0 | 10.000000 | 0.000000000 | 0.2500000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 15 | peace | 89.00000 | 8.525314 | 0 | 131245.8 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 16 | peace | 90.00000 | 8.581366 | 0 | 129057.9 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 17 | peace | 91.00000 | 8.562673 | 0 | 47000.0 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 18 | peace | 92.00000 | 8.560198 | 0 | 47000.0 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 19 | peace | 93.00000 | 8.541015 | 0 | 47000.0 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 20 | peace | 94.00000 | 61.100498 | 0 | 47000.0 | 0 | 10.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 21 | peace | 95.00000 | 54.437160 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 22 | peace | 96.00000 | 45.760941 | 0 | 47000.0 | -3 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 23 | peace | 97.00000 | 51.714909 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 24 | peace | 98.00000 | 46.355721 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 25 | peace | 99.00000 | 38.092049 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 26 | peace | 100.00000 | 38.050751 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0131874 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 27 | peace | 101.00000 | 35.823952 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0131874 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 28 | peace | 102.00000 | 45.656319 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0131874 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 29 | peace | 103.00000 | 33.937969 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0105231 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |

Drops year and country code. These are never used in the paper (as far as I can tell).

| | warstds | ager | agexp | anoc | army85 | autch98 | auto4 | autonomy | avgexp | centpol3 | coldwar | decade1 | decade2 | decade3 | decade4 | dem | dem4 | domch98 | dia... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | peace | 34.46177 | 8.510845 | 0 | 129472.8 | 0 | 3.925812 | 0.005150819 | 0.4329550 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 3.995929 | 0 | 70.0 |
| 2 | peace | 14.34635 | 8.478997 | 0 | 129413.0 | 0 | 10.000000 | 0.000000000 | 0.0450517 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 3 | peace | 77.00000 | 8.483015 | 0 | 130431.0 | 0 | 10.000000 | 0.000000000 | 0.0300345 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 4 | peace | 78.00000 | 8.451628 | 0 | 126781.7 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 5 | peace | 79.00000 | 8.500172 | 0 | 130079.2 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 6 | peace | 80.00000 | 8.528873 | 0 | 130616.5 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 7 | peace | 81.00000 | 8.540965 | 0 | 129142.7 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 8 | peace | 82.00000 | 8.550921 | 0 | 129067.3 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 9 | peace | 83.00000 | 8.530715 | 0 | 130133.3 | 0 | 10.000000 | 0.000000000 | 0.0225258 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 10 | peace | 84.00000 | 8.544636 | 0 | 131839.4 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 11 | peace | 85.00000 | 8.532668 | 0 | 129951.5 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 12 | peace | 86.00000 | 8.468276 | 0 | 129301.5 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 13 | peace | 87.00000 | 8.534172 | 0 | 130448.2 | 0 | 10.000000 | 0.000000000 | 0.2500000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 14 | peace | 88.00000 | 8.529870 | 0 | 131597.9 | 0 | 10.000000 | 0.000000000 | 0.2500000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 15 | peace | 89.00000 | 8.525314 | 0 | 131245.8 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 16 | peace | 90.00000 | 8.581366 | 0 | 129057.9 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 17 | peace | 91.00000 | 8.562673 | 0 | 47000.0 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 18 | peace | 92.00000 | 8.560198 | 0 | 47000.0 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 19 | peace | 93.00000 | 8.541015 | 0 | 47000.0 | 0 | 10.000000 | 0.000000000 | 0.0000000 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 20 | peace | 94.00000 | 61.100498 | 0 | 47000.0 | 0 | 10.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 21 | peace | 95.00000 | 54.437100 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 22 | peace | 96.00000 | 45.760941 | 0 | 47000.0 | -3 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 23 | peace | 97.00000 | 51.714909 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 24 | peace | 98.00000 | 46.355721 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 25 | peace | 99.00000 | 38.092049 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0078924 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 26 | peace | 100.00000 | 38.050751 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0131874 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 27 | peace | 101.00000 | 35.823952 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0131874 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 28 | peace | 102.00000 | 45.656319 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0131874 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |
| 29 | peace | 103.00000 | 33.932969 | 0 | 47000.0 | 0 | 7.000000 | 0.000000000 | 0.0105231 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0.000000 | 0 | 70.0 |

#Fearon and Laitin Model (2003) Specification###
model.fl.1<-
train(as.factor(warstds)~warhist+ln_gdpen+lpopns+lmtnest+ncontig+oil+nwstate
+inst3+pol4+ef+relfrac, #FL 2003 model spec
metric="ROC", method="glm", family="binomial",
trControl=tc, data=data.full)

Uses the present to predict the present

Assignment 2

# Class slides for Thursday, Sept 24: Armed conflict, part 2

Matthew J. Salganik

COS 597E/SOC 555 Limits to prediction
Fall 2020, Princeton University