

Bush 631-607: Quantitative Methods

Lecture 14 (11.30.2021): Review & Summary

Rotem Dvir

The Bush school of Government and Public Policy

Texas A&M University

Fall 2021

What is today's plan?

- ▶ Review course topics.
- ▶ Research designs.
- ▶ Predictions, probability and estimations.
- ▶ Social science and real-world politics.
- ▶ Data science in the real-world.

Course review: Causality

- ▶ Establish Cause → effect.
- ▶ Using (hypothetical) *Counterfactual*.



Establish causality

EXPERIMENTAL RESEARCH DESIGN

- ▶ Treatment and control groups.
- ▶ Same outcome measures.
- ▶ Gold standard → randomization.
- ▶ Calculate ATE over group of respondents.

$$SATE = \frac{1}{n} * \sum_{i=1}^n * Y_i(1) - Y_i(0)$$

Experiments: Example

► President's type and support for FP (2019)

```
# Diff-in-means = ATE of type  
mean(hawks$approve_b[hawks$hawk_t == 1], na.rm = T) -  
  mean(hawks$approve_b[hawks$hawk_t == 2], na.rm = T)
```

```
## [1] -0.1202408
```

```
# Also with tapply()  
tapply(hawks$approve_b, hawks$hawk_t, mean, na.rm = TRUE)
```

```
##          1          2  
## 0.5774336 0.6976744
```

```
# Can also use subsets and diff-in-means
```

Establish causality

OBSERVATIONAL RESEARCH DESIGN

- ▶ Using data to assess causality.
- ▶ Good for generalizing results.
- ▶ Not as good for randomization.
- ▶ Problem of pre-treatment variables (confounders).

Causality in observational data

- ▶ Survey: political polarization and views of China (2020)
- ▶ Party *Thermometer* → *Dems/Reps* > 50 support for party.

```
levels(as.factor(threat$china_frenemy))
```

```
## [1] ""          "Ally"       "Enemy"      "Friendly"   "Unfriendly"
```

Causality in observational data

► China as the enemy?

```
# Diff-in-means (Support Dems thermometer)
```

```
mean(threat$china[threat$affective_Dem < 50], na.rm = T) -  
mean(threat$china[threat$affective_Dem > 50], na.rm = T)
```

```
## [1] 0.2243064
```

```
# Using tapply() by political interest
```

```
app <- tapply(threat$china, threat$pol_interest, mean, na.rm = TRUE)  
sort(app)
```

```
## Not interested at all Moderately interested Slightly interested  
##          2.526570          2.580093          2.643333  
##      Very interested          Extremely interested  
##          2.650142          2.700000          2.713311
```

Assessing research designs

Strengths and Weaknesses

- ▶ Internal validity:
 - ▶ How does the design help answering the research Q?
 - ▶ Experiments → strong (randomization).
 - ▶ Observational → weak (confounders).
- ▶ External validity:
 - ▶ Can we generalize the results from sample?
 - ▶ Experiments → weak (hypothetical).
 - ▶ Observational → strong (real-world, cross-national).

Course review: Measurement

- ▶ Apply quant methods for social science.
- ▶ Measures → the context of concepts.
- ▶ Challenge: latent factors
 - ▶ What is ideology? How do we measure it?
 - ▶ Terrorism?
 - ▶ Democracy - polity vs. freedom house scores.

Measurement challenge

- ▶ Challenge of missing values (non-responses)
- ▶ No data collected / respondents refuse to answer (DKs).
- ▶ NAs in our data

```
# Solving with na.rm = TRUE  
mean(bushdata$Pol_survMusl)
```

```
## [1] NA
```

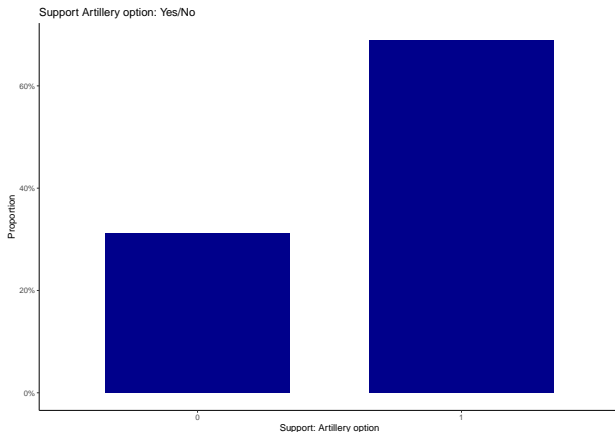
```
mean(bushdata$Pol_survMusl, na.rm = TRUE)
```

```
## [1] 2.067584
```

Visuals

- ▶ Barplot: counts/proportions for categories

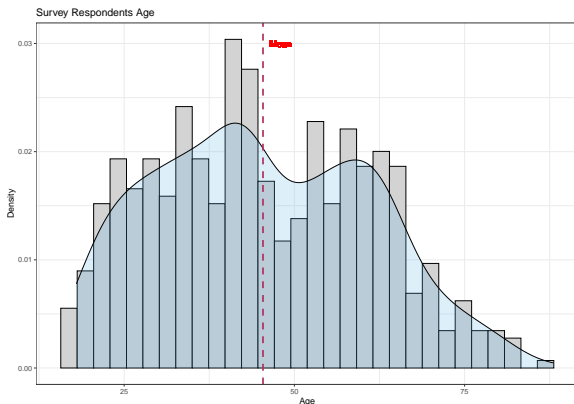
```
ggplot(wardata, aes(x=factor(prefer_artillery_dummy))) +  
  geom_bar(aes(y = (..count..)/sum(..count..)), width = 0.7, fill = "darkblue") +  
  xlab("Support: Artillery option") + ylab("Proportion") +  
  scale_y_continuous(labels=scales::percent) + ggtitle("Support Artillery option") +  
  theme_classic()
```



Visuals

► Histogram: distribution of numerical variable

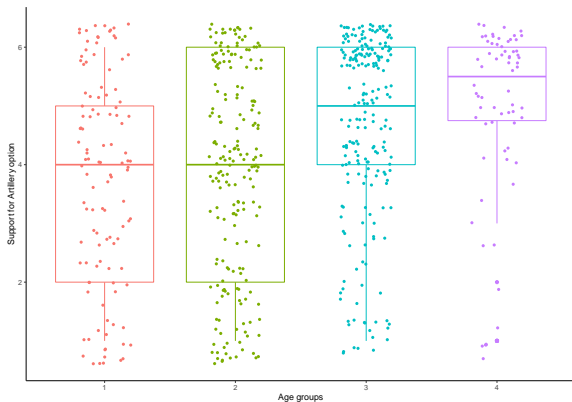
```
ggplot(wardata, aes(x=age)) +  
  geom_histogram(aes(y=..density..), colour="black", fill="lightgrey")+  
  geom_density(alpha=.2, fill="#56B4E9") +  
  xlab("Age") + ylab("Density") + theme_bw() + ggtitle("Survey Respondents Age") +  
  geom_vline(aes(xintercept=mean(age)),  
             color="maroon", linetype="dashed", size=1) +  
  geom_text(x = 48, y = 0.03, label = "Mean", col = "red")
```



Visuals

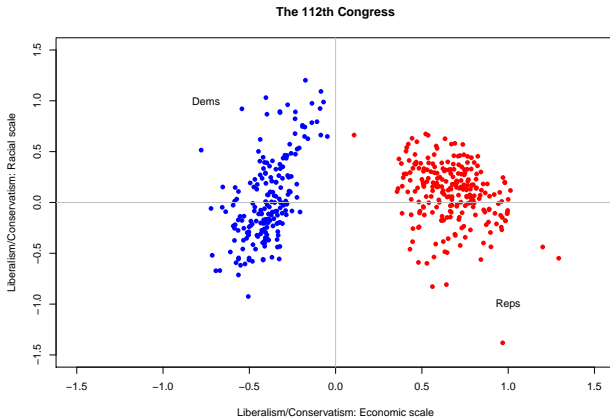
- ▶ Boxplot: Compare single variable distribution

```
ggplot(wardata, aes(x=factor(agegroup), y = artillery_approve,  
                    color = factor(agegroup))) +  
  geom_boxplot() +  
  geom_jitter(shape=16, position=position_jitter(0.2)) +  
  xlab("Age groups") + ylab("Support for Artillery option") +  
  theme_classic() + theme(legend.position = "none")
```



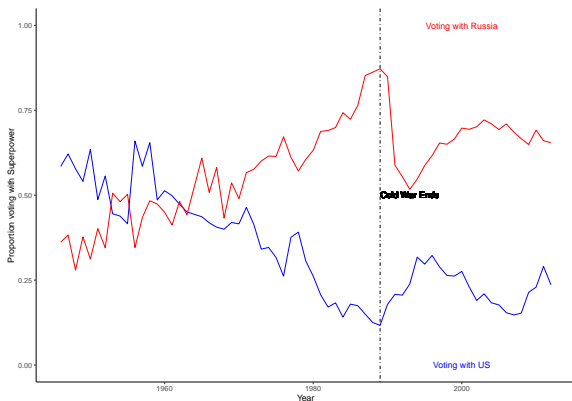
Visuals

- ▶ Scatterplot: Visualize bivariate relationship



Visuals

- ▶ Plot time trends (UN Voting data 1946-2012)



Predictions

- ▶ Predict with sample mean: using loops.
- ▶ Prediction error = actual outcome - predicted outcome.
- ▶ RMSE: average magnitude of prediction error.

- ▶ Correlations:
 - ▶ Summary of bivariate relationship.
 - ▶ How two factors 'move together' on average.
 - ▶ Always relative to mean value.

```
# Voting with US  
cor(unvoting$idealpoint, unvoting$PctAgreeUS, use = "pairwise")
```

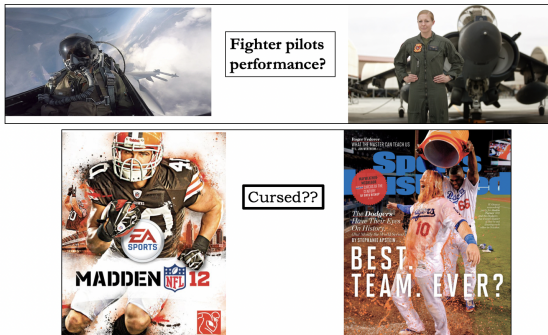
```
## [1] 0.7498446
```

Predictions

- ▶ The linear model: $Y = \alpha + \beta * X_i + \epsilon$
- ▶ Model elements:
 - ▶ Intercept (α): the average value of Y when X is zero.
 - ▶ Slope (β): the average increase in Y when X increases by 1 unit.
 - ▶ Error/disturbance term (ϵ): the deviation of an observation from a perfect linear relationship.
- ▶ Least squared:
 - ▶ How to estimate the regression line.
 - ▶ 'Select' $\hat{\alpha}, \hat{\beta}$ to minimize SSR.
 - ▶ R syntax: `lm(y ~ x, data = mydata)`

Regression to the mean

- ▶ High (low) observations are followed by low (high) observations.
- ▶ Observations 'regress' towards the average value of the data.



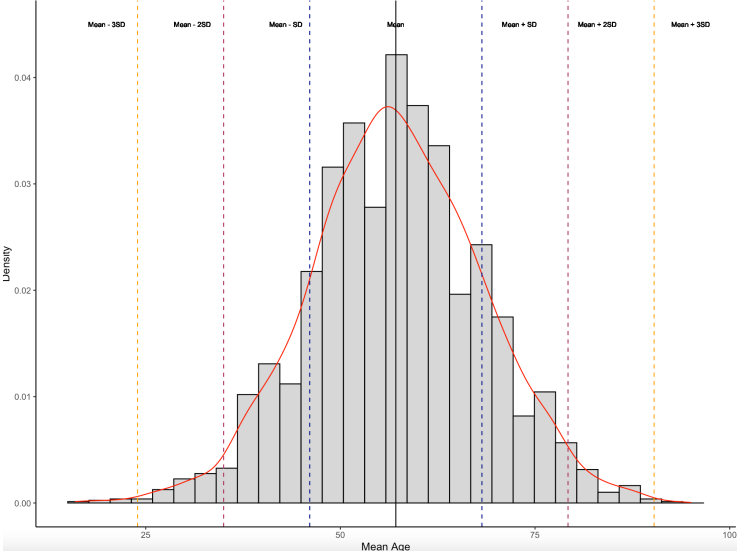
Probability

- ▶ Quantify uncertainty: step 1.
- ▶ Sample space, events (cards, coin toss)
- ▶ Probability: $P(A) = \frac{\text{Elements}(A)}{\text{Elements}(\Omega)}$
- ▶ Conditional probability = $P(A|B) = \frac{P(A \& B)}{P(B)}$
- ▶ *Monty Hall problem.*

Probability

- ▶ Random variables: from events to numbers.
- ▶ Uncertainty of sample means or sums.
- ▶ Probability distributions: Bernoulli (binary), Binomial (discrete).
- ▶ Expectations of r.v. \rightarrow population value.
- ▶ Variance of r.v. \rightarrow 'spread' of distribution.
- ▶ CLT and large samples.

Normal distribution

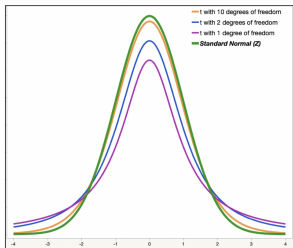


Uncertainty and estimation

- ▶ Estimate based on r.v.s.
- ▶ Quantity of interest: *point estimate* (mean / diff-in-means)
- ▶ How to learn of estimator distribution? simulations. . .
- ▶ Calculate SD, or SE in single sample.
- ▶ Construct 95% CIs - how to interpret?

t-distribution

- ▶ Small samples.
- ▶ Account for DOF.
- ▶ More conservative, why? 'fatter tails'.



Numbers in each row of the table are values on a t-distribution with (df) degrees of freedom for selected right-tail (greater-than) probabilities (α).

df \ α	0.40	0.25	0.10	0.05	0.025	0.01	0.005	0.0005
1	0.324988	1.000000	0.077084	0.317327	12.706193	31.820828	63.6574	308.4132
2	0.230602	0.691459	1.060718	2.353368	4.302672	6.964662	8.164564	31.526517
3	0.274621	0.764882	1.052784	2.353363	3.142670	4.540702	5.408191	13.28243
4	0.278222	0.740007	1.052066	2.131847	2.7945	3.74695	4.40425	8.6183
5	0.267181	0.729807	1.047084	2.015046	2.575926	3.36483	4.0214	6.8986
6	0.264253	0.717338	1.042676	1.941218	2.44681	3.14267	3.74742	5.9588
7	0.263167	0.711424	1.041824	1.894576	2.36482	2.99795	3.49248	5.4079
8	0.261921	0.706287	1.040815	1.860548	2.30600	2.89846	3.35528	5.0413
9	0.260954	0.702722	1.040029	1.831112	2.26276	2.81144	3.24984	4.7889
10	0.260185	0.699812	1.039384	1.804741	2.22814	2.73277	3.16227	4.5888
11	0.259568	0.697445	1.038800	1.780885	2.20009	2.71088	3.10071	4.4370
12	0.259023	0.695683	1.038271	1.762288	2.1781	2.68180	3.05454	4.3178
13	0.258531	0.694229	1.037771	1.748532	2.16027	2.65521	3.01228	4.2206
14	0.258102	0.693017	1.037299	1.738123	2.14639	2.64049	2.97684	4.1465
15	0.257738	0.691997	1.036866	1.730260	2.1345	2.62824	2.94671	4.0728
16	0.257408	0.691132	1.036452	1.724084	2.1181	2.61849	2.92078	4.0150
17	0.257147	0.689975	1.036037	1.719087	2.10962	2.60983	2.89822	3.9611
18	0.256912	0.689444	1.035691	1.714943	2.10267	2.60218	2.87944	3.9106
19	0.256742	0.689161	1.035278	1.711313	2.09632	2.59494	2.86305	3.8624
20	0.256614	0.688884	1.034931	1.707978	2.08956	2.57798	2.84834	3.8165
21	0.256508	0.688752	1.034748	1.70542	2.07981	2.57195	2.83126	3.7813
22	0.256421	0.688656	1.034597	1.703413	2.07197	2.56622	2.81675	3.7511
23	0.256343	0.688588	1.034469	1.701872	2.06508	2.56081	2.80374	3.7252
24	0.256273	0.688536	1.034356	1.700643	2.05898	2.55563	2.79344	3.7024
25	0.256212	0.688498	1.034256	1.700002	2.05362	2.55064	2.78544	3.6814
26	0.256158	0.688470	1.034167	1.699914	2.04894	2.54581	2.77971	3.6624
27	0.256108	0.688450	1.034087	1.699838	2.04481	2.54121	2.77528	3.6466
28	0.256061	0.688436	1.034017	1.699771	2.04111	2.53681	2.77114	3.6329
29	0.256018	0.688426	1.033954	1.699713	2.03781	2.53261	2.76724	3.6201
30	0.255978	0.688419	1.033897	1.699663	2.03481	2.52861	2.76354	3.6081
∞	0.255842	0.688314	1.033162	1.698484	2.01908	2.52028	2.75762	3.5768
α			80%	90%	95%	99%	99.5%	99.9%

Hypothesis tests

- ▶ Estimators: sample means / diff-in-means
- ▶ *Proof by contradiction.*
- ▶ Steps for testing:
 1. Define null and alternative hyps ($H_0; H_1$).
 2. Select *test statistic* and level of test (α).
 3. Derive reference distribution.
 4. Calculate p-values.
 5. Make a decision: reject/retain.
- ▶ **Decision rule:**
 - ▶ **Reject null** if p-value is *below* $\alpha = 0.05$
 - ▶ Otherwise **retain the null** or **fail to reject**.

Hypothesis test

- ▶ Run *Two-sample t-test* with `t.test()`

```
t.test(exp.dat$cont_cor1[exp.dat$trt1 == 0],  
       exp.dat$cont_cor1[exp.dat$trt1 == 1])
```

```
##
```

```
## Welch Two Sample t-test
```

```
##
```

```
## data: exp.dat$cont_cor1[exp.dat$trt1 == 0] and exp.dat$cont_cor1[exp.dat$trt1 == 1]
```

```
## t = -13.697, df = 993.53, p-value < 2.2e-16
```

```
## alternative hypothesis: true difference in means is not equal to 0
```

```
## 95 percent confidence interval:
```

```
## -23.59653 -17.68267
```

```
## sample estimates:
```

```
## mean of x mean of y
```

```
## 1489.333 1509.973
```

Least squared estimator

- ▶ Uncertainty in *least squared* estimator:
 - ▶ Generate reference distribution.
 - ▶ Calculate SEs.
 - ▶ Construct 95% CIs.
 - ▶ Run hypotheses tests.
 - ▶ Results are 'statistically significant', or not.

- ▶ Assumptions for regression estimates:

(1) Exogeneity: mean of ϵ_i does not depend on X_i

$$E(\epsilon_i|X_i) = E(\epsilon_i) = 0$$

(2) Homoskedasticity: variance of ϵ_i does not depend on X_i

$$V(\epsilon_i|X_i) = V(\epsilon_i) = \sigma^2$$

Putting everything together

- ▶ Hypotheses:

- ▶ $H_0 : \beta_1 = 0$

- ▶ $H_a : \beta_1 \neq 0$

- ▶ Our estimators: $\hat{\beta}_0, \hat{\beta}_1$

- ▶ SE and CIs:

- ▶ $\hat{\beta}_0 \pm 1.96 * \hat{SE}(\hat{\beta}_0)$

- ▶ $\hat{\beta}_1 \pm 1.96 * \hat{SE}(\hat{\beta}_1)$

- ▶ Hypotheses test:

- ▶ Test statistic: $\frac{\hat{\beta}_1 - \beta_1^*}{\hat{SE}(\hat{\beta}_1)} \sim N(0,1)$

- ▶ $\hat{\beta}_1$ is **statistically significant** if $p < 0.05$.

Rebels and Nukes (2015)

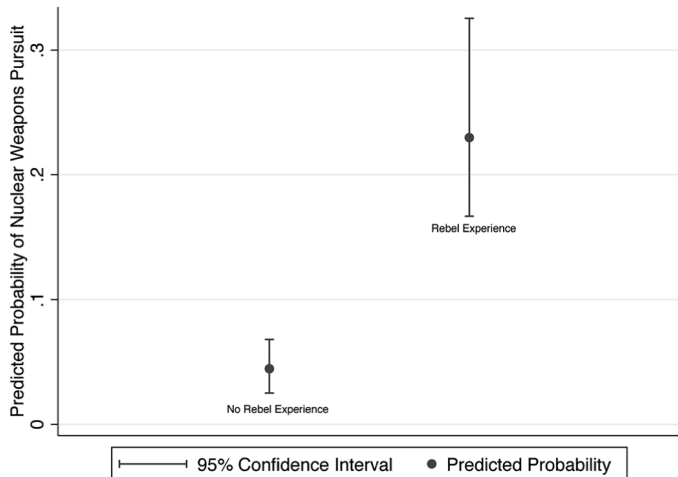
- ▶ Multivariate regression: account for confounders

```
summary(lm(pursuit ~ rebel + milservice + polity2, data = nukes))
```

```
##
## Call:
## lm(formula = pursuit ~ rebel + milservice + polity2, data = nukes)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.06587 -0.04408 -0.02544 -0.01020  0.99682
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.0073899  0.0027782   2.660  0.00783 **
## rebel        0.0320096  0.0044238   7.236  5.08e-13 ***
## milservice   0.0217914  0.0045106   4.831  1.38e-06 ***
## polity2     0.0004679  0.0002801   1.670  0.09489 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1672 on 7684 degrees of freedom
## (1164 observations deleted due to missingness)
## Multiple R-squared:  0.01596,    Adjusted R-squared:  0.01558
## F-statistic: 41.54 on 3 and 7684 DF,  p-value: < 2.2e-16
```

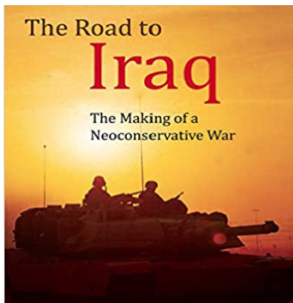
OLS coefficient interpretation

- ▶ Rebel experience and nuclear technology (2015)



How to use our research?

Applying theories and IR research



Applying IR research in global affairs

- ▶ The motivation:
 - ▶ We study IR or social science dynamics.
 - ▶ Do policymakers use? or even care about this knowledge?

The New York Times | <http://nyti.ms/1f4Huyy>

SundayReview | OP-ED COLUMNIST

Professors, We Need You!

Nicholas Kristof FEB. 15, 2014

SOME of the smartest thinkers on problems at home and around the world are university professors, but most of them just don't matter in today's great debates.

The most stinging dismissal of a point is to say: "That's academic." In other words, to be a scholar is, often, to be irrelevant.

 **David Rothkopf** 
@djrothkopf

Kristof gets why we at FP are dialing back academic contributions--too many are opaque, abstract, incremental, dull. nyti.ms/1fpsd84

6:30 AM · Feb 16, 2014 · Twitter for Webaltes

50 Retweets 31 Likes

Applying IR research

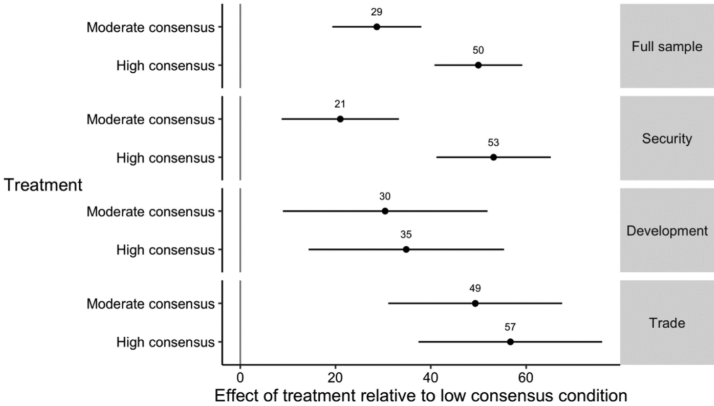
- ▶ Main challenge → **Time** and reading research article.
- ▶ Scholars adapt:
 - ▶ Joint forums: *bridging the gap*.
 - ▶ Policy-focused writing (Lawfare blog, War on the rocks, Monkey cage).
- ▶ Is it working?
- ▶ Ask policymakers. . .
- ▶ Previous work (2014):
 - ▶ Not really.
 - ▶ Academic work not aimed to 'close the gap'.

Applying IR research

- ▶ Recent evidence (2021) → replicate 2014 survey.
- ▶ Innovations for better insights:
 - ▶ Diverse sample - more areas FP.
 - ▶ Both high and low-ranking officials.
 - ▶ Embedded experiment for direct effects.
 - ▶ Broader conception for engagement (social media).
- ▶ Sample: 616 officials (Clinton, Bush, Obama administration).

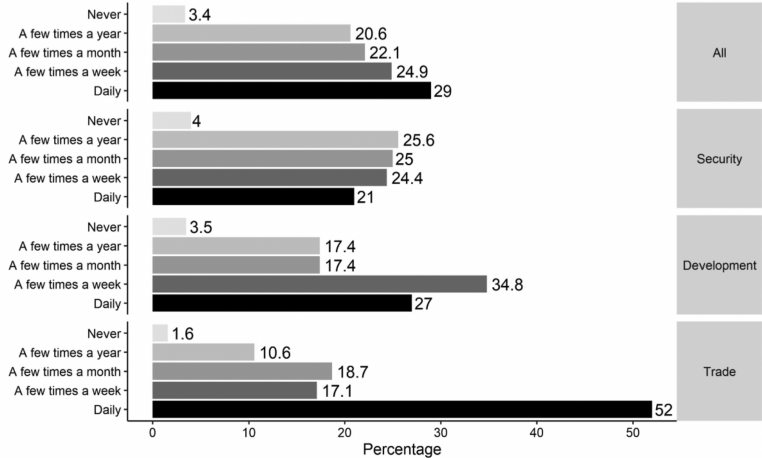
Useful political research?

- ▶ Research consensus and updating policy views.



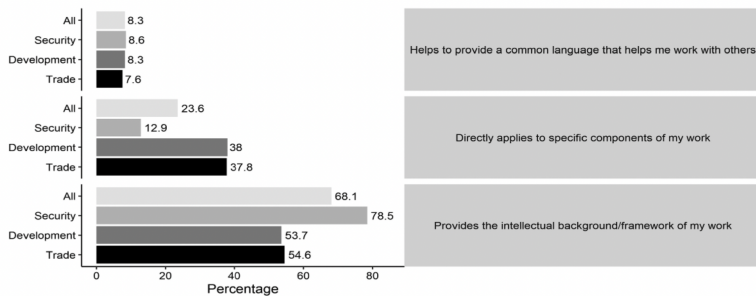
Useful political research?

► Frequency of using research into government work.



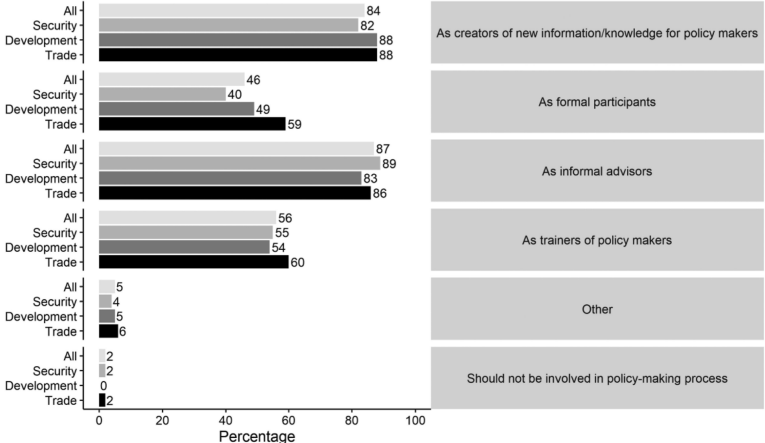
Useful political research?

► How do you use academic research?

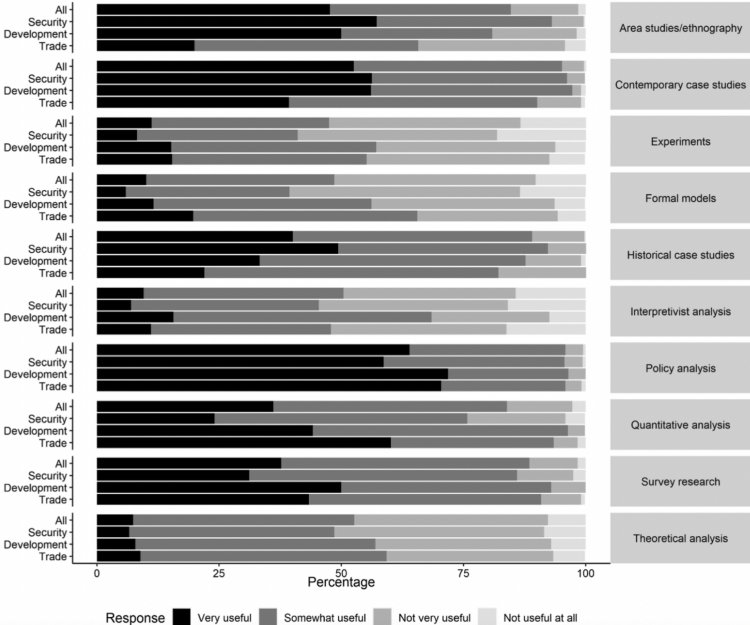


Useful researchers

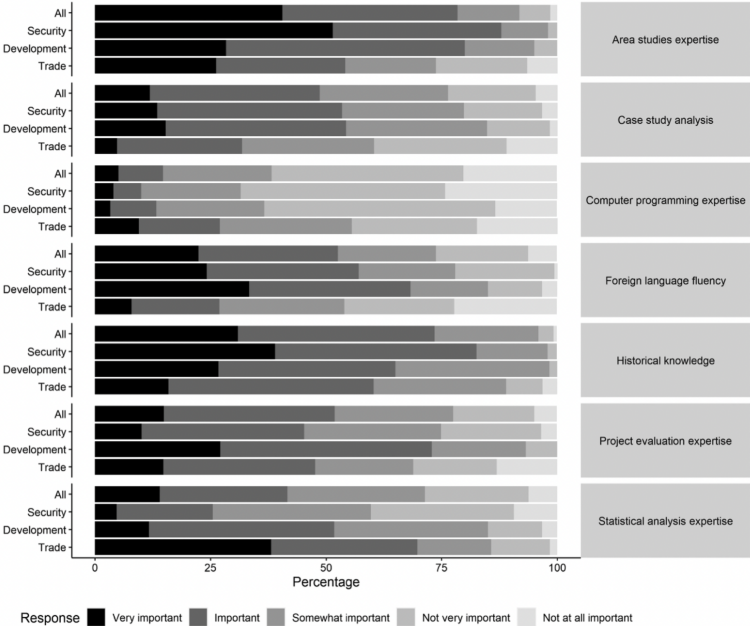
► How can researchers contribute the most?



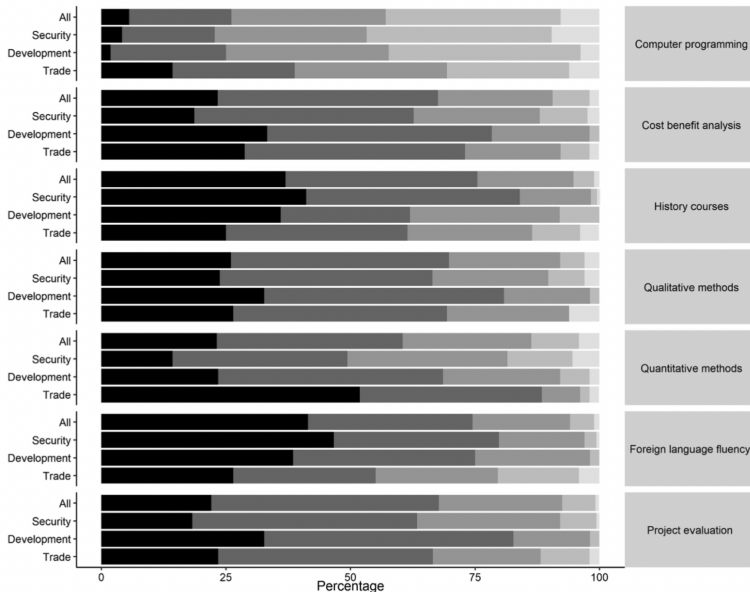
Useful skills...



Useful skills...



Working for the government...



Response **Very valuable** Valuable Somewhat valuable Not very valuable Not valuable at all

Using political research

- ▶ How can research be useful (Byman & Kroenig 2016)
 - ▶ Practical and useful recommendations.
 - ▶ Focus - clarify complex situations.
 - ▶ Time relevant research.
- ▶ Scenarios for applying academic insights:
 - ▶ Challenge existing government knowledge - shocks (9.11, Soviet collapse).
 - ▶ Policy failures (Iraq insurgency outbreak).
 - ▶ Missing baseline knowledge (Somalia intervention).

Becoming useful political advisor

- ▶ Concrete steps:
 - ▶ Networking and personal connections.
 - ▶ 'Inject' research into bureaucracy.
 - ▶ Concise and clear reports in nonacademic outlets.
- ▶ Tamper expectations:
 - ▶ What is being relevant?
 - ▶ Not likely to drastically shape policy.
 - ▶ Influence the **deliberation** process.
- ▶ What's in it for policymakers?
 - ▶ Offer *contrarian arguments* to accepted view.

Data science in the real world

- ▶ Data analysis → set of tools to understand the world.
- ▶ The core role of probability.
- ▶ Apply complex concepts like repeated sampling.
- ▶ Bayesian logic and saving lives.
- ▶ Movies:
 - ▶ Prediction by the numbers: ([Link](#))
 - ▶ Tails you win: ([Link](#))