

# PREDICTIVE ELECTION MODELS

Data Analysis for Journalism and Political Communication  
(Fall 2024)

Prof. Bell

# DART-THROWING CHIMPANZEES



# PREDICTION

- Prediction is not the same thing as inference

# PREDICTION

- Prediction is not the same thing as inference
- Prediction uses **training data** to make predictions on **test data**

# PREDICTION

- Prediction is not the same thing as inference
- Prediction uses **training data** to make predictions on **test data**
  - ▶ With inference, we are generating an estimate from a sample that is representative of the population

# PREDICTION

- Prediction is not the same thing as inference
- Prediction uses **training data** to make predictions on **test data**
  - ▶ With inference, we are generating an estimate from a sample that is representative of the population
  - ▶ Prediction focuses on an “out-of-sample” population

# PREDICTION

- Prediction is not the same thing as inference
- Prediction uses **training data** to make predictions on **test data**
  - ▶ With inference, we are generating an estimate from a sample that is representative of the population
  - ▶ Prediction focuses on an “out-of-sample” population
- Biggest risk in prediction is **overfitting**: generating a model that relies too heavily on the training data and is not able to make good predictions on the test data

# PREDICTION

- Prediction is not the same thing as inference
- Prediction uses **training data** to make predictions on **test data**
  - ▶ With inference, we are generating an estimate from a sample that is representative of the population
  - ▶ Prediction focuses on an “out-of-sample” population
- Biggest risk in prediction is **overfitting**: generating a model that relies too heavily on the training data and is not able to make good predictions on the test data
  - ▶ With inference, the larger the sample, the better



# PREDICTION

- Prediction is not the same thing as inference
- Prediction uses **training data** to make predictions on **test data**
  - ▶ With inference, we are generating an estimate from a sample that is representative of the population
  - ▶ Prediction focuses on an “out-of-sample” population
- Biggest risk in prediction is **overfitting**: generating a model that relies too heavily on the training data and is not able to make good predictions on the test data
  - ▶ With inference, the larger the sample, the better
- By design, the **prediction interval** is wider than the confidence interval (e.g., MOE)

# MACHINE LEARNING AND AI

- Machine learning is a wide spectrum ranging from simple decision trees to large language models

# MACHINE LEARNING AND AI

- Machine learning is a wide spectrum ranging from simple decision trees to large language models
- The “learning” is simply that the data draws on patterns in the training data to make predictions on the test data

# MACHINE LEARNING AND AI

- Machine learning is a wide spectrum ranging from simple decision trees to large language models
- The “learning” is simply that the data draws on patterns in the training data to make predictions on the test data
- It is “machine” learning because computation takes over some of the decisions that are usually left to the researcher

# MACHINE LEARNING AND AI

- Machine learning is a wide spectrum ranging from simple decision trees to large language models
- The “learning” is simply that the data draws on patterns in the training data to make predictions on the test data
- It is “machine” learning because computation takes over some of the decisions that are usually left to the researcher
  - ▶ **Supervised learning:** The researcher provides the machine with a target (e.g., an election outcome) and the machine determines the features of the test data that best predict the target
  - ▶ **Unsupervised learning:** What we typically think of as “artificial intelligence,” the machine is not given a target and uncovers new patterns in the data

# MACHINE LEARNING AND AI

- Machine learning is a wide spectrum ranging from simple decision trees to large language models
- The “learning” is simply that the data draws on patterns in the training data to make predictions on the test data
- It is “machine” learning because computation takes over some of the decisions that are usually left to the researcher
  - ▶ **Supervised learning:** The researcher provides the machine with a target (e.g., an election outcome) and the machine determines the features of the test data that best predict the target
  - ▶ **Unsupervised learning:** What we typically think of as “artificial intelligence,” the machine is not given a target and uncovers new patterns in the data
- Because computation replaces some researcher decisions, many machine learning/AI models are considered “black boxes” where it is difficult to decipher why the machine makes the predictions that it does (called **explainability**)

# EXPERT POLITICAL JUDGMENT



VS.



# WISDOM OF THE CROWDS

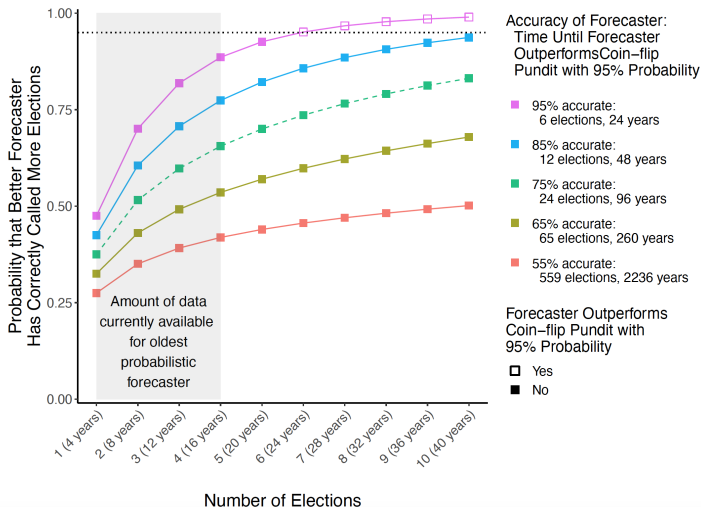




# WISDOM OF THE CROWDS

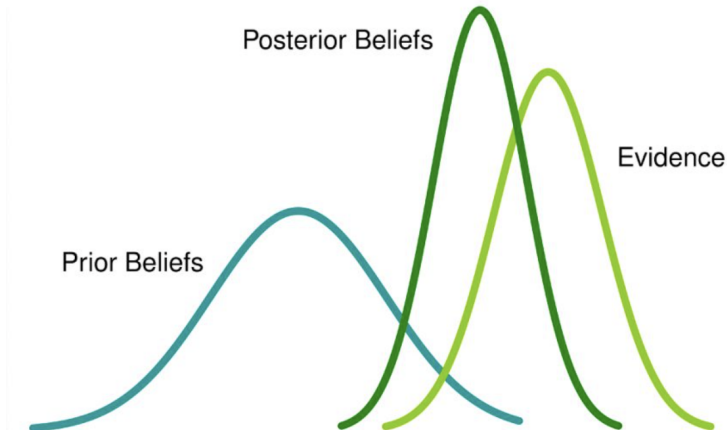
Polling aggregation	Clinton	Trump	Clinton	Trump	Clinton	Trump
	Probabilities		Electoral college		Vote share	
FiveThirtyEight	70.3	29.6	299	238	48.6	45.1
The Upshot	84.0	16.0	322	216		
RCP average of polls			301	237	47.2	44.3
The Daily Kos	88.0	12.0	313	225		
Princeton EC	99.0	1.0	312	226	51.3	48.8
HuffPost	98.1	1.6	323	215		
PollyVote			323	215	52.6	47.4
<b>Mean</b>	<b>87.9</b>	<b>12.0</b>	<b>313.3</b>	<b>224.6</b>	<b>49.0</b>	<b>46.1</b>

# WISDOM OF THE CROWDS



Source: Grimmer, Justin, Dean Knox, and Sean Westwood. 2024. "Assessing the Reliability of Probabilistic US Presidential Election Forecasts May Take Decades." OSF Preprints. August 26. doi:10.31219/osf.io/6g5zq.

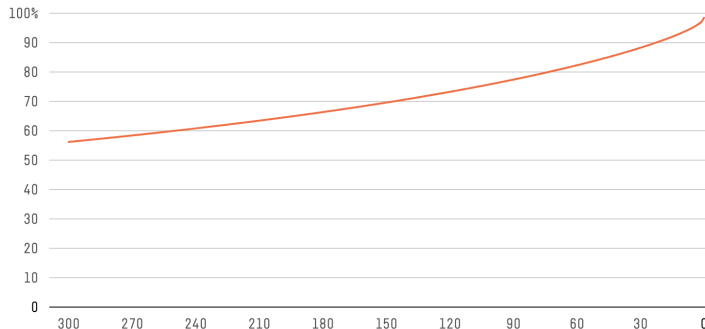
# HOW PREDICTION MODELS WORK



# HOW PREDICTION MODELS WORK

## We put more weight on the polls as Election Day approaches

Estimated\* share of the overall 538 presidential forecast that is based on polls (as opposed to non-polling historical "fundamentals"), by day before the election



\*As of Aug. 23, this estimate uses a standard deviation of 9 percentage points in 538's fundamentals model and a daily standard deviation of 0.35 points in polls, plus overall uncertainty of 1 point about the polling average. Real values will depend on how many polls we have and can differ from this estimate by a few points.

538

# HOW PREDICTION MODELS WORK

Election prediction models must decide:

- What polls to include (and how much to weight them):  
quality, quantity, sample size, time, etc.

# HOW PREDICTION MODELS WORK

Election prediction models must decide:

- What polls to include (and how much to weight them): quality, quantity, sample size, time, etc.
- How to adjust polls for house effects and/or mode effects

# HOW PREDICTION MODELS WORK

Election prediction models must decide:

- What polls to include (and how much to weight them): quality, quantity, sample size, time, etc.
- How to adjust polls for house effects and/or mode effects
- How to quantify the uncertainty in poll results

# HOW PREDICTION MODELS WORK

$$Y_{0, Support_t} \sim N(\mu_t, \sigma)$$
 WHERE  $\mu_t \sim N(\mu_{t-1}, T_t)$  &  $T_t = N(t^{\frac{1}{\alpha}}, (t+1))$   
 or  $T_t = N(0, \sigma_{true})$

AND "X..."

Full obs.  $i \sim N(\mu_i + (\text{pollster} + \text{mode} + \text{third party} + \text{bullet}), \text{obs. sig.} + \sigma_{US})$

where  $X_{i...} \sim N(0, 1) \times \sigma_{X...}$  &  $\sigma_{X...} \sim N(0, 1)$

obs. sig. =  $\frac{1}{\sqrt{N}}$

$\sigma_{US} \sim N(0, 1), [0, 2]$

one party, one geo

2 parties to go

And with the following trend expansion for states  $S=1 \dots S$  and parties  $P=1 \dots P$ :

$M[P, P, S, S] \sim MN(\mu_{P, P, S, S}, \Sigma_{P, P, S, S})$

&  $\mu_{P, P, S, S}$  for  $P, P, S$

$\mu_{P, P, S, S} \sim MN(\mu_{P, P, S, S}, \Sigma_{P, P, S, S})$

For pulling together  $\Sigma_{P, P, S, S}$  as a prior of  $\Sigma_{P, P, S, S}$  and random  $\Sigma_{P, P, S, S}$  with  $\Sigma_{P, P, S, S} \sim \text{diag}(S)$

M.U. Sign.  $\sigma$  added ahead of the prior band.

$P_1/P_2$	$P_3$
1	0.5
0.5	1
0.1	0.1
0.1	0.1
1	1
$P_2$	$P_3$

$S=1 \dots S$



# HOW PREDICTION MODELS WORK

Election prediction models must decide:

- What polls to include (and how much to weight them): quality, quantity, sample size, time, etc.
- How to adjust polls for house effects and/or mode effects
- How to quantify the uncertainty in poll results
- How to model election outcomes (e.g, intra-state correlation)

# HOW PREDICTION MODELS WORK

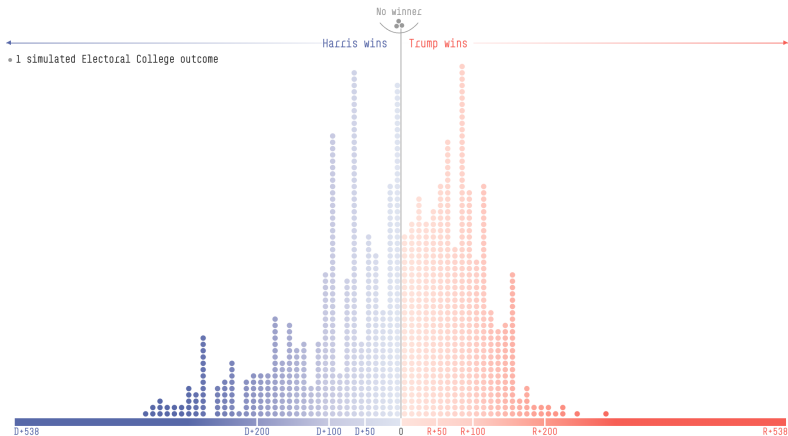
Election prediction models must decide:

- What polls to include (and how much to weight them): quality, quantity, sample size, time, etc.
- How to adjust polls for house effects and/or mode effects
- How to quantify the uncertainty in poll results
- How to model election outcomes (e.g, intra-state correlation)
- How to communicate probabilities

# HOW PREDICTION MODELS WORK

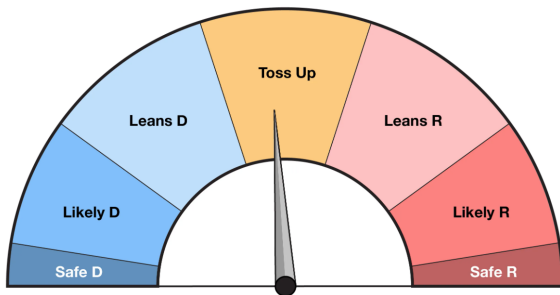


Trump	509
Harris	488
No winner	3
<hr/>	
1,000 simulations	



# HOW PREDICTION MODELS WORK

## 2024 Presidential Forecast



**President**

**Harris 52%**

**241 Harris, 246 Trump, 51 Toss Up**

# HOW PREDICTION MODELS WORK

