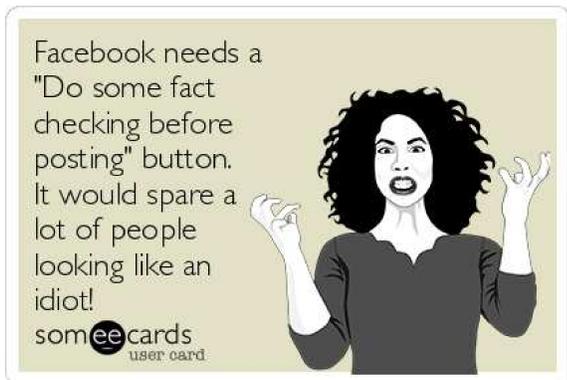


**Motivation** Fact checking is one of the main tasks performed by journalists. Automating it is one of the key goals in *computational journalism*.



**Task definition** Given a sentence:

1. **Identify** the components of a simple statistical claim
2. **Verify** the claim against a knowledge base

Similar to **knowledge base population**, if we assume that all claims are true.

**Desiderata** No labeled data for identification or verification, just a knowledge base and raw text

**Distant supervision?** Yes, but matching numerical values is not the same as named entity linking:

- Same number can be the value of unemployment, inflation, etc.
- Values are often approximated in text

**Dataset construction**

- Chose 16 statistical properties of approx. 175 countries from Freebase: *population*, *inflation rate*, *gross domestic product*, *life expectancy*, etc.
- Collected 100K pages from the Web querying with Bing for each country and property combination
- Converted HTML to text (BoilerPipe)
- Processed text with Stanford CoreNLP
- Extracted dependency and surface patterns between countries and values mentioned in the same sentence

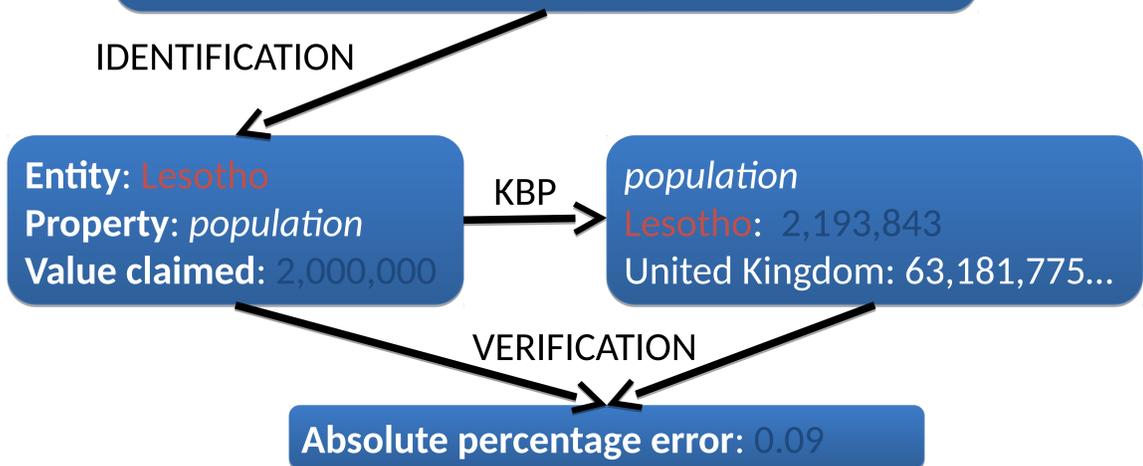
	France	Russia	Iceland
<b>population</b>	66,028,467	143,700,000	325,600
the population of X is _	66,000,000	140,000,000	325,000
X's population is estimated at _	66,030,000	145,000,000	
X's inflation rate is _	0.9		4.0
X has _ inhabitants		140,000,000	300,000
_ tourists visited X	68,000,000		

**Mean Absolute Percentage Error** measures the accuracy of real-valued predictions:

$$MAPE(V, \hat{V}) = \frac{1}{|V|} \sum \frac{|V - \hat{V}|}{|V|}$$

- scale-free (unlike RMSE)
- lower is better
- predicting 0s gives MAPE of 1

Lesotho, a landlocked enclave of South Africa, has a population of nearly 2 million and covers an area slightly smaller than the U.S. state of Maryland.



**Claim identification learning algorithm**

**Input:** *property* with entity-value pairs  $EV$ , patterns  $p_1, p_2, \dots$  with entity-value pairs  $EV_{p_1}, EV_{p_2}, \dots$

**Output:** Selected patterns  $P_{sel}$

$P_{sel} = \text{emptySet}$

calculate *propertyDefault*  
 $\text{currentMAPE} = MAPE(E, P_{sel})$

order patterns according to  $MAPE(EV, EV_p)$

**while** more patterns left:

add top pattern  $p$  to  $P_{sel}$

$\text{newMAPE} = MAPE(EV, PREDICT(E, P_{sel}))$

**if**  $\text{newMAPE} > \text{currentMAPE}$  **break**

**else**  $\text{currentMAPE} = \text{newMAPE}$

**function**  $PREDICT(\text{entities } E, \text{patterns } P)$

**foreach** entity  $e$  in  $E$ :

collect all values  $V$  for  $e$  in  $P$

**if**  $V$  not empty:

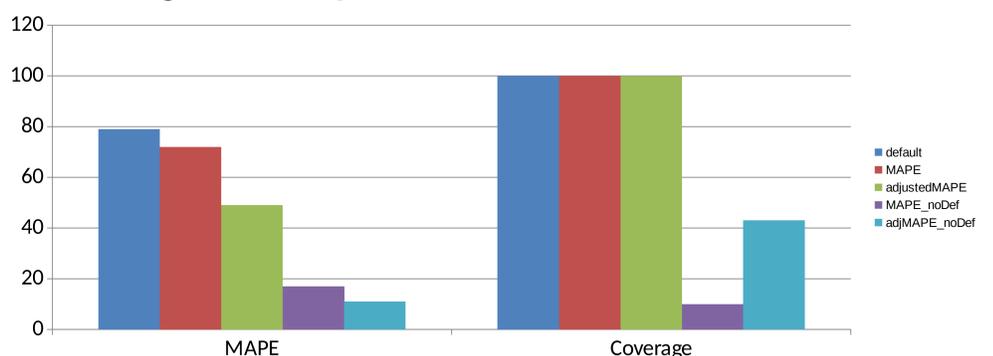
return  $MEAN(V)$

**else** return *propertyDefault*

Pick among the mean, median and 0

If the new pattern worsened the performance, stop

**Knowledge Base Population evaluation**



**Claim Identification** For each property, applied the selected patterns to all 100K pages obtained from Bing. The claims identified were assessed manually:

- Overall precision: 60% over 7,092 claims
- Precision was better on properties with distinct value range
- Approach is limited due to looking at patterns between the country and the number, other context matters too.

**Claim Verification** No corpus of "lies" available. In our qualitative analysis we often found out-of-date statistics used.

Code and data: <https://github.com/uclmr/simpleNumericalFactChecker>