

Applied AI & The Fake News Problem

Mike Tamir Ph.D.
mike.tamir@berkeley.edu



Agenda

1. The Fake news Problem
 - a. Solving the right problem
2. Solving the Problem
 - a. Deep Learning Representation: Latest Techniques
 - b. How it works: Applying the Algorithm
3. The Algorithm in Action

The Fake News Problem

Fake News: The Problem

- Large players are having difficulty preventing malicious use of their services to propagate “fake news”
- Why?
 - **Solving the wrong problem!**
 - For Bigger Social Media and Search enterprises:
Stopping Fake News amounts to Censorship
- But Censorship is very serious:
 - Requires **Truth Detection**
 - Not just about the article
 - About **article ↔ world** correspondance

Fake News: The *Right* Problem

Image Source

Strong scientific evidence* that

Increased emotional content

**negatively impacts cognitive
function**



Practically perfect people never permit sentiment to muddle their thinking.

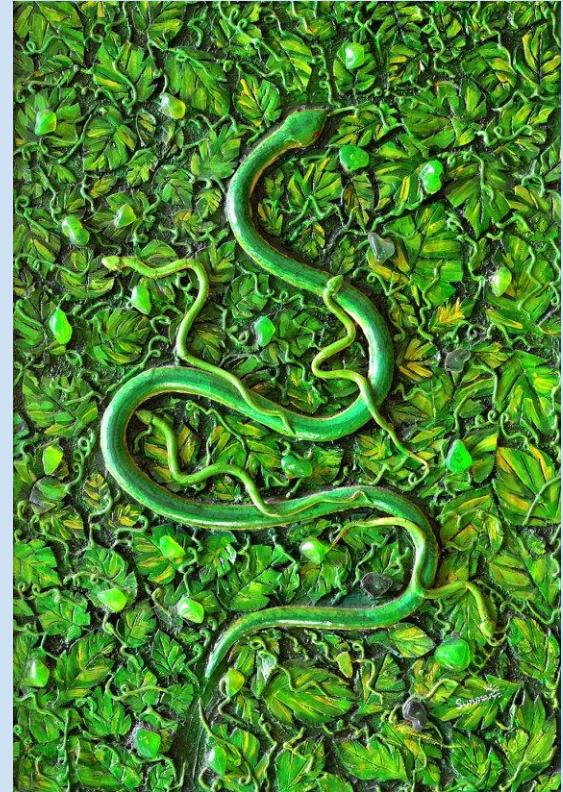
* Studies in the field of cognitive control emotion interactions in healthy adults.

Solving “Fake News” means...

Enable readers to detect when
the article's goal is

credible information sharing
vs.
emotional manipulation.

Image Source

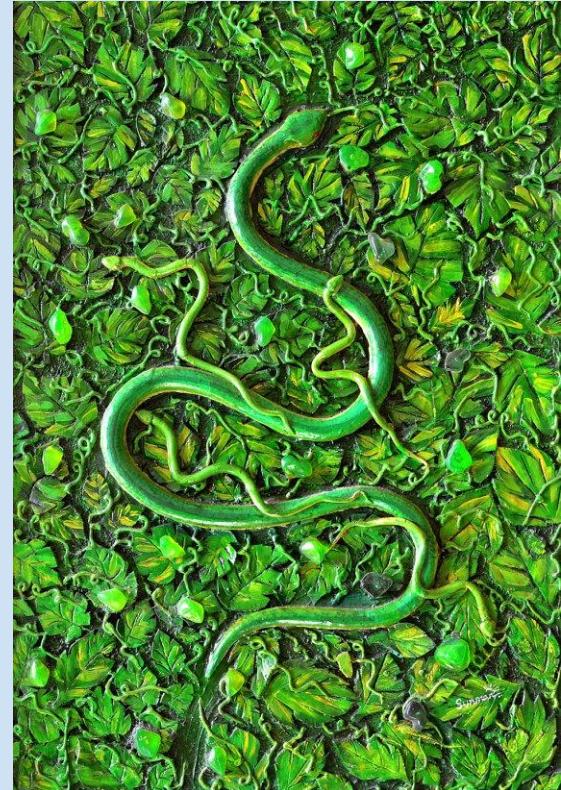


Solving “Fake News” means...

Image Source

Enable **cognitive engagement**

Alert and prevent
**manipulated emotional
reaction**

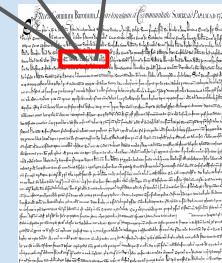
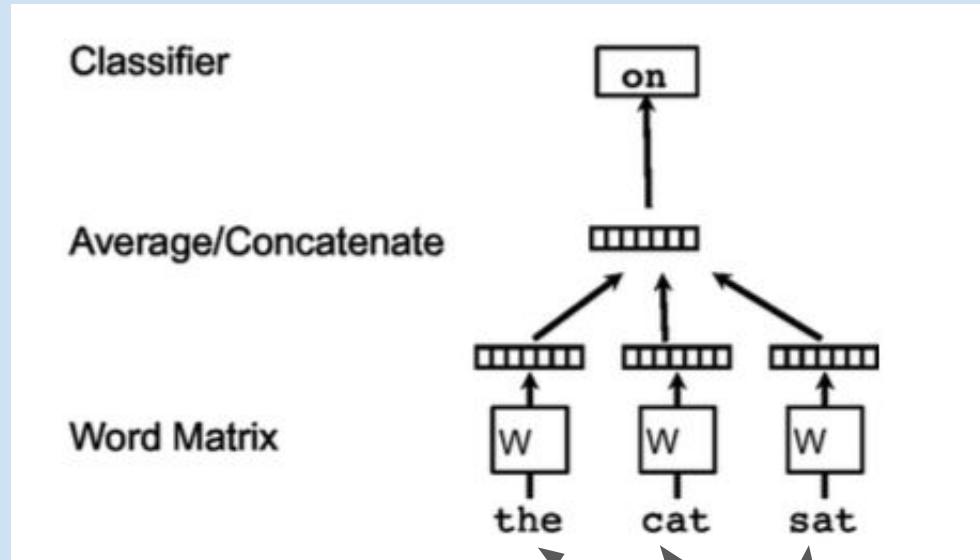


Solving the Problem: Deep Learning Techniques

Word2Vec: Algorithm

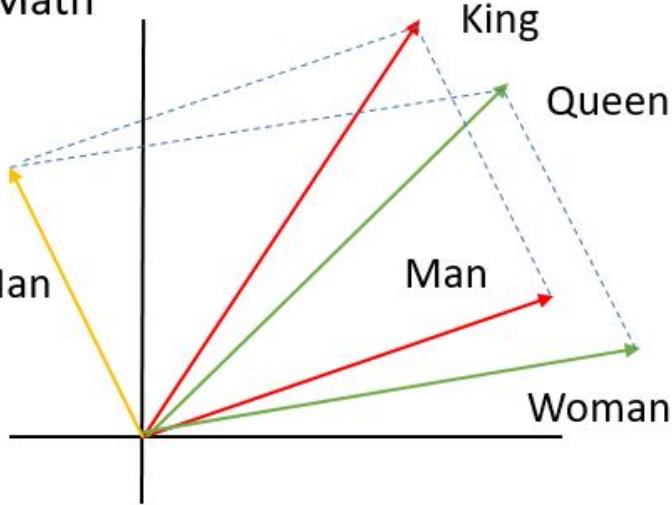
N-gram ANN classifier:

1. Project the “context” h_j
(w_{j-n+1} to w_{j-1})
2. Soft-max predictor for output layer
3. Use BackProp algorithm to execute gradient descent.



Word2Vec: Semantics in the Geometry

Vector Math



[Image from Mathworks](#)

- $|king\rangle - |man\rangle + |woman\rangle = |queen\rangle$
- $|England\rangle - |London\rangle + |Paris\rangle = |France\rangle$
- $|Microsoft\rangle - |Nadella\rangle + |Zuckerberg\rangle = |Facebook\rangle$
- ...

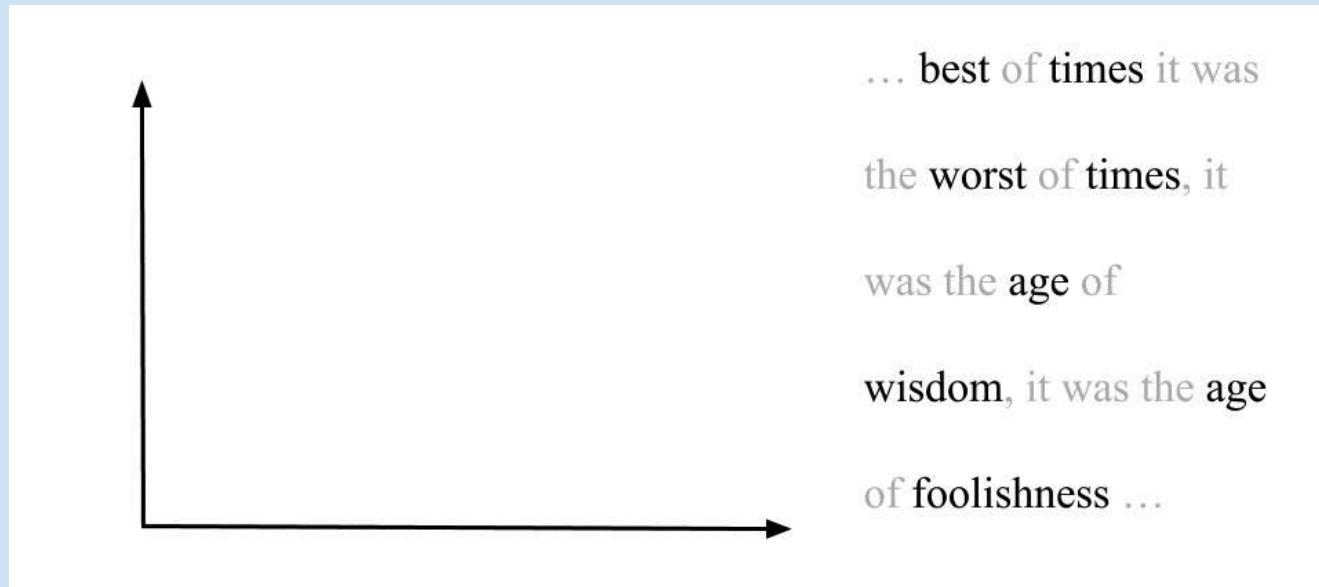
* While W2V was originally achieved with (shallow) Neural Networks, it has been shown that PMI based matrix factorization will generate similar results.

Popular Word Embedding Options:

- Word2Vec
 - [Mikolov et al \(2013\) Efficient Estimation of Word Representations in Vector Space](#)
 - [Google's Gensim Implementation](#)
- Glove
 - [Pennington et al \(2014\) - GloVe: Global Vectors for Word Representation](#)
 - [Stanford NLP implementation](#)
- FastText
 - [Joulin et al \(2016\) - Bag of Tricks for Efficient Text Classification](#)
 - [FacebookResearch implementation](#)

Naive Doc2Vec

Take the L2 normalization of the vector sum of terms in the document creates a **random walk** with **bias** process in the w2v space.

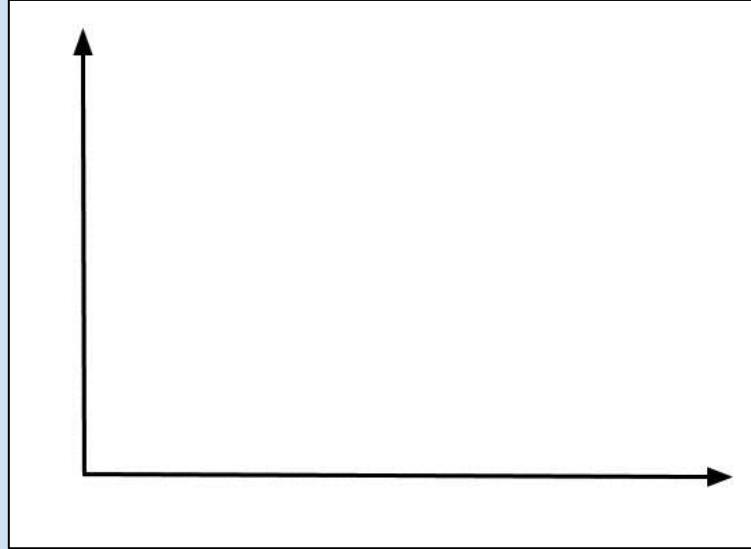


Naive Doc2Vec

ND2V: Easy to implement and captures large scale aggregate signal.

Drawbacks

1. High sensitivity to document length.
 - o Long document sums wash out stronger signal points in FN classification.
2. “Nativity” Assumption in raw Sum:
 - o Summation is commutative
 - o Cannot capture word order or negation context:
 - i. “Jill broke her crown but did not fall down.”
 - ii. “Jill fell down but did not break her crown.”



P2V (Non-Naive Doc2Vec)

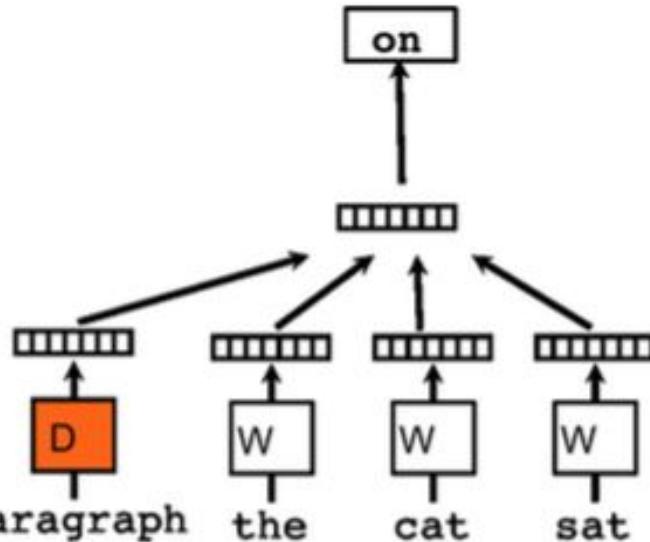
- P2V embed the full “paragraph” (doc) vector.
- Provides context for the whole paragraph.
- Tend to capture the semantics of the entire paragraph (doc).
- Some sensitivity to signal from word order.

Classifier

Average/Concatenate

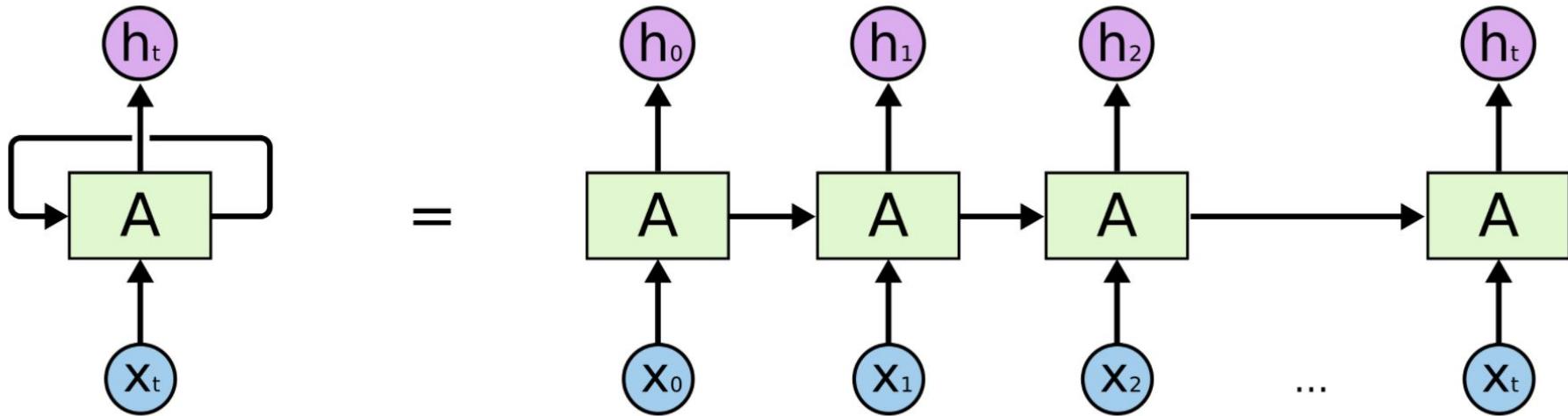
Paragraph Matrix----->

Paragraph id



Recurrent Neural Networks (RNNs)

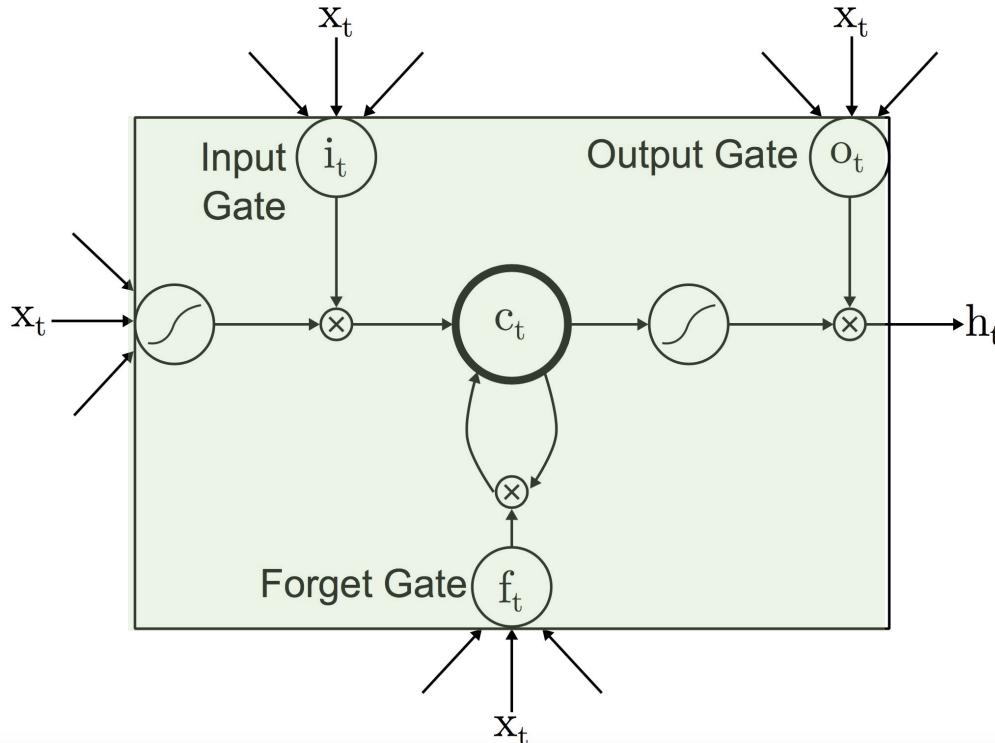
RNNs look at data points **in order**



[Image by Christopher Ohla](#)

Long Short Term Memory (LSTM)

- LSTM RNNs are designed to capture context and track what it should “remember” and when it should “forget”.



$$i_t = \sigma(x_t W_{xi} + h_{t-1} W_{hi} + b_i) \quad (1)$$

$$f_t = \sigma(x_t W_{xf} + h_{t-1} W_{hf} + b_f) \quad (2)$$

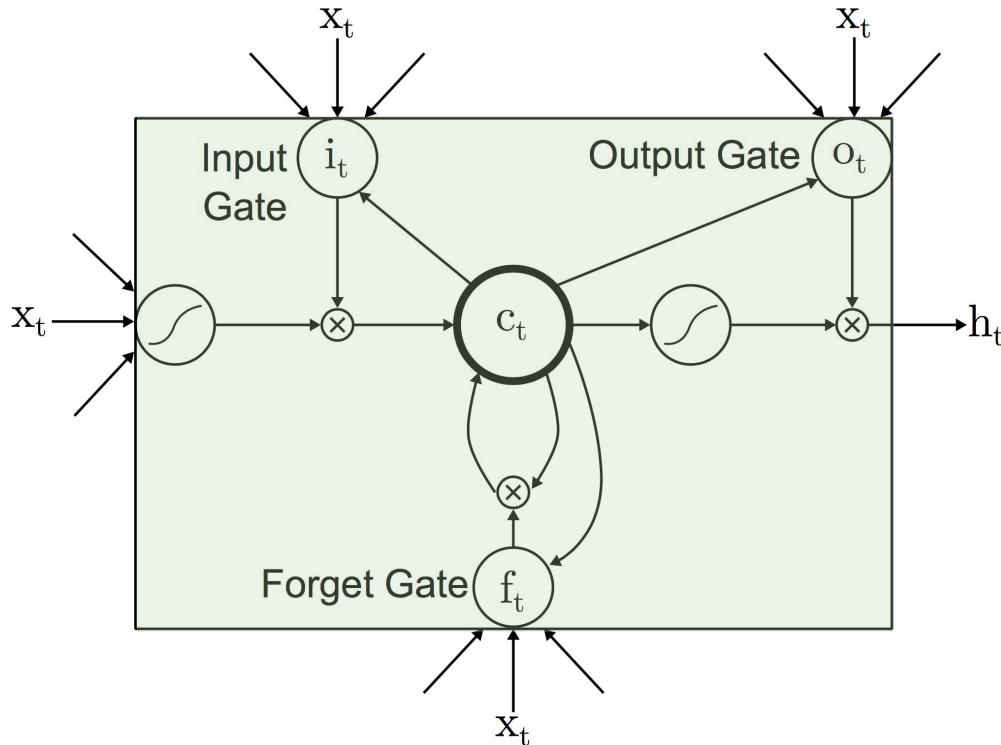
$$o_t = \sigma(x_t W_{xo} + h_{t-1} W_{ho} + b_o) \quad (3)$$

$$c_t = f_t \cdot c_{t-1} + i_t \sigma(x_t W_{xc} + h_{t-1} W_{hc} + b_c) \quad (4)$$

$$h_t = o_t \cdot \tanh(c_t) \quad (5)$$

LSTM with Peep-holes

- Many variants of LSTMs offer high flexibility in model development.



$$i_t = \sigma(x_t W_{xi} + h_{t-1} W_{hi} + b_i + w_{ci} c_{t-1}) \quad (1)$$

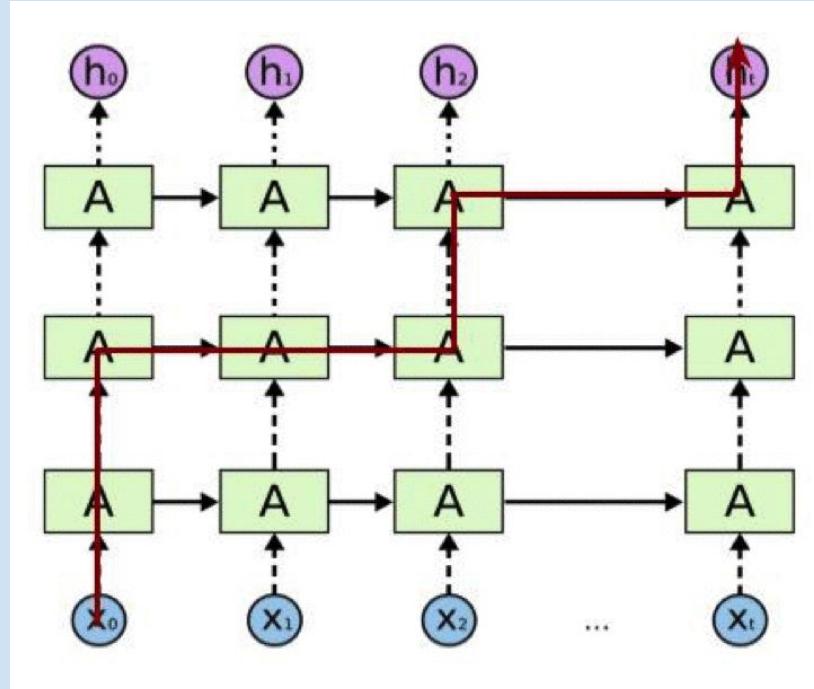
$$f_t = \sigma(x_t W_{xf} + h_{t-1} W_{hf} + b_f + w_{cf} c_{t-1}) \quad (2)$$

$$o_t = \sigma(x_t W_{xo} + h_{t-1} W_{ho} + b_o + w_{co} c_{t-1}) \quad (3)$$

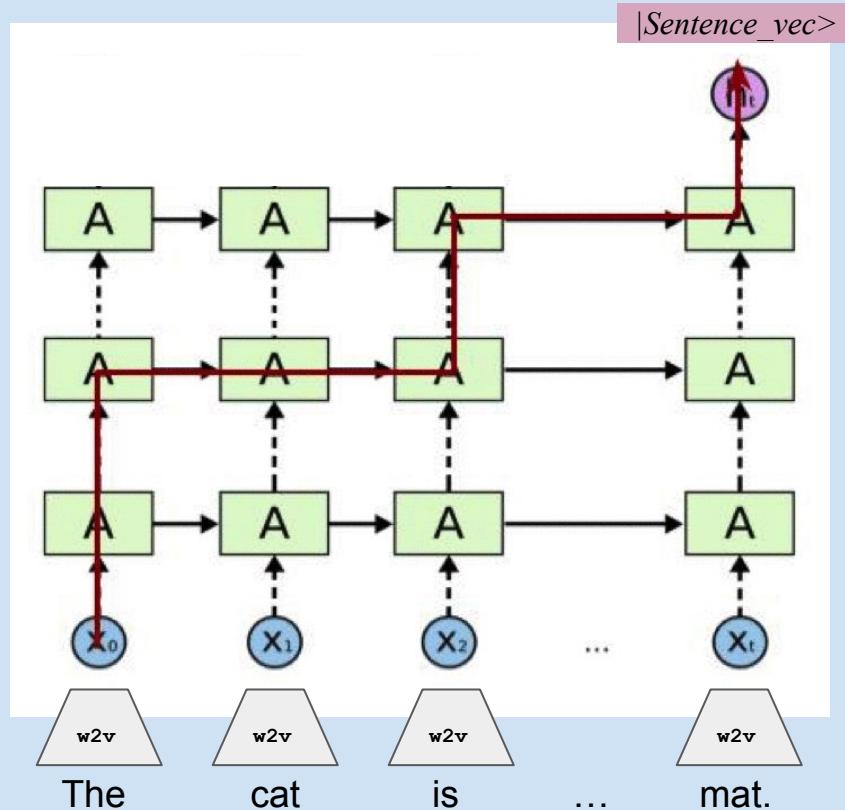
$$c_t = f_t \cdot c_{t-1} + i_t \sigma(x_t W_{xc} + h_{t-1} W_{hc} + b_c) \quad (4)$$

$$h_t = o_t \cdot \tanh(c_t) \quad (5)$$

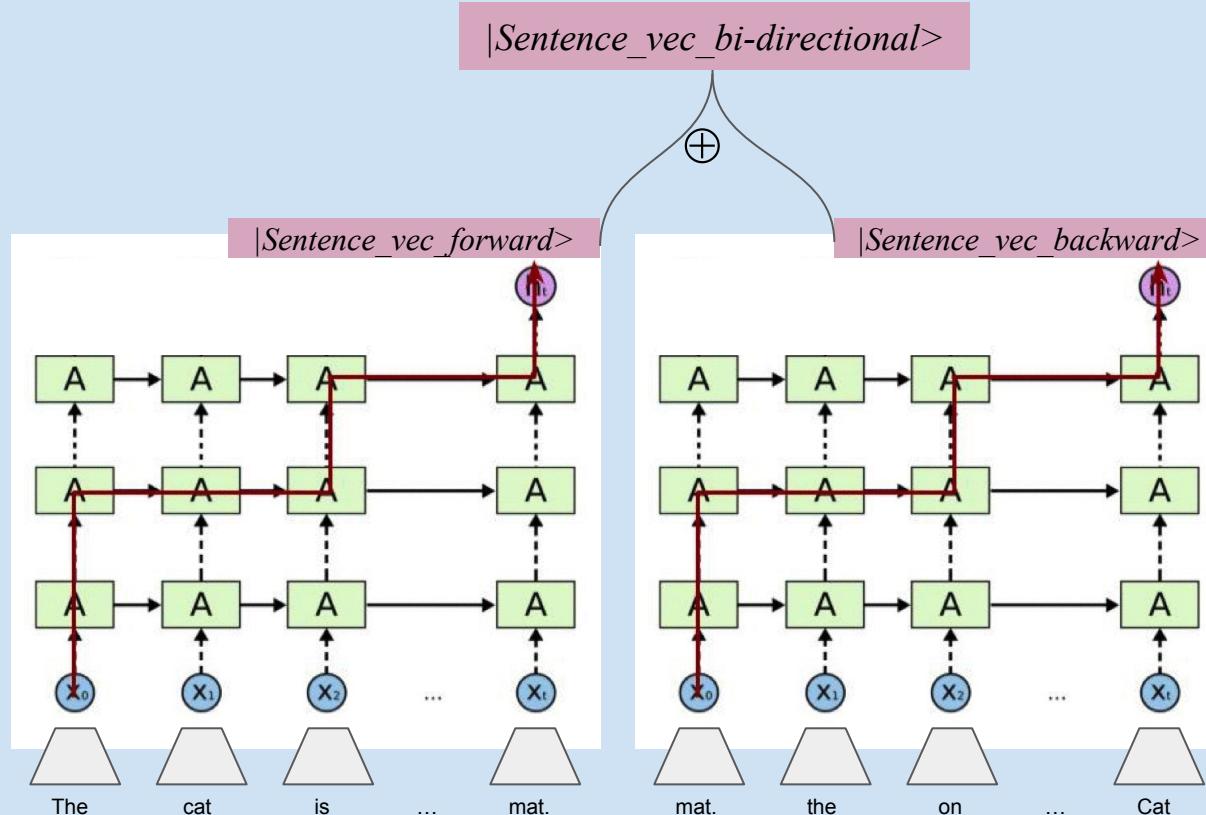
LSTM Stacking



LSTM Sentence-Vec Encoding



Bi-LSTM Encoding



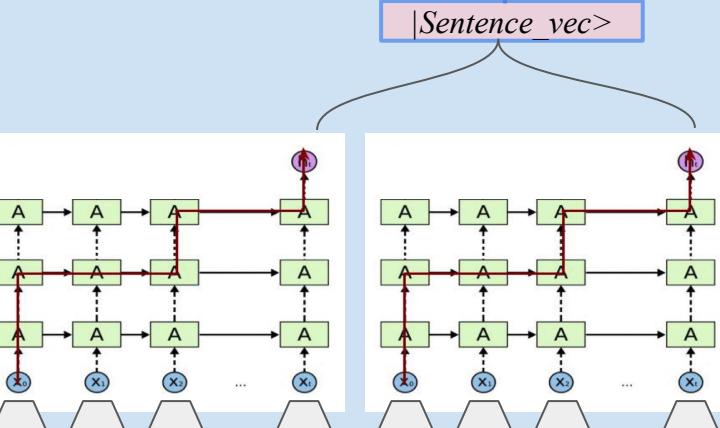
seq2seq Encoding-Decoding (Text Generation)

Bi-LSTMs encode semantics of the entire sentence.

Strong performance in several Use cases including...

- **Text generation**
- Translation
- Summarization...

Decoder LSTM



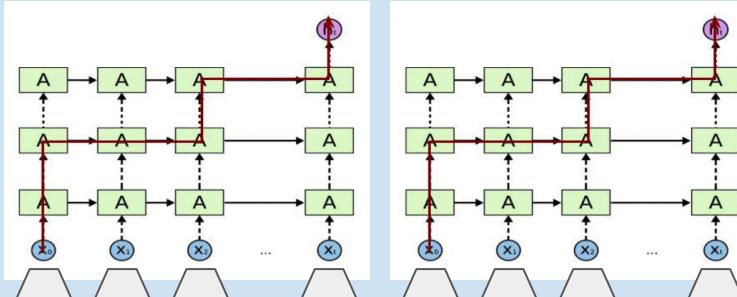
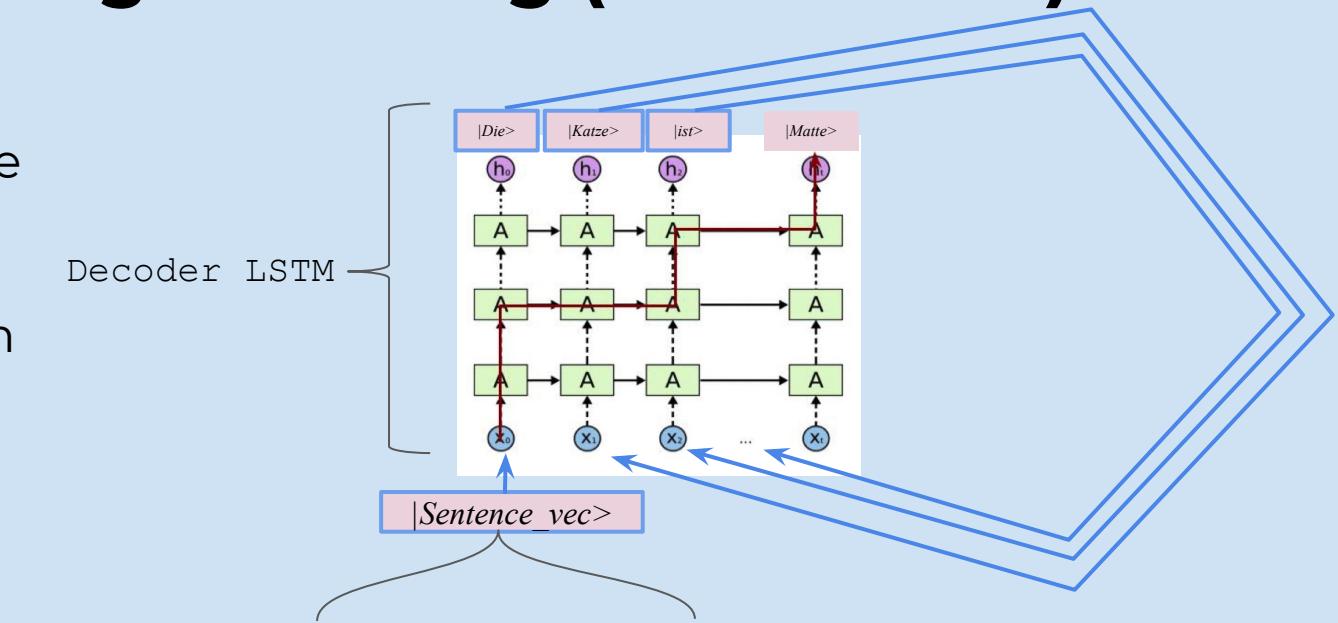
seq2seq Encoding-Decoding (Translation)

Bi-LSTMs encode semantics of the entire sentence.

Strong performance in several Use cases including...

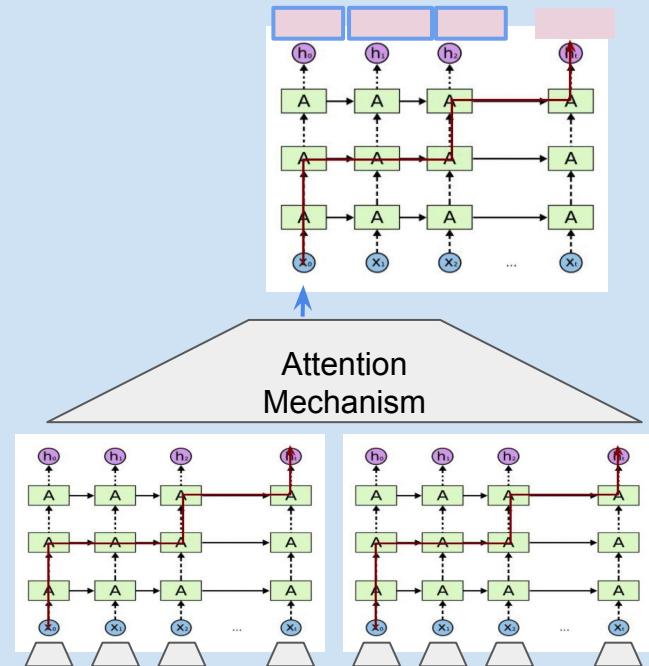
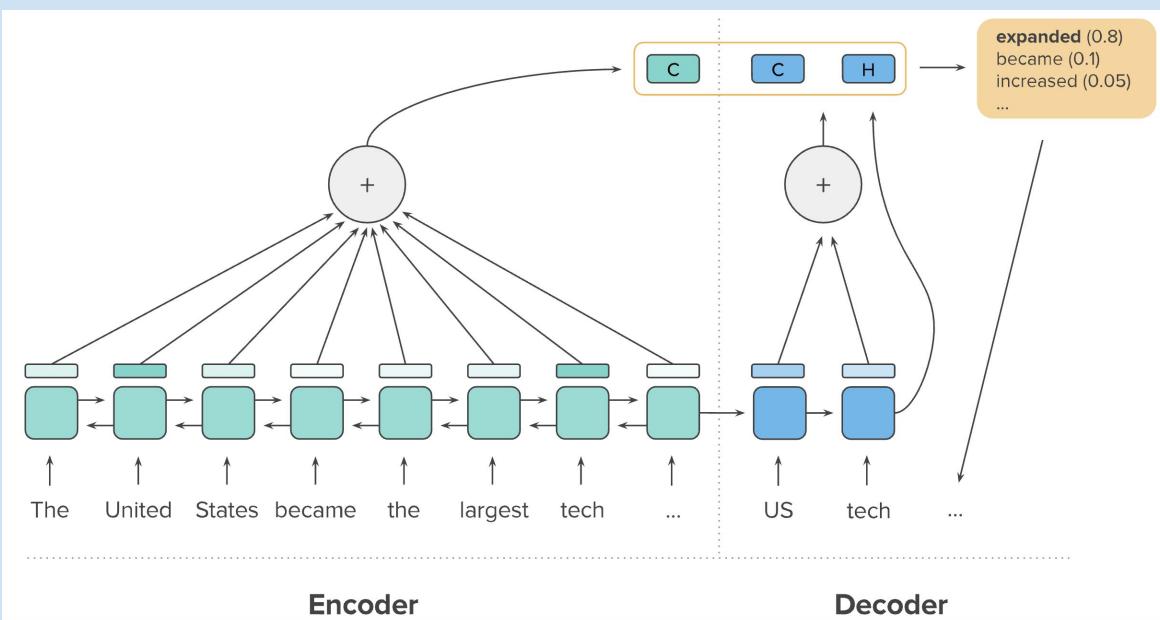
- Text generation
- **Translation**
- Summarization...

Decoder LSTM

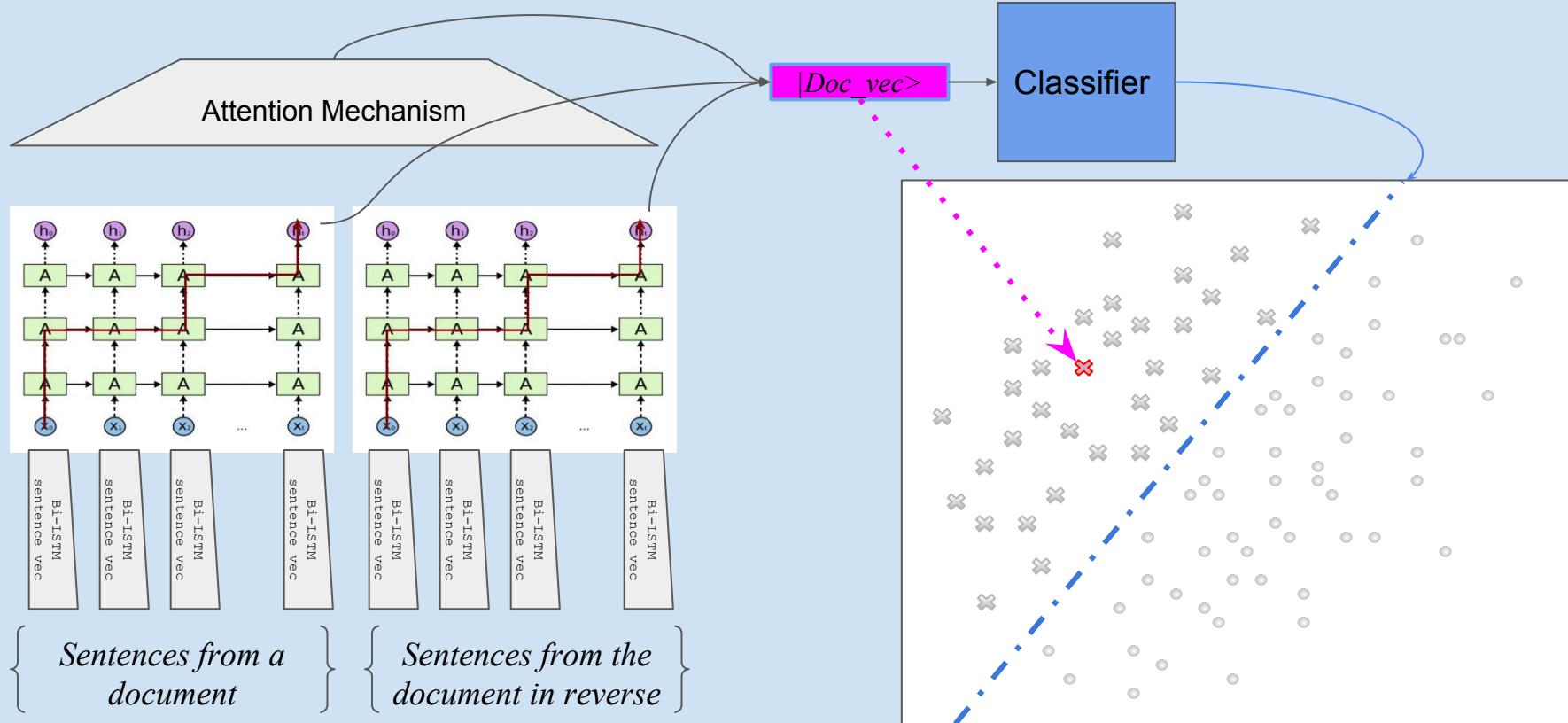


seq2seq Attention (Summarization)

- Attention mechanisms look at the outputs from each step of the Bi-LSTM network and learn to leverage the entire text sequence.



Bi-LSTM doc2vec Classification



Solving the Problem: Applying the Algorithm

(Remember) Solving “Fake News” means...

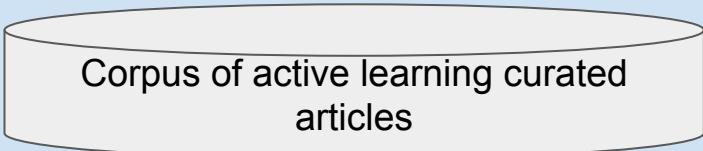
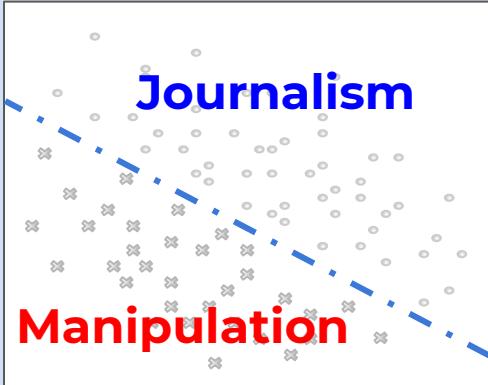
Image Source

Enable readers to detect when
the article's goal is

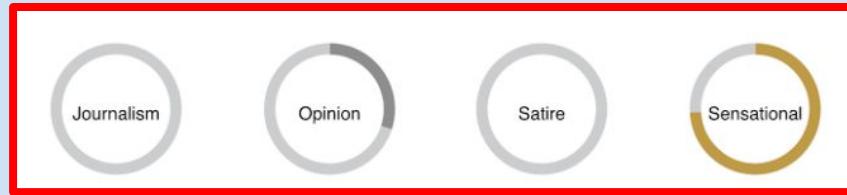
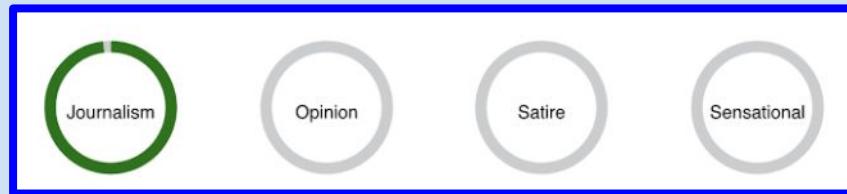
credible information sharing
vs.
emotional manipulation.



Applying the Algorithm



1. **Training:** $O(10^6)$ of active learning curated articles for clear signal.
2. Use DL techniques to generate rich vector representation of articles.
3. Classification:



Algorithm in Action



FakerFact

Secure | https://www.fakerfact.com

FakerFact

ABOUT LOGOUT

Fake News? Find out!

URL

CHECK Q

Download the Browser Extensions
Let FakerFact be there for you wherever you go on the web

<https://chrome.google.com/webstore/detail/fakerfact-plugin/hmcmekfmgfmlmmnicpmkfkcggnfgeef>



Screenshot of the Chrome Web Store search results for "fakerfact".

The search bar at the top shows "Secure | https://chrome.google.com/webstore/search/fakerfact".

The left sidebar contains a search input field with "fakerfact" and navigation links like "Home", "Extensions", "Themes", "FEATURES", and "RATINGS".

The main content area is titled "Extensions" and shows "1 of 1 Extension Results".

The result is the "FakerFact" extension by jeff. It has a 5-star rating and the description: "This extension will open up a new tab on www.fakerfact.com to check your current tab's content for indications of fake news".

A red circle highlights the "ADD TO CHROME" button, which is located next to the extension's name and rating.



Trump Readies Sweeping Tariff x

Secure | https://www.nytimes.com/2018/03/15/us/politics/trump-china-trade-measures.html?hp&action=click&pgtype=Homepage&clickSource=story-heading&... Mike

POLITICS Trump Readies Sweeping Tariffs and Investment Restrictions on China

By ANA SWANSON MARCH 15, 2018

Trump Readies Sweeping Tariffs and Investment Restrictions on China

432

f t e m r b 432



A factory for the smartphone maker Oppo in Dongguan, China. The tariffs could extend to more mundane products, including consumer electronics, apparel and even shoes. Nicolas Asfouri/Getty Images

The Trump White House
The historic moments, head-spinning developments and inside-the-White House intrigue.

Newly Emboldened, Trump Says What He Really Feels MAR 18

Before Saudi's Visit, Administration Implores Congress Not to Block Military Aid MAR 18

Army, Struggling to Get Technology in Soldiers' Hands, Tries the Unconventional MAR 18

Andrew McCabe, Fired F.B.I. Deputy, Is Said to Have Kept Memos on Trump MAR 17

Smuggling of U.S. Technology Is Outpacing Cold War Levels, Experts Say MAR 17

See More »



Trump Readies Sweeping Tariff × FakerFact

Secure | https://www.fakerfact.com/walt-says/5686812383117312

Mike

ABOUT LOGOUT

URL https://www.nytimes.com/2018/03/15/us/politics/trump-china-trade-measures.html?hp&action=click

CHECK Q

Walt says this one could be credible. But Walt always likes to verify the facts with more than one source.

Journalism Opinion Satire Sensational Outlandish Hate

> View Details



**Every result has
a “Walt says”
summary of the
scores.**

Secure | https://www.fakerfact.com/walt-says/5686812383117312

FakerFact

ABOUT LOGOUT

URL https://www.nytimes.com/2018/03/15/us/politics/trump-china-trade-measures.html?hp&action=click

CHECK Q

Walt says this one could be credible. But Walt always likes to verify the facts with more than one source.

Journalism Opinion Satire Sensational Outlandish Hate

View Details

Detailed analysis
of how the
article scores in
several
dimensions of
news credibility.

Trump Readies Sweeping Tariff × FakerFact x Mike

Secure | https://www.fakerfact.com/walt-says/5686812383117312

FakerFact

ABOUT LOGOUT

URL https://www.nytimes.com/2018/03/15/us/politics/trump-china-trade-measures.html?hp&action=click CHECK Q

Walt says this one could be credible. But Walt always likes to verify the facts with more than one source.

Journalism Opinion Satire Sensational Outlandish Hate

▶ View Details

A screenshot of a web browser displaying the FakerFact website. The URL in the address bar is https://www.fakerfact.com/walt-says/5686812383117312. The page title is "FakerFact". On the right side, there are links for "ABOUT" and "LOGOUT". Below the title, there is a search bar with the URL and a blue "CHECK Q" button. The main content area contains a statement: "Walt says this one could be credible. But Walt always likes to verify the facts with more than one source." Below this statement are six circular icons, each representing a dimension of news credibility: Journalism (green), Opinion (grey), Satire (grey), Sensational (grey), Outlandish (grey), and Hate (grey). Each icon has a thumbs up and a thumbs down icon underneath it. At the bottom of the content area is a link labeled "View Details". A large white box on the left side contains the text "Detailed analysis of how the article scores in several dimensions of news credibility.", with a thin black arrow pointing from the word "Journalism" towards this text.

Trump Readies Sweeping Tariff × FakerFact x Mike

Secure | https://www.fakerfact.com/walt-says/5686812383117312

FakerFact ABOUT LOGOUT

URL https://www.nytimes.com/2018/03/15/us/politics/trump-china-trade-measures.html?hp&action=click CHECK Q

Walt says this one could be credible. But Walt always likes to verify the facts with more than one source.

Journalism Opinion Satire Sensational Outlandish Hate

Like Dislike Like Dislike Like Dislike Like Dislike Like Dislike

I agree with Walt's opinion about whether this article is Journalism

Opportunity for
feedback on any
results for model
fine tuning.



Stop Apologizing for Being Elite x

Secure | https://www.nytimes.com/2018/03/16/opinion/sunday/elites-stop-apologizing.html?action=click&pgtype=Homepage&clickSource=story-heading&module...

Mike

SECTIONS HOME SEARCH

The New York Times

SundayReview OPINION

Stop Apologizing for Being Elite

By SUSAN JACOBY MARCH 16, 2018

f t e m r



RECENT COMMENTS

MidwesternReader 6 hours ago
Susan Jacoby's column disappoints in her failure to address the following: Elites and their views are not getting rejected. The elites are...

GT 6 hours ago
I have met many well educated -- very well paid individuals in NYC who are completely clueless of what goes on in the rest of the USA. ...

David 6 hours ago
"endless self-flagellation among well-educated liberals — "the elites," in pejorative parlance — about their failure to "get" the concerns...

SEE ALL COMMENTS



Stop Apologizing for Being Elite | FakerFact

Secure | https://www.fakerfact.com/walt-says/5740240702537728

Mike

ABOUT LOGOUT

URL https://www.nytimes.com/2018/03/16/opinion/sunday/elites-stop-apologizing.html?action=click&pgt

CHECK Q

Walt thinks this one is more about opinions than Journalism.

Journalism Opinion Satire Sensational Outlandish Hate

> View Details



BREAKING: Male Democratic C x

Secure | https://www.dailystrike.com/news/28350/breaking-male-democratic-operative-arrested-ryan-saavedra

Mike

DAILYWIRE News Podcasts Login SUBSCRIBE

DONALD TRUMP • HOLLYWOOD • BARACK OBAMA • GUNS • GUN CONTROL • DIANNE FEINSTEIN

BREAKING: Male Democratic Operative Arrested For Assaulding Female Trump Official

By RYAN SAAVEDRA
@RealSaavedra

March 16, 2018
70.7k views



Tony Williams/CQ Roll Call via Getty Images

HOT WIRE

- 1 Cate Blanchett: Sandra Bullock And I Got 'Penis Facials'
By HANK BERRIEN
- 2 BOMBSHELL: Obama's DOJ Forced Deletion Of 500,000 Fugitives From Gun Background Check System
By RYAN SAAVEDRA
- 3 Trump Challenges Mueller On Twitter: End The Probe
By EMILY ZANOTTI
- 4 NYT Columnist Pens Stupidest Column In Recent History: 'Go Ahead, Millennials, Destroy Us!'
By BEN SHAPIRO
- 5 The Women's March Is Losing Members After Aligning With Anti-Semitic Louis Farrakhan

A red circle highlights the "Advertisements OFF" button in the top right corner of the browser window.



A screenshot of a web browser window showing the FakerFact website. The browser's address bar displays a secure connection to <https://www.fakerfact.com/walt-says/5646748928180224>. The FakerFact logo is at the top left, and the user "Mike" is logged in at the top right. Below the header is a dark blue navigation bar with "ABOUT" and "LOGOUT" links. A URL input field contains the link from the address bar, and a blue "CHECK" button is to its right. The main content area features a bold, centered text: "Walt says this one grabs your attention, but other articles are probably better to get just the facts." Below this text are six circular icons representing different article types: Journalism, Opinion, Satire, Sensational, Outlandish, and Hate. Each icon has a "like" and "dislike" button underneath it. At the bottom left, there is a link to "View Details".