



INFORMATION EXTRACTION ENGINE FOR SENTIMENT-TOPIC MATCHING

IN PRODUCT INTELLIGENCE APPLICATIONS

CORNELIA FERNER | INTERNATIONAL DATA SCIENCE CONFERENCE | SALZBURG 2017

WERNER POMWENGER | MARTIN SCHNÖLL | VERONIKA HAAF | ARNOLD KELLER | STEFAN WEGENKITTL

MOTIVATION

Lenovo ThinkPad X1 Carbon 20A7003DUS Ultrabook

★★★★☆ PCWORLD RATING

Almost everything about the new X1 Carbon is better than the original. But the company's engineers should have left most of the keyboard alone. We can also do without its speech- and gesture-recognition.

6 COMMENTS



Michael Brown | Executive Editor, PCWorld | Apr 2, 2014 3:00 AM

meta data

Update, 4/2/2014: Lenovo provided us with incorrect pricing on this eval unit. As configured, the correct price is \$1609, not \$1529 as originally reported.

Lenovo took its sweet time cooking up Haswell versions of its flagship business Ultrabook line. While there are signs the 2014 ThinkPad X1 Carbon spent too much time in the test kitchen, it's still the best notebook I've laid hands on.

This new X1 Carbon is thinner and lighter than the original, which was first introduced in 2012 and updated (with a touchscreen) in early 2013. Yet it has a higher-resolution display, more I/O ports, an improved docking-station option, and an entirely new "adaptive function key" row. Lenovo's cooks should have served it up that way, instead of going on to mess with its keyboard layout and lard it with half-baked voice- and gesture-control features.

Lenovo is justifiably proud of its innovative "adaptive function row."

topics

Let's go over the new machine's many positive attributes first, starting with its display. Lenovo sent an eval unit equipped with an Intel Core i5-4300U processor, 4GB of low-power DDR3/1600 memory, and a skimpy 180GB SSD. Its 14-inch IPS touchscreen is delightful, boasting a resolution of 2560 by 1440 pixels. Pixel density is 210 pixels per inch (PPI), just shy of Apple's 15-inch MacBook Pro with Retina display (which has a resolution of 2560 by 1800 and a pixel density of 220 PPI).

facts

The screen is very bright, responsive to touch, and it delivers excellent contrast in all lighting situations. As with the previous model, its hinges allow for 180 degrees of rotation. Lenovo also offers less-expensive non-touch configurations and models with 1600-by-900-pixel displays.

sentiment

AGENDA

- ARIE – article and review information extraction engine
- Topic classification
- Sentiment analysis
- Recap

ARIE

Update, 4/2/2014: Lenovo provided us with incorrect pricing on this eval unit. As configured, the correct price is \$1

Lenovo took its sweet time cooking up Haswell versions of its flagship business Ultrabook line. While there are signs the 2014 ThinkPa

This new X1 Carbon is thinner and lighter than the original, which was first introduced in 2012 and updated (with a touchscreen) in ea

“adaptive function key” row. Lenovo’s cooks should have served it up that way, instead of going on to mes

ard layout an

uters/the-bes

Lenovo is justifiably proud of its innovative “adaptive function row.”

Let’s go over the new machine’s many positive attributes first, starting with its display. Lenovo sent

delightful, boasting a resolution of 2560 by 1440 pixels. Pixel density is 210 pixels per inch (PPI),

ped with an I

15-inch MacB

The screen is very bright, responsive to touch, and it delivers excellent contrast in all lighting situ

le previous mo

1600-by-900-pixel displays.

2014” width=“580” height=“388”/”><a class=“zoom” href=“http://images.techhive.com/images/article/2014/04/leno

2014” width=“580” height=“388”/”> <figcaption>

The 2014 X1 Carbon is thinner and lighter, has a higher-resolution touchscreen

Robust, but lightweight

The top half of the X1 Carbon is fabricated from carbon fiber, hence its na

inches) and it weighs only a bit more (3.15 pounds compared to 2.99 pounds)

More importantly, Lenovo managed to cram more I/O ports into the new design

docking station, either of which would consume one of its USB ports.

2014” width=“580” height=“388”/”><a class=“zoom” href=“http://images.techhive.com

2014” width=“580” height=“388”/”> <figcaption>

Can’t decide if you like DisplayPort or HDMI? The new X1 Carbon has both!

The 2014 model sacrifices the SD card slot, but it has two USB 3.0 ports, b

href=“http://www.pcworld.com/article/2042360/wi-fi-alliance-announces-802-1

next to the right arrow key.

A proprietary port links the notebook to Lenovo’s new OneLink docking station, so you won’t need to gi

up a USB port when you’re des

has four USB 3.0 ports (one with always-on charging), two USB 2.0 ports, gigabit ethernet, and a stere

mic/headphone jack.

Lenovo is justifiably proud of its replacement for the row of keys at the top of the typical keyboard. Its “adaptive function row” cha

In place of mechanical keys that serve multiple functions depending on which other key you hold down, both the function keys

keys when you’re using a word processor, a different set when you’re using a web browser, and so on. According to Lenovo, this is made

electroluminescent layer beneath a Gorilla Glass panel to deter scratches.

2014” width=“580” height=“388”/”><a class=“zoom” href=“http://images.techhive.com/images/article/2014/04/lenovocarbonx1_5-100258974-orig.jpg

2014” width=“580” height=“388”/”> <figcaption>

Lenovo’s innovative “adaptive function row” changes according to the application you’re using.

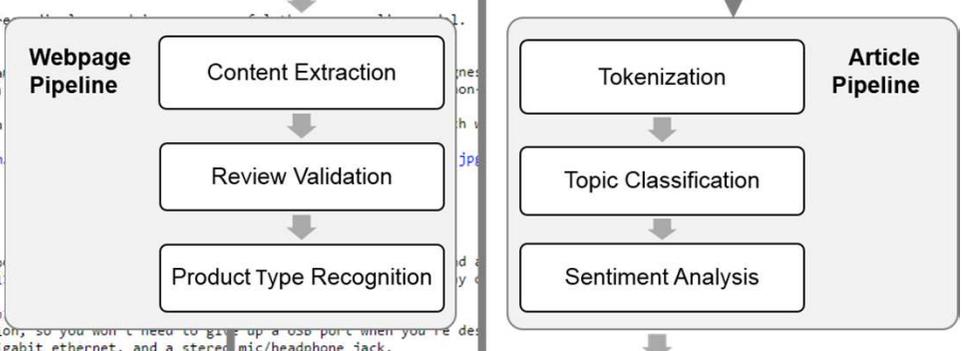
Lenovo’s

height=“505”/”><img sr

height=“505”/”> <figcaption>

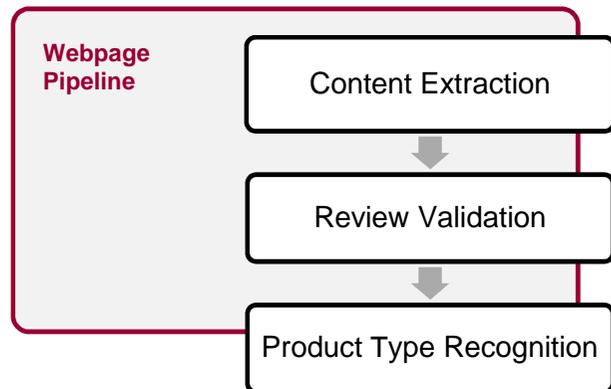
This mock-up shows all four of the function rows and which applications activate them.

Target Websites



Annotated Articles

WEBPAGE PIPELINE



■ Content extraction:

- Boilerplate removal (comments, ads, teasers etc.)
- Raw text extraction (without html tags)
- Store meta data

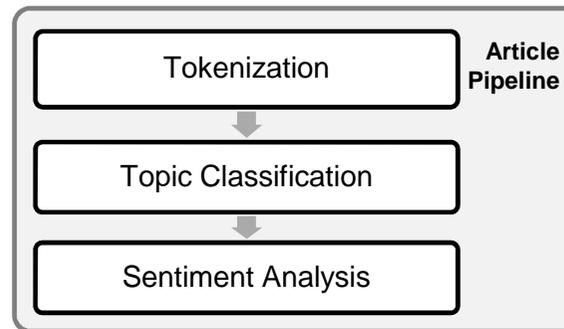
■ Review validation:

- Only expert reviews are needed
- Sort out ads, comparisons etc.
- Latent Dirichlet Allocation + SVM

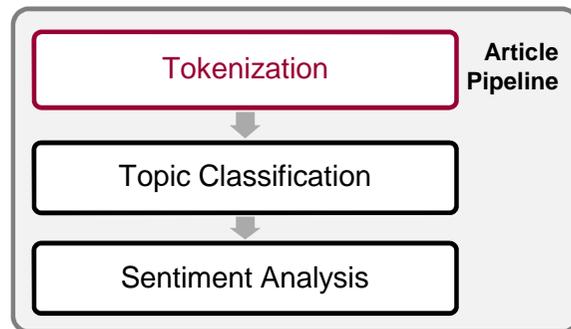
■ Product type recognition:

- Only laptops (e.g. ultrabooks, convertibles) are needed
- Sort out reviews on displays, speakers etc.
- Maximum Entropy (logistic regression for multiclass problems)

ARTICLE PIPELINE



ARTICLE PIPELINE



■ Tokenization:

- Words and sentences
- Sentence-level annotations for topic classification and sentiment analysis

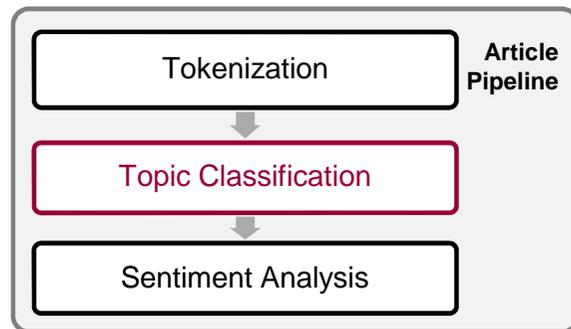
■ Preprocessing:

- Lowercasing
- No stopwords, no stemming, no removal of infrequent/frequent words

■ Features:

- Bag-of-words (also with bigrams)
- Word2Vec

ARTICLE PIPELINE



- $D = \{1, \dots, d\}$... dictionary; set of all words
- $C = \{1, \dots, c\}$... set of topics
- $S \in C$ and $W_t \in D$... topics and sequence of words, respectively

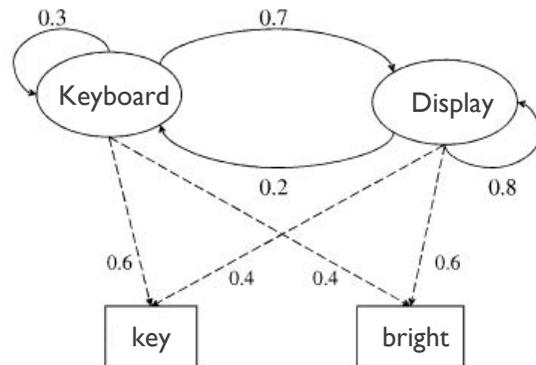
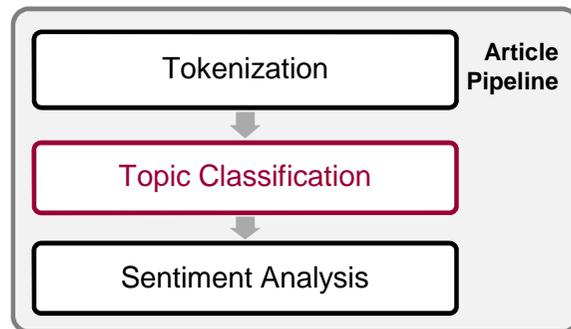
- MaxEnt (logistic regression, softmax):

- $\vec{V} = (v_1, \dots, v_d)$... count vector (absolute word counts in a sequence) with $v_i = \sum_{t=1}^n \mathbf{1}(W_t = i)$

- $$P[S = j | \vec{V}] = \frac{e^{\sum_{i \in D} v_i \cdot \lambda_{ij} + \mu_j}}{Z} = \prod_{t=1}^n \frac{e^{\lambda_{W_t j} + \mu_j}}{\sum_{l=1}^c e^{\lambda_{W_t l} + \mu_l}}$$

$$= \prod_{t=1}^n P[S = j | W_t]$$

ARTICLE PIPELINE



■ Hidden Markov Model (HMM):

- decoding: find the most probable sequence of hidden states (topics) given the model and a sequence of observations (words)?

- $A = a_{ij} = P[S_{t+1} = j | S_t = i]$... transition probabilities

- $B = b_{jk} = P[W_t = k | S_t = j]$... emission probabilities

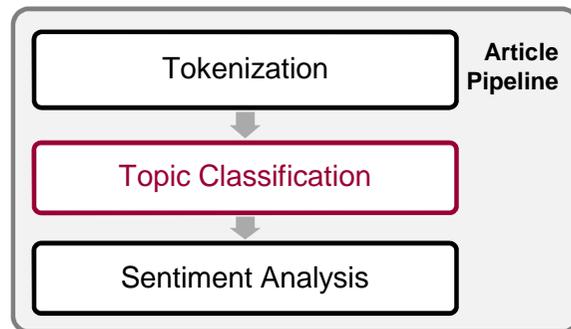
- $\pi_i = P[S_1 = i]$... initial state probabilities

- $M = (C, D, A, B, \pi)$

■ Combining MaxEnt and HMM (Bayes):

- $$b_{jk} = P[W_t = k | S_t = j] = \frac{P[S=j | W=k] \cdot P[W=k]}{\sum_{r=1}^d P[S=j | W=r] \cdot P[W=r]}$$

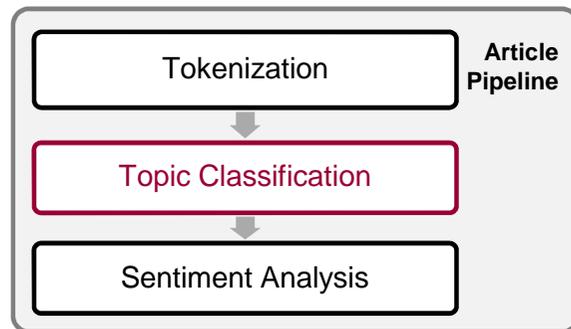
ARTICLE PIPELINE



■ Combining MaxEnt and HMM (Bayes):

$$■ \quad b_{jk} = \frac{P[S=j | W=k] \cdot P[W=k]}{\sum_{r=1}^d P[S=j | W=r] \cdot P[W=r]}$$

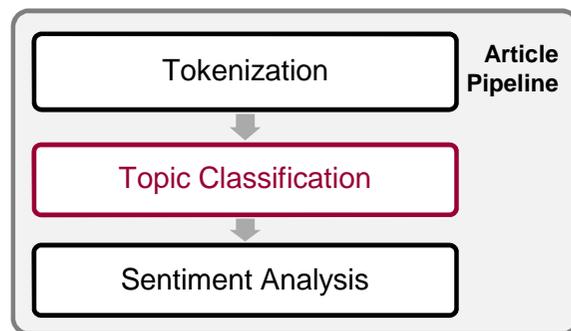
ARTICLE PIPELINE



■ Combining MaxEnt and HMM (Bayes):

$$■ \quad b_{jk} = \frac{P[S=j | W=k] \cdot \hat{p}_k}{\sum_{r=1}^d P[S=j | W=r] \cdot P[W=r]} \quad \hat{p}_k \cong P[W = k]$$

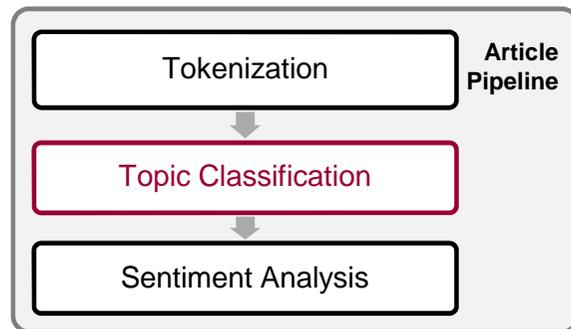
ARTICLE PIPELINE



■ Combining MaxEnt and HMM (Bayes):

$$■ \quad b_{jk} = \frac{P[S=j | W=k] \cdot \hat{p}_k}{f^{-1}(j)} \quad \hat{p}_k \cong P[W = k]$$

ARTICLE PIPELINE



■ Combining MaxEnt and HMM (Bayes):

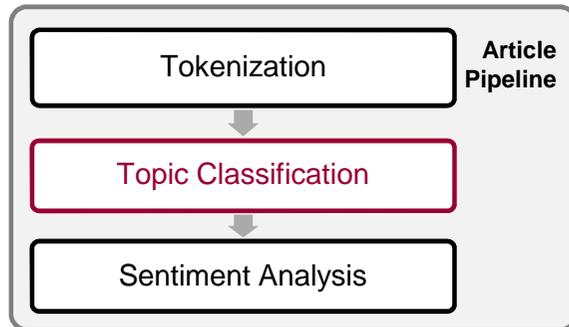
$$■ \quad b_{jk} = \frac{P[S=j | W=k] \cdot \hat{p}_k}{f^{-1}(j)}$$

$$\hat{p}_k \cong P[W = k]$$

$$■ \quad b_{jk} = \frac{e^{\lambda_{kj} + \mu_j} \cdot \hat{p}_k}{\sum_{l=1}^c e^{\lambda_{kl} + \mu_l}} \cdot f(j)$$

$$P[S = j | W_t] = \frac{e^{\lambda_{wtj} + \mu_j}}{\sum_{l=1}^c e^{\lambda_{wtl} + \mu_l}}$$

ARTICLE PIPELINE



■ Combining MaxEnt and HMM (Bayes):

$$■ \quad b_{jk} = \frac{P[S=j | W=k] \cdot \hat{p}_k}{f^{-1}(j)}$$

$$\hat{p}_k \cong P[W = k]$$

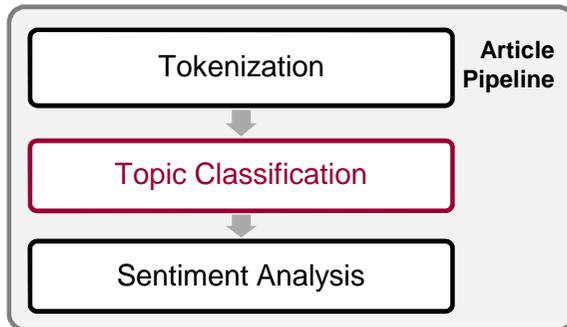
$$■ \quad b_{jk} = \frac{e^{\lambda_{kj} + \mu_j} \cdot \hat{p}_k}{\sum_{l=1}^c e^{\lambda_{kl} + \mu_l}} \cdot f(j)$$

$$P[S = j | W_t] = \frac{e^{\lambda_{wtj} + \mu_j}}{\sum_{l=1}^c e^{\lambda_{wtl} + \mu_l}}$$

$$■ \quad b_{jk} = g(j) \cdot f(j)$$

$$f(j) = \frac{1}{\sum_{i=1}^d g(i)}$$

ARTICLE PIPELINE



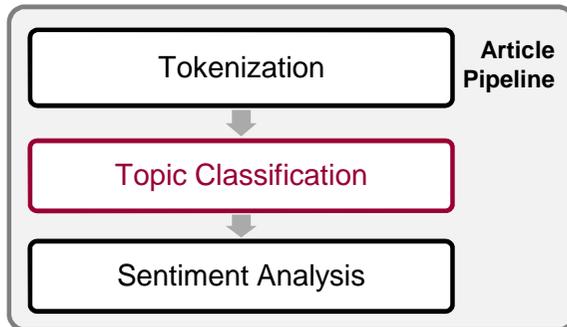
■ Combining MaxEnt and HMM (Bayes):

$$\blacksquare b_{jk} = \frac{P[S=j | W=k] \cdot \hat{p}_k}{f^{-1}(j)} \quad \hat{p}_k \cong P[W = k]$$

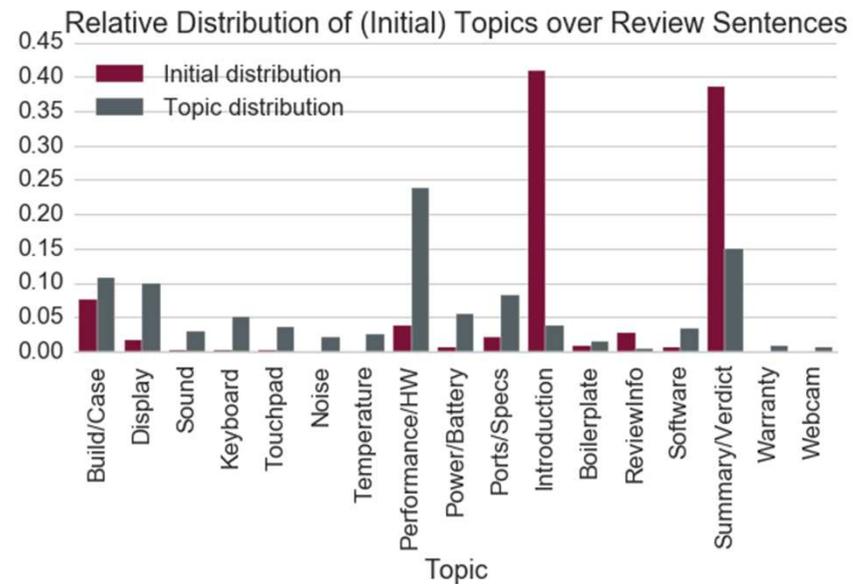
$$\blacksquare b_{jk} = \frac{e^{\lambda_{kj} + \mu_j} \cdot \hat{p}_k}{\sum_{l=1}^c e^{\lambda_{kl} + \mu_l}} \cdot f(j) \quad P[S = j | W_t] = \frac{e^{\lambda_{wtj} + \mu_j}}{\sum_{l=1}^c e^{\lambda_{wtl} + \mu_l}}$$

$$\blacksquare b_{jk} = g(j) \cdot f(j) \quad f(j) = \frac{1}{\sum_{i=1}^d g(i)}$$

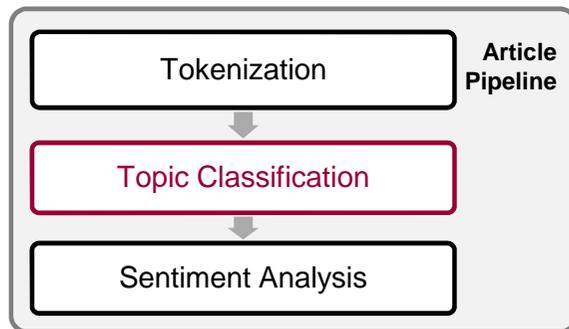
ARTICLE PIPELINE



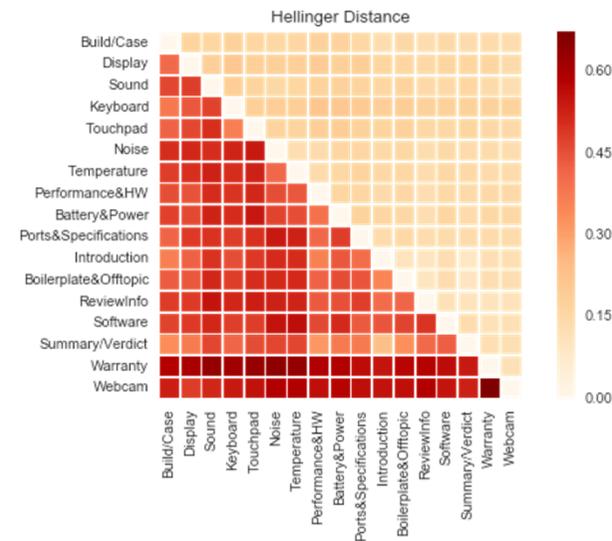
- 3152 reviews
- 240220 manually labelled sentences
- 17 predefined topics



ARTICLE PIPELINE



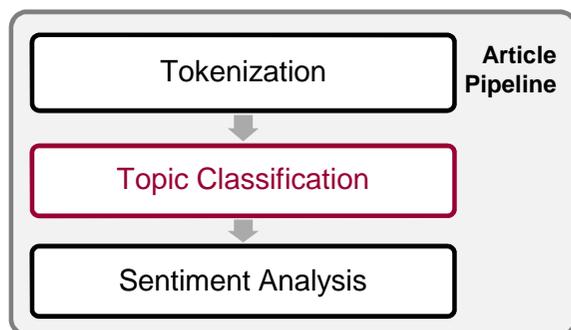
- Distance of the topics' MaxEnt distributions (bottom left) and the topics' word frequencies (top right).



- Hellinger distance:

$$H(S_1, S_2)^2 = \frac{1}{\sqrt{2}} \sqrt{\sum (\sqrt{S_1} - \sqrt{S_2})^2}$$

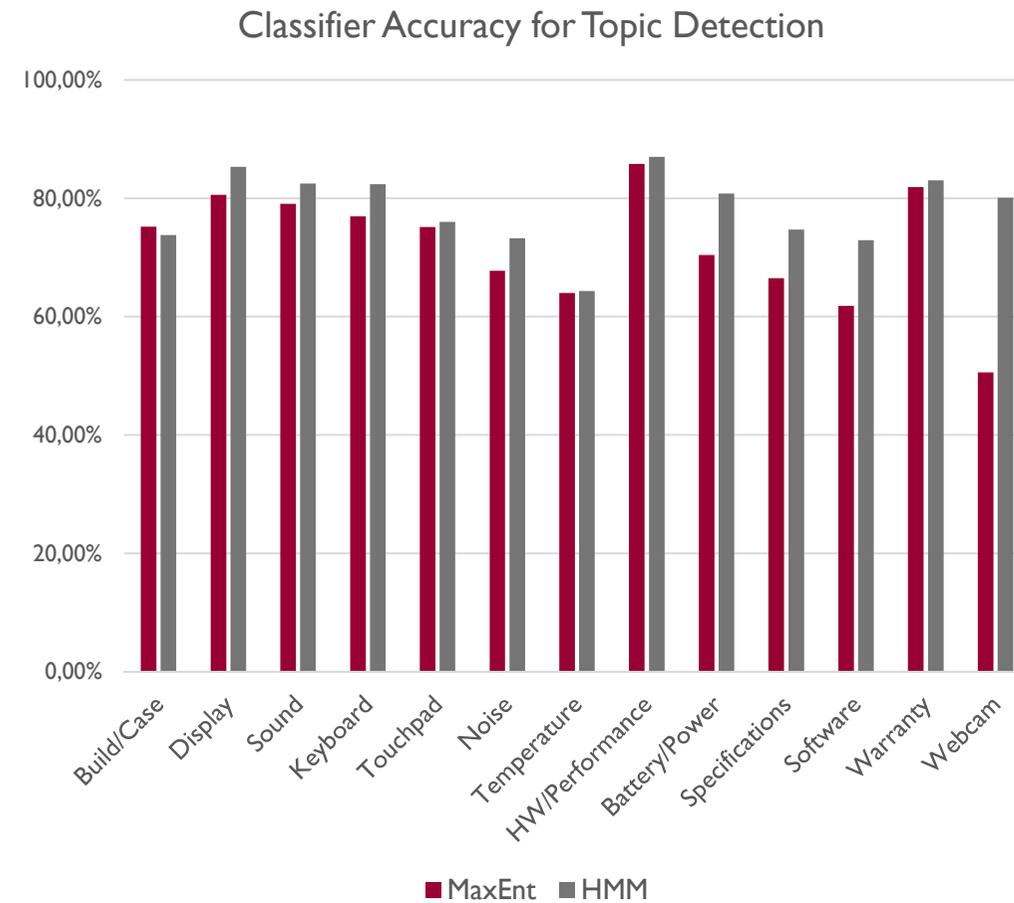
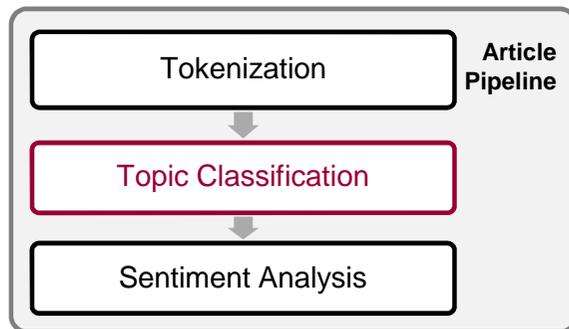
ARTICLE PIPELINE



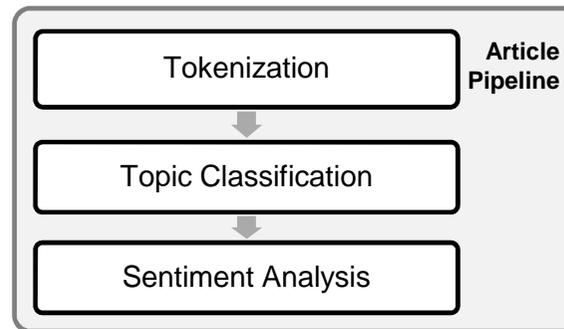
Sound	Noise	Temperature	Summary
sound	db	cool	verdict
music	quiet	heat	price
bass	noise	cooling	lasted
audio	fan	lap	watts
speakers	audible	temperature	read
and	db	and	of
the	to	to	and
speakers	and	degrees	for
sound	the	the	the
to	is	is	it

Table 1: Words with highest weights as trained by the MaxEnt classifier (above) and with maximum frequency (below) in four exemplary topics.

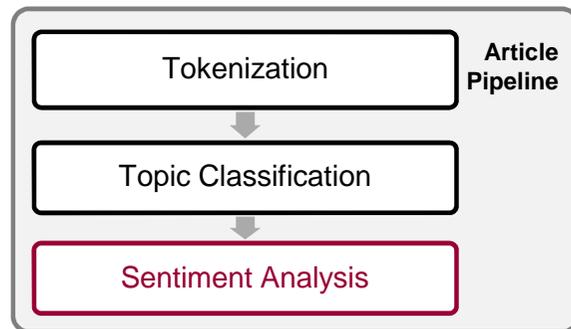
ARTICLE PIPELINE



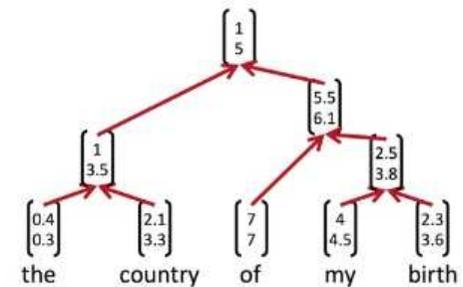
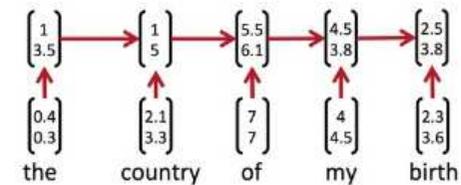
ARTICLE PIPELINE



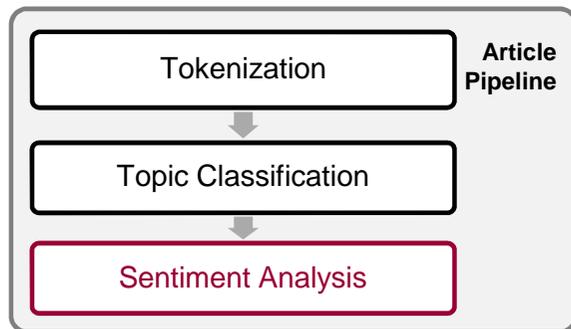
ARTICLE PIPELINE



- MaxEnt
 - baseline
- Recurrent neural network (RNN)
 - Best performance on sentence level
- Recursive neural tensor network (RNTN)
 - Require parsed syntax tree
 - Good on word level

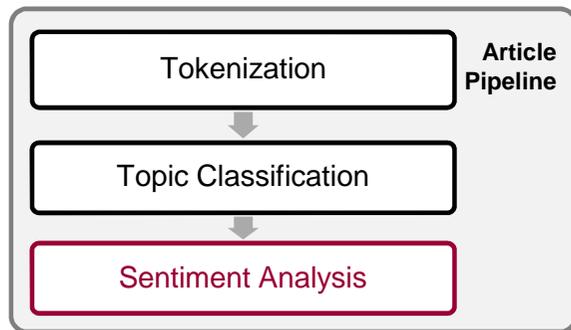


ARTICLE PIPELINE

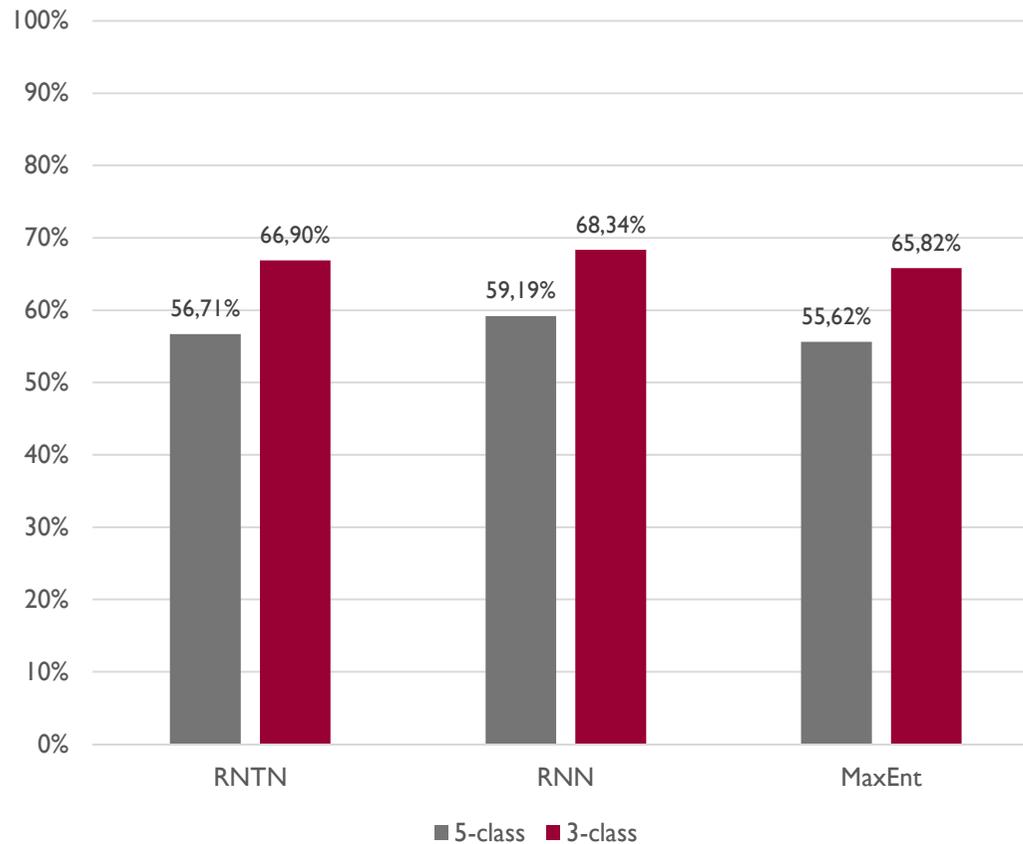


- 21695 manually labelled sentences
- 5-class analysis:
 - very positive
 - positive
 - neutral
 - negative
 - very negative
- 3-class analysis:
 - positive
 - neutral
 - negative

ARTICLE PIPELINE



Classifier Accuracy for Sentiment Analysis

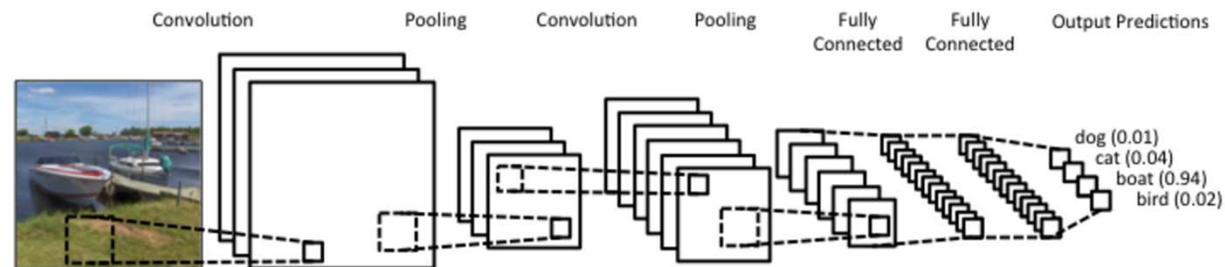


LESSONS LEARNT

- Language is ambiguous.
- Lack of standardized features or off-the-shelf (preprocessing) methods.
- There isn't only English.

ONGOING RESEARCH

- Representation learning with CNNs (deep neural network)
- Change fully connected layers to adapt net to new tasks
- Multiple languages: reuse network architecture
- Character-wise approach for reusability of representations



LITERATURE

- [MaxEnt] A. Berger, S. Della Pietra and V. Della Pietra, “A maximum entropy approach to natural language processing,” *Computational Linguistics*, vol. 22(1), pp. 39-71, 1996.
- [HMM] Cappé, O., Moulines, E., and Ryden, T., “Inference in Hidden Markov Models,” New York, Springer, 2005.
- [RNN] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9(8), pp. 1735-1780, 1997.
- [RNTN] Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C., Ng, A., and Potts, C., “Recursive deep models for semantic compositionality over a sentiment treebank,” *Conference on Empirical Methods in Natural Language Processing*, 2013.
- [UIMA] D. Ferrucci and A. Lally, “UIMA: An architectural approach to unstructured information processing in the corporate research environment,” *Natural Language Engineering*, 10(3), pp. 237-348, 2004.