

Website Categorisation

Joseph Atkins-Turkish Jake Hodges James Patrick-Evans

May 5, 2015

The Idea

The aim is to take a website and automatically assign it to a descriptive category label.

- ▶ Original idea to sort bookmarks into folders
- ▶ Categorisation of websites
 - ▶ Organising bookmarks
 - ▶ Automated SEO
 - ▶ Content filtering
- ▶ Supervised multi-class classification problem
 - ▶ Unsupervised element

Approach

1. Data Collection

- ▶ Compile corpus¹ of websites with pre-existing class labels

2. Data Sanitisation

- ▶ Process and sanitise raw HTML

3. Feature Extraction

- ▶ Find key features in corpus (LDA or tf-idf)

4. Classification

- ▶ Use generated features to categorise websites

¹a collection of documents

Approach

1. Data Collection

- ▶ Compile corpus¹ of websites with pre-existing class labels

2. Data Sanitisation

- ▶ Process and sanitise raw HTML

3. Feature Extraction

- ▶ Find key features in corpus (LDA or tf-idf)

4. Classification

- ▶ Use generated features to categorise websites

¹a collection of documents

Approach

1. Data Collection

- ▶ Compile corpus¹ of websites with pre-existing class labels

2. Data Sanitisation

- ▶ Process and sanitise raw HTML

3. Feature Extraction

- ▶ Find key features in corpus (LDA or tf-idf)

4. Classification

- ▶ Use generated features to categorise websites

¹a collection of documents

Approach

1. Data Collection

- ▶ Compile corpus¹ of websites with pre-existing class labels

2. Data Sanitisation

- ▶ Process and sanitise raw HTML

3. Feature Extraction

- ▶ Find key features in corpus (LDA or tf-idf)

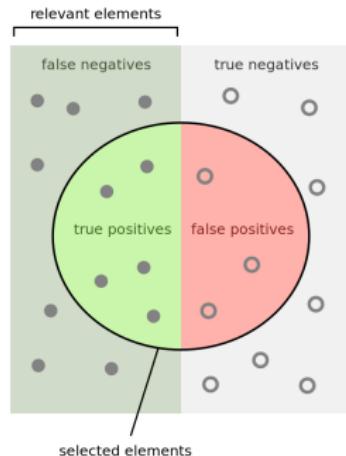
4. Classification

- ▶ Use generated features to categorise websites

¹a collection of documents

F1 Metric

- ▶ Need an objective metric to evaluate results of classification
- ▶ Removes biases due to non-uniform class distribution
- ▶ A weighted average of the *precision* and *recall*
- ▶ Also average of F1 score of each class, weighted by class proportions



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Compiling a corpus (1)

Data Collection

Require corpus of websites arranged into hierarchical categories

- ▶ 20 Newsgroups
- ▶ dmoz.org
- ▶ Reuters-21578

Compiling a corpus (2)

Data Collection

dmoz is a large, hierarchical web directory curated by humans.

The screenshot shows the homepage of the dmoz website. At the top is a green header bar with the 'dmoz' logo on the left and the AOL logo with a 'partner' badge on the right. Below the header is a navigation bar with links for 'about dmoz', 'dmoz blog', 'suggest URL', 'help', 'link', and 'editor login'. To the left of the main content area is a 'Follow @dmoz' button with a Twitter icon. In the center is a search bar with a 'Search' button and a link to 'advanced' search options. The main content area is divided into several columns of category links:

Arts	Business	Computers
Movies , Television , Music ...	Jobs , Real Estate , Investing ...	Internet , Software , Hardware ...
Games	Health	Home
Video Games , RPGs , Gambling ...	Fitness , Medicine , Alternative ...	Family , Consumers , Cooking ...
Kids and Teens	News	Recreation
Arts , School Time , Teen Life ...	Media , Newspapers , Weather ...	Travel , Food , Outdoors , Humor ...
Reference	Regional	Science
Maps , Education , Libraries ...	US , Canada , UK , Europe ...	Biology , Psychology , Physics ...
Shopping	Society	Sports
Clothing , Food , Gifts ...	People , Religion , Issues ...	Baseball , Soccer , Basketball ...

Compiling a corpus (3)

Data Collection

We built a scraper using Scrapy² and used it to scrape 2000 websites from 15 top-level categories.

[Top: Sports: Winter Sports: Skiing \(1,572\)](#)

-
- [Associations](#) (405)
 - [Consumer Information Ski Resorts](#) (2)
 - [Disabled](#) (21)
 - [Equipment](#) (8)
 - [Fantasy](#) (0)
 - [For Kids and Teens](#) (6)
 - [Gay, Lesbian, and Bisexual](#) (4)
 - [Guides](#) (69)
 - [Personal Pages](#) (13)
 - [Regional](#) (430)
 - [Shopping](#) (67)
 - [Snow and Ski Forecasts](#) (24)
 - [Travel Programs](#) (98)
 - [Video Games](#) (9)
-

- [Alpine](#) (237)
- [Backcountry](#) (106)
- [Monoskiing](#) (3)
- [Nordic](#) (148)
- [Ski Joring](#) (2)

See also:

- [News: Weather: Snow and Ski Forecasts](#) (24)
 - [Sports: Winter Sports: Snowboarding](#) (98)
-

- [About.com Skiing](#) - Guide to ski resorts, skiing vacations and trip planning, ski clothing and equipment, conditioning tips for downhill and cross country skiers.
- [Aviemore Medical Practice](#) - Supplying general practice information and details relating to sports injuries particularly related to ski accidents.
- [Born2ski.com](#) - Snow reports for over 180 resorts worldwide, plus a searchable archive of historical data, and online Ski Shops.
- [Doglotion.com](#) - A freeskiing web magazine with videos, pictures, news, stories, and reviews.
- [New England Ski Museum](#) - The New England Museum houses the most extensive collection of historical ski equipment, clothing, film, photographs, literature
- [New Zealand Ski Net](#) - Resource for New Zealand skiing and boarding information with links to ski areas, snow and weather reports, and photos.
- [Ski Law](#) - Provides information on ski litigation, including ski cases of note, current US state ski law and national trends in ski litigation.
- [Ski Net](#) - Skiing resource for gear, travel, resorts, snow conditions, ski instruction, news, and racing results.
- [Ski Safety](#) - Promotes knowledge of skiers' rights and responsibilities, with discussion of recent ski accidents, legal cases and national trends. Updated routinely
- [Ski-Ski.Ski](#) - A ski-and-snowboard-portal featuring worldwide resorts, travel, shopping, weather, tips, news and club sites for skiing and snowboarding.



Storing & Querying Documents

Data Collection

- ▶ MongoDB, JSON document store
- ▶ Apache Mahout
 - ▶ LDA
 - ▶ Hadoop
 - ▶ Scala + Spark
- ▶ Hadoop-style Map Reduce
- ▶ Process terabytes of data
- ▶ REST API for querying corpus

```
intopics=[]
alltopics=[]
hierarchy='^sports.foot(.*)$'
projection=[]
```



¹<https://mahout.apache.org/users/clustering/latent-dirichlet-allocation.html>

Website Content & Metadata

Data Collection

We take the body text...

The image contains two side-by-side screenshots of the Skinect website. The left screenshot shows the 'skiing' magazine page, which features a large red 'skiing' logo at the top. Below it, there's a paragraph about serious skiers looking for gear and powder, followed by a section titled 'What's new on Skiing:' with a list of trip ideas: Cannon, New Hampshire; Jay Peak, Vermont; Sun Valley, Idaho; and Mt. Baker, Washington. At the bottom is a blue button labeled 'GO TO SKIING MAGAZINE'. The right screenshot shows the 'SKI' magazine page, featuring a large 'SKI' logo at the top. It has a similar layout with a paragraph about luxury travel and a section titled 'What's new on Ski:' with a list of items: Used and Abused: Weekly Gear Reviews; Want to be a SKI Magazine Intern?; We Hell Ski with CMH; Ski Resorts That Are Still Open: 2015-16; and Feel Like A Ski Tester: Win The Ski... At the bottom is a blue button labeled 'GO TO SKI MAGAZINE'.

And the metadata tags

```
<html xmlns="http://www.w3.org/1999/xhtml"
      xml:lang="en" lang="en">
  <head>
    <title>SkiNet.com: Home of SKI Magazine, Skiing Magazine, Warren Miller and Nastar</title>
    <meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
    <meta name="keywords" content="ski magazine, skiing magazine, skinet, warren miller, nastar, ski resorts, ski gear, skiing videos, skiing photos" />
    <meta name="description" content="Ski resorts, ski gear reviews, lift ticket deals, skiing videos and photos and more at SkiNet, home of SKI Magazine, Skiing Magazine, Warren Miller and Nastar." />
```

Website Content & Metadata

Data Collection

We take the body text...

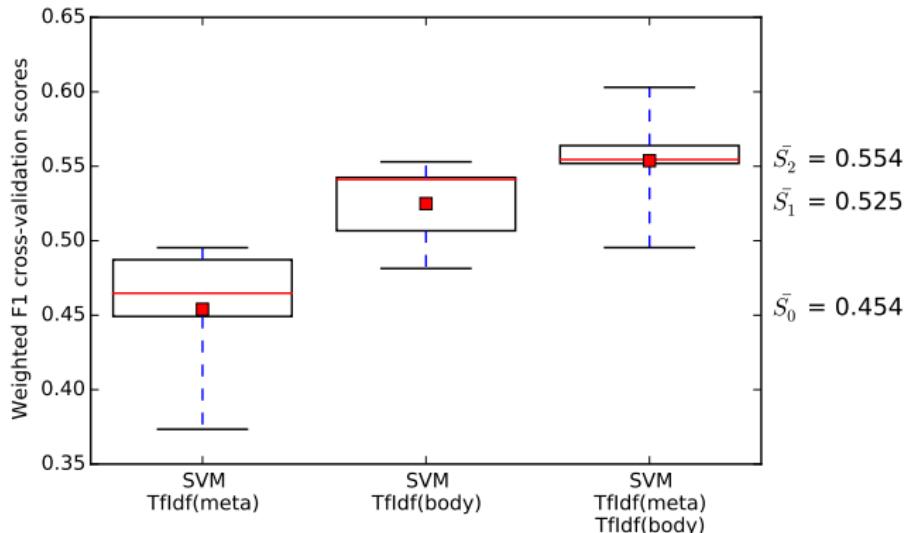


And the metadata tags

```
<html xmlns="http://www.w3.org/1999/xhtml"
      xml:lang="en" lang="en">
  <head>
    <title>SkiNet.com: Home of SKI Magazine, Skiing Magazine, Warren Miller and Nastar</title>
    <meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
    <meta name="keywords" content="ski magazine, skiing magazine, skinet, warren miller, nastar, ski resorts, ski gear, skiing videos, skiing photos" />
    <meta name="description" content="Ski resorts, ski gear reviews, lift ticket deals, skiing videos and photos and more at SkiNet, home of SKI Magazine, Skiing Magazine, Warren Miller and Nastar." />
```

Metadata vs. Body Text

Data Collection



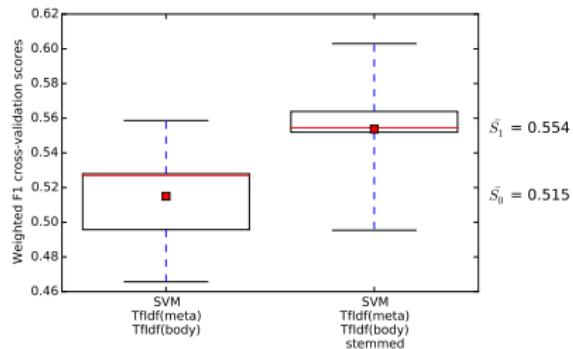
- ▶ Body text beats metadata by 0.07 (p-value: 0.028)
- ▶ Combining both is an improvement, but only of 0.029 (p-value: 0.22)

Cleaning and Preprocessing

Many words have multiple forms, e.g. plurals. Is it worth counting them as identical features?

- ▶ site (2), sites (1), ski (8), skiers (1), skiing (4), skis (1), slog (1) ...
- ▶ Use stemming to remove suffixes
 - ▶ Ski, **Skiing**, **Skis** → Ski
 - ▶ **Skiers** → Skier
- ▶ site (3), ski (13), skier (1), slog (1), ...

F1-score improvement is marginal (0.009), but the algorithm runs 28% quicker.

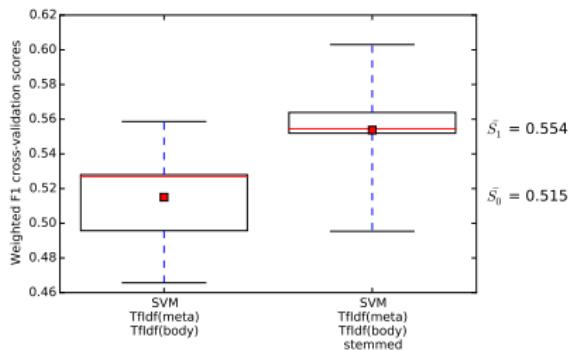


Cleaning and Preprocessing

Many words have multiple forms, e.g. plurals. Is it worth counting them as identical features?

- ▶ site (2), sites (1), **ski (8)**, **skiers (1)**, **skiing (4)**, **skis (1)**, slog (1) ...
- ▶ Use stemming to remove suffixes
 - ▶ Ski, **Skiing**, **Skis** → Ski
 - ▶ **Skiers** → Skier
- ▶ site (3), ski (13), skier (1), slog (1),
...

F1-score improvement is marginal (0.009), but the algorithm runs 28% quicker.

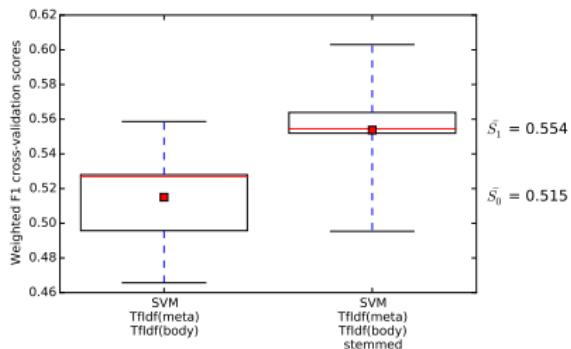


Cleaning and Preprocessing

Many words have multiple forms, e.g. plurals. Is it worth counting them as identical features?

- ▶ site (2), sites (1), **ski (8)**, **skiers (1)**, **skiing (4)**, **skis (1)**, **slog (1)** ...
- ▶ Use stemming to remove suffixes
 - ▶ Ski, **Skiing**, **Skis** → **Ski**
 - ▶ **Skiers** → **Skier**
- ▶ site (3), **ski (13)**, **skier (1)**, **slog (1)**, ...

F1-score improvement is marginal (0.009), but the algorithm runs 28% quicker.



Feature Extraction Methods

Feature Extraction

There are many ways to extract useful numerical features from text documents:

- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Latent Semantic Analysis (LSA), (and pLSA, LSI, . . .)
- ▶ String Kernels
- ▶ Term-Frequency Inverse-Document-Frequency (tf-idf)
- ▶ Word Counts

Feature Extraction Methods

Feature Extraction

There are many ways to extract useful numerical features from text documents:

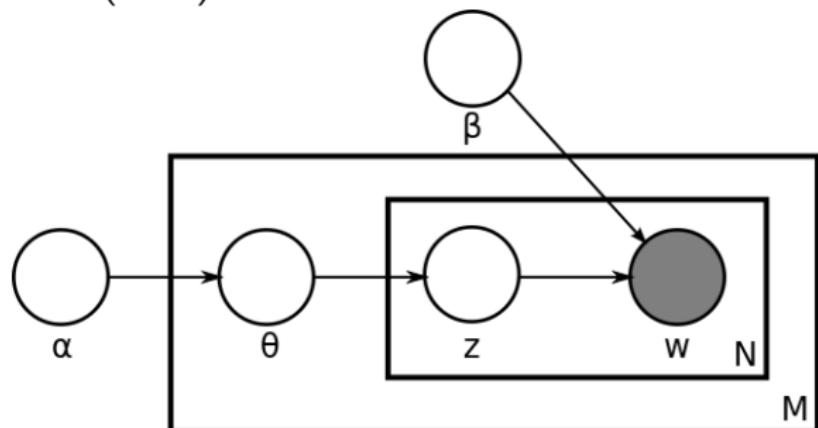
- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Latent Semantic Analysis (LSA), (and pLSA, LSI, . . .)
- ▶ String Kernels
- ▶ Term-Frequency Inverse-Document-Frequency (tf-idf)
- ▶ Word Counts

Topic Modelling

Feature Extraction

Latent Dirichlet Allocation (LDA)

- ▶ Bag of words model
- ▶ Per document topic distributions
 $\sim Dir(\alpha)$
- ▶ Per topic word distribution
 $\sim Dir(\beta)$
- ▶ $Doc\{Topic\{W\}\}$



Term-frequency inverse-document-frequency (tf-idf)

Feature Extraction

$$\begin{aligned} tfidf(t, d) &= \frac{\text{Term Frequency}}{\text{Document Frequency}} \\ &= \frac{0.5 \times f(t, d)}{\max \{f(w, d) : w \in d\}} \times \log \left(\frac{N}{|\{d \in D : t \in d\}|} \right) \end{aligned}$$

- ▶ Importance of word with relation to corpus
- ▶ Word count only or just *TermFrequency*, { "mountain" : 9 }
- ▶ Stop words

Classifiers

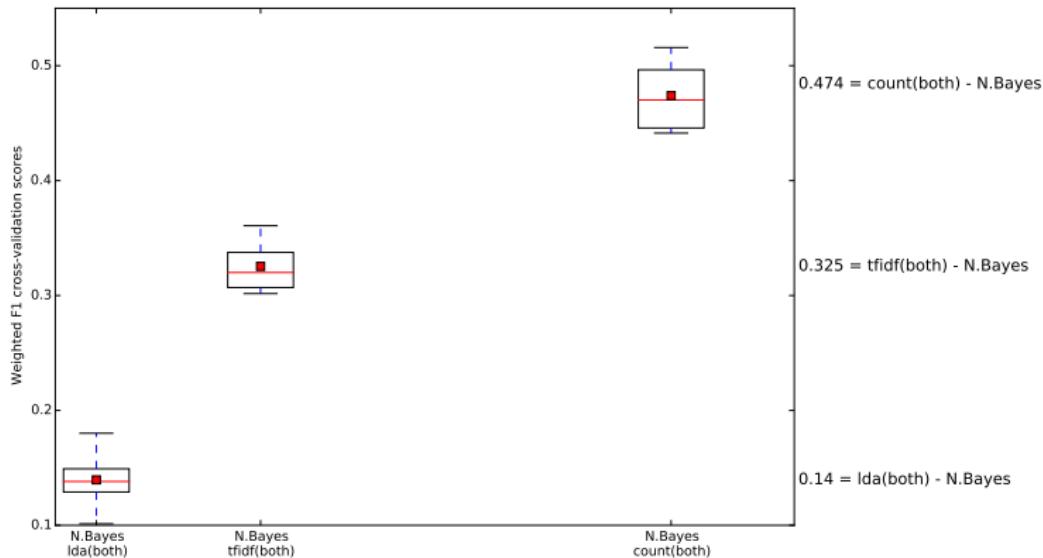
Comparing Classifiers

- ▶ Naive Bayes
 - ▶ Models classes as Gaussian distributions and classifies based on posterior probability.
- ▶ Randomised Trees
 - ▶ Ensemble method that fits extremely randomised trees on subsets of the data and averages the results.
- ▶ Support Vector Machine Classification
 - ▶ Calculates optimal class boundaries by maximising margins between classes.
- ▶ Stochastic Gradient Descent optimisation
 - ▶ Can be trained ‘online’ with new data as it comes in
 - ▶ Can be used to approximate SVMs, or algorithms which give actual probability estimates.

Results

Classifier comparison

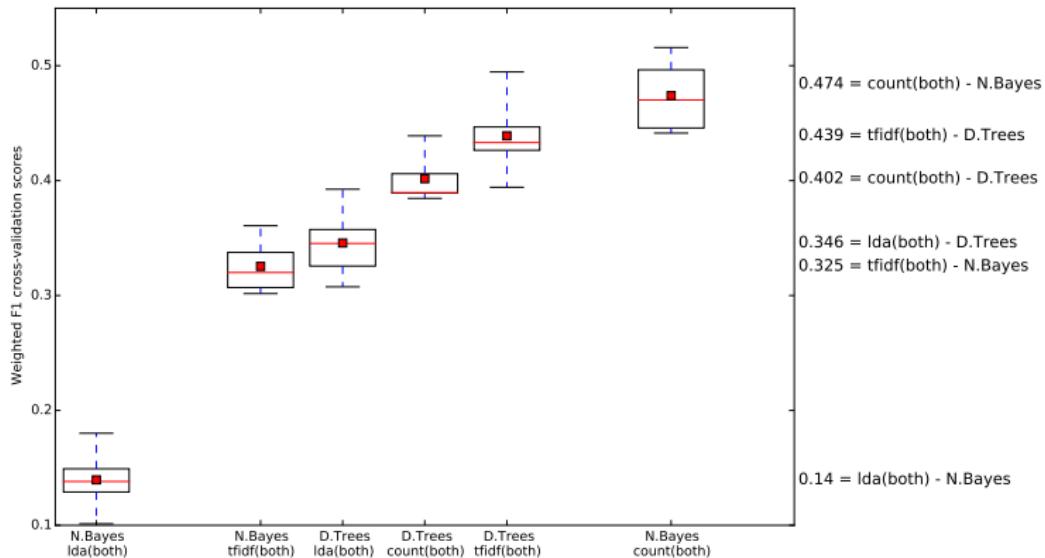
Different classifiers work best with different representations



Results

Classifier comparison

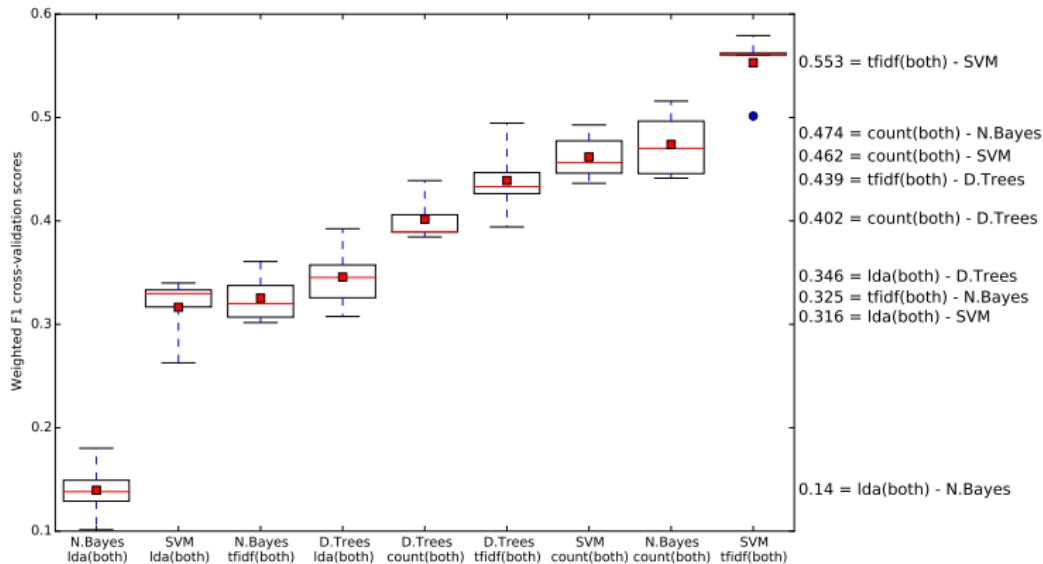
Different classifiers work best with different representations



Results

Classifier comparison

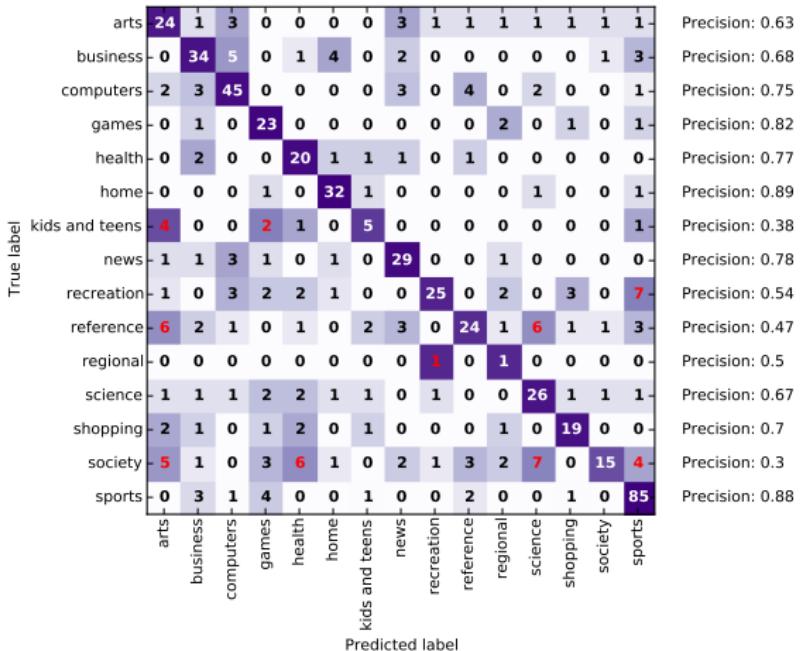
Different classifiers work best with different representations



Confusion Matrices

Which things got mistaken for what?

- ▶ Confusion matrices let you see where misclassifications happen
- ▶ This one is for all data, tf-idf and SVM Classification
- ▶ But you can also look at the confusion of *second* guesses...



Confusion Matrices

Second guesses when first guesses were wrong

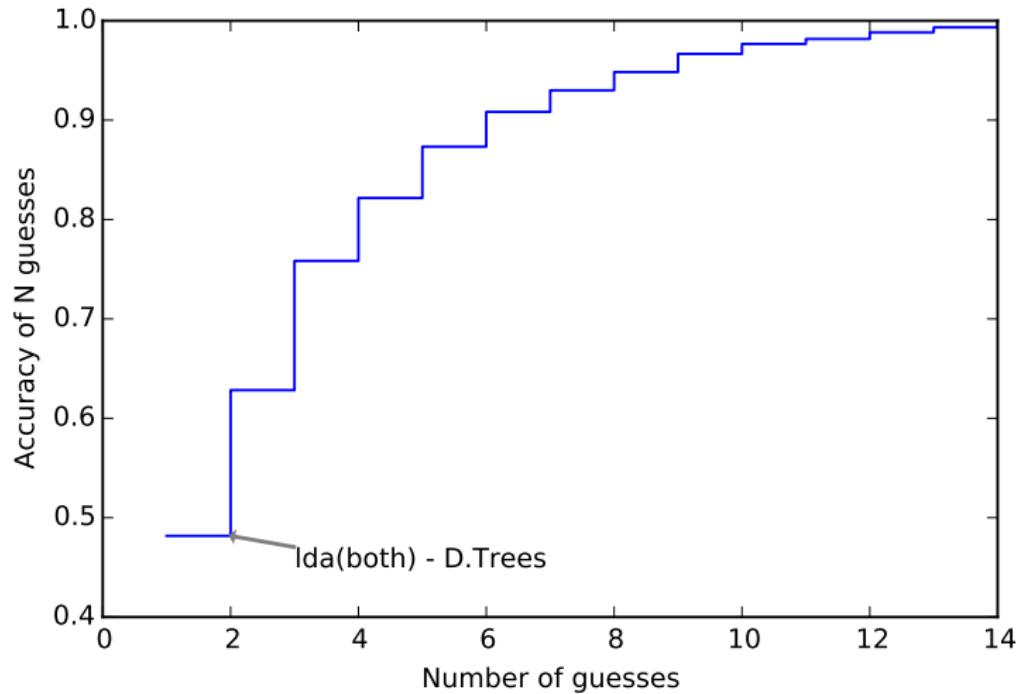
True label	arts	business	computers	games	health	home	kids and teens	news	recreation	reference	regional	science	shopping	society	sports	Precision
arts	8	1	1	0	0	0	1	1	0	1	0	0	0	1	0	0.57
business	0	7	0	1	1	0	0	0	1	0	0	2	3	0	1	0.44
computers	0	0	9	0	0	0	0	0	1	1	0	3	1	0	0	0.6
games	0	0	1	3	0	0	0	1	0	0	0	0	0	0	0	0.6
health	1	0	1	0	3	1	0	0	0	0	0	0	0	0	0	0.5
home	0	3	0	0	0	1	0	0	0	0	0	0	0	0	0	0.25
kids and teens	0	0	0	0	0	0	5	1	1	1	0	0	0	0	0	0.62
news	0	0	1	0	0	0	0	5	0	1	0	0	0	0	1	0.62
recreation	2	1	0	1	0	1	0	3	8	2	1	1	0	0	1	0.38
reference	0	1	3	1	1	1	0	1	1	11	1	1	1	2	2	0.41
regional	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1.0
science	1	2	1	0	0	0	0	1	0	2	0	3	0	2	1	0.23
shopping	0	0	0	0	1	0	0	0	0	0	0	6	1	0	0	0.75
society	1	2	4	0	3	1	3	1	1	3	3	1	2	10	0	0.29
sports	1	0	1	0	0	0	0	0	0	1	0	0	1	0	8	0.67

Confusion Matrices

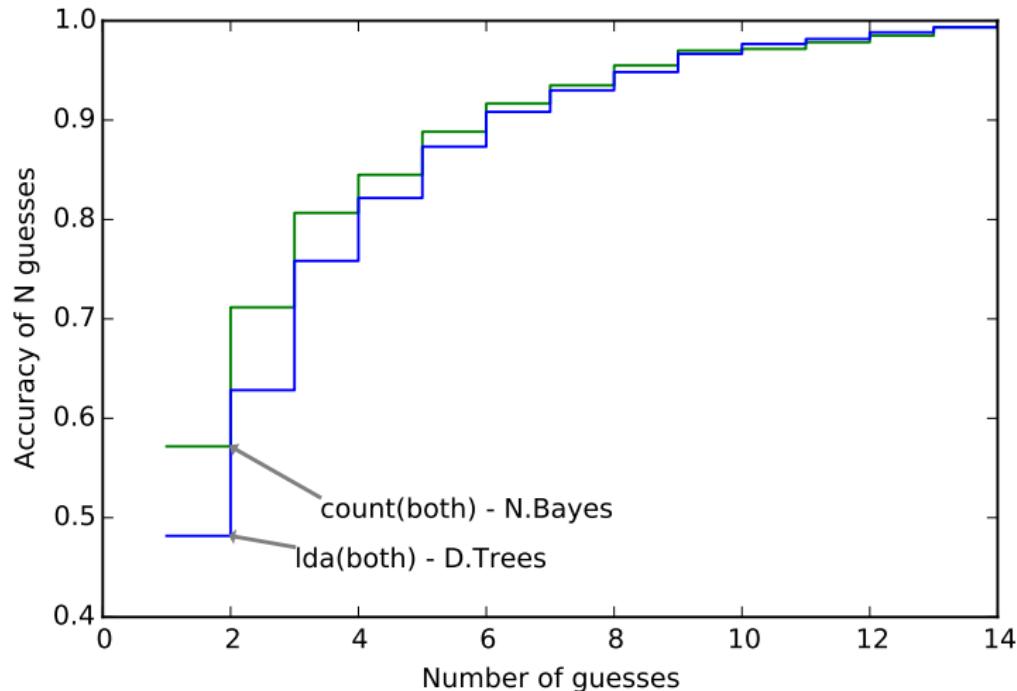
Second guesses when first guesses were right

True label	arts	business	computers	games	health	home	kids and teens	news	recreation	reference	regional	science	shopping	society	sports	Precision
arts	0	2	2	2	2	1	2	1	0	6	0	2	5	4	1	0.0
business	1	0	8	0	4	2	1	3	0	0	1	4	2	0	3	0.0
computers	2	4	0	1	1	0	0	1	0	3	0	9	2	3	2	0.0
games	1	0	3	0	0	0	4	0	1	0	0	0	0	0	4	0.0
health	0	7	0	0	0	0	3	2	0	1	0	5	2	5	0	0.0
home	1	3	0	1	0	0	2	5	2	0	2	2	4	1	1	0.0
kids and teens	0	0	0	2	3	0	0	2	0	2	0	1	0	0	1	0.0
news	4	3	1	0	2	1	0	0	0	2	1	1	0	5	4	0.0
recreation	2	1	3	1	1	3	0	2	0	1	0	0	5	2	9	0.0
reference	7	1	2	1	0	1	2	2	3	0	2	6	3	4	3	0.0
regional	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0.0
science	0	3	7	0	3	1	0	2	2	10	1	0	1	1	1	0.0
shopping	2	3	1	1	1	9	0	0	5	1	0	0	0	0	4	0.0
society	0	1	2	1	4	1	1	4	0	1	1	1	0	0	1	0.0
sports	3	3	7	25	4	2	3	10	17	5	3	5	2	8	0	0.0

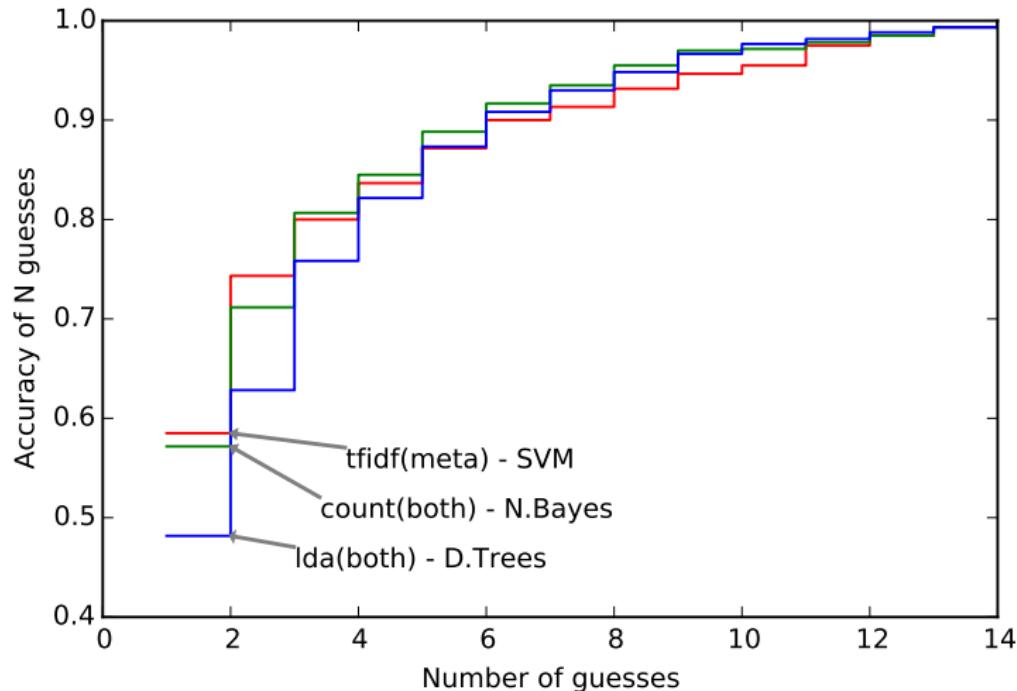
How many guesses does it take to get it right?



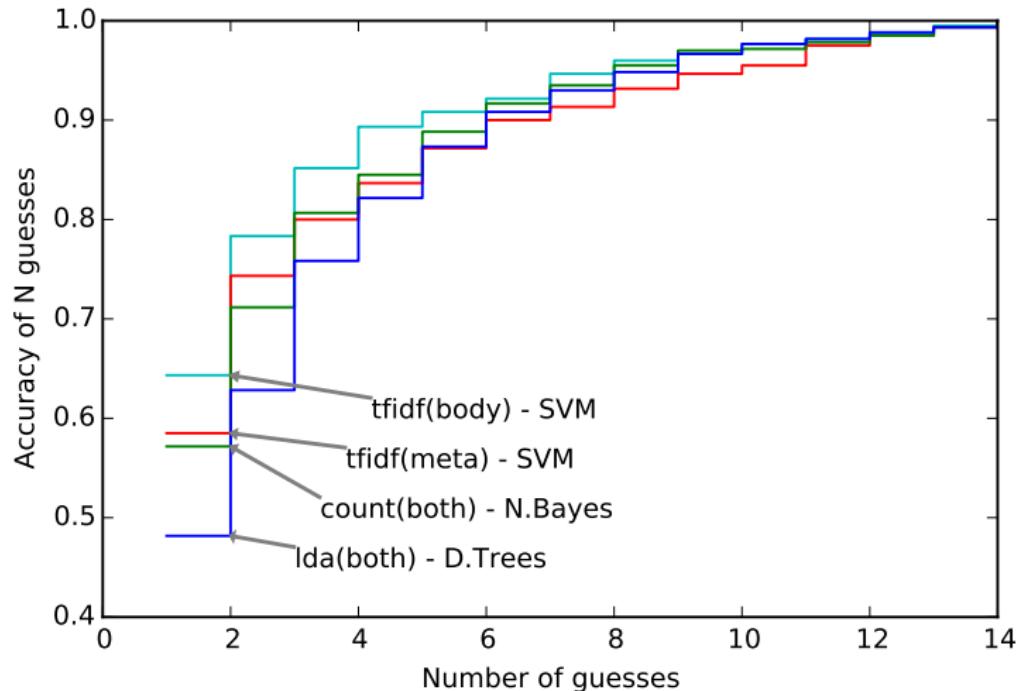
How many guesses does it take to get it right?



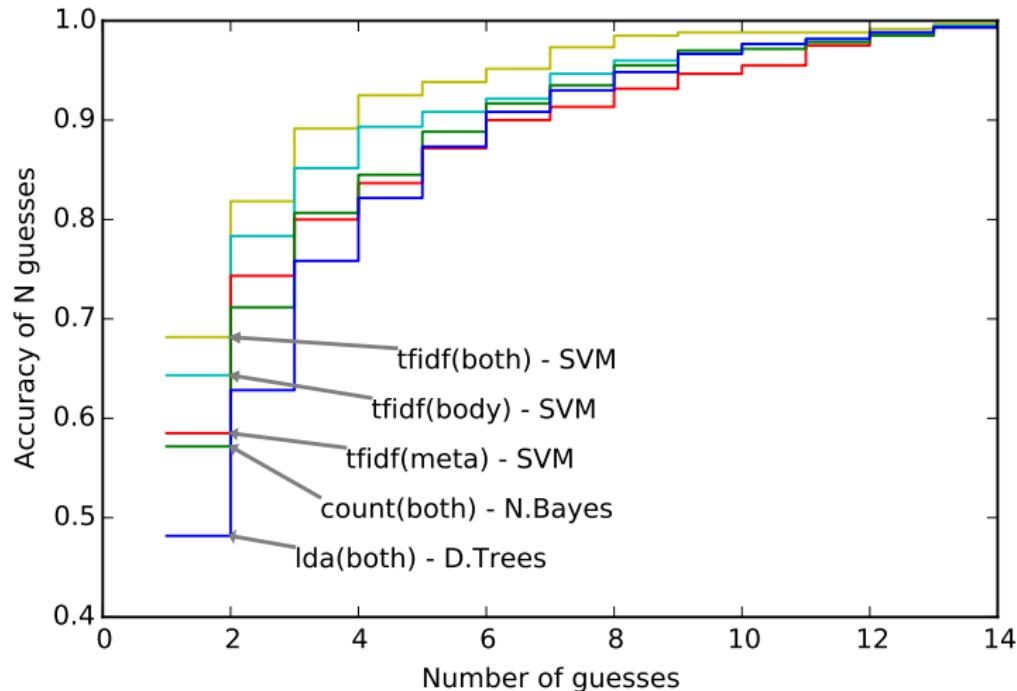
How many guesses does it take to get it right?



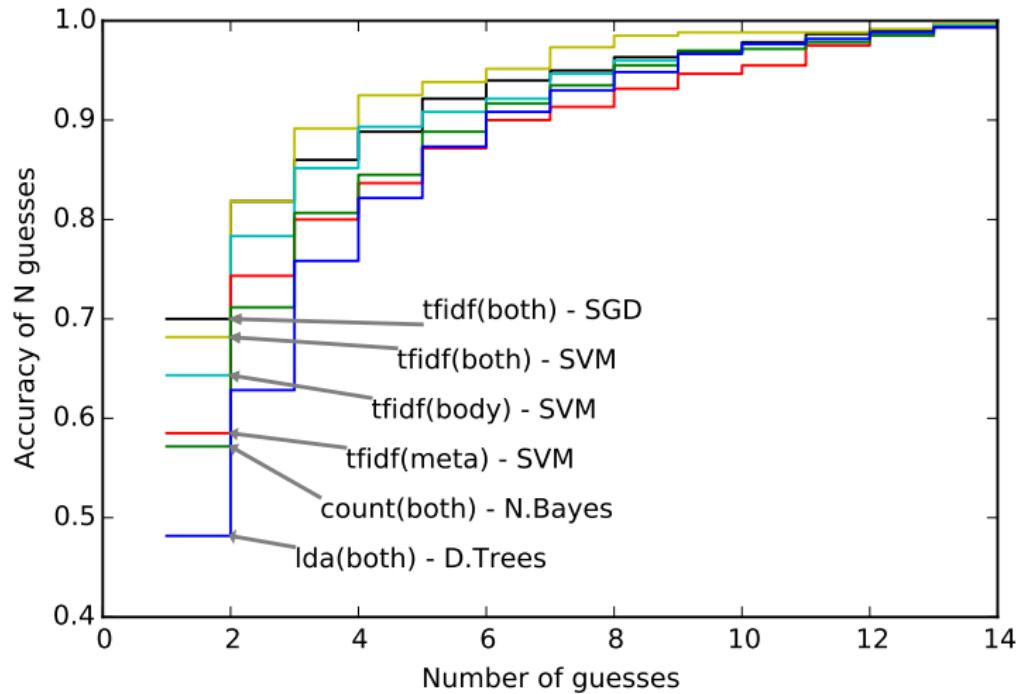
How many guesses does it take to get it right?



How many guesses does it take to get it right?



How many guesses does it take to get it right?



Top Words

SVM feature weightings for website body

1. **arts**: art, artist, award, design, costum, font, film, tattoo
2. **business**: busi, compani, product, investor, drill, sustain, financi, industri
3. **computers**: comput, file, usenet, softwar, bb, algorithm, robot, window
4. **games**: game, dice, puzzl, rpg, wargam, miniatur, stack, diplomaci
5. **health**: health, pharmaci, addict, diseas, drug, osha, age, veterinari
6. **home**: garden, diy, apart, toy, consum, clean, move, kitchen
7. **kids and teens**: kid, parent, teacher, clipart, fun, homeschool, children, disney
8. **news**: newspap, netanyahu, weather, obama, news, radio, iran, polit
9. **recreation**: hike, lock, climb, bird, cave, travel, camp, pet
10. **reference**: librari, knot, dictionari, isbn, thesauru, biographi, word, encyclopedia
11. **regional**: popul, countri, cia, island, nat, territori, geograph, gdp
12. **science**: scienc, astronomi, physic, scientif, econom, laboratori, mathemat, chemistri
13. **shopping**: furnitur, accessori, 00, 99, price, shop, camera, pipe
14. **society**: astrolog, sex, aesthet, horoscop, crime, legend, folklor, gay

Top Words

SVM feature weightings for website body

1. **arts**: art, artist, award, design, costum, font, film, tattoo
2. **business**: busi, compani, product, investor, drill, sustain, financi, industri
3. **computers**: comput, file, usenet, softwar, bb, algorithm, robot, window
4. **games**: game, dice, puzzl, rpg, wargam, miniatur, stack, diplomaci
5. **health**: health, pharmaci, addict, diseas, drug, osha, age, veterinari
6. **home**: garden, diy, apart, toy, consum, clean, move, kitchen
7. **kids and teens**: kid, parent, teacher, clipart, fun, homeschool, children, disney
8. **news**: newspap, **netanyahu**, weather, **obama**, news, radio, iran, polit
9. **recreation**: hike, lock, climb, bird, cave, travel, camp, pet
10. **reference**: librari, knot, dictionari, isbn, thesauru, biographi, word, encyclopedia
11. **regional**: popul, countri, cia, island, nat, territori, geograph, gdp
12. **science**: scienc, astronomi, physic, scientif, econom, laboratori, mathemat, chemistri
13. **shopping**: furnitur, accessori, **00**, **99**, price, shop, camera, pipe
14. **society**: astrolog, sex, aesthet, horoscop, crime, legend, folklor, gay

Pipeline & Tools

- ▶ Scikit-learn
 - ▶ Python machine learning library
 - ▶ Range of classifiers and other processing tools
 - ▶ Built-in pipeline framework for cascading algorithms
- ▶ Gensim, a python topic modelling toolbox used for LDA
- ▶ NTLK, python natural language library used for word stemming

Final Pipeline

1. Sanitise
 - ▶ Strip tags
 - ▶ Remove stopwords
 - ▶ Stem words
2. Tf-idf, both metadata and body
3. Support Vector Classification

Cross-validated accuracy of 56%

NB. a random assignment would give an accuracy of 6%

Limitations

Issues & Difficulties

- ▶ Polysemy; multiple meanings of the same word.
- ▶ Bag of words model, order is disregarded.
- ▶ Different results each iteration due to randomness in Gibbs Sampling.
- ▶ Stemming, capturing the meaning of words. Literary tenses.
- ▶ Topics in a document can be subjective not objective.
- ▶ Quantifying performance, critical evaluation.
- ▶ Storing and processing large quantities of corpora.
- ▶ Currently using old websites, which tend to lack content and metadata.

Multiple correct answers

Misclassifications

DMOZ class was Sports, top two guesses were Shopping, Sports

SALE ITEMS

<>

Deluxe Gift Pen \$14.99 \$2.99	Under Armour Women's HeatGear \$39.99 \$27.99	Velocity Shorts - Neon Orange/Stripe \$10.99 \$14.99	Under Armour Women's \$49.99 \$34.99
-----------------------------------	---	---	--

NEW PRODUCTS

<>

Molten MS500 NCAA Volleyball \$15.95	Under Armour Women's Great \$24.99	Under Armour Play Up Shorts - Pink \$24.99
--	---------------------------------------	---

THOUSANDS OF VOLLEYBALL PRODUCTS

Volleyball.Com is proud to be celebrating 20 years of bringing athletes what they need to stay on top of their game! With brands like Asics, Mizuno, Under Armour, Tachikara and many more we have everything you need to hit the court in style and comfort. Don't forget to visit our Volleyball 101 section while you're here. It's full of information on how to keep you on the court in top form! At Volleyball.Com we're committed to helping you enjoy every minute you play! Volleyball.Com - Love the Game - Get the Gear.

This is clearly a sports shopping website!

Training Bias

Misclassifications

DMOZ class was Business, misclassified as Sports



The website has practically no text or metadata.

Objective nature of the problem

Misclassifications

DMOZ class was “Society”, our algorithm says Science.

The screenshot shows the homepage of THE WIY FILES, THE SCIENCE BEHIND THE NEWS. The header features the site's logo and the date MONDAY, 4 MAY, 2015. Below the header is a navigation bar with links to HOME, ARCHIVES, BOOK REVIEWS, COOL SCIENCE IMAGES, TEACHING, ABOUT, and a search bar. The main content area has a large image of a dog's eyes. Below the image are several news articles:

- Bombardier beetle spray-bottle explained at last!** (Thumbnail: A bombardier beetle spraying a toxic spray.)

The pulse jet – biology had it first!
Tiny, fast explosions inside the bombardier beetle make a toxic spray, then heat and shoot it at the target. Here's how!
[More ▾](#)
- Dogs and their owners: A chemical bond** (Thumbnail: A close-up of a dog's face.)

Cortisol, the "birth hormone," works on both sides! [More ▾](#)
- Continental connection: North, South America linked much earlier than thought** (Thumbnail: A person standing on a bridge over a river, with mountains in the background.)

Ancient river moved from Argentina to Colombia, proving early land bridge
[More ▾](#)
- LAPD shooting: neighborhood on edge of police body-worn cameras** (Thumbnail: A view of a city street with police cars in the background.)

In the aftermath of the video-taped death of Eric Garner at the hands of police in New York City, we look at the effect of cameras on policing. [More ▾](#)
- Surge in MERS cases reported in Saudi Arabia** (Thumbnail: A view of a desert landscape with mountains in the background.)

How do scientists predict an eruption? Watch volcanologists track a massive lava flow moving toward a town in Chile. [More ▾](#)
- Villarrica Volcano erupts in Chile, thousands flee** (Thumbnail: A view of a volcano with smoke rising from its peak.)

How do scientists predict an eruption? Watch volcanologists track a massive lava flow moving toward a town in Chile. [More ▾](#)

We think DMOZ is wrong here! Many ‘society’ websites on DMOZ seem to be about other things.

Extensions

Future work ...

- ▶ Learn and classify a topic hierarchy
- ▶ Deal with (almost) empty websites
 - ▶ Filter them out entirely?
 - ▶ Train a 'null' category?
- ▶ Recognise websites not in any existing category
- ▶ Try other classifiers (e.g. Neural Networks)
- ▶ Make a website service

Extensions

Future work ...

- ▶ Learn and classify a topic heirarchy
- ▶ Deal with (almost) empty websites
 - ▶ Filter them out entirely?
 - ▶ Train a 'null' category?
- ▶ Recognise websites not in any existing category
- ▶ Try other classifiers (e.g. Neural Networks)
- ▶ Make a website service
 - ▶ Oh wait, we did that!

Web service!

<http://autodmoz.apps.veryjoe.com>

The screenshot shows a web browser window with the URL autodmoz.apps.veryjoe.com/?url=ecs.soton.ac.uk. The page title is "DMOZ Auto-Categoriser". The main content area displays the results for the URL "ecs.soton.ac.uk". It lists the primary category as "Computers, Science" and provides a detailed breakdown of sub-categories: Sports, Reference, Society, Recreation, Business, Shopping, Arts, Health, News, Home, Games, Kids and teens, Regional. Below this, two sections of keywords are listed: "computers keywords on this site" and "science keywords on this site".

DMOZ Auto-Categoriser

The DMOZ Auto-Categoriser analyses a website and decides which category on [DMOZ.org](#) it would be placed in.

ecs.soton.ac.uk

Computers, Science, Sports, Reference, Society, Recreation, Business, Shopping, Arts, Health, News, Home, Games, Kids and teens, Regional

computers keywords on this site: comput, intranet, ec, univers, electron, engin, acm, user, search, access, download, cancel, open, southampton, volum

science keywords on this site: scienc, univers, engin, faculti, phd, undergradu, career, uk, studi, seminar, editor, electron, outreach, electr, impact

The End

Any questions?

