

Intro to Machine Learning

Crowdsourcing and Human Computation
Lecture 12

Instructor: Chris Callison-Burch
TA: Ellie Pavlick

Website: crowdsourcing-class.org

What is Machine Learning?

- How can we build computer systems that automatically improve with experience and what are the fundamental laws that govern all learning processes?
–Tom Mitchell, CMU

The age of machine learning is now

- In the 90s we rediscovered probabilistic models and statistics and applied it to CS and machine learning
- We now have access to much better computing infrastructure
- We have so much data that we can barely store it, and it provides great opportunities for analysis

What can we do with machine learning?

- Find groups of related things via “clustering”. Used for recommendations by Amazon, Netflix, etc
- Are two items the same? Named entity disambiguation
- Classification: Is this email spam? What language is this web page written in? Whose face is shown in a picture?

Shop by
Department ▾

Search

All ▾

Go

Hello, **Chris**
Your Account ▾Your
Prime ▾

Cart ▾

Wish
List ▾

Your Amazon.com Your Browsing History Recommended For You Rate Items You Like Improve Your Recommendations Your Profile Learn More

Your Amazon.com > Recommended for You

(If you're not Chris Callison-Burch, click here.)

Just For Today

[Browse Recommended](#)These recommendations are based on [items you own](#) and more.view: [All](#) | [New Releases](#) | [Coming Soon](#)[More results](#) ▾

Recommendations

[Amazon Instant Video](#)[Amazon MP3 Store](#)[Appliances](#)[Appstore for Android](#)[Arts, Crafts & Sewing](#)[Automotive](#)[Baby](#)[Beauty](#)[Books](#)[Books on Kindle](#)[Camera & Photo](#)[Cell Phones & Accessories](#)[Clothing & Accessories](#)[Computers](#)[Electronics](#)[Grocery & Gourmet Food](#)[Health & Personal Care](#)[Home & Kitchen](#)[Home Improvement](#)[Industrial & Scientific](#)[Jewelry](#)[Kitchen & Dining](#)[Magazine Subscriptions](#)[Magazines on Kindle](#)

[Carter's Keep Me Dry Waterproof Fitted Quilted Crib Pad, White](#)

by Kids Line (December 11, 2009)

Average Customer Review: ★★★★★ (383)

In Stock

List Price: \$12.99**Price:** \$12.79[17 used & new from \\$11.02](#)[Add to Cart](#)[Add to Wish List](#) I own it Not interested Rate this itemRecommended because you added **Summer Infant Contoured Changing Pad Amazon Frustration F...** to your Shopping Cart and more ([Fix this](#))

[Nosefrida The Snotsucker Nasal Aspirator](#)

by FridaBaby (April 1, 2010)

Average Customer Review: ★★★★★ (1,859)

In Stock

List Price: \$15.99**Price:** \$14.78[43 new from \\$9.86](#)[Add to Cart](#)[Add to Wish List](#) I own it Not interested Rate this itemRecommended because you added **Summer Infant Infant Character Change Pad Cover, Safari S...** to your Shopping Cart and more ([Fix this](#))

[Safety 1st Heavenly Dreams White Crib Mattress](#)

by Dorel Home Products (December 11, 2010)

Average Customer Review: ★★★★★ (627)

In Stock

List Price: \$54.99**Price:** \$52.99[Add to Cart](#)[Add to Wish List](#)

Shop by
Department ▾

Search

All ▾

Your Amazon.com

Your Browsing History

[Your Amazon.com](#) > **Recommended for You**
(If you're not Chris Callison-Burch, click here.)**Just For Today**[Browse Recommended](#)**Recommendations**[Amazon Instant Video](#)[Amazon MP3 Store](#)[Appliances](#)[Appstore for Android](#)[Arts, Crafts & Sewing](#)[Automotive](#)[Baby](#)[Beauty](#)[Books](#)[Books on Kindle](#)[Camera & Photo](#)[Cell Phones & Accessories](#)[Clothing & Accessories](#)[Computers](#)[Electronics](#)[Grocery & Gourmet Food](#)[Health & Personal Care](#)[Home & Kitchen](#)[Home Improvement](#)[Industrial & Scientific](#)[Jewelry](#)[Kitchen & Dining](#)[Magazine Subscriptions](#)[Magazines on Kindle](#)

Amazon.com: Why is this recommended for you?

amazon.com

[Help](#) | [Close window](#)

Recommended for You



[Nosefrida The Snotsucker Nasal Aspirator](#)
by FridaBaby (April 1, 2010)
In Stock
List Price: \$15.99
Price: \$14.78
[43 new from \\$9.86](#)

[Add to Cart](#) [Add to Wish List](#)

Rate this item ★★★★★

I own it
 Not interested

Because you purchased...



[GE 51386 Metal Shade With Flower Design Incandescent Night Light](#)
★★★★★
 This was a gift
 Don't use for recommendations



[Munchkin Arm & Hammer Diaper Pail Refill Bags, 30 Count](#)
by Munchkin
★★★★★
 This was a gift
 Don't use for recommendations



[My Brest Friend Original Pillow, Bluebells](#)
by Zenoff Products
★★★★★
 This was a gift
 Don't use for recommendations

Because your Wish List includes...

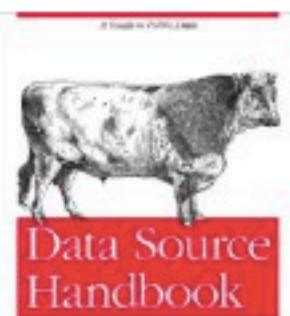


[WiFi Baby 2.0 \(2013 Model\) - iPhone, iPad, Android, Baby Monitor & Nanny Cam DVR. Video, Audio, Recording. Anywhere. Same Look, New Features \(WFB2013\)](#)
by WiFi Baby
★★★★★
 Don't use for recommendations

[More results](#) [Wish List](#)[Shopping Cart and more \(Fix this\)](#)[Wish List](#)[Shopping Cart and more \(Fix this\)](#)

More Items to Consider

You viewed



Data Source Handbook

O'REILLY™
Data Source Handbook

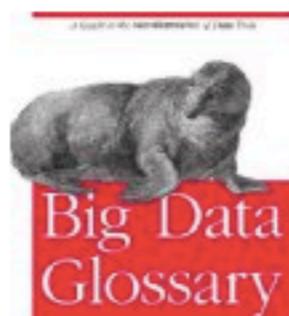
› Pete Warden

Paperback

★★★★★ (8)

\$29.99 \$20.78

Customers who viewed this also viewed



Big Data Glossary

O'REILLY™
Big Data Glossary

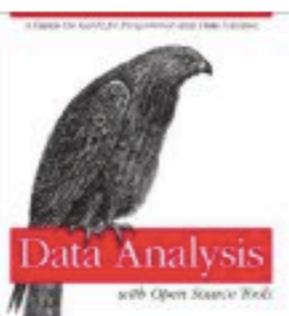
› Pete Warden

Paperback

★★★★★ (5)

\$19.99 \$14.78

[View or edit your browsing history](#) ›



Data Analysis
with Open Source Tools

O'REILLY™
Data Analysis with Open Source Tools

› Philipp K. Janert

Paperback

★★★★★ (35)

\$39.99 \$24.22



1. ASICS Women's GT-1000
Running Shoe
ASICS



2. ASICS Women's GEL-Noosa Tri
Running Shoe
ASICS



3. ASICS Women's Gel-Nimbus 14
Running Shoe
ASICS



4. ASICS Women's Gel-Kayano 19
Running Shoe
ASICS



5. ASICS Women's GEL-Cumulus
Running Shoe
ASICS

[See all Best Sellers in Women's Running](#)

Related to Items You've Viewed

You viewed



Bose QuietComfort 20i
Acoustic Noise...
★★★★★ (144)

Customers who viewed this also viewed



Bose IE2 Audio Headphones
★★★★★ (1,153)
\$99.95

[View or edit your browsing history](#) ›



Bose® MIE2i Mobile Headset
★★★★★ (632)
\$129.95



1. Amazon Gift Card - E-mail
Amazon
\$50.00



2. Amazon Gift Card - E-mail - Happy Birthday (Candles)
Amazon
\$50.00



3. Amazon Gift Card - E-mail - Thank You (Note)
Amazon
\$50.00



4. Amazon Gift Card Upload Your Photo - Gift for You
Amazon
\$50.00



5. Amazon.com Gift Cards - E-mail

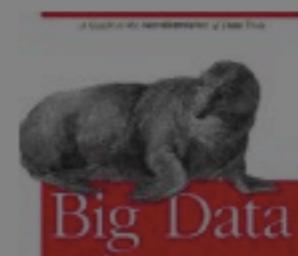
More Items to Consider

You viewed



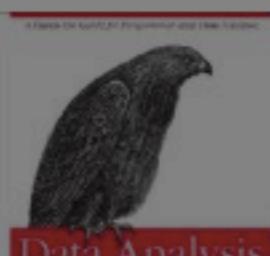
Data Source

Customers who viewed this also viewed



Big Data

[View or edit your browsing history](#) ▾



Data Analysis



1. ASICS Women's GT-1000 Running Shoe
ASICS



2. ASICS Women's GEL-Noosa Tri Running Shoe
ASICS



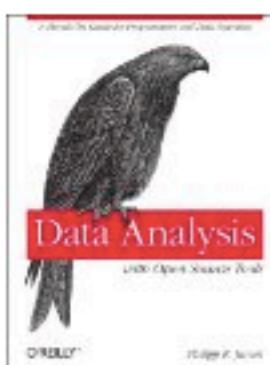
3. ASICS Women's Gel-Nimbus 17 Running Shoe
ASICS



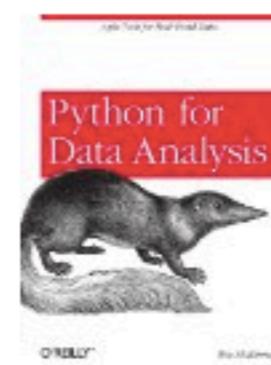
4. ASICS Women's Gel-Kayano 19

Continue Shopping: Customers Who Bought Items in Your Recent History Also Bought

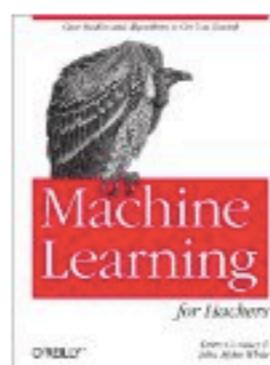
Page 1 of 13



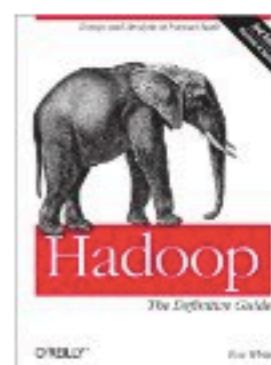
[Data Analysis with Open Source Tools](#)
► Philipp K. Janert
 (35)
Paperback
\$24.22 



[Python for Data Analysis](#)
► Wes McKinney
 (34)
Paperback
\$24.68 



[Machine Learning for Hackers](#)
► Drew Conway
 (22)
Paperback
\$33.66 



[Hadoop: The Definitive Guide](#)
► Tom White
 (36)
Paperback
\$29.99 



Bose QuietComfort 20i
Acoustic Noise...
 (144)



Bose IE2 Audio Headphones
 (1,153)
\$99.95



Bose® MIE2i Mobile Headset
 (632)
\$129.95



- Amazon Gift Card - E-mail - Thank You (Note)
Amazon
\$50.00



3. Amazon Gift Card Upload Your Photo - Gift for You
Amazon
\$50.00



4. Amazon.com Gift Cards - E-mail

Does Anne Hathaway News Drive Berkshire Hathaway's Stock?

ALEXIS C. MADRIGAL | MAR 18 2011, 10:50 AM ET

1

Like

483

Tweet

14

+1

258

Share

More ▾



A couple weeks ago, Huffington Post blogger Dan Mirvish noted a funny trend: when Anne Hathaway was in the news, Warren Buffett's Berkshire Hathaway's shares went up. He pointed to [six dates going back to 2008](#) to show the correlation. Mirvish then suggested a mechanism to explain the trend: "automated, robotic trading programming are picking up the same chatter on the Internet about 'Hathaway' as the IMDb's StarMeter, and they're applying it to the stock market."

The idea seems ridiculous. But the more I thought about the strange behavior of algorithmic trading systems and the news that [Twitter sentiment analysis could be used](#) by stock market analysts and the fact that many computer programs are simply looking for tradeable correlations, I really started to

VIDEO



Advice to a Younger Me: Michelle Peluso

The CEO of Gilt reflects on leadership and work-life balance

WRITERS



Molly Ball

Republicans Shut Down the Government for Nothing 11:17 AM ET

Julie Beck

Study: People Love to Cheat 11:10 AM ET

Matthew O'Brien

The Undeadline: Why Halloween Could (Seriously) Be Our Last Day to Save America From Default 10:56 AM ET

Kasia Cieplak-Mayr von Baldegg

The Highs (and Lows) of Starting Your Own Company 10:53 AM ET

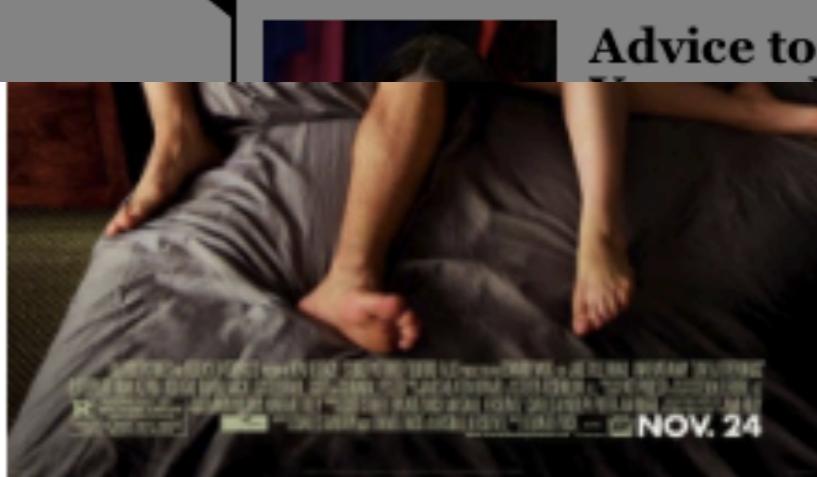
Derek Thompson

American Politics Deserves to Be Downgraded 10:20 AM ET

Does Anne Hathaway News Drive Berkshire Hathaway's Stock?

ALEXIS C.

VIDEO



Advice to a

use:
luso
ership

On the Friday before the Oscars, Berkshire shares rose a whopping 2.02%. And on the Monday just after the Academy Awards, they rose again, this time 2.94%. But it's not just an Oscar bounce, or something Warren Buffett may have said in the newspaper, or even necessarily something the company itself is doing (i.e. rumors afoot to buy [Costco](#)). Just look back at some other landmark dates in Anne Hathaway's still young career:

Oct. 3, 2008 - *Rachel Getting Married* opens: BRK.A up .44%

Jan. 5, 2009 - *Bride Wars* opens: BRK.A up 2.61%

Feb. 8, 2010 - *Valentine's Day* opens: BRK.A up 1.01%

March 5, 2010 - *Alice in Wonderland* opens: BRK.A up .74%

Nov. 24, 2010 - *Love and Other Drugs* opens: BRK.A up 1.62%

Nov. 29, 2010 - Anne announced as co-host of the Oscars: BRK.A up .25%

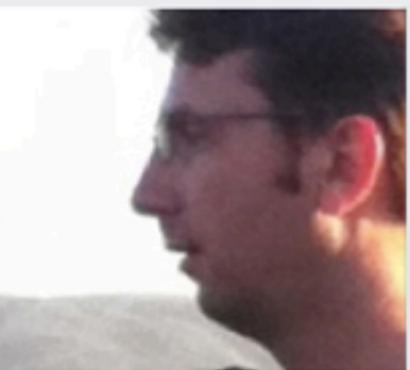
My guess is that all those automated, robotic trading programming are picking up the same chatter on the internet about "Hathaway" as the IMDb's StarMeter, and they're applying it to the stock market. Of course, this isn't necessarily bad news for the investor. After all, can you imagine what might have happened to Berkshire stock if Warren Buffett had appeared nude in *Love and Other Drugs* rather than Anne Hathaway? Perhaps it's best if we don't think about it.

of algorithmic trading systems and the news that Twitter sentiment analysis could be used by stock market analysts and the fact that many computer programs are simply looking for tradeable correlations, I really started to

Derek Thompson

American Politics Deserves to Be
Downgraded 10:20 AM ET

LIBRARY
Events
Photos
Faces
Places
RECENT
Oct 3, 2009
Last 12 Months
Last Import
Flagged
Trash
WEB
Photo Stream
callison-burch
DEVICES
ccb's iPhone
EVENTS
ALBUMS
Last 12 Months
2nd floor bathroom
3029 Elm Ave



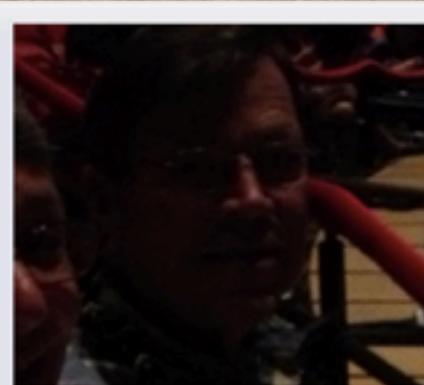
Alex Klementiev



Alexa Callison-Burch



Ann Irvine



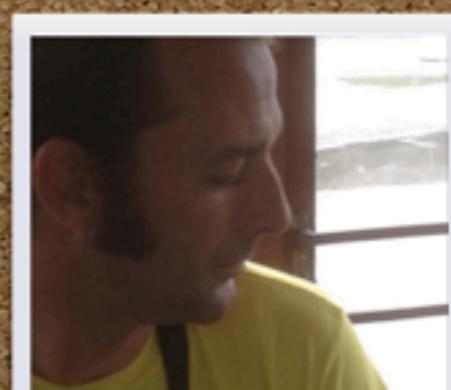
Archie Burch



Charley Chan



Chris Callison-Burch



Chris Landers



Colin Bannard



Danielle Matthews



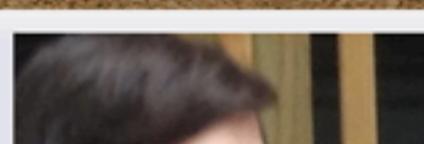
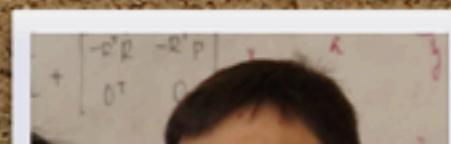
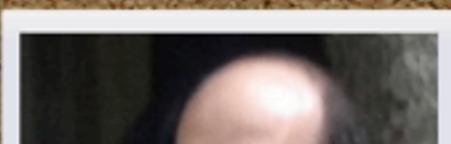
Dave Hawkey



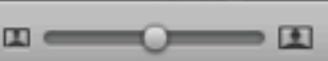
Dawn Mercurio



Diane Callison



Search



Zoom



Slideshow



Find Faces

LIBRARY

- Events
- Photos
- Faces
- Places

RECENT

- Oct 3, 2009
- Last 12 Months
- Last Import
- Flagged
- Trash

WEB

- Photo Stream
- callison-burch

DEVICES

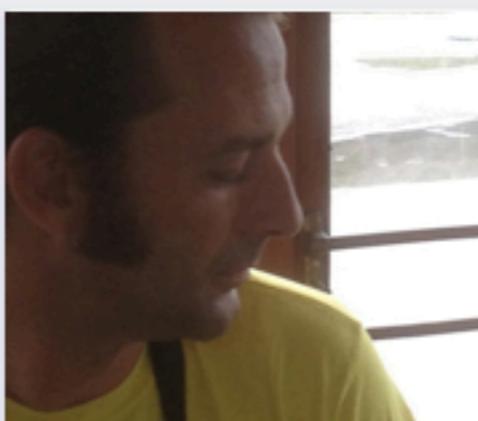
- ccb's iPhone

EVENTS**ALBUMS**

- Last 12 Months
- 2nd floor bathroom
- 3029 Elm Ave



Chris Callison-Burch



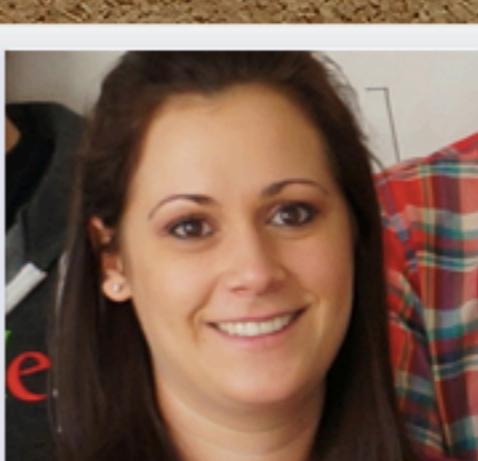
Chris Landers



Colin Bannard



Colin Bannard



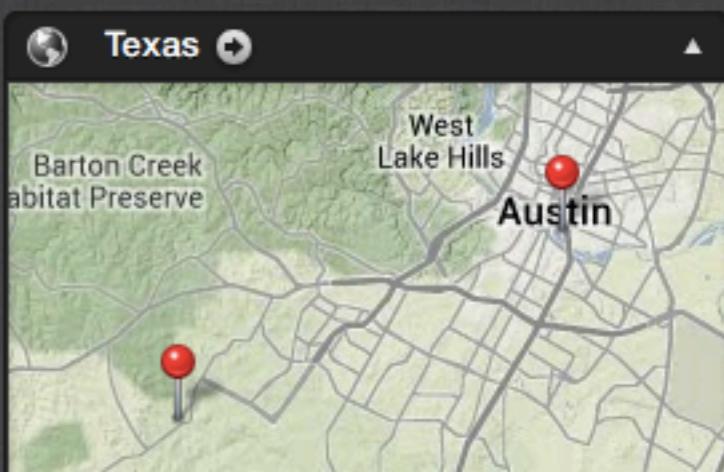
Courtney Napoles



Danielle Matthews

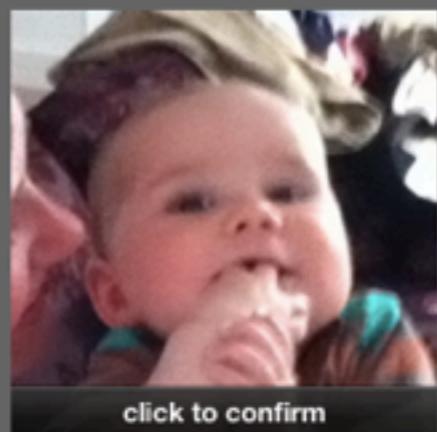


Dave Hawkey

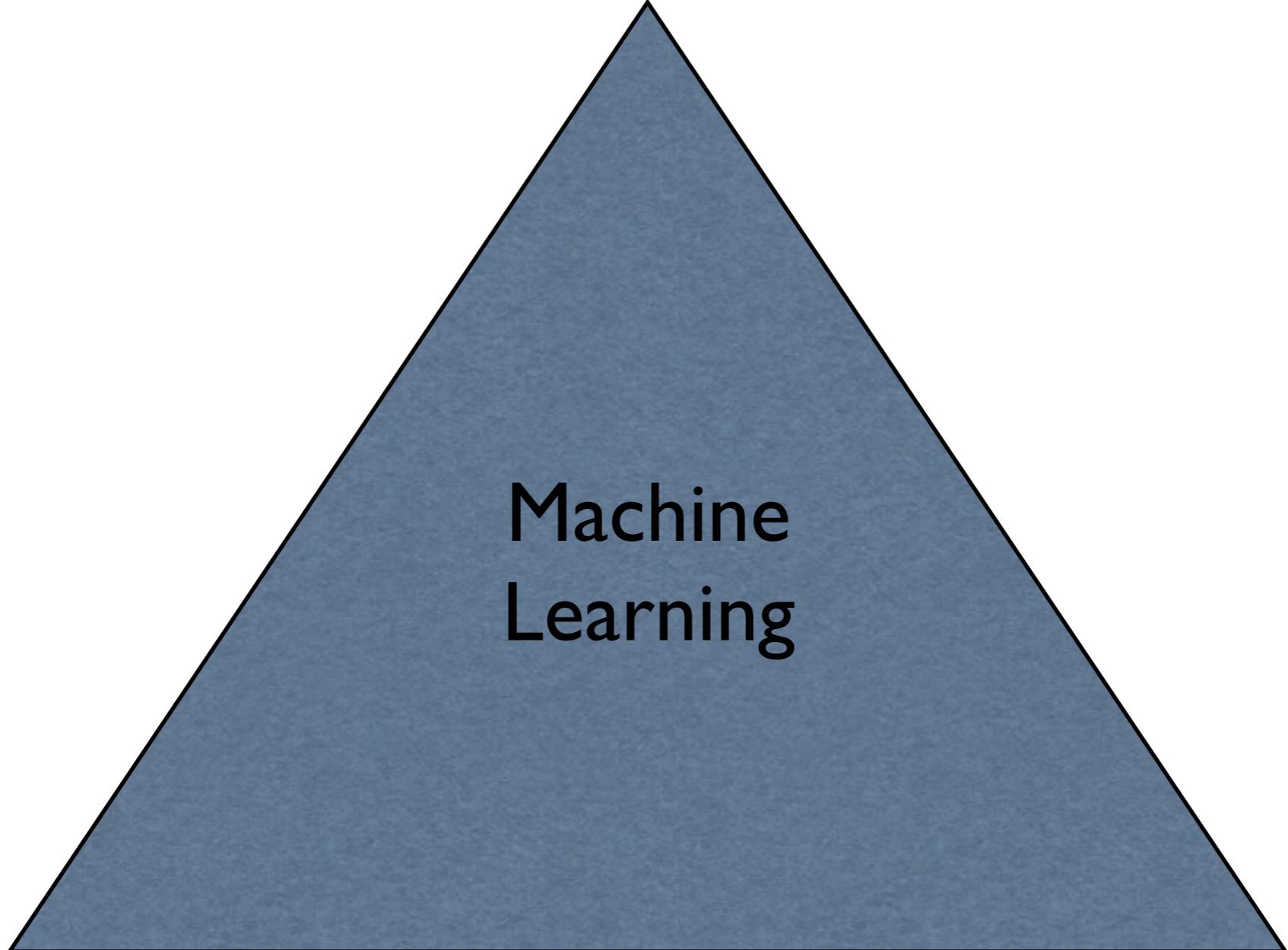




Unconfirmed Faces



Data



Machine
Learning

Model

Algorithm

Supervised v. Unsupervised Learning

- In supervised learning you are starting with a labeled training set of data
- In unsupervised learning you don't (yet) have that data

Kinds of data?

- Text and speech
- Images and video
- Geographic information
- Time series information
- Transaction data from customers
- Climate data
- Census data

Where does data come from?

- Some datasets are available for free:
<http://delicious.com/pskomoroch/dataset>
- Some is owned by companies
- Sometimes you can assemble it yourself
- Crowdsourcing!

Classification

- Classification is the assignment of a label to unlabeled input based on previously seen data
- Learn $f(x)$...
- that outputs a label ...
- along with a probability that that label is true

Example classification tasks

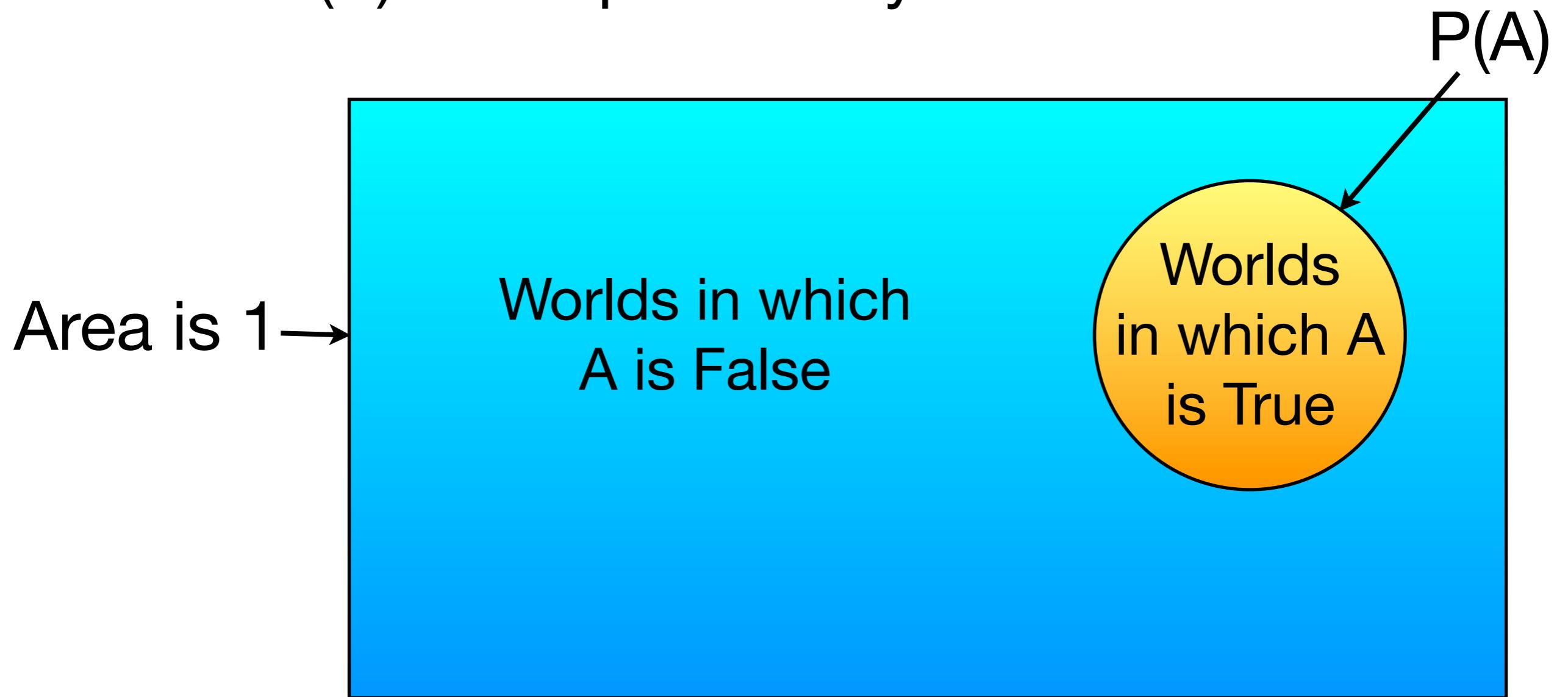
- Spam v. not Spam
- Face detection
- Language Identification

Naive Bayes

- First used for Spam filters in the 1990s
- Math is just counting, dividing and multiplying
- No calculus

Probability

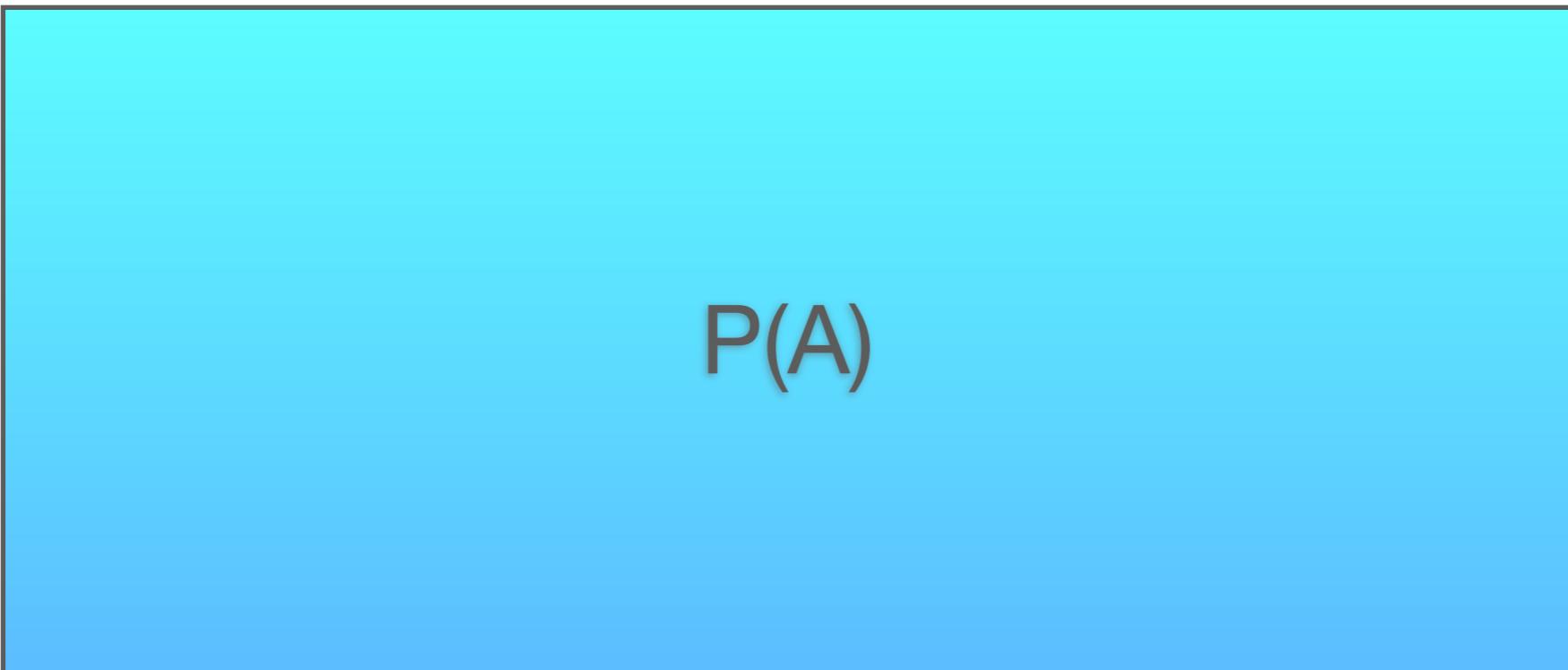
$P(A)$ is the probability that A is true



Axioms of Probability

- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $0 \leq P(A) \leq 1$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

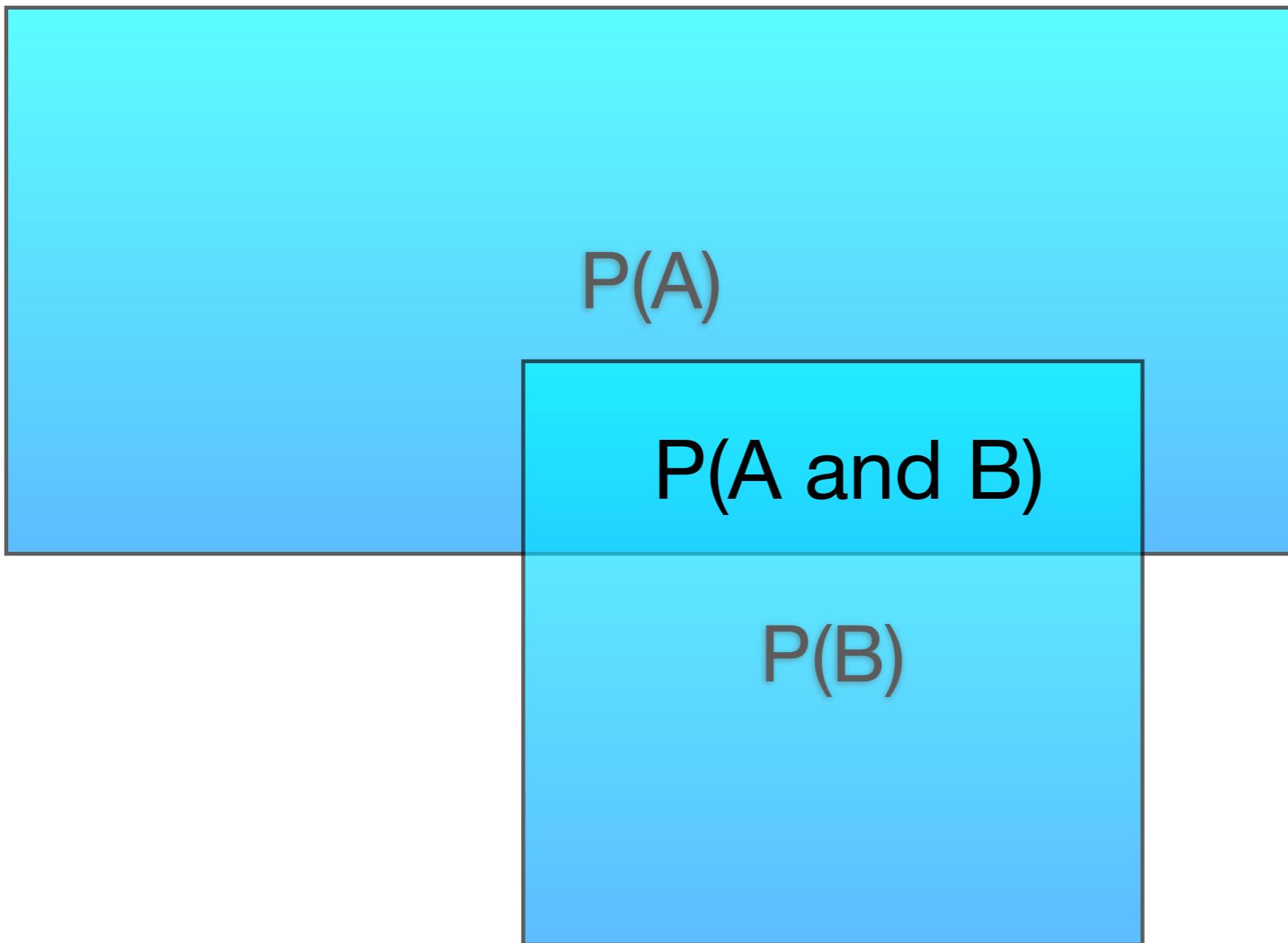


P(A)



P(B)

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$



Bayes Law

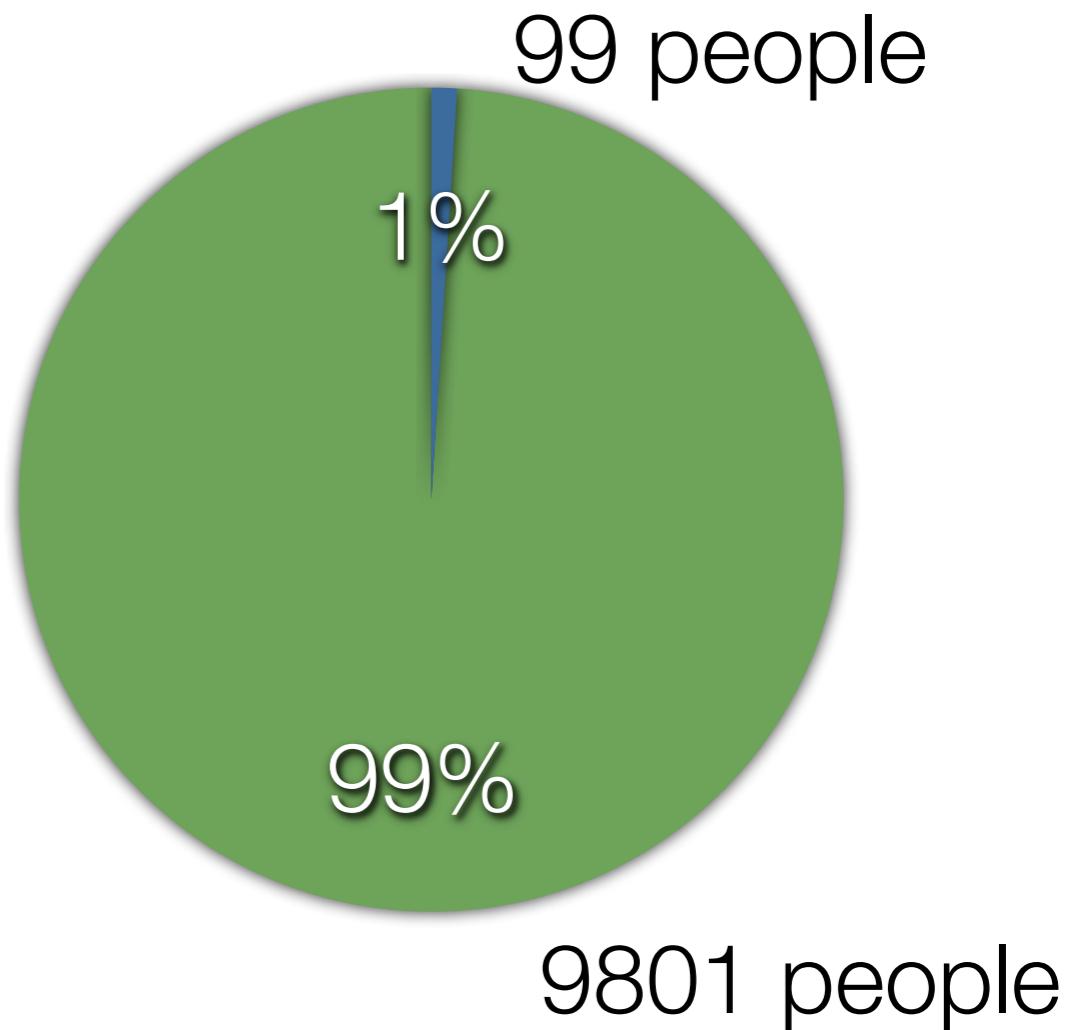
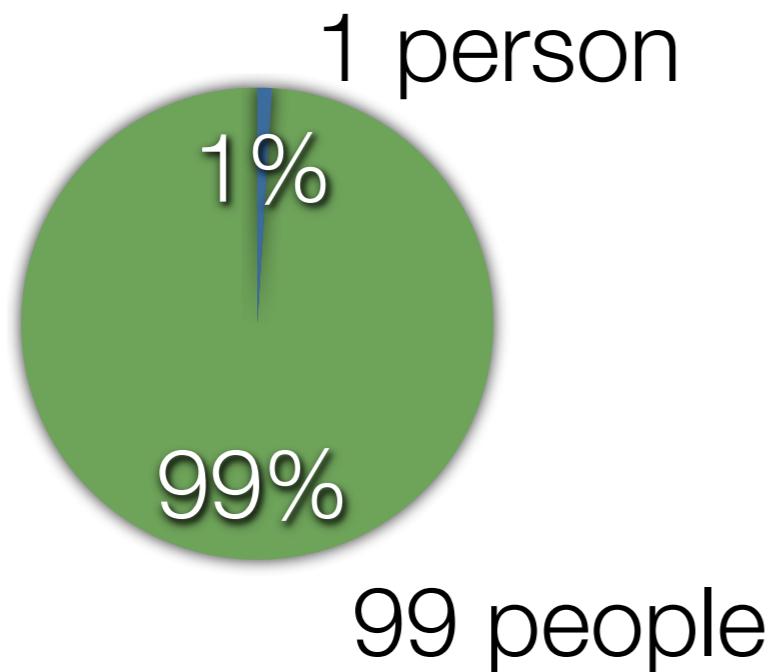


$$P(B | A) = P(A | B) * P(B) / P(A)$$

Example

- In a group of 10,000 people
- 1% of them have a rare disease
- There is a test that is 99% effective
 - 99% of sick patients test positive
 - 99% of healthy patients test negative
- Given a positive result, what is the probability that the patient are sick?

Sick population (100 people) Healthy population (9900 people)



- False positives
- True positives

- False negative
- True negatives

Disease diagnosis

- 99 sick patients test positive, 99 healthy patients test positive
- Give a positive test result, there is a 50% chance that a patient has the disease

Bayesian disease

$$p(\text{sick} \mid \text{test_pos}) = \frac{\frac{99}{100} * \frac{1}{100}}{p(\text{test_pos})} = \frac{99}{198} = \frac{1}{2}$$

$$\frac{99}{10,000} + \frac{99}{10,000} = \frac{198}{10,000}$$

(true positives + false positives)

We know the effectiveness of test (probability of testing positive given being sick), and the prior probability of being sick.

How do we use it for
machine learning?

Naive Bayes Classifier

$$p(C|F_1, \dots, F_n) = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$
$$= \frac{p(C) p(F_1, \dots, F_n|C)}{\sum_{C'} p(C') p(F_1, \dots, F_n|C')}.$$

posterior

evidence

$$p(C) p(F_1, \dots, F_n|C) = p(C, F_1, \dots, F_n)$$

Simplifying assumption

assume each
feature is
independent of
the others

$$\begin{aligned} p(C|F_1, \dots, F_n) &\propto p(C, F_1, \dots, F_n) \\ &\propto p(C) p(F_1|C) p(F_2|C) p(F_3|C) \dots \\ &\propto p(C) \prod_{i=1}^n p(F_i|C). \end{aligned}$$

Text classification example

- NYTimes labels the type of each of its articles
- Has hand curated its categorization for the past 100 years
- Provides a Developer API
- You can get an API key from developer.nytimes.com

API Documentation and Tools

You already know that NYTimes.com is an unparalleled source of news and information. But now it's a premier source of data, too — why just read the news when you can hack it?

Overview

[APIs](#)[Keys](#)[Tools](#)[Forum](#)[Gallery](#)[API Console](#)

Getting Started

The Times Developer Network is our API clearinghouse and community. Here's how to get started:

- Request an [API key](#)
- Read the [API documentation](#), [FAQ](#) and [Terms of Use](#)
- Use the [API Tool](#) to experiment without writing code
- Browse the [application gallery](#)
- Connect with other developers in the [forum](#)

To see your API keys and rate limits, visit the [Keys page](#).

News

New York State Legislature API Deprecated

As of April 15, the [New York State Legislature API](#) is no longer supported. Learn more on the [Open blog](#).

Article Search API: Full Service Restored

Jan. 30: After an outage over the weekend, the [Article Search API](#) is now working properly. Thanks for your patience.

API Logos and Branding Guide

Get Times API logos for your apps: review The New York Times Branding Guide and the Attribution Guidelines.

Featured Applications and Projects

OPEN

A blog written by New York Times developers.



[Go to Open Blog »](#)

Blogroll

NOTABLE SITES

- [Amazon Web Services](#)
The AWS blog
- [BBC Backstage](#)
The BBC's developer network
- [Daylife](#)
- [Google Code](#)
- [Mashable](#)
- [Programmable Web](#)
- [ReadWriteWeb](#)
- [Reuters Labs](#)
- [The Sunlight Foundation](#)
- [Yahoo Developer](#)

OPEN-SOURCE SOFTWARE

- [The Apache Software Foundation](#)
- [Django](#)
- [MySQL](#)
- [PHP](#)
- [Ruby on Rails](#)

JOURNALISM & TECHNOLOGY

- [Adrian Holovaty](#)
Creator of Django
- [NICAR](#)
The National Institute for Computer-Assisted Reporting

Get the Silverlight kit for New York Times APIs.





```
#!/usr/bin/env python
# encoding: utf-8
"""

nytimes_pull.py

Created by Hilary Mason on 2011-02-17.
Copyright (c) 2011 Hilary Mason. All rights reserved.
"""

import urllib
import json

def main(api_key, category, label):

    content = []
    for i in range(0,5):
        # print "http://api.nytimes.com/svc/search/v1/article?query=classifiers_facet:%s&api-key"
        offset=%s" % (category, api_key, i)
        h = urllib.urlopen("http://api.nytimes.com/svc/search/v1/article?query=classifiers_facet"
                           api-key=%s&offset=%s" % (category, api_key, i))
        data = json.loads(h.read())
        for result in data['results']:
            content.append(result['body'])

    f = open(label, 'w')
    for line in content:
        f.write('%s\n' % line.encode('utf-8'))

    f.close()

if __name__ == '__main__':
    main("e561d1c48263398319384f3e1343fe2c:19:68267067", "[Top/Features/Arts]", "arts")
    main("e561d1c48263398319384f3e1343fe2c:19:68267067", "[Top/News/Sports]", "sports")
```

Derian Hatcher scored midway through the third period and Antero Niittymaki stopped 38 shots, helping the Philadelphia Flyers beat the Florida Panthers, 1-0, on Sunday for their third straight victory. It was the third career shutout and first of the season for Niittymaki, who has appeared in 111 games over parts of four seasons with the Flyers.

The University of Hawaii's star-kissed run to the Sugar Bowl has so many improbable story lines that it practically reads like a movie script. The Warriors had four fourth-quarter comeback victories, they compete on a shoestring budget and their two best players, quarterback Colt Brennan and wide receiver Davone Bess, each spent time in jail before

Kerry Collins waited all season to prove his value in Tennessee. The Titans would not have been celebrating Sunday night without him. Collins rescued his teammates in the most precarious of circumstances, leading them to three second-half field goals in relief of the injured Vince Young and helping Tennessee rally for a 16-10 victory over the

One of the risks the NFL Network faced by having CBS and NBC simulcast the Patriots-Giants game Saturday night was the exposure to a wider world it would give Bryant Gumbel, the channel's play-by-play announcer. The NFL Network reaches 43 million cable and satellite subscribers, but the simulcast increased the viewing universe to 113 million TV

It almost sounds like a joke. Florida State's football team will be short-handed and vulnerable in the Music City Bowl against Kentucky on Monday because the Seminoles will be missing players who were suspended for cheating on, of all things, a music test. Ten to 15 years ago, at the height of its prominence in college football, Florida State was

Rutgers seemed refreshed after a nearly three-week layoff. Epiphanny Prince scored 20 points to lead the No. 6 Scarlet Knights to a 70-34 victory against visiting Temple on Sunday. Rutgers had been off for the last 18 days for exams, last playing Dec. 12 against Princeton. "We had a chance to recharge ourselves after playing a brutal schedule

Whether slumped at his locker in a dress shirt and tie, or slumped on the bench in his blue warm-up, Eddy Curry presented an apt portrait Sunday of all that has gone horribly wrong for the Knicks. The day began with Curry out of the starting lineup for the first time in two seasons, the latest desperation move by Coach Isiah Thomas to break his

For a split second in overtime Sunday night, Rangers center Chris Drury and forward Brendan Shanahan spotted the same sliver of space that makes the difference between an inconsequential play and a game-winning goal. Drury snagged the puck at mid-ice after a Montreal turnover and saw Shanahan -- 38-year-old legs and all -- skating past the

This was at an N.B.A. Board of Governors meeting in New York not long ago, James L. Dolan facing his colleagues-in-ownership for the first time since the sexual harassment trial brought by Anucha Browne Sanders had slimed Madison Square Garden and by extension the league. According to a person who attended



Google Prediction API



907

Feedback on this document

▶ Getting Started

▶ Using Prediction

▶ Reference

▶ Tools and Resources



What is the Google Prediction API?

▶ Feedback and Discussion



Customer sentiment analysis

▶ Related APIs



Message routing decisions

FAQ



Document and email classification

Terms of Service



Churn analysis



Recommendation systems



Spam detection



Upsell opportunity analysis



Diagnostics



Suspicious activity identification

*And much more...*

RESTful API

Asynchronous cloud-based training, automatic model selection and tuning, and the ability to add training data on the fly.

Flexible Input

Numeric or text input that can output hundreds of discrete categories or continuous values.

Multi-Platform Support

Google App Engine, Apps Script (Google Docs), web & desktop apps, and the command line.

Open Source Code

- We will look at open source code in Python for doing Machine Learning
- scikit-learn: classification, regression, clustering, etc.
- Hilary Mason's video: Introduction to Machine Learning with Web Data
- https://github.com/hmason/ml_class

```
#!/usr/bin/env python
# encoding: utf-8
"""
classify.py

Created by Hilary Mason on 2011-02-17.
Copyright (c) 2011 Hilary Mason. All rights reserved.
"""

import re, string

from nltk import FreqDist
from nltk.tokenize import word_tokenize
from nltk.stem.porter import PorterStemmer

class NaiveBayesClassifier(object):

    def __init__(self):
        self.feature_count = {}
        self.category_count = {}

    def probability(self, item, category):
        """
        probability: prob that an item is in a category
        """
        category_prob = self.get_category_count(category) / sum(self.category_count.values())
        return self.document_probability(item, category) * category_prob

    def document_probability(self, item, category):
        features = self.get_features(item)

        p = 1
        for feature in features:
            print "%s - %s - %s" % (feature, category, self.weighted_prob(feature, category))

-uuu:---F1  classify.py      Top L1      (Python)-----
```

```
w_prob = ((weight*ap) + (totals * basic_prob)) / (weight + totals)
return w_prob

def train(self, item, category):
    features = self.get_features(item)

    for f in features:
        self.increment_feature(f, category)

    self.increment_cat(category)

if __name__ == '__main__':
    labels = ['arts', 'sports'] # these are the categories we want
    data = {}
    for label in labels:
        f = open(label, 'r')
        data[label] = f.readlines()
        # print len(data[label])
        f.close()

    nb = NaiveBayesClassifier()
    nb.train_from_data(data)
    print nb.probability("Early Friday afternoon, the lead negotiators for the N.B.A. and the playe
s union will hold a bargaining session in Beverly Hills – the latest attempt to break a 12-month st
lēmate on a new labor deal.", 'arts')
    print nb.probability("Early Friday afternoon, the lead negotiators for the N.B.A. and the playe
s union will hold a bargaining session in Beverly Hills – the latest attempt to break a 12-month st
lēmate on a new labor deal.", 'sports')
```

the - arts - 0.915670103093
a - arts - 0.618356164384
, - arts - 0.975051546392
. - arts - 0.154545454545
12-month - arts - 0.5
Beverly - arts - 0.193333333333
Early - arts - 0.5
Friday - arts - 0.132
Hills - arts - 0.5
N.B.A. - arts - 0.25
afternoon - arts - 0.17
and - arts - 0.7368
attempt - arts - 0.14
bargaining - arts - 0.5
break - arts - 0.25
deal - arts - 0.5
for - arts - 0.304444444444
hold - arts - 0.26
in - arts - 0.657837837838
labor - arts - 0.5
latest - arts - 0.25
lead - arts - 0.1
negotiators - arts - 0.5
new - arts - 0.155
on - arts - 0.460888888889
players - arts - 0.125
session - arts - 0.5
stalemate - arts - 0.5
to - arts - 0.598461538462
union - arts - 0.5
will - arts - 0.151666666667
-- arts - 0.5
6.1474279934e-16

the - sports - 0.994845360825
a - sports - 0.815616438356
, - sports - 0.935463917526
. - sports - 0.118181818182
12-month - sports - 0.5
Beverly - sports - 0.166666666667
Early - sports - 0.5
Friday - sports - 0.132
Hills - sports - 0.5
N.B.A. - sports - 0.26
afternoon - sports - 0.125
and - sports - 0.7368
attempt - sports - 0.155
bargaining - sports - 0.5
break - sports - 0.26
deal - sports - 0.5
for - sports - 0.578222222222
hold - sports - 0.25
in - sports - 0.795945945946
labor - sports - 0.5
latest - sports - 0.26
lead - sports - 0.164
negotiators - sports - 0.5
new - sports - 0.14
on - sports - 0.421777777778
players - sports - 0.17
session - sports - 0.5
stalemate - sports - 0.5
to - sports - 0.677230769231
union - sports - 0.5
will - sports - 0.133333333333
-- sports - 0.5
2.06771958824e-15