

Machine Learning The Art and Science of Algorithms that Make Sense of Data

Machine Learning for SA

AND MACHINE LEARNING CHRISTOPHER M. BISHOP



DEEP LEARNING Grans Carrier

Lyle Ungar Computer and information Science University of Pennsylvania



Tutorial outline

Core sentiment analysis (SA) methods

- Simple: using lexica (dictionaries)
- Aspect-based: using information extraction

Machine learning for SA

- Unsupervised: open language SA
- Supervised: regression and deep learning

SA extensions

- Post, person and community
- Multi-media





SA Goals and methods

- Prediction vs. Insight
 - Most SA is for insight
 - Not "how well did people like this product?"
 - But "what did they like or dislike about it?"
- Closed vs. Open vocabulary
 - Closed: lexica selected a priori
 - Open: look to see what words correlate



Outline

Unsupervised learning: Open language SA

- Differential language analysis
- LDA
- Visualization
- Supervised learning
 - Regression and regularization
 - Semi-supervised learning
 - Deep learning
 - CNN, GRNN, LSTM
 - SA deep learning example: Socher



SA for hospitals from Yelp

US hospitals are rated by formal surveys (HCAHPS) and on Yelp



Differential Language Analysis

- To get insight into what drives the difference in sentiment, compare the words in 1 and 5 star reviews
 - Find the words most correlated with high/low reviews
 - Sometimes controlling for demographics
 - Display words based on correlation (size) and frequency (color)



Yelp - positive

SINE MORIBUS



7

Yelp - negative



Yelp finds 'missing' topics

| Topic name | Correlation, Pearson's r | Covered by HCAHPS |
|---|---|------------------------|
| YELP TOPICS MOST CORRELATED WITH PO | SITIVE YELP RATINGS | |
| Caring doctors, nurses, and staff Comforting Clean, private, nice hospital rooms Surgery/procedure and peri-op Labor and delivery | 0.46 0.29 0.25 0.23 0.20 | No No No No |
| YELP TOPICS MOST CORRELATED WITH NE | GATIVE YELP RATINGS | |
| Horrible hospital Rude doctor/nurse communication Pain control Insurance and billing Cost of hospital visit | -0.33 -0.29 -0.28 -0.26 -0.26 | Yes Yes No No |



LDA topics

SINE MORIBUS

| | Pediatric (son, daughter, hospital, child, kids, time, children's) | " My goddaughter was treated at [hospital] for a congenital heart defect[The hospital] has gone above and beyond for my family in every respect possible, and its doctors and nurses have the utmost care and professionalism." |
|-------------------------------------|--|--|
| | Labor and delivery (baby, birth, nurses, labor, delivery, experience, nurse) | "Excellent for having natural childbirth, they respected all my wishes. I had a midwife and my own doula." |
| Specific type of medical care | Surgical management and physical therapy (dr, physical, back, surgery, mri, therapy, year) | "This is one of the best rehab hospitals in the world. First class service by therapists, doctors, nurses, aids dietary, housekeeping. They taught me to walk again after total knee replacement." |
| | Wait time in ED (room, emergency, er, waiting, hours, wait, time) | "They suck!!!!! They will make u wait for 3 hrs in the emergency waiting room even though the waiting room is empty!!!!!" |
| | ED imaging after injury (er, broken, back, accident, x-ray, foot, car) | "Their ER is awesome! I came in with a possible broken toe. They took me in right away, x-rayed my foot and diagnosed it (yup, broken toe) 1.5 hours from arrival to dismissal. Best ER service Ive ever experienced." |
| | | |

| Yelp Domains | Yelp Topic (topic terms) | Example Quote | |
|-------------------------------|---|---|----|
| Cost of hospital visit | Cost of hospital visit (insurance, bill, pay, cost, charge, visit, hospital) | "Ask what you will be charged before you allow this place to touch you! [T]heir charges are excessive. In hospital 18 hours and got a 36k bill for knee surgery that cost 1/3 the amount in Los Angeles." | |
| Insurance and Billing | Insurance and billing (insurance, billing, bill, hospital, department, company, paid) | "Billing department is awful. They billed the wrong insurance company twice and then finally had the correct one but billed it incorrectly Instead of contacting me they sent the bill to a collection agency" | |
| Ancillary testing | Ancillary testing and results (blood, test, doctor, results, tests, lab, work) | "Service way too slowfor relatively standard testing90 mins so far and sample hasn't even been collected because doc needs lab to tell him what tests need to be ordered." | |
| | Facility (place, big, hand, time, TV, watch, cool) | "This place is awesomeseek out the quite courtyards for some solitude." | |
| Facilities | Security and front desk (front, security, desk, room, friend, back, waiting) | "There was no sign that said to check-in, so I didn't bother with going to the information deskThe security guard who told us to check-in was rude" | |
| Amonities | Parking (parking, hospital, building, free, lot, nice, valet) | "2 hours free parking with validation I like it how they are building a garden between the parking structure and the hospital." | |
| Ameniues | Hospital food (food, cafeteria, hospital, good, eat, coffee, order) | "This review is for the cafeteria, prices are super cheap and the burger and nachos were terrific. Chili cheese fries were soggy" | 11 |



"Bag of Words" Models

Assume that all the words within a document are exchangeable.

• The order of the words doesn't matter, just the count





Mixture of Unigrams = Naïve Bayes





Plate Model

(equivalent)

Mixture of Unigrams = Naïve Bayes

Model: For each document:

- **Choose a topic** z_d with $p(topic_i) = \theta$
- Choose N words w_n by drawing each one independently from a multinomial conditioned on z_d with $p(w_n = word_i | topic_i = z) = \beta_z$
 - Multinomial: take a (non-uniform prior) dice with a word on each side; roll the dice N times and count how often each word comes up
- In NB, we have exactly one topic per document

4

LDA: Each doc is a mixture of topics

LDA lets each document be a (different) mixture of topics

- Naïve Bayes assumes each *document* is on a single topic
- LDA lets each *word* be on a different topic
- For each document,
 - Choose a multinomial distribution θ_d over topics for that document
 - For each of the *N* words w_n in the document
 - -Choose a topic z_n with $p(topic) = \theta_d$
 - -Choose a word w_n from a multinomial conditioned on z_n with $p(w=w_j|topic=z_n)$
 - —Note: each topic has a different probability of generating each word



Dirichlet Distributions

- In the LDA model, we want the *topic mixture proportions* for each document to be drawn from some *distribution*.
 - *distribution* = "probability distribution", so it sums to one
- So, we want to put a prior distribution on multinomials. That is, ktuples of non-negative numbers that sum to one.
 - We want probabilities of probabilities
 - These multinomials lie in a (k-1)-*simplex*
 - *Simplex* = generalization of a triangle to (k-1) dimensions.

• Our prior:

- Needs to be defined for a (k-1)-simplex.
- Conjugate to the multinomial



3 Dirichlet Examples (over 3 topics)



Corners: only one topic Center: uniform mixture of topics

SINE MORIDUS

Dirichlet Distribution

$$p(\theta|\alpha) = \frac{\Gamma(\sum_{i=1}^{k} \alpha_i)}{\prod_{i=1}^{k} \Gamma(\alpha_i)} \prod_{i=1}^{k} \theta_i^{\alpha_i - 1}$$

Dirichlet distribution

- is defined over a (k-1)-simplex. I.e., it takes k non-negative arguments which sum to one.
- is the conjugate prior to the multinomial distribution.
 - I.e. if our likelihood is multinomial with a Dirichlet prior, then the posterior is also Dirichlet
- The Dirichlet parameter $\alpha_{\rm i}$ can be thought of as the prior count of the i^{\rm th} class.

• For LDA, we often use a "symmetric Dirichlet" where all the α are equal;

• α is then a "concentration parameter"

Effect of α

- When α < 1.0, the majority of the probability mass is in the "corners" of the simplex, generating mostly documents that have a small number of topics.
- When $\alpha > 1.0$, the most documents contain most of the topics.



The LDA Model



- Then choose a word $w_n \sim Multinomial(\beta_z)$
 - Where each topic has a different parameter vector β for the words



Lyle H Ungar, University of Pennsylvania

ZU

The LDA Model: "Plate representation"



For each of M documents

- Choose the topic distribution $\theta \sim \text{Dirichlet}(\alpha)$
- For each of the N words w_n:
 - Choose a topic z ~ Multinomial(θ)
 - Choose a word $w_n \sim \text{Multinomial}(\beta_z)$



Parameter Estimation

Given a corpus of documents, find the parameters α and β which maximize the likelihood of the observed data (words in documents), marginalizing over the hidden variables θ , z

- E-step:
 - Compute $p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta)$, the posterior of the hidden variables (θ, \mathbf{z}) given a document \mathbf{w} , and hyper-parameters α and β .
- ♦ M-step
- Estimate parameters α and β given the current hidden variable distribution estimates

 θ: topic distribution for the document,
 z: topic for each word in the document

 Unfortunately, the E-step cannot be solved in a closed form

So people use a "variational" approximation Lyle H Ungar, University of Pennsylvania

Variational Inference





•In variational inference, we consider a simplified graphical model with variational parameters γ , ϕ and minimize the KL Divergence between the variational and posterior distributions.

• q approximates p

 $(\gamma^*, \phi^*) = \arg\min_{(\gamma, \phi)} KL(q(\theta, z | \gamma, \phi) || p(\theta, z | w, \alpha, \beta))$



Parameter Estimation: Variational EM

 Given a corpus of documents, find the parameters α and β which maximize the likelihood of the observed data.

• E-step:

 Estimate the variational parameters γ and φ in q(γ,φ;α,β) by minimizing the KL-divergence to p (with α and β fixed)

M-step

 Maximize (over α and β) the lower bound on the log likelihood obtained using q in place of p (with γ and φ fixed)



| "Arts" | ${\bf ``Budgets''}$ | "Children" | "Education" |
|---------|---------------------|------------|-------------|
| | | | |
| NEW | MILLION | CHILDREN | SCHOOL |
| FILM | TAX | WOMEN | STUDENTS |
| SHOW | PROGRAM | PEOPLE | SCHOOLS |
| MUSIC | BUDGET | CHILD | EDUCATION |
| MOVIE | BILLION | YEARS | TEACHERS |
| PLAY | FEDERAL | FAMILIES | HIGH |
| MUSICAL | YEAR | WORK | PUBLIC |
| BEST | SPENDING | PARENTS | TEACHER |
| ACTOR | NEW | SAYS | BENNETT |
| FIRST | STATE | FAMILY | MANIGAT |
| YORK | PLAN | WELFARE | NAMPHY |
| OPERA | MONEY | MEN | STATE |
| THEATER | PROGRAMS | PERCENT | PRESIDENT |
| ACTRESS | GOVERNMENT | CARE | ELEMENTARY |
| LOVE | CONGRESS | LIFE | HAITI |

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

a and a second

LDA requires fewer topics than NB





Lyle H Ungar, University of Pennsylvania

20

There are many LDA extensions

The author-topic model





Ailment Topic Aspect Model



Observed

word w

SINE MORIBUS

aspect y = symptom, treatment or other Hidden

topic type: background? (I), non-ailment (x)

- Set the background switching binomial λ
- Draw an ailment distribution $\eta \sim \text{Dir}(\sigma)$
- Draw word multinomials $\phi \sim \text{Dir}(\beta)$ for the topic, ailment, and background distributions
- For each message $1 \le m \le D$:
 - Draw a switching distribution $\pi \sim \text{Beta}(\gamma_0, \gamma_1)$
 - Draw an ailment $a \sim \operatorname{Mult}(\eta)$
 - Draw a topic distribution $\theta \sim \text{Dir}(\alpha_a)$
 - For each word $w_i \in N_m$
 - Draw aspect $y_i \in \{0, 1, 2\}$ (observed)
 - Draw background switcher $\ell \in \{0,1\} \sim \operatorname{Bi}(\lambda)$
 - If $\ell == 0$:
 - Draw $w_i \sim \text{Mult}(\phi_{B,y})$ (a background)
 - Else:
 - Draw $x_i \in \{0,1\} \sim \operatorname{Bi}(\pi)$
 - If $x_i == 0$: (draw word from topic z)
 - Draw topic $z_i \sim \text{Mult}(\theta)$
 - Draw $w_i \sim \operatorname{Mult}(\phi_z)$
 - Else: (draw word from ailment a aspect y)
 - Draw $w_i \sim \operatorname{Mult}(\phi_{a,y})$
- Lyle H Ungar, University of Pennsylvania

Paul & Dredze

LDA generation - example

- Topics = {sports, politics}
- Words = {football, baseball, TV, win, president}
- α = (0.8,0.2)
- β=

| | sports | politics |
|-----------|--------|----------|
| football | 0.3 | 0.01 |
| baseball | 0.25 | 0.01 |
| TV | 0.1 | 0.15 |
| win | 0.3 | 0.25 |
| president | 0.01 | 0.2 |
| | 0.04 | 0.38 |

Lyle H Ungar, University of Pennsylvania

ZЭ

LDA generation - example

♦ For each document, d

- Pick a topic distribution, θ_d using α
- For each word in the document
 - pick a topic, z
 - given that topic, pick a word using β

α = (0.8,0.2)

| ♦ β = | sports politics |
|-----------|---|
| football | 0.3 0.01 |
| baseball | 0.25 0.01 |
| TV | 0.1 0.15 |
| win | 0.3 0.25 |
| president | 0.01 Lyle HOundar, University of Pennsylvania |

LDA in practice

♦ Be careful what your unit of analysis is

- Obvious: use each review as a document
 - Reviews are mostly Zipfian
 - Most book reviews are about Harry Potter
 - So most of the topics end up capturing product features
 - ◆ Harry, Voldemort, Hermione, ...
- Often better: group reviews for each product (book, hospital, ..)



LDA in practice

- Use a good package
 - Mallet or Factorie

(Usually) build topics for the whole corpus

• And observe how topic frequencies differ for different products or sentiments

Pick an appropriate number of topics

- 2,000 for all of Facebook or Twitter
- 200 for a specialized domain like hospital reviews



Yields 2 happy topics

happy, Christmas, year, family, friends, hope, merry, Thanksgi happy, birthday, day, love, wishes, mom, miss, wonderful, dad

Kern, M. L., Park, G., Eichstaedt, J. C., Schwartz, H. A., Sap, M., Smith, L. K., & Ungar, L. H. (2016, August 8). Gaining Insights From Social Media Language: Methodologies and Challenges. *Psychological Methods*. Advance online publication. http://dx.doi.org/10.1037/met0000091



Yields 8 happy topics

day, happy, mothers, mother's, mom, mother, wonderful, moms, mommy, mama day, happy, valentines, fathers, valentine's, father's, dad, independence, dads, sing birthday, happy, wishes, bday, wished, b-day, birthdays, present, celebrate, cake year, happy, 2010, 2011, joy, wishing, bring, happiness, safe, diwali happy, 4th, July, Halloween, year, fireworks, safe, fourth, holiday, holidays happy, thanksgiving, Easter, family, thankful, turkey, holiday, bunny, enjoy, eggs happy, birthday, anniversary, wishing, brother, son, bday, daddy, mommy, celebra happy, makes, sooo, soo, soooo, Easter, Thanksgiving, camper, ending, sooooo



Yields 20 happy topics

happy, birthday, mommy, daddy, mama, momma, dearest, bestest, 21st, 18th happy, birthday, sis, lil, bday, b-day, luv, cousin, 21st, nephew happy, mothers, mother's, mom, moms, mother, mommy, mom's, mama, mommies happy, makes, camper, unhappy, extremely, happier, smiling, satisfied, contented, content happy, diwali, wishing, eid, happiness, mubarak, holi, festival, prosperous, gibran Easter, happy, bunny, eggs, egg, hunt, holidays, risen, candy, basket happy, birthday, brother, wishing, 18th, 21st, xxxx, 16th, monthsary, nephew year, happy, 2010, 2011, chinese, 2009, cheers, prosperous, tiger, rabbit happy, independence, friendship, valentines, Canada, valentine's, republic, memorial, Aust year, happy, joy, happiness, bring, 2010, 2011, health, wishing, brings happy, fathers, father's, dad, dads, father, daddy, dad's, mothers, papa 4th, July, happy, fireworks, fourth, safe, independence, bbq, 5th, quarter happy, birthday, celebrate, anniversary, celebrating, birthdays, dad's, b-day, b'day, mom's happy, valentines, valentine's, single, valentine, hump, pi, awareness, singles, v-day happy, birthday, grandma, mama, aunt, Beth, Mary, anniversary, papa, grandpa

Yields 1 play topic

game, play, win, playing, football, team, won, games, beat, lets



Yields 5 play topics

guitar, play, playing, music, piano, band, bass, hero, practice, played game, football, play, soccer, basketball, playing, games, team, practice, bas place, chuck, find, meet, play, birth, Norris, interesting, babies, profile play, playing, game, games, xbox, halo, Wii, video, Mario, 360 play, playing, game, ball, games, played, golf, tennis, poker, cards



Yields 25 play topics

golf, played, ultimate, frisbee, mini, ball, balls, golfing, tennis, disc play, game, let's, role, sims, rules, chess, basketball, plays, poker words, comment, note, play, wake, jail, copy, paste, sport, fair black, cod, ops, playing, play, mw2, modern, warfare, ps3, online game, team, won, win, played, boys, soccer, season, proud, football soccer, football, game, play, team, basketball, playing, ball, practice, field kids, park, playing, boys, played, pool, blast, playground, swimming, toys sand, beach, water, toes, carl, grain, playin, mountain, rocks, desert music, band, playing, piano, guitar, songs, sound, metal, bass, played na, stuck, everyday, ki, replay, melody, ami, er, ta, singin http://www.youtube.com, feature, related, =p, marcus, channel, double, user, nr, youtube gdat guitar, bass, drum, playing, amp, drums, string, strings, electric, acoustic play, guitar, learn, piano, learning, playing, learned, lessons, songs, rules games, play, playing, game, video, played, card, board, begin, playin play, playing, starcraft, warcraft, sims, ii, beta, online, nerds, nerd watchin, sittin, chillin, waitin, doin, havin, gettin, eatin, playin, drinkin pokemon, playing, mon, shiny, version, pikachu, pok, cards, ds, ash player, dvd, cd, record, printer, bought, set, mp3, ink, borrow anima manga nariita blaach anisada sarias coonlaw anisadas alchamist Iananasa

38

Find how topics correlate with ...

Positive/negative sentiment

• For different product categories



Yelp hospital reviews again

| | | Cost of hospital visit | "Ask what you will be charged before you allow this place | - |
|-------------|-----------------------------|---|---|-------------|
| - | Cost of | (insurance, bill, pay, | to touch you! [T]heir charges are excessive. In hospital 18 | |
| | hospital visit | cost, charge, visit, | hours and got a 36k bill for knee surgery that cost 1/3 the | |
| | | hospital) | amount in Los Angeles." | |
| | | Insurance and billing | "Billing department is awful. They billed the wrong | |
| | Insurance and | (insurance, billing, bill, | insurance company twice and then finally had the correct | - |
| - | Billing | hospital, department, | one but billed it incorrectly Instead of contacting me they | |
| | | company, paid) | sent the bill to a collection agency" | |
| | | Ancillary testing and | "Service way too slowfor relatively standard testing90 | |
| | Ancillary | results (blood, test, | mins so far and sample hasn't even been collected because | |
| | testing | doctor, results, tests, lab, | doc needs lab to tell him what tests need to be ordered." | - |
| | | work) | | |
| | | Eagility (place hig | "This place is average a solt out the quite courtwards for | |
| | | Facility (place, big, | This place is awesomeseek out the quite courtyards for | |
| + | | hand, time, TV, watch, | some solitude." | + |
| + | | hand, time, TV, watch, cool) | some solitude." | + |
| + | Facilities | hand, time, TV, watch, cool) Security and front | "There was no sign that said to check-in, so I didn't bother | + |
| + | Facilities | hand, time, TV, watch, cool) Security and front desk (front, security, | "There was no sign that said to check-in, so I didn't bother with going to the information deskThe security guard | ÷ |
| + | Facilities | hand, time, TV, watch, cool) Security and front desk (front, security, desk, room, friend, back, | "There was no sign that said to check-in, so I didn't bother with going to the information deskThe security guard who told us to check-in was rude" | + |
| + | Facilities | hand, time, TV, watch, cool) Security and front desk (front, security, desk, room, friend, back, waiting) | "There was no sign that said to check-in, so I didn't bother with going to the information deskThe security guard who told us to check-in was rude" | + |
| + | Facilities | hand, time, TV, watch, cool) Security and front desk (front, security, desk, room, friend, back, waiting) Parking (parking, | "This place is awesomeseek out the quite courtyards for some solitude." "There was no sign that said to check-in, so I didn't bother with going to the information deskThe security guard who told us to check-in was rude" "2 hours free parking with validation I like it how they are | + - + |
| + + + | Facilities | hand, time, TV, watch, cool) Security and front desk (front, security, desk, room, friend, back, waiting) Parking (parking, hospital, building, free, | This place is awesomeseek out the quite courtyards for some solitude." "There was no sign that said to check-in, so I didn't bother with going to the information deskThe security guard who told us to check-in was rude" "2 hours free parking with validation I like it how they are building a garden between the parking structure and the | + - + |
| + - + | Facilities | Facility (place, big, hand, time, TV, watch, cool)Security and front desk (front, security, desk, room, friend, back, waiting)Parking (parking, hospital, building, free, lot, nice, valet) | This place is awesomeseek out the quite courtyards for some solitude." "There was no sign that said to check-in, so I didn't bother with going to the information deskThe security guard who told us to check-in was rude" "2 hours free parking with validation I like it how they are building a garden between the parking structure and the hospital." | + - + |
| + + | Facilities Amenities | hand, time, TV, watch, cool) Security and front desk (front, security, desk, room, friend, back, waiting) Parking (parking, hospital, building, free, lot, nice, valet) Hospital food (food, | This place is awesomeseek out the quite courtyards for some solitude." "There was no sign that said to check-in, so I didn't bother with going to the information deskThe security guard who told us to check-in was rude" "2 hours free parking with validation I like it how they are building a garden between the parking structure and the hospital." "This review is for the cafeteria, prices are super cheap and | + + + |



Unsupervised learning

Differential language analysis

- Using words
- Using topics

LDA

- Each document is a mixture over topics
- Each topic looks like a Naïve Bayes model
 - It produces words with some probability
- Estimation of LDA is messy
 - Requires variational EM or Gibbs sampling
- Lots of extensions and variations





Key methods

Linear and logistic regression

• Requires regularization

Deep Learning

- Static convolutional neural nets
- **Dynamic** GRNNs and LSTMs



Regression for sentiment analysis

Given a set of observations with labels, y

- Observations
 - Sentences or reviews labeled with a sentiment score
- ♦ Generate features, x, for each observation
 - E.g. presence or count of a 'word'
- Learn a regression model to predict y
 - $y = f(x) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 \dots$
 - Most of the w_i are zero.



Overfitting and Regularization

 With 10,000 labeled tweets and a vocabulary of 40,000 words, you can predict sentiment *really* well



Regularization approaches

Use PCA or LDA to to dimensionality reduction

- 40,000 words -> 200 topics or factors
 - Often semi-supervised

Do feature selection

Zero out some words

Use Ridge or Lasso or other penalty to "shrink" weights



A word on words

- "I'm sick of dos movies! :("
- How many words?
 - 5? 6? 7? More?
- I'm sick_of dos movies !:(
- Remove punctuation and stopwords?
- Normalize spelling?
 - Depends on the application
 - often yes for product names
 - No for other words



To find multi-word expressions

Use Pointwise mutual information

p(w₁,w₂) / p(w₁)p(w₂)



Transform word counts

- Linear regression assumes linearity (duh!)
 - $y = f(x) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 \dots$
 - But seeing a word 10 vs. 11 times is not the same as seeing it never vs. once.

So transform word counts

- Easy: square root of count
- Often better: Anscombe transformation

 $p(phrase) + \frac{3}{2}$

Then normalize for document length to get word frequency



Now we're ready for regression

Linear regression if y is real-valued

- But may require transformation
- Logistic regression if y is binary



Feature selection for regression

- Goal: minimize error on a test set
- Approximation: minimize a penalized training set error
 - Argmin_w (Err + $\lambda |w|_p^p$) where Err = $\sum_i (y_i \sum_j w_j x_{ij})^2 = ||y w^T x||^2$
 - Different norms
 - p = 2 "ridge regression"
 - Makes all the w's a little smaller
 - p = 1 "LASSO" or "LARS" (least angle regression)
 - drives some w's to zero; makes all smaller

p = 0 – "stepwise regression" – the number of features

 Zeros some w's; leaves rest untouched the confusion in the names of the of optimization method with

♦ How to pick λ?

Artificial Intelligence

the objective function

Solving with regularization penalties

- Argmin_w $|\mathbf{y} \mathbf{w} \cdot \mathbf{x}|_2^2 + \lambda |\mathbf{w}|_p^p$
- ♦ p = 2 (Ridge)
 - Closed form: (X'X + λ I)⁻¹ X'y
- ♦ p = 1 (Lasso)
 - Convex optimization
- - Search

"Elastic net" combines Ridge and Lasso widely used convex



Copyright © Andrew W. Moore

Regularization Summary

- Penalized error approximates test error
 - L_2 , L_1 , L_0 penalties
 - Convex and non-convex
 - Use tests set/cross validation to evaluate methods
- ◆ PCA gives roughly similar results to Ridge (L₂)
- Highest accuracy from hybrid methods
 - But requires careful regularization

If you have many more features than observations you need to zero out some of the features