

Computer-assisted content analysis 2

Will Lowe Princeton University

James Lo University of Southern California

Session 1: Classical Content Analysis

Session 2:

- Document classification

 - The direct approach to classification

 - The indirect approach to classification

 - Evaluation and interpretation

 - Topic models

Session 3: Scaling Models

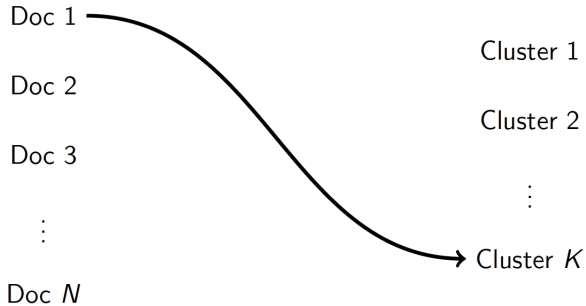
Classification Approaches when Categories are Known

We often want to read text to put them into categories

Examples:

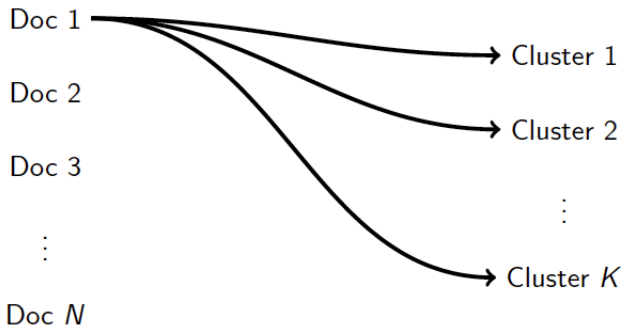
1. Are campaign advertisements positive or negative?
2. What policy areas do newspaper editorials cover?
3. Are international statements belligerent or peaceful?
4. Do court letters liberal or conservative?
5. Is this email spam?

Classification



For classification, each document belongs to a single cluster

Topic Models



In LDA, each document can contain multiple topics

Classification Approaches when Categories are Known

Yesterday: Classification using dictionary approach.

An alternative is **supervised machine learning** methods:

1. coders categorize a set of documents by hand
2. the algorithm “learns” how to sort the documents in categories
3. characteristics of training set are used to assign new documents to categories.

Naive Bayes, Maximum Entropy, Support Vector Machines, Neural Networks, Bagging, Boosting, ...

Classification Approaches when Categories are Known

Let θ_k be the *probability* that $Z=k$ for each document

Words are denoted $\{W\}$

Each document has a *single* topic Z

Some topic labels are observed

We essentially want $P(Z | \{W\}) = \theta_{z|w}$

Supervised machine learning

Tries to give good predictions of observed Y given X

Naive Bayes, Maximum Entropy, Support Vector
Machines, Neural Networks, Bagging, Boosting, ...

Classification approaches

As before, we have two approaches

Discriminative: Model $P(Z | \{W\}) = \theta_{z|w}$ directly

Generative: Model $P(\{W\} | Z)$ and $P(Z) = \theta_z$, then get $P(Z | \{W\}) = \theta_{z|w}$ via Bayes theorem

Bayes theorem is about inverse probabilities

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

If $P(\{W\} | Z)$ known, estimate $P(Z | \{W\})$

Either way...

Desirable classification outcome:

	$P(Z = \text{'Domestic'} \mid \{W\}_d)$	$P(Z = \text{'Foreign'} \mid \{W\}_d)$
D_1	0.75	0.25
D_2	0.82	0.18
...
D_{N-1}	0.02	0.98
D_N	0.45	0.55

where $\theta_{z|w} = P(Z = z \mid \{W\})$

The Basic Steps

1. Construct a training set (a) create a coding scheme
(b) select documents (ideally randomly sampled)
2. Apply the supervised learning method to learn features
of a training set and infer labels
3. Validate and classify remaining documents

We will do this in the lab together with NY Times titles

An Applied Example: Affirmative Action

Evans et al. (2007) apply naive Bayes to legal text

Affirmative action programs at university level

Examine Amicus Curiae briefs sent to appellate court

Can words tell us direction of AC brief?

Are the briefs liberal or conservative?

Naive Bayes

Naive Bayes is a relatively old (~ 1975) classification method

Want to guess if document j is liberal/conservative based on its word profile $\{W\}_j$.

Probability can be derived by applying Bayes theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$
$$P(Z = \text{'Lib'} \mid \{W\}_j) = \frac{P(\{W\}_j \mid Z = \text{'Lib'}) P(Z = \text{'Lib'})}{P(\{W\}_j)}$$

Naive Bayes

$$P(Z = \text{'Lib'} \mid \{W\}_j) \propto P(\{W\}_j \mid Z = \text{'Lib'}) P(Z = \text{'Lib'})$$

We can drop $P(\{W\}_j)$ since it is constant across categories.

Given a representative training set, estimating $P(Z = L)$ is easy:

$$\hat{P}(Z = \text{'Lib'}) = \frac{\text{\# training docs that are liberal}}{\text{\# of training docs}}$$

Naive Bayes

Estimating the probability that a word profile $\{W\}_j$ occurs given that the document is liberal $P(\{W\}_j | Z = \text{'Lib'})$ is more challenging, because any one word profile is likely to occur only once.

Assumption: words are assumed to be generated *independently* given the category Z

$$P(\{W\}_j | Z = \text{'Lib'}) = \prod_i P(W_i | Z = \text{'Lib'})$$

$$P(\text{'Affirmative Action'} | Z = \text{'Lib'}) = P(\text{'Affirmative'} | Z = \text{'Lib'}) \cdot P(\text{'Action'} | Z = \text{'Lib'})$$

Naive Bayes

With this assumption, we can estimate the probability of observing a word i given that the document is liberal: proportion of word i in liberal training set.

The classifier then chooses the class Z (Liberal or Conservative) with the highest aggregate probability.

Note that every new word adds a bit of information that re-adjusts the conditional probabilities.

Naive Bayes

Note that with two classes (here: liberal and conservative) this has a rather neat interpretation:

$$\frac{P(Z = \text{'Lib'} \mid \{W\}_j)}{P(Z = \text{'Con'} \mid \{W\}_j)} = \prod_i \frac{P(W_i \mid Z = \text{'Lib'})}{P(W_i \mid Z = \text{'Con'})} \times \frac{P(Z = \text{'Lib'})}{P(Z = \text{'Con'})}$$

Logging this probability ratio, every new word *adds* a bit of information that pushes the ratio above or below 0

Naive Bayes

Example: Naive Bayes with only word class 'discriminat*'.

$$P(W = \text{'discriminat*'} \mid Z = \text{'Lib'}) = (26 + 13)/(20002 + 18722) \approx 0.001$$

$$P(W = \text{'discriminat*'} \mid Z = \text{'Con'}) = (70 + 48)/(17368 + 17698) \approx 0.003$$

Assume that liberal and conservative supporting briefs are equally likely (true in the training set)

$$\frac{P(Z = \text{'Lib'})}{P(Z = \text{'Con'})} = 1$$

Last step: calculate posterior classification probabilities for a new document (based on occurrence of this word).

Naive Bayes

Amicus brief from 'King County Bar Association' containing 3667 words and 4 matches to discriminat*.

that "the state shall not [discriminate] against, or grant preferential treatment the lingering effects of racial [discrimination] against minority groups in this remedy the effects of societal [discrimination]. Another four Justices (Stevens that "the state shall not [discriminate] against, or grant preferential treatment

Naive Bayes

A priori, the probabilities are...

Probability that we observe the word discriminat* 4 out of 3667 times if the document is liberal:

```
> dbinom(4, size=3667, prob=0.001007127)
[1] 0.1930602
```

Probability that we observe the word discriminat* 4 out of 3667 times if the document is conservative:

```
> dbinom(4, size=3667, prob=0.003365083)
[1] 0.004188261
```

Logged probability ratio = 3.83

Naive Bayes

Conclusion: Seeing 4 instances of discriminat* gives the posterior classification probabilities

$$\theta_{\text{liberal}} = \frac{0.193}{0.193+0.004} = 0.979$$

$$\theta_{\text{conservative}} = 1-0.979=0.021$$

This is *quite* confident

...but other words will be less loaded or push the other way

Evaluating Classifiers: Accuracy, Precision and Recall

		Machine	
		Liberal	Conservative
Training Data	Liberal	40	10
	Conservative	40	60

$$\text{Accuracy} = (40+60)/150 = .66$$

$$\text{Precision } (Z_{\text{machine}}=\text{Liberal}) = 40/(40+40) = 0.5$$

$$\text{Precision } (Z_{\text{machine}}=\text{Cons}) = 60/(10+60) = 0.86$$

$$\text{Recall } (Z_{\text{training}}=\text{Liberal}) = 40/(40+10) = 0.80$$

$$\text{Recall } (Z_{\text{training}}=\text{Cons}) = 60/(40+60) = 0.60$$

Latent Dirichlet Allocation

What if documents don't belong to only one category?

- Each document is a mixture over topics

- Each topic is a mixture over words

Using LDA gives us:

- Distribution of words for each topic (β)

- Proportion of a document in each topic (θ)

Setting up LDA

- Still uses bag of words

- Number of topics fixed ex ante

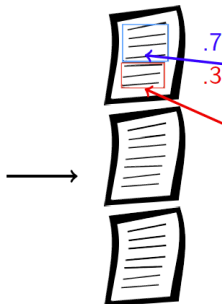
- Topics examined ex post – unsupervised learning

Visualizing Topic Models

Say you have
a lot of people.



Each writes
some texts



that discuss a few
different topics

Politics

congress, nations,
power, votes, agree-
ment, bargaining

Statistics

estimator, data, anal-
ysis, variance, model,
inference

The Latent Dirichlet Allocation estimates:

- ① The topics- each is a distribution over words
- ② The proportion of each document in each topic

Courtesy of Brandon Stewart

Latent Dirichlet Allocation Components

Given dimensions:

- **N** documents, **J** different topics, **K** unique words

We want the following matrices:

- $X = N \times K$ document-term matrix (observed)
- $\theta = N \times K$ matrix with row $\theta_i = (\theta_{i1}, \dots, \theta_{iK})$
 - θ_{ik} : Proportion of document i allocated to topic k
- $\beta = K \times J$ matrix with row $\beta_k = (\beta_{k1}, \dots, \beta_{kJ})$
 - β_{kj} : Probability of using word J , if topic k is chosen
- $\alpha = K$ length population prior for θ

Objective function: $f(X, \beta, \theta, \alpha)$

LDA Generative Model

When writing a word (m) for document i :

$$p(\theta_i | \alpha) = \text{Dirichlet}(\alpha) \quad (\text{Pick potential topics } \theta_i)$$

$$p(z_{im} | \theta_i) = \text{Multinomial}(1, \theta_i) \quad (\text{Pick the topic } z \text{ to discuss})$$

$$x_{im} | \beta_k, z_{im} = \text{Multinomial}(1, \beta_k) \quad (\text{Pick a word from topic } z_{im})$$

$$p(\beta_k) = \text{Dirichlet}(\eta) \quad (\text{Prior on topics})$$

with

$$p(\beta, \theta, Z, \alpha | X) \propto p(\beta | \eta) \cdot p(\theta | \alpha) \cdot p(Z | \theta) \cdot p(X | \beta, Z)$$

Intuition on LDA

$$p(\beta, \theta, Z, \alpha | X) \propto p(\beta | \eta) \cdot p(\theta | \alpha) \cdot p(Z | \theta) \cdot p(X | \beta, Z)$$

Remember rows of β and θ must sum to 1!

$p(X | \beta, Z)$: Favors co-occurring words, segregated topics

- $\beta = (0.5, 0.5, 0, 0)$ vs $\beta = (0.25, 0.25, 0.25, 0.25)$

$p(Z | \theta)$: Higher if θ is concentrated (i.e. fewer potential topics)

- $\theta = (0.5, 0.5, 0, 0)$ vs. $\theta = (0.25, 0.25, 0.25, 0.25)$

Joint distribution favors sparse topics, small topic clusters

But data often need a larger number of topics to assign small topic clusters to data

Japanese Campaign Manifestos (Catalinac 2016)

Examines the effect of electoral reform

Before 1994: SNTV-MMD (more intraparty competition)

After 1994: Mixed member majoritarian (PR + SMD)

LDA on N=7497 candidate manifestos for Diet 1950-2009

Hand transcribed, OCR failed

More complicated in Japanese

Characterizing campaigns across 50+ years

What do candidates talk about?

How did electoral reform change incentives?

Why increasing interest in militaristic foreign action?

Japanese Manifesto Topics

Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6
改革	年金	推進	区	政治	日本
郵政	円	郵政	政策	改革	国
民営	廃止	国	地域	国民	外交
小泉	改革	つとめる	まち	企業	国家
構造	光	社会	鹿児島	自民党	社会
政府	実現	対策	全力	日本	国民
官	無駄	振興	選挙	共産党	保障
推進	日本	充実	国政	郵政	安全
民	増税	促進	作り	金権	地域
自民党	削減	安定	塊	党	拉致
日本	一元化	確立	対策	選挙	経済
制度	設備	企業	中小	禁止	守る
民間	子供	実現	発電	憲法	問題
年金	地域	中小	推進	腐敗	北朝鮮
実現	ひと	育成	エネルギー	団体	教育
進める	サラリーマン	制度	企業	区	責任
断行	制度	政治	声	ソ連	力
地方	議員	地域	実現	守る	創る
止める	区	福祉	活性	平和	安心
保障	民主党	自民党	円	反	日露
財政	年寄	改革	地方	反	防
作る	一掃	確保	尽くす	真	憲法
策定	郵政	強化	商店	是正	可能
社会	道路	教育	いかす	一掃	道
国民	交代	施設	全国	憲法	未来
公務員	社会保障庁	生活	政党	抜本	ひと
力	月額	支援	ひと	定数	再生
経済	手前	環境	支援	政治	将来
国	談合	発展	経済	金丸	解決
安心	支援	産業	福祉	改革	基本
Postal privatization	Reducing Wasteful Public Spending	Pork for the District	Policies for the district	Political Reform	National Security Policy

Manifesto Content over Time

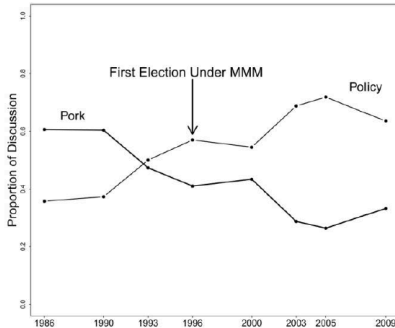


Figure 1. LDP candidates switched to more policy and less pork in the 1993 election and continued with this strategy under MMM. This figure plots the mean proportions of discussion devoted to pork and policy, respectively, in the 2,355 manifestos produced by LDP candidates in these eight elections.

Foreign Policy Content

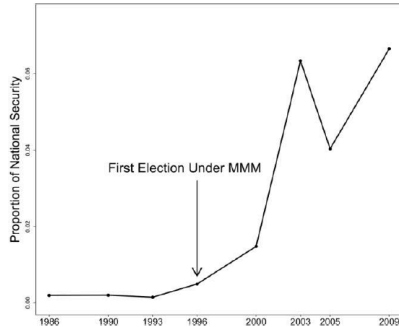


Figure 2. LDP candidates increased their discussion of national security in the first election under MMM. This figure plots the mean proportion of discussion comprised of national security in the 2,355 manifestos produced by LDP candidates in these eight elections.

Wrapping up

For single category classification

Discriminative: Model $P(Z | \{W\}, \beta) = \theta_{z|w}$ directly

Generative: Model $P(\{W\} | Z, \beta)$ and $P(Z) = \theta_z$, then get $P(Z | \{W\}, \beta) = \theta_{z|w}$ via Bayes theorem

For multiple categories

Latent Dirichlet Allocation

More generally used for discovery of topics

Tomorrow's Lab

Classifying NY Times articles by topic