Using Text as Data

# Which pages are those of job ads?



Display Ad 133 -- No Title
*Boston Globe (1960-1985);* Nov 4, 1979; ProQuest Historical Newspapers: The Boston Globe
pg. E51

# Which pages are those of job ads?

# Topic classification

- There are a number of possible types of pages of advertisements: for jobs, for real estate, for cars, in retail, etc...
- Within each type of advertisements certain types of words are more likely to appear:
  - employment, responsibility, manager in job ads
  - foot, loft, architectural in real estate ads
- If we knew the predictive words for each type of ad, we could estimate the likelihood for each topic for each page.
- Goal: Figure out the "predictive words" of each type of ad.
- Idea: Words come from different *topics*. Within a given page of ads, certain topics are more or less represented.

# Latent Dirichlet Allocation

Let

- $V$ = size of vocabulary; $K$ = number of topics
- $\alpha$ ($K \times 1$) and $\beta$ ($K \times V$) parameters to be estimated
- $w$ (a word) is a vector of length $V$. $w_i = 1$ for one $i$ and 0 elsewhere.
- $d_m$ is a particular page of ads.

# Latent Dirichlet Allocation

Each page of ads comes from a mixture of topics:

- For page $m$, random variables $\theta_{mk}$ are drawn from a (Dirichlet-$\alpha$) distribution.
- For each word within a page, the probability that a word is from topic $z = z_k$ is $\theta_{mk}$
- Conditional on the topic, $\beta_{ki} = \Pr(w_i = 1 | z = z_k)$
- We can write the likelihood of each doc:

$$\Pr(d_m | \alpha, \beta) = \int p(\theta | \alpha) \cdot \times$$
$$\left( \prod_{n \in \text{Words in } d_m} \sum_{k \in \text{Topics}} \Pr(z_k | \theta) \Pr(w_n | z_k, \beta) \right) d\theta$$

# Latent Dirichlet Allocation

- Need to reduce the dimensionality of the maximization problem.
    - Stem
    - Drop very frequent and infrequent words.

Results for $K = 3$ (highest $\beta$ word stems for each topic):

- Topic 1: new, home, owner, acr, call, car, hous, den, area, ask
- Topic 2: resum, seek, call must, work, exp, excel, new, salari, send.
- Topic 3: build, ave, new, park, call, studio, east, avail, fee, firm

# Similarity

- Do the words/phrases lpn, rn, registered nurse, nurse represent the same occupation?
  - What about finance professor and econometrician?
  - What about finance professor and nurse?
- Do the words/phrases delegating and directing refer to the same work activity?
  - What about monitor and motivating?
  - What about scheduling and coordinating?

Idea: Similar words appear in similar contexts:

- "Limits how many patients can be assigned to each *registered nurse* in Massachusetts hospitals and certain other healthcare facilities."
- "Newton-Wellesley Hospital *registered nurse* Betty Sparks is a member of the Committee to Ensure Safe Patient Care, which advocates for increased staffing levels."
- "With a busy hospital job as a *licensed practical nurse*, it was hard not to miss some of those frequent injections."

# Continuous Bag of Words

- Suppose we see a chunk of text that is of the form " each registered X Massachusetts hospitals" What is the probability that...
  - X="nurse"
  - X="econometrician"
  - X="programmer" ?

- Would like representation of words that assigns a high probability to X=nurse.

- Model to estimate $\mathcal{C}$, $\tilde{\mathcal{C}}$ each which is of dimension $V$ (number of words) $\times$ $n$ (number of "features")

- As before: $w$ (a word) is a vector of length $V$. $w_i = 1$ for one $i$ and 0 elsewhere.

# Continuous Bag of Words

1. Take word vectors for each of the words in the context of $\left(\text{e.g., } w_{\text{nurse}}, w_{\text{registered}}, w_{\text{Massachusetts}}, w_{\text{hospitals}}\right)$
2. Recover the embedded word vectors for a given context: $\tilde{\mathcal{C}}' w_{\text{nurse}}, \tilde{\mathcal{C}}' w_{\text{registered}}, \tilde{\mathcal{C}}' w_{\text{Massachusetts}}, \tilde{\mathcal{C}}' w_{\text{hospitals}}.$
3. Average these to get a measure of the context: $\tilde{v} = \frac{\tilde{\mathcal{C}}' w_{\text{nurse}} + \tilde{\mathcal{C}}' w_{\text{registered}} + \tilde{\mathcal{C}}' w_{\text{Massachusetts}} + \mathcal{C}' w_{\text{hospitals}} \cdot}{4}.$
4. Generate a probability of this context: $p_c = \frac{C\tilde{v}}{\sum C\tilde{v}}$

Log likelihood of a data set is given by comparing the actual words in the "X" position to what would be predicted out of step 4

$$\log \mathcal{L} = \sum_c y_c \cdot \log p_c$$

# Examples Elsewhere

- Gentzkow and Shapiro (2010): Measuring media slant.
  - Link Congressional record (words/phrases used by Dem./Repub. congresspeople) to newspapers.
    - Democrats mention "estate tax," while Republicans mention "death tax"
    - Washington Post: 13.7 ratio mentions of estate tax to death tax.
    - Washington Times: 1.3 ratio mentions of estate tax to death tax.
- Baker, Bloom, Davis (2016): Measuring policy uncertainty.
- In the optional papers.
  - Hoberg and Philips: SEC Filings to learn about firms' products/competitors
  - Mann and Puttmann: Google text of all US patents to learn about new technologies which lead to automation.
- For other examples, see Gentzkow, Kelly, Taddy (2017, JEL)