

Neural Models for Documents with Metadata

Dallas Card, Chenhao Tan, Noah A. Smith

July 18, 2018

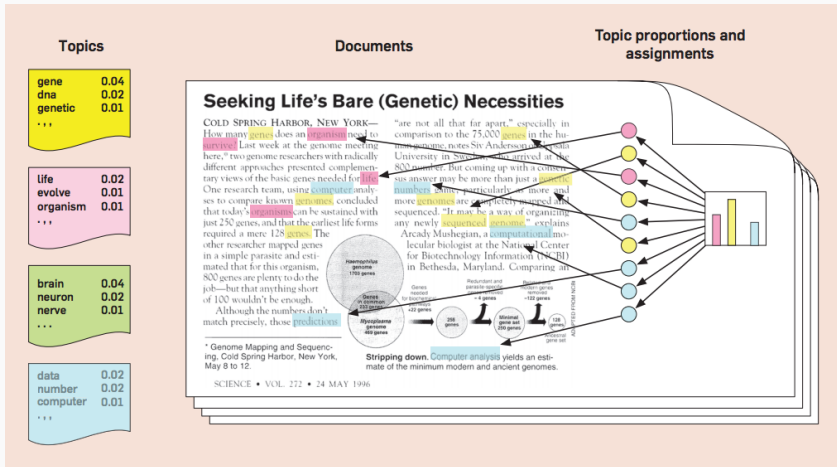


Main points of this talk:

1. Introducing *Scholar*¹: a neural model for documents with metadata
 - Background (LDA, SAGE, SLDA, etc.)
 - Model and related work
 - Experiments and Results
2. Power of neural variational inference for interactive modeling

¹Sparse Contextual Hidden and Observed Language Autoencoder

Latent Dirichlet Allocation



Blei, Ng, and Jordan. *Latent Dirichlet Allocation*. JMLR. 2003.

David Blei. *Probabilistic topic models*. Comm. ACM. 2012

Types of metadata

- Date or time
- Author(s)
- Rating
- Sentiment
- Ideology
- etc.

Variations and extensions

- Author topic model (Rosen-Zvi et al 2004)
- Supervised LDA (SLDA; McAuliffe and Blei, 2008)
- Dirichlet multinomial regression (Mimno and McCallum, 2008)
- Sparse additive generative models (SAGE; Eisenstein et al, 2011)
- Structural topic model (Roberts et al, 2014)
- ...

Desired features of model

- Fast, scalable inference.
- Easy modification by end-users.

Desired features of model

- Fast, scalable inference.
- Easy modification by end-users.
- Incorporation of metadata:
 - Covariates: features which influences text (as in SAGE).
 - Labels: features to be predicted along with text (as in SLDA).

Desired features of model

- Fast, scalable inference.
- Easy modification by end-users.
- Incorporation of metadata:
 - Covariates: features which influences text (as in SAGE).
 - Labels: features to be predicted along with text (as in SLDA).
- Possibility of sparse topics.

Desired features of model

- Fast, scalable inference.
- Easy modification by end-users.
- Incorporation of metadata:
 - Covariates: features which influences text (as in SAGE).
 - Labels: features to be predicted along with text (as in SLDA).
- Possibility of sparse topics.
- Incorporate additional prior knowledge.

Desired features of model

- Fast, scalable inference.
 - Easy modification by end-users.
 - Incorporation of metadata:
 - Covariates: features which influences text (as in SAGE).
 - Labels: features to be predicted along with text (as in SLDA).
 - Possibility of sparse topics.
 - Incorporate additional prior knowledge.
- Use variational autoencoder (VAE) style of inference (Kingma and Welling, 2014)

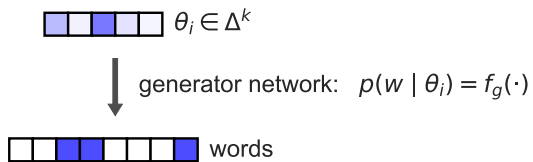
- Coherent groupings of words (something like topics), with offsets for observed metadata

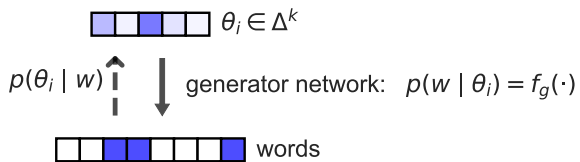
Desired outcome

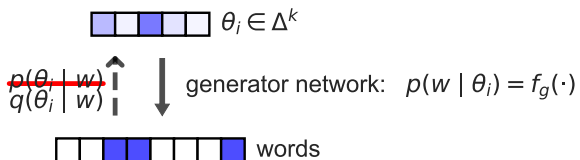
- Coherent groupings of words (something like topics), with offsets for observed metadata
- Encoder to map from documents to latent representations

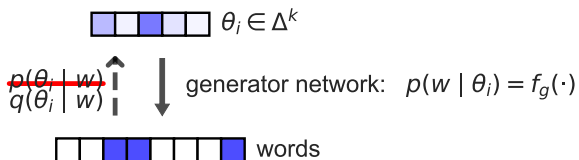
Desired outcome

- Coherent groupings of words (something like topics), with offsets for observed metadata
- Encoder to map from documents to latent representations
- Classifier to predict labels from from latent representation

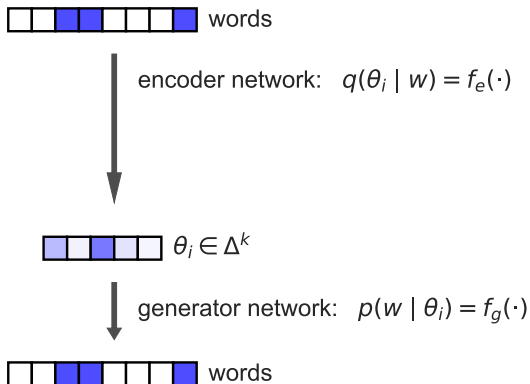




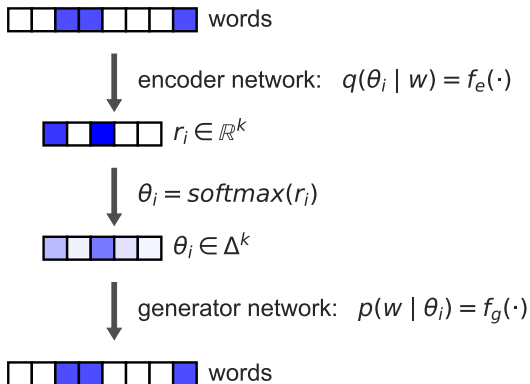




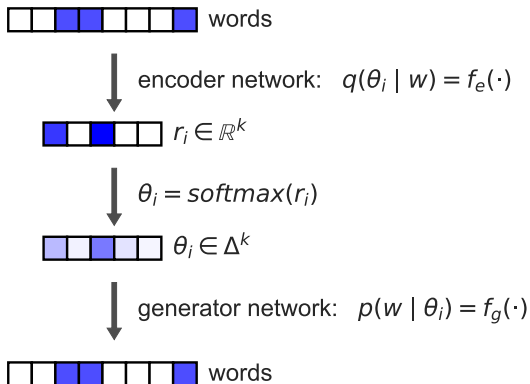
$$\text{ELBO} = \mathbb{E}_q[\log p(\text{words} | \theta_i)] - D_{\text{KL}}[q(\theta_i | \text{words}) || p(\theta_i)]$$



$$\text{ELBO} = \mathbb{E}_q[\log p(\text{words} | \theta_i)] - D_{\text{KL}}[q(\theta_i | \text{words}) || p(\theta_i)]$$

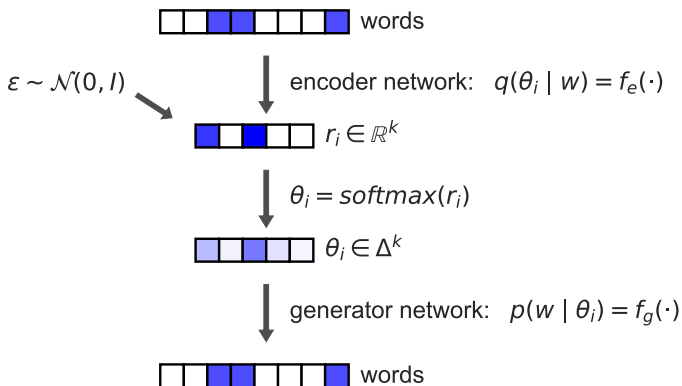


$$\text{ELBO} = \mathbb{E}_q[\log p(\text{words} | r_i)] - D_{\text{KL}}[q(r_i | \text{words}) || p(r_i)]$$

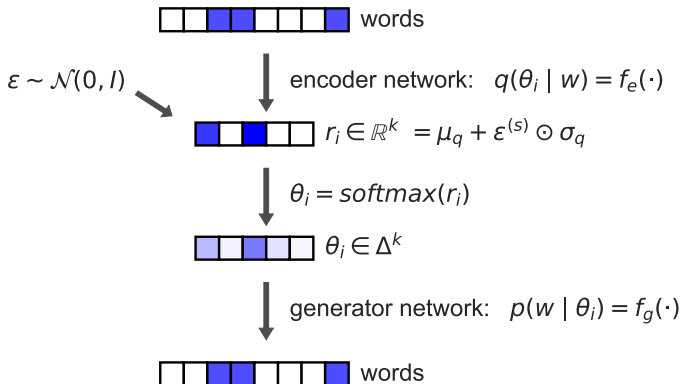


$$\text{ELBO} \approx \frac{1}{S} \sum_{s=1}^S [\log p(\text{words} | r_i^{(s)})] - D_{\text{KL}}[q(r_i | \text{words}) || p(r_i)]$$

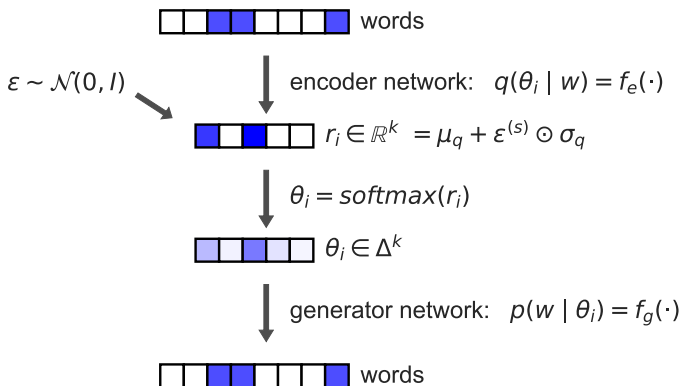
Model



$$\text{ELBO} \approx \frac{1}{S} \sum_{s=1}^S [\log p(\text{words} | r_i^{(s)})] - D_{\text{KL}}[q(r_i | \text{words}) || p(r_i)]$$

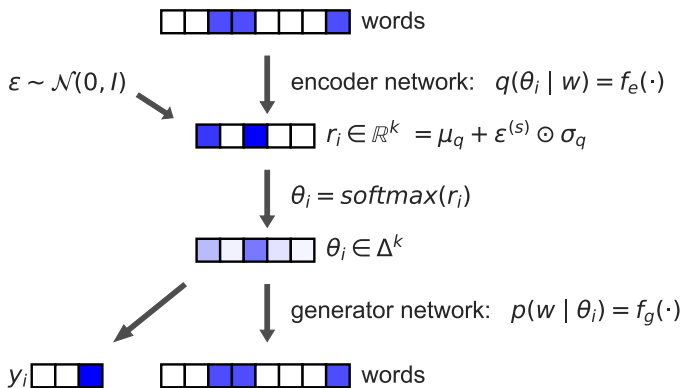


$$\text{ELBO} \approx \frac{1}{S} \sum_{s=1}^S [\log p(\text{words} | r_i^{(s)})] - D_{\text{KL}}[q(r_i | \text{words}) || p(r_i)]$$

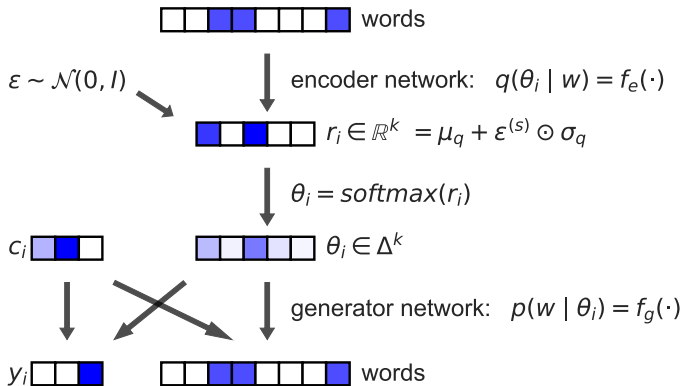


Srivastava and Sutton, 2017, Miao et al, 2016

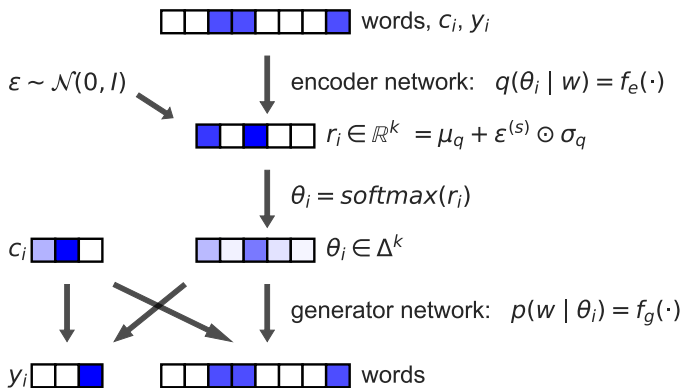
Model



Model



Model



Generator network:

- $p(\text{word} \mid \theta_i, c_i) = \text{softmax}(d + \theta_i^T B^{(\text{topic})} + c_i^T B^{(\text{cov})})$

Generator network:

- $p(\text{word} \mid \theta_i, c_i) = \text{softmax}(d + \theta_i^T B^{(\text{topic})} + c_i^T B^{(\text{cov})})$
- Optionally include interactions between topics and covariates

Generator network:

- $p(\text{word} \mid \theta_i, c_i) = \text{softmax}(d + \theta_i^T B^{(\text{topic})} + c_i^T B^{(\text{cov})})$
- Optionally include interactions between topics and covariates
- $p(y_i \mid \theta_i, c_i) = f_y(\theta_i, c_i)$

Generator network:

- $p(\text{word} \mid \theta_i, c_i) = \text{softmax}(d + \theta_i^T B^{(\text{topic})} + c_i^T B^{(\text{cov})})$
- Optionally include interactions between topics and covariates
- $p(y_i \mid \theta_i, c_i) = f_y(\theta_i, c_i)$

Encoder:

- $\mu_i = f_\mu(\text{words}, c_i, y_i)$
- $\log \sigma_i = f_\sigma(\text{words}, c_i, y_i)$
- Optional incorporation of word vectors to embed input

- Stochastic optimization using mini-batches of documents
- Tricks from Srivastava and Sutton, 2017:
 - Adam optimizer with high-learning rate to bypass mode collapse
 - Batch-norm layers to avoid divergence
- Annealing away from batch-norm output to keep results interpretable

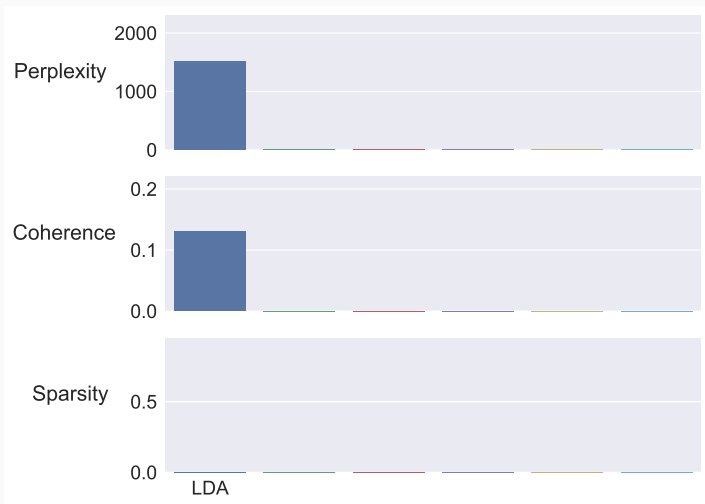
- $B^{(topic)}, B^{(cov)}$: Coherent groupings of positive and negative deviations from background (\sim topics)

- $B^{(topic)}, B^{(cov)}$: Coherent groupings of positive and negative deviations from background (\sim topics)
- f_{μ}, f_{σ} : Encoder network: mapping from words to topics:
 $\hat{\theta}_i = \text{softmax}(f_e(\text{words}, c_i, y_i, \epsilon))$

- $B^{(topic)}, B^{(cov)}$: Coherent groupings of positive and negative deviations from background (\sim topics)
- f_{μ}, f_{σ} : Encoder network: mapping from words to topics:
 $\hat{\theta}_i = \text{softmax}(f_e(\text{words}, c_i, y_i, \epsilon))$
- f_y : Classifier mapping from $\hat{\theta}_i$ to labels: $\hat{y} = f_y(\theta_i, c_i)$

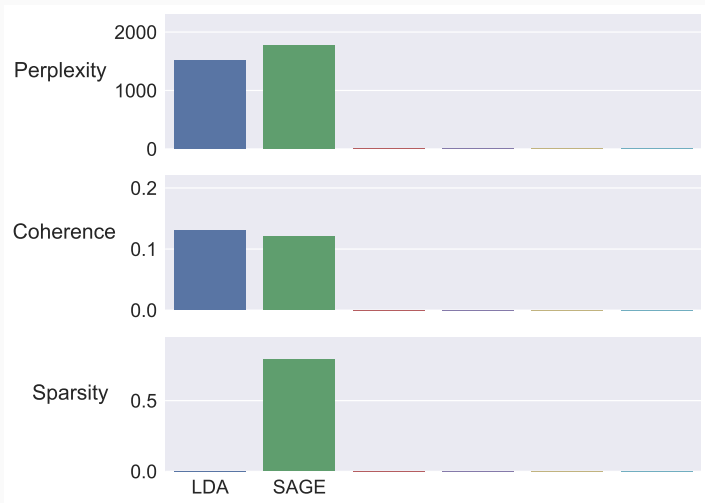
1. Performance as a topic model, without metadata (perplexity, coherence)
2. Performance as a classifier, compared to SLDA
3. Exploratory data analysis

Quantitative results: basic model



IMDB dataset (Maas, 2011)

Quantitative results: basic model



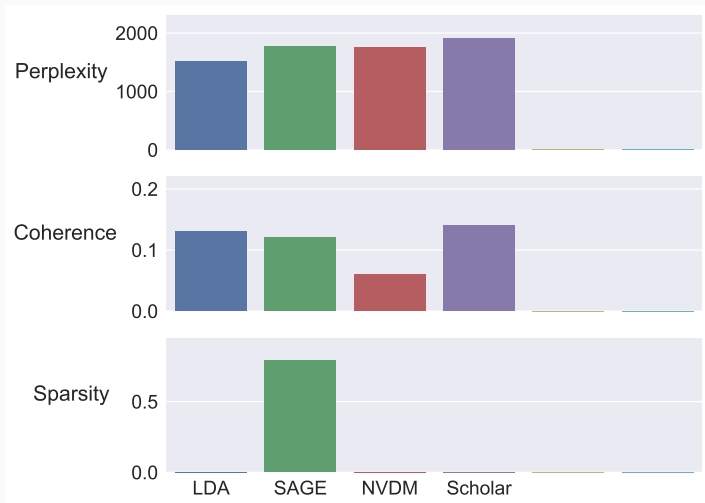
IMDB dataset (Maas, 2011)

Quantitative results: basic model



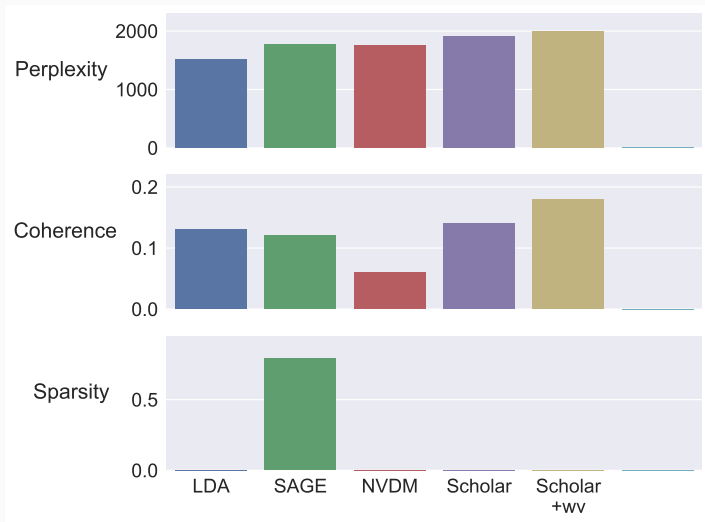
IMDB dataset (Maas, 2011)

Quantitative results: basic model



IMDB dataset (Maas, 2011)

Quantitative results: basic model



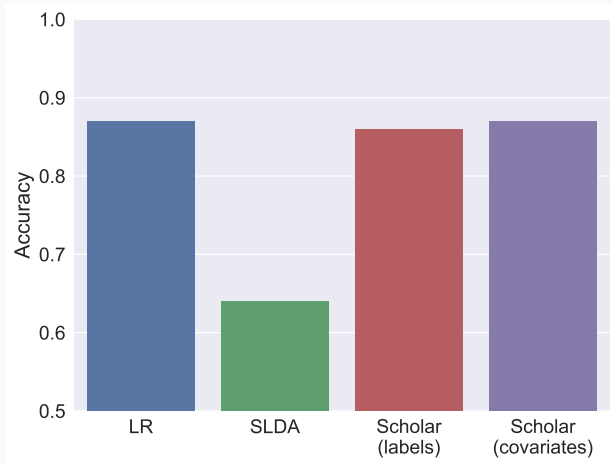
IMDB dataset (Maas, 2011)

Quantitative results: basic model



IMDB dataset (Maas, 2011)

Classification results

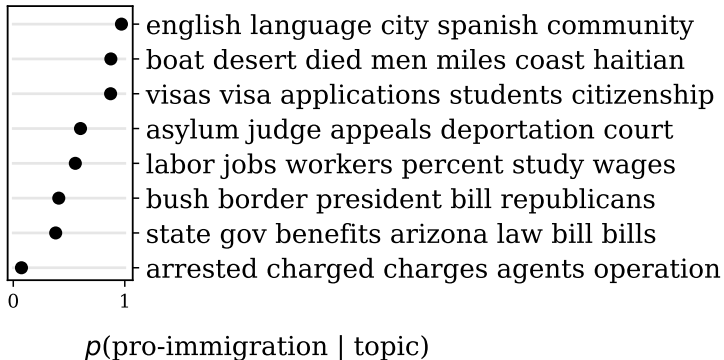


IMDB dataset (Maas, 2011)

Data: Media Frames Corpus (Card et al, 2015)

- Collection of thousands of news articles annotated in terms of tone and framing
- Relevant metadata: year of publication, newspaper, etc.

Tone as a label



Tone as a covariate, with interactions

Base topics

ice customs agency
population born percent
judge case court guilty
patrol border miles
licenses drivers card
island story chinese
guest worker workers
benefits bill welfare

Anti-immigration

criminal customs
jobs million **illegals**
guilty charges man
patrol border
foreign sept visas
smuggling federal
bill border house
republican california

Pro-immigration

detainees detention
english **newcomers**
asylum court judge
died authorities desert
green citizenship card
island school ellis
workers tech skilled
law welfare students

Conclusions

- Variational autoencoders (VAEs) provide a powerful framework for latent variable modeling
- We use the VAE framework to create a customizable model for documents with metadata
- We obtain comparable performance with enhanced flexibility and scalability
- Code is available: www.github.com/dallascard/scholar