

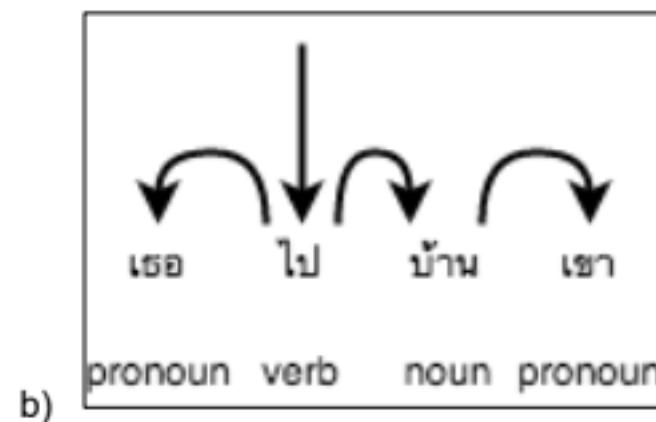
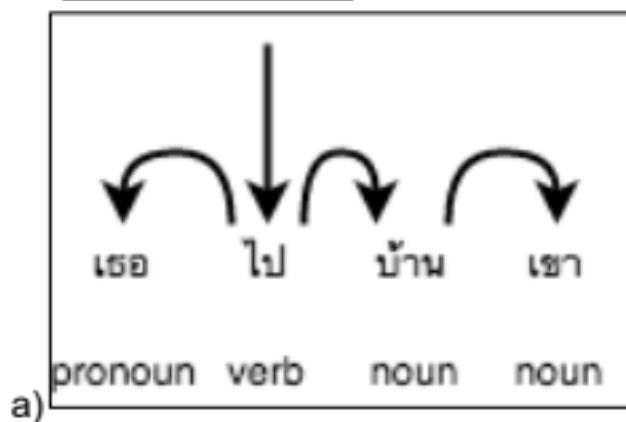
TEXT CLASSIFICATION

Intent, topic, sentiment, etc.

Midterm

- Ambiguity

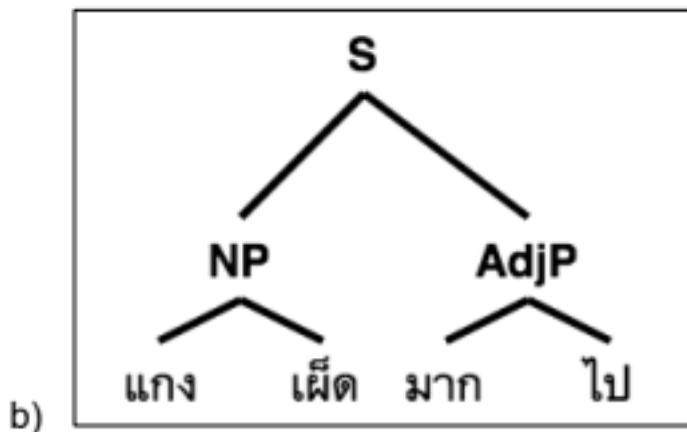
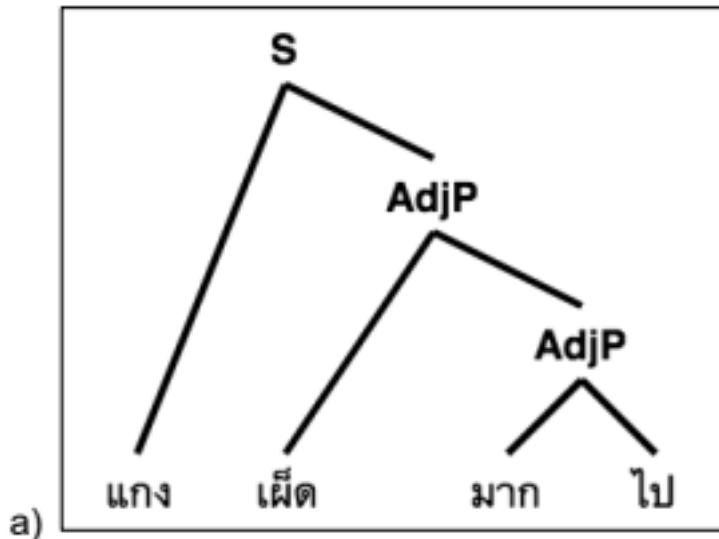
15. ເຮືອ ໄປ ບ້ານ ແຂງ



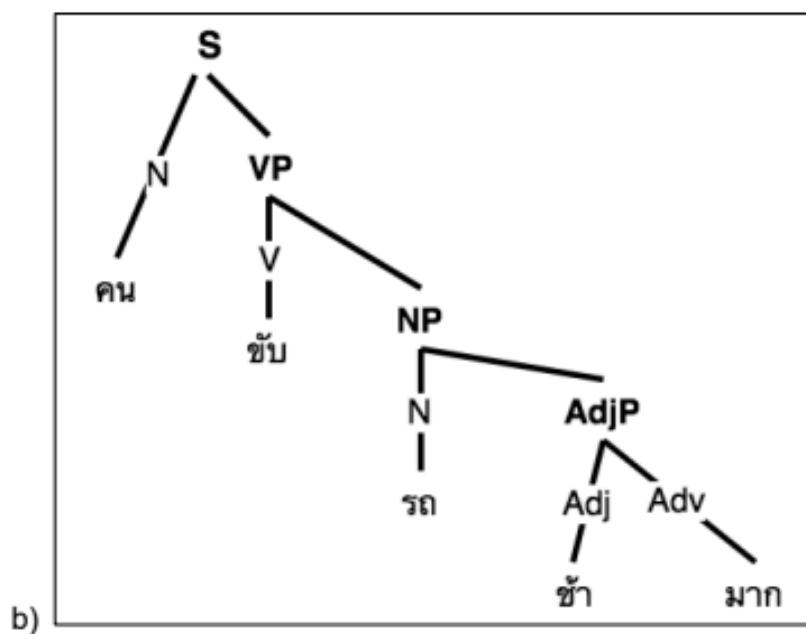
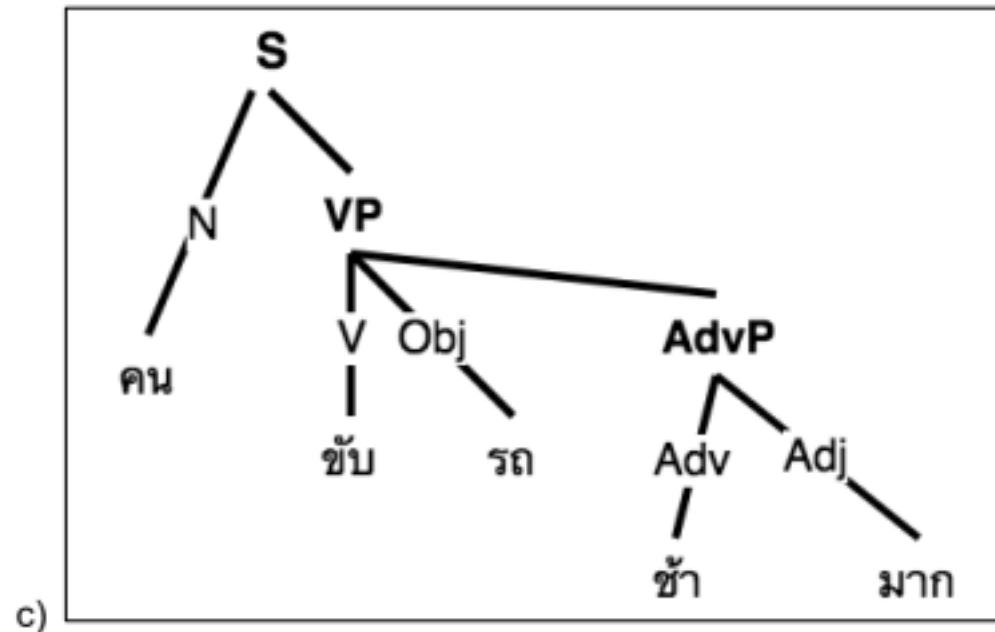
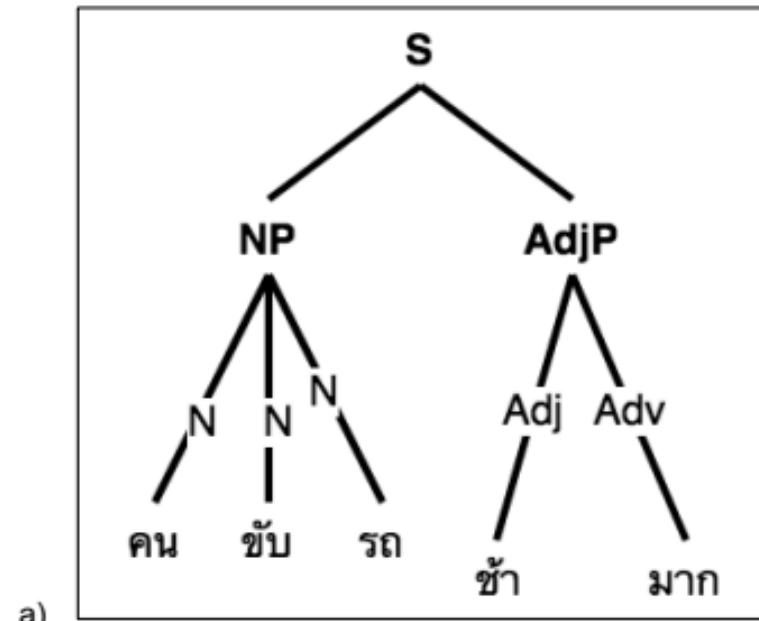
Midterm

- Ambiguity

16. แกง เพ็ด มาก ไป



17. คน ขับ รถ ช้า มาก



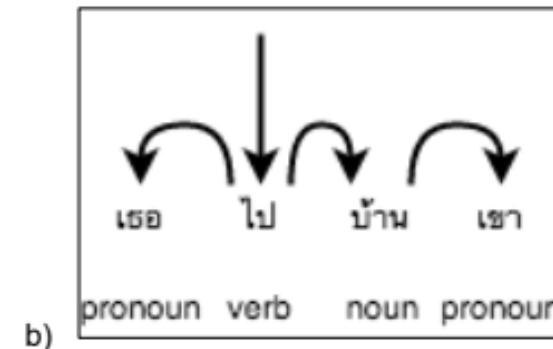
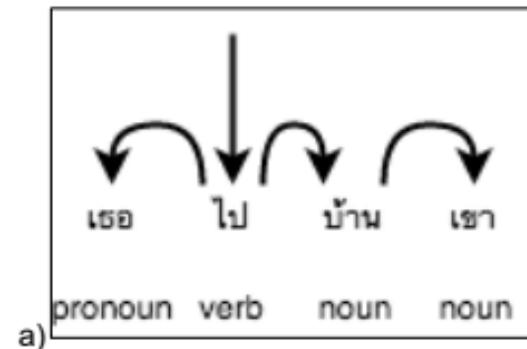
Error propagation and NLP systems

- Typical flow of NLP systems



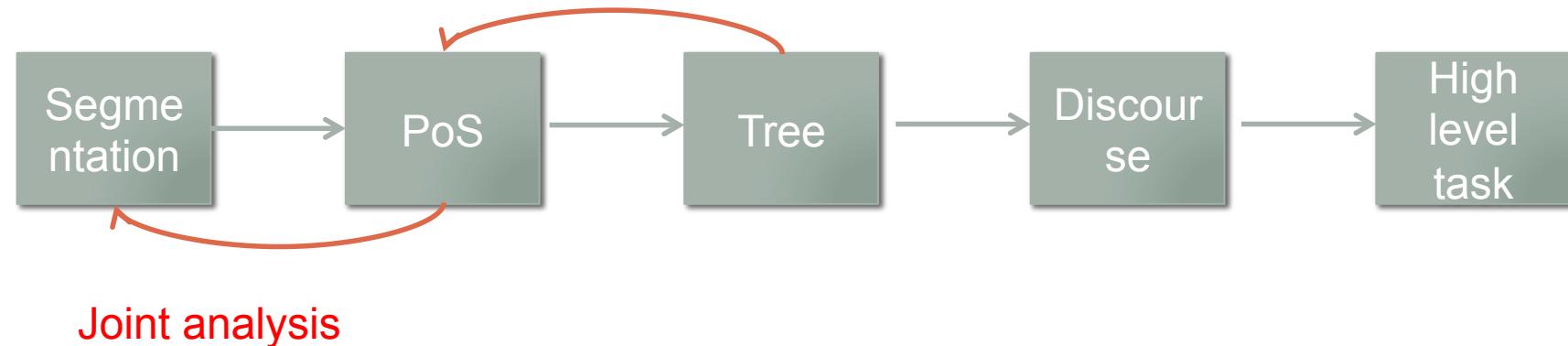
Pipeline analysis

15. เธอไปบ้านเขา



Error propagation and NLP systems

- Use higher level info to help lower level
 - Iterative process



Error propagation and NLP systems

- Do the task together



Joint analysis

Error propagation and NLP systems

- One single task
 - Input closest to the input form – characters, words

High level task with low level structures as input

End-to-end analysis

Dropped words recovery

- Some sentence are ambiguous because of possibility of dropped words

เธอไปบ้าน [ของ] เข้า
แทง [ที่] เพ็ดมากไป
คนขับรถ [มา] ช้ำมาก

- One solution (Pipeline processing)
 - Separate into two parts:
 - Dropped word detection – determine the location
 - Dropped word prediction – determine what is missing

Dropped word detection

- เธอ ไป บ้าน เข้า
- 0 0 1 0
- Similar to word segmentation. A detection task.
 - Some recurrent network should be good for this
 - Things to consider
 - Are the words already segmented in the corpus?
 - Words representation?
 - Outside corpora?
 - Use PoS tags?
 - Does the training sentence need PoS tags?
 - PoS embeddings?
- Special mention: use tree structure to calculate the probability. If there are two trees that have similar probability there should be dropped words

Dropped word prediction

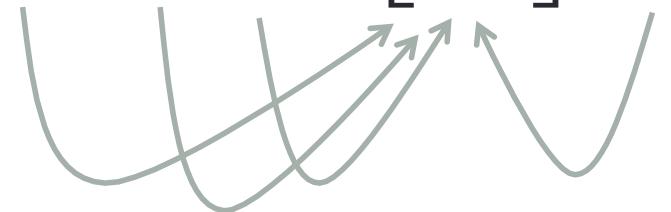
- ເຮືອ ໄປ ບ້ານ [ຂອງ] ເຂົາ



- Context predict the middle word -> CBOW structure
- Or previous words predict next word -> n-gram
 - CBOW should be more powerful (bigger context), but both answers are accepted

End-to-end analysis

- ເຮືອ ໄປ ບ້ານ [-X-] ເຂົາ



- X can be a proper word or blank
 - Vocabulary = N+1
 - Skip the detection task

Assumptions

- This only handle one missing word
 - How to handle multiple missing words in a row?
 - Predict phrases?
- Discourse analysis
 - Previous solution does not allow discourse
 - Can expand to discourse by expanding the context to be bigger than a sentence
- Types of dropped words considered

Corpora

- What kind of text?
 - Pantip, twitter, wikipedia, textbook, ...
- Extra: How to actually construct the training corpora?
 - Painfully look at the sentences and determine dropped words
 - or
 - Look at standard sentences and determine which words can be dropped

Wongnai Challenge



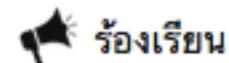
- Predict star rating from review text

output



ก๋วยเตี๋ยวอร่อย ราคาถูก นั่งทานในห้องแอร์

เมนูเด็ด: บะหมี่ต้มยำหมูแดง



input

ส่วนตัวชอบก๋วยเตี๋ยวต้มยำมากค่ะ รสชาติอร่อยไม่ต้องป魯งเลย แต่ชุปเปอร์ตินไก่ใส่ถั่วงอกมาด้วย เหมือนเป็น
เกาเหลาตินไก่มากกว่าชุปเปอร์ตินไก่ รสชาติก็ยังไม่กลมกล่อมเท่าก๋วยเตี๋ยวต้มยำ ตินไก่เปื่อยดี ทานง่าย

#	△priv	Team Name	Kernel	Team Members	Score ⓘ	Entries	Last
1	—	blead			0.59139	30	9h
2	—	Danupat Kns			0.58333	21	6h
3	▼ 1	Phizaz			0.58172	4	9d
4	▲ 1	Attawit Chaiyaroj			0.57848	14	3h
5	▼ 4	Chal			0.57633	42	1h
6	▼ 5	James R.			0.57419	32	13m
7	▼ 6	Khajornal Anantanit			0.57258	60	7h
8	▲ 1	burin			0.57258	44	14m
9	▼ 1	Nattasit			0.56774	4	5h
10	▲ 5	Wannita			0.56666	18	6h
1	—	blead			0.58138	30	9h
2	—	Danupat Kns			0.56919	21	6h
3	▲ 1	Attawit Chaiyaroj			0.56297	14	3h
4	▼ 1	Phizaz			0.56182	4	9d
5	▲ 5	Wannita			0.55789	18	6h
6	▲ 6	ChanatipSaetia			0.55306	5	1h
7	▲ 1	burin			0.55284	44	14m
8	▲ 6	AtomicBoyZ			0.55192	11	17h
9	▼ 4	Chal			0.55192	42	1h
10	▼ 1	Nattasit			0.55123	4	5h

Public

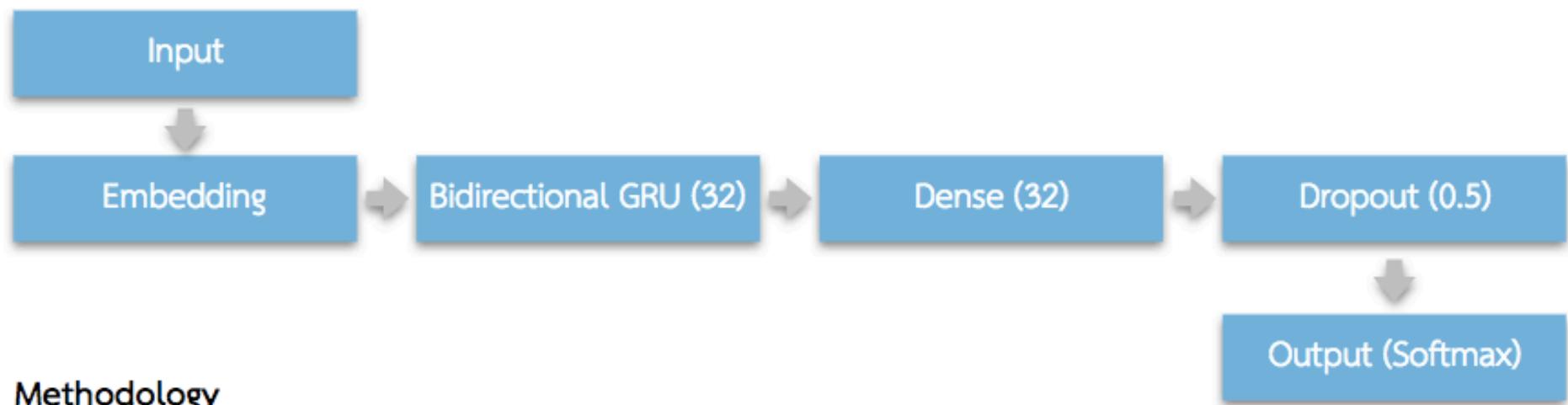
Private

Cute ideas

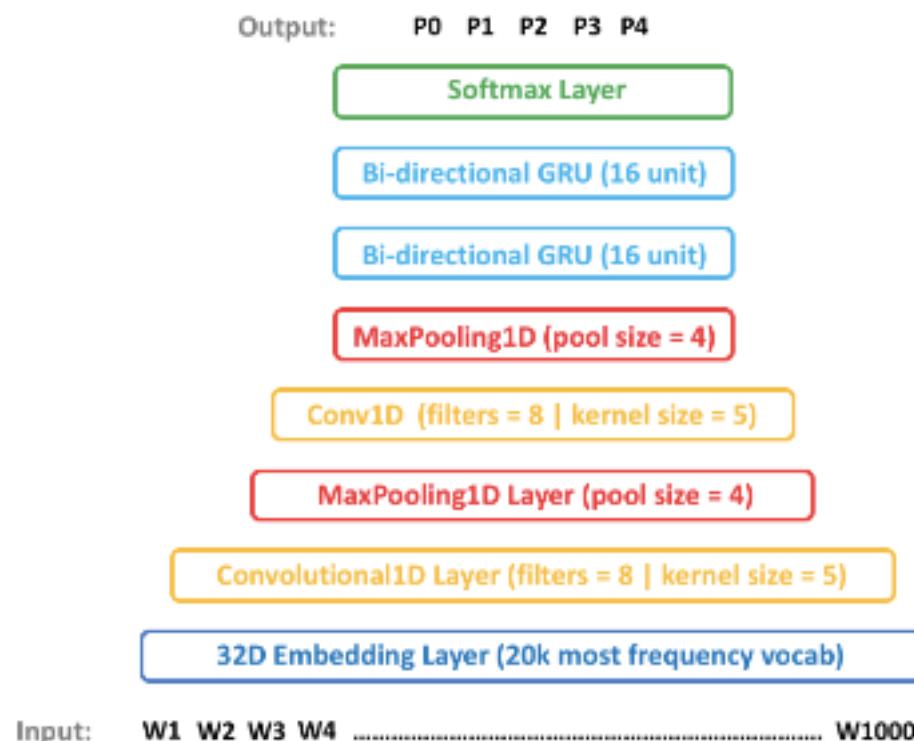
- Post-processing
 - ໃໝ່ດາວ
- Fasttext just got better
 - Trained on webcrawl instead of wiki
 - Tokenized using ICU
 - Character n-gram model to handle OOV
- Thai2vec
 - trained on thai wiki <https://github.com/cstorm125/thai2vec>
 - thaiNLP things <https://github.com/PyThaiNLP/>
- One person tried both with no significant difference (same amount of OOV)

3rd best model

- 8 words threshold for unk
- Pre-train weights from homework but adapted to task



2nd best model



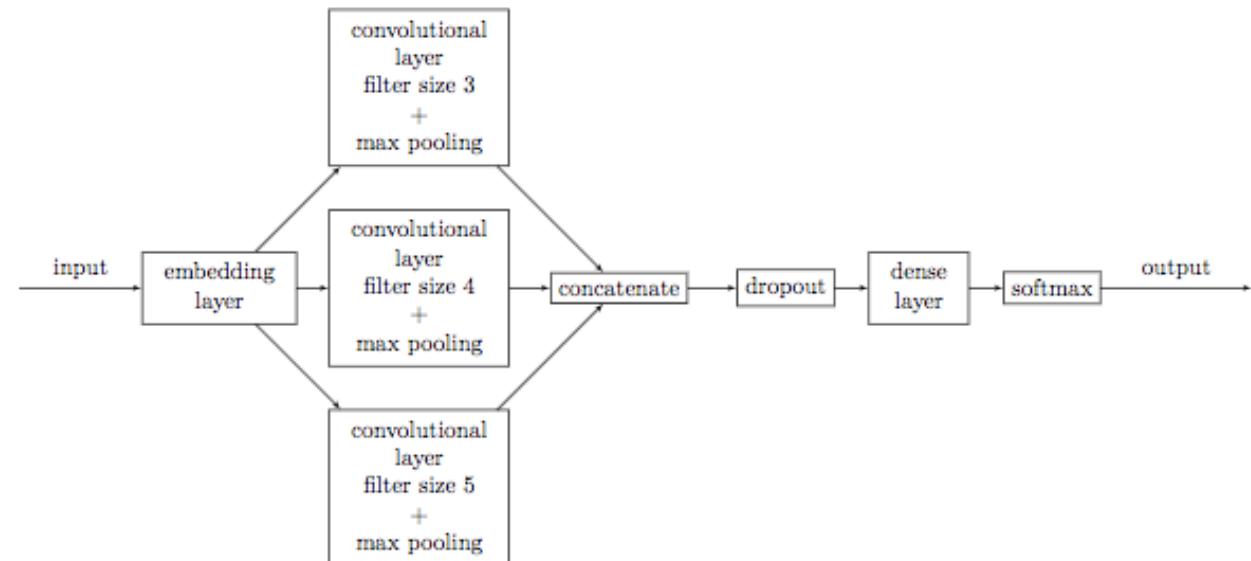
Filter unknowns with freq = 5
20% dropout every layer

No pre-trained vectors (?)

Neural Network Diagram

Best model

- shallow-and-wide CNN (<https://arxiv.org/abs/1510.03820>)
 - People should read papers!
- Thai2vec
 - Further adapted to the task



Character model (end-to-end)

- Char embedding
- 3 GRU layers
- 1000 character segment
 - Vote if longer than 1000

1	—	blead	0.58138
2	—	Danupat Kns	0.56919
3	▲ 1	Attawit Chaiyaroj	0.56297
4	▼ 1	Phizaz	0.56182
5	▲ 5	Wannita	0.55789
6	▲ 6	ChanatipSaetia	0.55306
7	▲ 1	burin	0.55284
8	▲ 6	AtomicBoyZ	0.55192
9	▼ 4	Chal	0.55192
10	▼ 1	Nattasit	0.55123

- Seems like a lucky random seed...

Yelp reviews

Duyu Tang Document Modeling with Gated Recurrent Neural Network
for Sentiment Classification, 2015 <http://aclweb.org/anthology/D15-1167>

Wongnai challenge top model Accuracy 0.5844

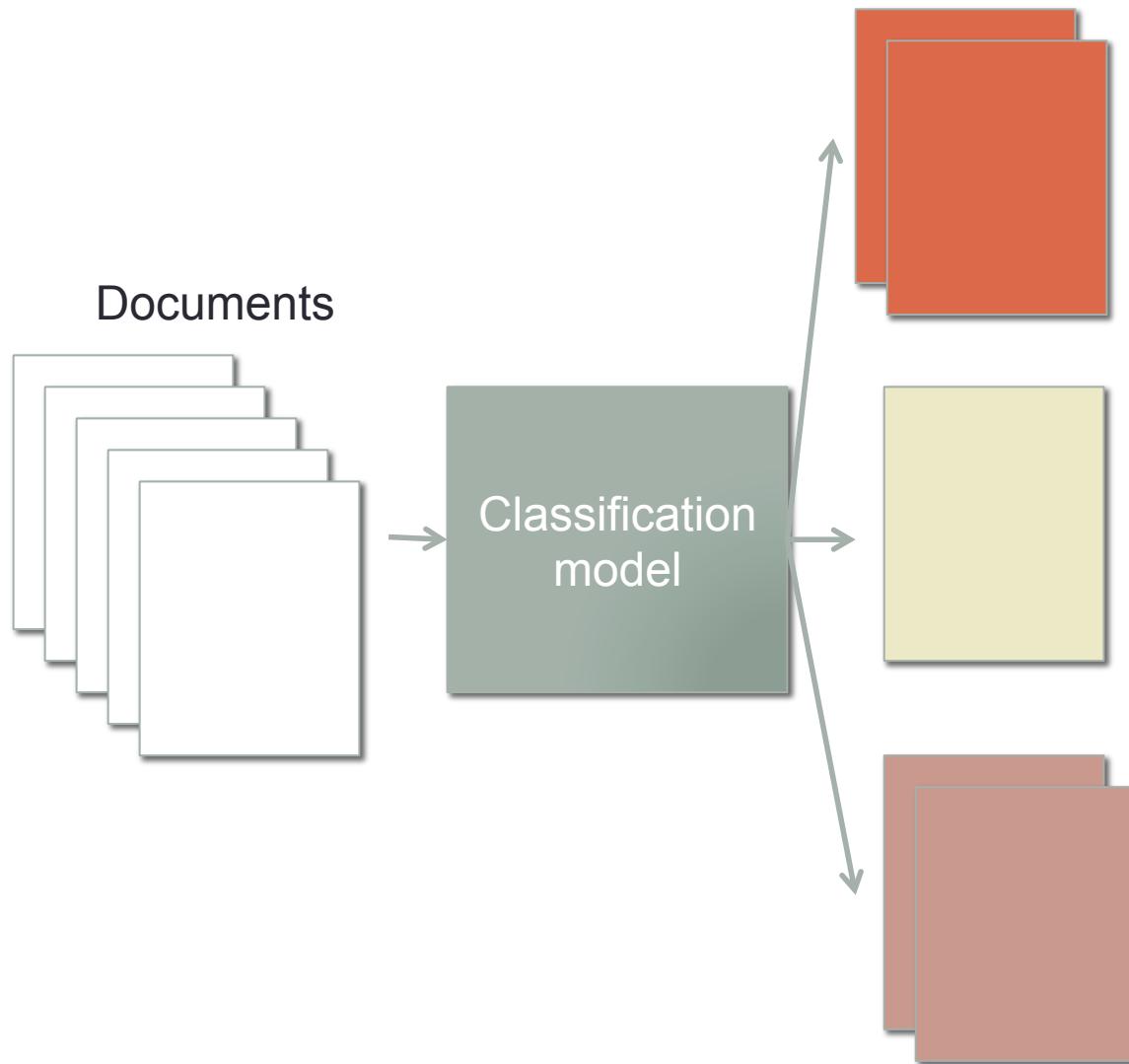
Corpus	#docs	#s/d	#w/d	V	#class	Class Distribution
Yelp 2013	335,018	8.90	151.6	211,245	5	.09/.09/.14/.33/.36
Yelp 2014	1,125,457	9.22	156.9	476,191	5	.10/.09/.15/.30/.36
Yelp 2015	1,569,264	8.97	151.9	612,636	5	.10/.09/.14/.30/.37
IMDB	348,415	14.02	325.6	115,831	10	.07/.04/.05/.05/.08/.11/.15/.17/.12/.18

	Yelp 2013		Yelp 2014		Yelp 2015		IMDB	
	Accuracy	MSE	Accuracy	MSE	Accuracy	MSE	Accuracy	MSE
Majority	0.356	3.06	0.361	3.28	0.369	3.30	0.179	17.46
SVM + Unigrams	0.589	0.79	0.600	0.78	0.611	0.75	0.399	4.23
SVM + Bigrams	0.576	0.75	0.616	0.65	0.624	0.63	0.409	3.74
SVM + TextFeatures	0.598	0.68	0.618	0.63	0.624	0.60	0.405	3.56
SVM + AverageSG	0.543	1.11	0.557	1.08	0.568	1.04	0.319	5.57
SVM + SSWE	0.535	1.12	0.543	1.13	0.554	1.11	0.262	9.16
JMARS	N/A	–	N/A	–	N/A	–	N/A	4.97
Paragraph Vector	0.577	0.86	0.592	0.70	0.605	0.61	0.341	4.69
Convolutional NN	0.597	0.76	0.610	0.68	0.615	0.68	0.376	3.30
Conv-GRNN	0.637	0.56	0.655	0.51	0.660	0.50	0.425	2.71
LSTM-GRNN	0.651	0.50	0.671	0.48	0.676	0.49	0.453	3.00

TEXT CLASSIFICATION

Intent, topic, sentiment

Text/document classification



Document classification

Type	Focus	Example
Topic	Subject matter	Sports vs Technology
Sentiment/opinion	Emotion (current state)	Negative vs Positive
Intent	Action (future state)	Order vs Inquiry

คินนี่จะได้ดูตอนใหม่แล้ว #ออกเจ้า #อดใจไม่ไหว

Topic: บุพเพสันนิวาส

Sentiment: positive

Action: watch

อยากจะสั่งพิชช่าหน้า纱瓦ยเอื้อนหน่อยครับ

Action: order_hawaiian

Does Anne Hathaway News Drive Berkshire Hathaway's Stock?

ALEXIS C. MADRIGAL | MAR 18, 2011 | TECHNOLOGY

[Share](#) [Tweet](#) [...](#)

TEXT SIZE
-

Like The Atlantic? Subscribe to [The Atlantic Daily](#), our free weekday email newsletter.

Email

[SIGN UP](#)

Given the awesome correlating powers of today's stock trading computers, the idea may not be as far-fetched as you think.



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."

<https://www.theatlantic.com/technology/archive/2011/03/does-anne-hathaway-news-drive-berkshire-hathaways-stock/72661/>

Other classification application

- Spam filtering
- Authorship id
 - Age id
 - Gender id
 - Native speaker id
 - First language id
- Auto tagging (information retrieval)
- Trend analysis
 - Patent landscape

Text classification definition

- Input
 - Set of documents: $D = \{d_1, d_2, d_3, \dots, d_M\}$
 - Each document is composed of words
 - $d_1 = [w_{11}, w_{12}, \dots w_{1N}]$
 - Set of classes: $C = \{c_1, c_2, c_3, \dots, c_J\}$
- Output
 - The predicted class c from the set C

Rule-based classification

- Rules based on phrases or other features
 - Wongnai Rating
 - แมลงสาบ -> 2 ดาว
 - อร่อย -> 4 ดาว
 - ไม่อร่อย -> 2 ดาว
 - ...
 - What if the phrase is ไม่ค่อยอร่อย
 - New rule
 - อร่อย โดยที่ไม่มีคำว่า ไม่อร่อย แล้วนั้น -> 4 ดาว
 - What if the phrase is ไม่ถูกแต่อร่อย
- This can yield very good results but...
- Building and maintaining rules is expensive!
- Or keep a word list of positive and negative words

Supervised text classification definition

- Input
 - Set of documents: $D = \{d_1, d_2, d_3, \dots, d_M\}$
 - And labels $Y = \{y_1, y_2, y_3, \dots, y_M\}$
 - Each document is composed of words
 - $d_1 = [w_{11}, w_{12}, \dots w_{1N}]$
 - Set of classes: $C = \{c_1, c_2, c_3, \dots, c_J\}$
- Output
 - A classifier $H: d \rightarrow c$

What classifier?

- Any classifier you like
- k-NN
- Naïve Bayes
- Logistic regression
- SVM
- Neural networks ←

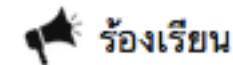
We use this kind of classifier
before, in homework and midterm

Outline

- Naïve Bayes
- Topic Models
 - Latent topic models (LDA)
- LDA2vec

Bag of words representation

★★★☆☆ 1 check-in



กิวยเตี๋ยวอร่อย ราคาถูก นั่งทานในห้องแอร์

เมนูเด็ด: บะหมี่ต้มยำหมูแดง

ส่วนตัวชอบกิวยเตี๋ยวต้มยำมากค่ะ รสชาติอร่อยไม่ต้องป魯เงyley แต่ชุปเปอร์ตินไก่ใส่ถั่วงอกมาด้วย เหมือนเป็นเกาเหลาตินไก่มากกว่าชุปเปอร์ตินไก่ รสชาติก็ยังไม่กลมกล่อมเท่ากิวยเตี๋ยวต้มยำ ตินไก่เปื่อยดี ทานง่าย

H

ส่วนตัวชอบกิวยเตี๋ยวต้มยำมากค่ะ รสชาติอร่อยไม่ต้องป魯เงyley แต่ชุปเปอร์ตินไก่ใส่ถั่วงอกมาด้วย เหมือนเป็นเกาเหลาตินไก่มากกว่าชุปเปอร์ตินไก่ รสชาติก็ยังไม่กลมกล่อมเท่ากิวยเตี๋ยวต้มยำ ตินไก่เปื่อยดี ทานง่าย

= 3

Bag of words representation

H [ส่วนตัวชอบกุ้ยเดียวต้มยำมากค่ะ รสชาติอร่อยไม่ต้องปูนเลย แต่ชุปเปอร์ตินไก่ใส่ถั่งอกมาด้วย เหมือนเป็น
เกาเหลาตินไก่มากกว่าชุปเปอร์ตินไก่ รสชาติก็ยังไม่กลมกล่อมเท่ากุ้ยเดียวต้มยำ ตินไก่เปื่อยดี ทานง่าย] = 3

Bag of words on care the presence of words or features but ignore word position and context

H [ชอบ, อร่อย, ไม่, ไม่, กลมกล่อม, ทานง่าย] = 3

Bag of words representation

H [ส่วนตัวชอบกิจกรรมยามากค่ะ รสนชาตior้อยไม่ต้องปูรงเลย แต่ชุปเปอร์ตินไก่ใส่ถังอกมาด้วย เหมือนเป็น
เงาเหลาตินไก่มากกว่าชุปเปอร์ตินไก่ รสชาติก็ยังไม่กลมกล่อมเท่ากิจกรรมยามาก ตินไก่นี่อยู่ดี ทานง่าย] = 3

Bag of words on care the presence of words or features but ignore word position and context

H [

Word	Count
ชอบ	1
อร่อย	1
ไม่	2
กลมกล่อม	1
ทานง่าย	1

] = 3

Bag of words for classification intuition

Test review

แมงสาบ
ใช้ได้
ถูก
อร่อย

1star

แมงสาบ
สกปรก
ແຍ່ງ

3star

ใช้ได้
ถูก
อร่อย
คาดฝัน

5star

ເຫດ
ເຂົ້າ
ອຮັບຍາຍ
ຍອດ

Bag of words for classification intuition



Figure 15.3 A spatial visualization of the document vectors for the four Shakespeare play documents, showing just two of the dimensions, corresponding to the words *battle* and *fool*. The comedies have high values for the *fool* dimension and low values for the *battle* dimension.

Reference: Jurafsky, Dan, and James H. Martin. Speech and language processing. 3rd edition draft, <https://web.stanford.edu/~jurafsky/slp3/>, August 2017

Bayes' Rule for classification

- A simple classification model
- Given document d , find the class c
 - $\underset{c}{\operatorname{Argmax}} P(c|d)$

$$= \underset{c}{\operatorname{Argmax}} \frac{P(d|c) P(c)}{P(d)}$$

Bayes' Rule

$$= \underset{c}{\operatorname{Argmax}} P(d|c) P(c)$$

$P(d)$ is constant wrt to c

$$= \underset{c}{\operatorname{Argmax}} P(x_1, x_2, \dots, x_n | c) P(c)$$

The document is represented by features
 x_1, x_2, \dots, x_n

Bayes' Rule for classification

- A simple classification model
- Given document d , find the class c

$$\underset{c}{\operatorname{Argmax}} P(c|d)$$

$$= \underset{c}{\operatorname{Argmax}} \frac{P(d|c) P(c)}{P(d)}$$

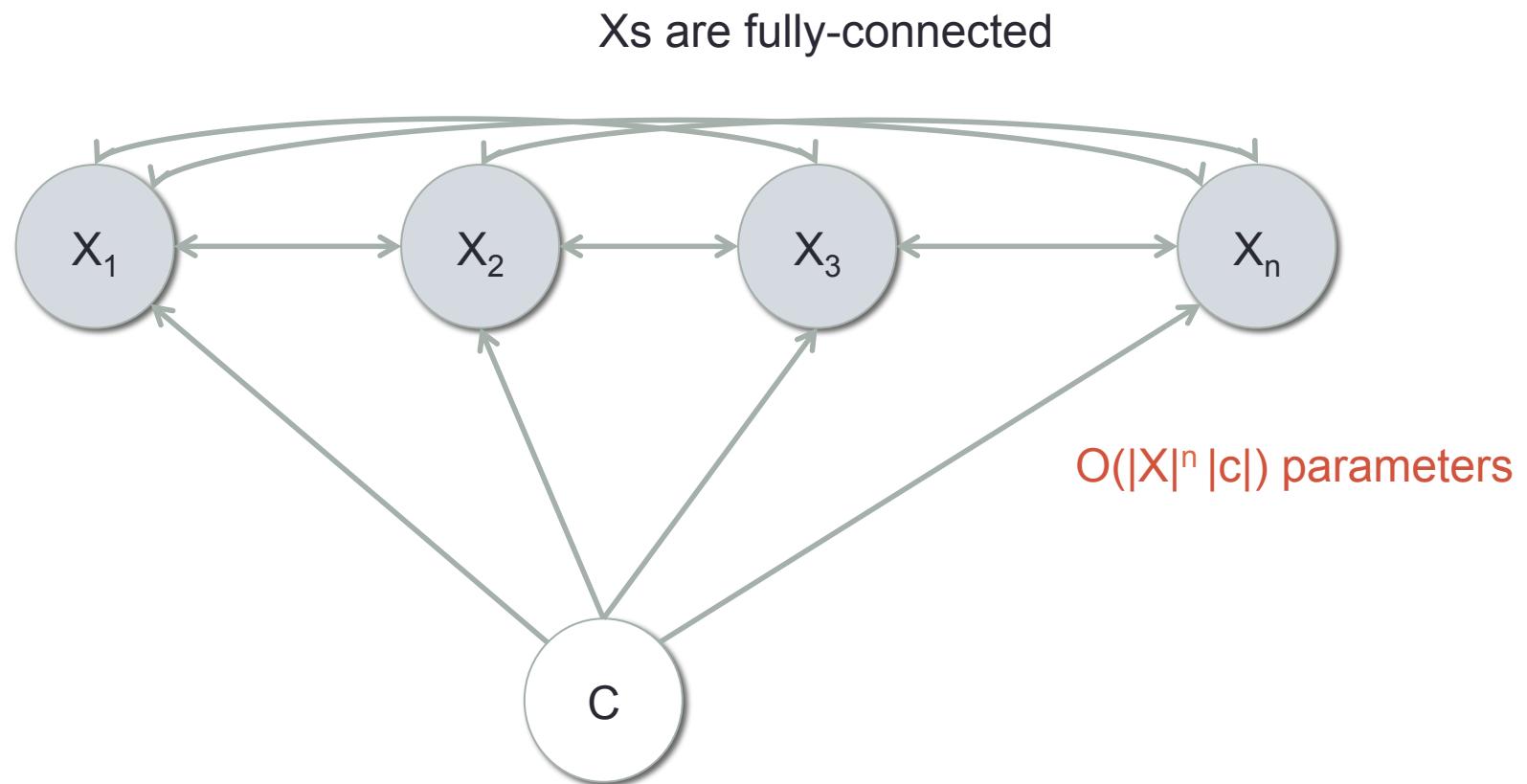
$$= \underset{c}{\operatorname{Argmax}} P(d|c) P(c)$$

$$= \underset{c}{\operatorname{Argmax}} P(x_1, x_2, \dots, x_n | c) P(c)$$

$P(x_1, x_2, \dots, x_n | c)$ requires $O(|X|^n |c|)$ parameters. Cannot train

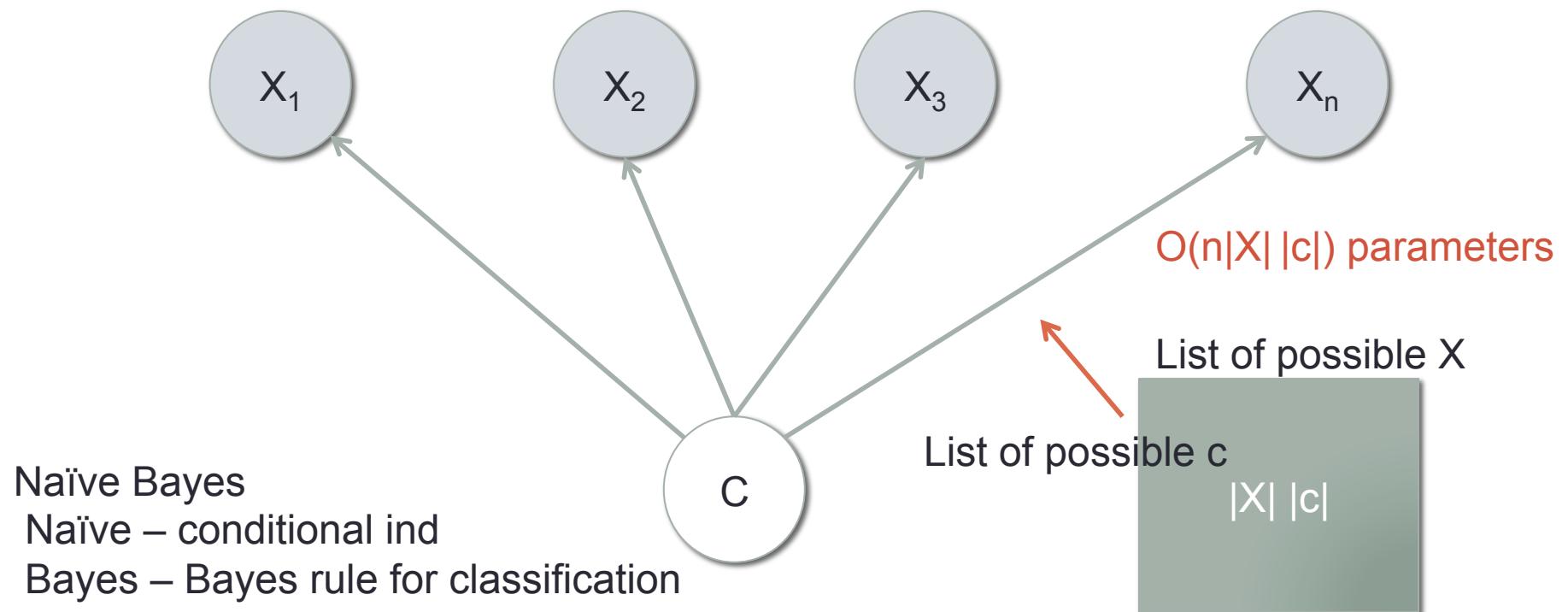
Graphical view

- $P(x_1, x_2, \dots, x_n | c) P(c)$



Bag of words assumption

- $P(x_1, x_2, \dots, x_n | c) P(c) = P(x_1|c) P(x_2|c) P(x_3|c) \dots P(x_n|c) P(c)$
 - Conditional independence



Learning the Naïve Bayes model

- As usual counts x_n is a feature that counts word occurrence
- $P(x_1|c)$ x_1 how many times ຂອດ appear
- $P(x_1 = 1| c=5) = \frac{\text{count}(x_1 = 1 , c = 5)}{\text{count}(c = 5)}$ List of possible counts
- $P(c)$ List of classes $|X| |c|$
- $P(c = 5) = \frac{\text{count} (c = 5)}{\text{count} (\text{all reviews})}$
- This is the Maximum Likelihood Estimate (MLE)

Learning the Naïve Bayes model

- What if we encounter zeroes in our table
- $P(x_1|c)$
 - x_n is a feature that counts word occurrence
 - x_1 how many times ยอด appear
- $P(x_1 = 1 | c=2) = \frac{\text{count}(x_1 = 1, c = 2)}{\text{count}(c = 2)}$
 - List of possible counts
 - $= 0$
 - List of classes $|X| |c|$

$$\begin{aligned} & P(\text{'ร้าน นี่ ราด หน้า ยอด ผัก ไม่ อร่อย'} | c = 2) \\ & = P(x_1=1|c=2)*... \\ & = 0 \end{aligned}$$

Zero probability regardless of other words

One solution: add-1 smoothing

Naives Bayes

- Can use other features beside word counts
 - Feature engineering – restaurant name, location, price range, reviewer id, date of review
 - Tedious but very powerful
 - Features > 10000
- Pros: very fast, very small model
- Robust especially for small training data (hand-crafted rules)
- A good fast baseline. Always try Naive Bayes or logistic regression in model search.
- Even with lots of data and rich features, Naives Bayes can be very competitive and fast!

Fast and good classification using n-grams

- Features: n-grams (bag of phrases)
- Model: logistic regression
- Very competitive results

Model	Yelp'13	Yelp'14	Yelp'15	IMDB
SVM+TF	59.8	61.8	62.4	40.5
CNN	59.7	61.0	61.5	37.5
Conv-GRNN	63.7	65.5	66.0	42.5
LSTM-GRNN	65.1	67.1	67.6	45.3
fastText	64.2	66.2	66.6	45.2

Table 3: Comparision with Tang et al. (2015). The hyperparameters are chosen on the validation set. We report the test accuracy.

Fast and good classification using n-grams

- Features: n-grams (bag of phrases)
- Model: logistic regression
- Very competitive results

	Zhang and LeCun (2015)		Conneau et al. (2016)			fastText
	small char-CNN	big char-CNN	depth=9	depth=17	depth=29	$h = 10, \text{bigram}$
AG	1h	3h	24m	37m	51m	1s
Sogou	-	-	25m	41m	56m	7s
DBpedia	2h	5h	27m	44m	1h	2s
Yelp P.	-	-	28m	43m	1h09	3s
Yelp F.	-	-	29m	45m	1h12	4s
Yah. A.	8h	1d	1h	1h33	2h	5s
Amz. F.	2d	5d	2h45	4h20	7h	9s
Amz. P.	2d	5d	2h45	4h25	7h	10s

Table 2: Training time for a single epoch on sentiment analysis datasets compared to char-CNN and VDCNN.

Tag prediction

Model	prec@1	Running time	
		Train	Test
Freq. baseline	2.2	-	-
Tagspace, $h = 50$	30.1	3h8	6h
Tagspace, $h = 200$	35.6	5h32	15h
fastText, $h = 50$	31.2	6m40	48s
fastText, $h = 50$, bigram	36.7	7m47	50s
fastText, $h = 200$	41.1	10m34	1m29
fastText, $h = 200$, bigram	46.1	13m38	1m37

Table 5: Prec@1 on the test set for tag prediction on YFCC100M. We also report the training time and test time. Test time is reported for a single thread, while training uses 20 threads for both models.

Naïve Bayes tricks for text classification

• Domain specific features

- Count words after “not” as a different word
 - I don’t go there. -> I don’t go_not there_not
- Upweighting: double counting words at important locations

- Words in titles
- First sentence of each paragraph
- Sentences that contain title words

Automatic text categorization using the importance of sentences
<https://dl.acm.org/citation.cfm?id=1072331>

Context-Sensitive Learning Methods for Text Categorization
https://www.researchgate.net/publication/2478208_Context-Sensitive_Learning_Methods_for_Text_Categorization

Information retrieval using location and category information
https://www.jstage.jst.go.jp/article/jnlp1994/7/2/7_2_141/_article

Relationship to language modeling

- .
- $P(x_1|c)$ x_1 is a feature that looks at the first word
- $P(x_1 = \text{ยอด} | c=5) = \frac{\text{count}(x_1 = \text{ยอด}, c = 5)}{\text{count}(c = 5)}$ List of words
- $P(c)$ List of classes
- $P(c = 5) = \frac{\text{count}(c = 5)}{\text{count}(\text{all reviews})}$ |X| |c|

List of words
List of classes
 $|X| |c|$

This looks like... n-grams, but instead of conditioning on the past, we condition on the topic –
bag of words model for topic modeling (unigram with topic)

Language modeling view

- Which class is this review
- $P(w|c)$

อร่อย แต่ ไม่ ถูก

Class= 1
อร่อย 0.01
แต่ 0.4
ไม่ 0.4
ถูก 0.03
...

Class= 5
อร่อย 0.4
แต่ 0.05
ไม่ 0.25
ถูก 0.15
...

$$P(s|c=1) = 0.01 * 0.4 * 0.4 * 0.03 = 0.000048$$
$$P(s|c=5) = 0.4 * 0.05 * 0.25 * 0.15 = 0.00075$$

Topic modeling

- Sometimes you want to model the topic of a document

Class=	
บรรยายการคุย	
อร่อย	0.01
แต่	0.4
ไม่	0.4
ถูก	0.03
...	
Class=	
อาหาร	
อร่อย	0.4
แต่	0.05
ไม่	0.25
ถูก	0.15
...	

อาหารที่นี่ไม่ค่อยอร่อย แต่ขนม
ใช้ได้เลย ถ้าว่างอาจจะกลับมา
กินอีก และนำให้สั่งเค้กใบเตย

ด้านบรรยายการคุย มีเสียงก่อสร้าง
มาจากตึกข้างๆ แต่นอกนั้นตก
แต่งโวโว แต่ยังขาดอะไรไป
หลายอย่าง

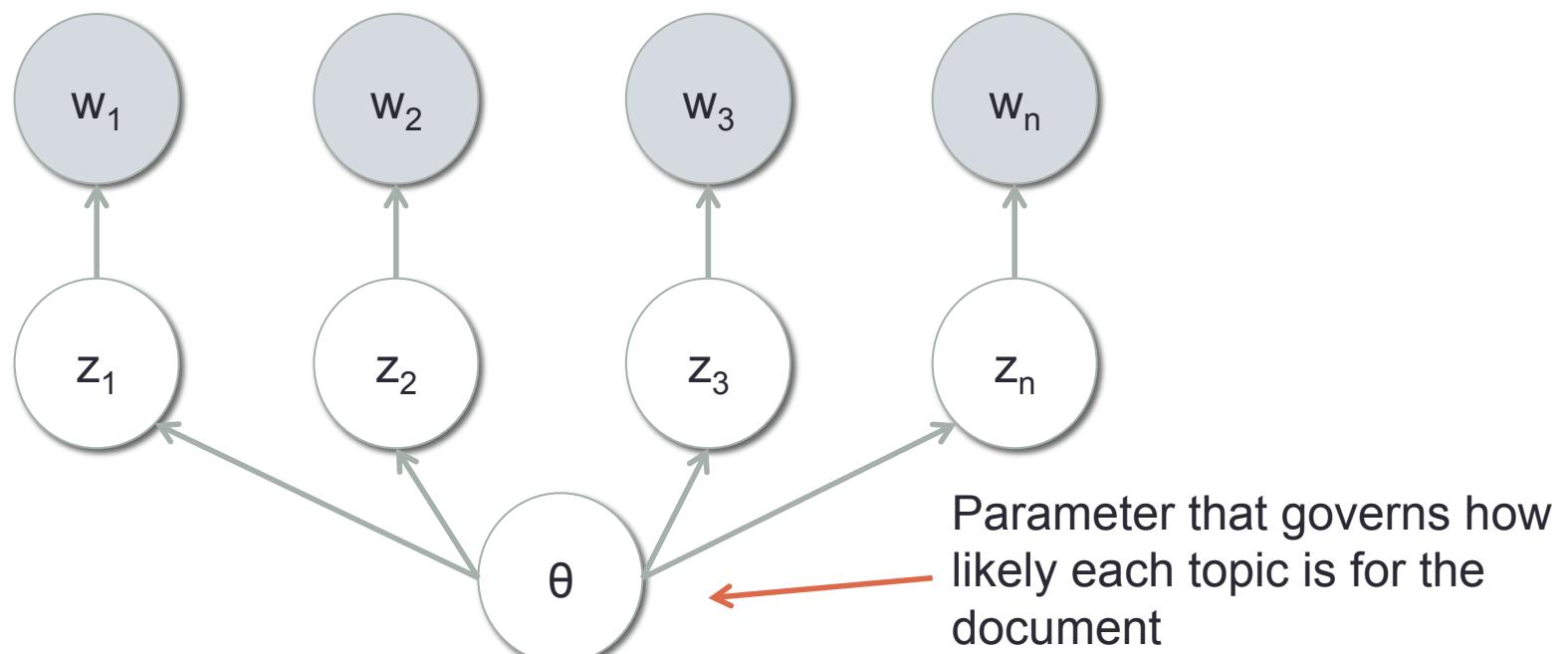
$$P(s|c=\text{บรรยายการคุย}) = ?$$
$$P(s|c=\text{อาหาร}) = ?$$

Naïve Bayes Topic modeling issues

- Most documents have multiple topics (multi-label).
 - Our model assumes 1 document 1 topic.
- Solution: Let a document be a mixture of topics (language model interpolation). Each word has its own topic, z .
 - $P(w) = P(w \text{ is topic A}) P(w | \text{topicA}) + P(w \text{ is topic B}) P(w | \text{topicB})$
 - $P(w \text{ is topic A}) + P(w \text{ is topic B}) = 1$

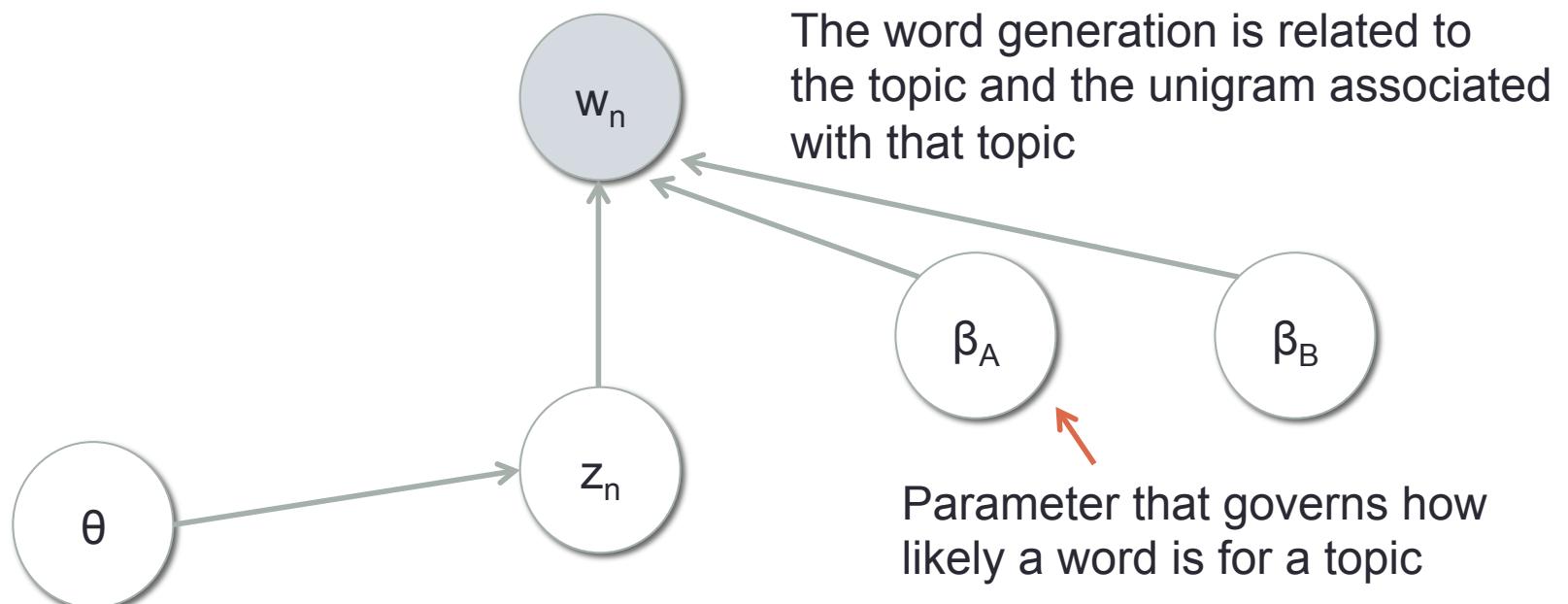
Naïve Bayes Topic modeling issues

- Most document have multiple topics. Our model assumes 1 document 1 topic.
 - Let a document be a mixture of topics (language model interpolation). Each word has its own topic, z .
 - $P(w) = P(z = A) P(w | z = A) + P(z = B) P(w | z = B)$
 - $P(z = A) + P(z = B) = 1, \quad \theta = P(z = A)$



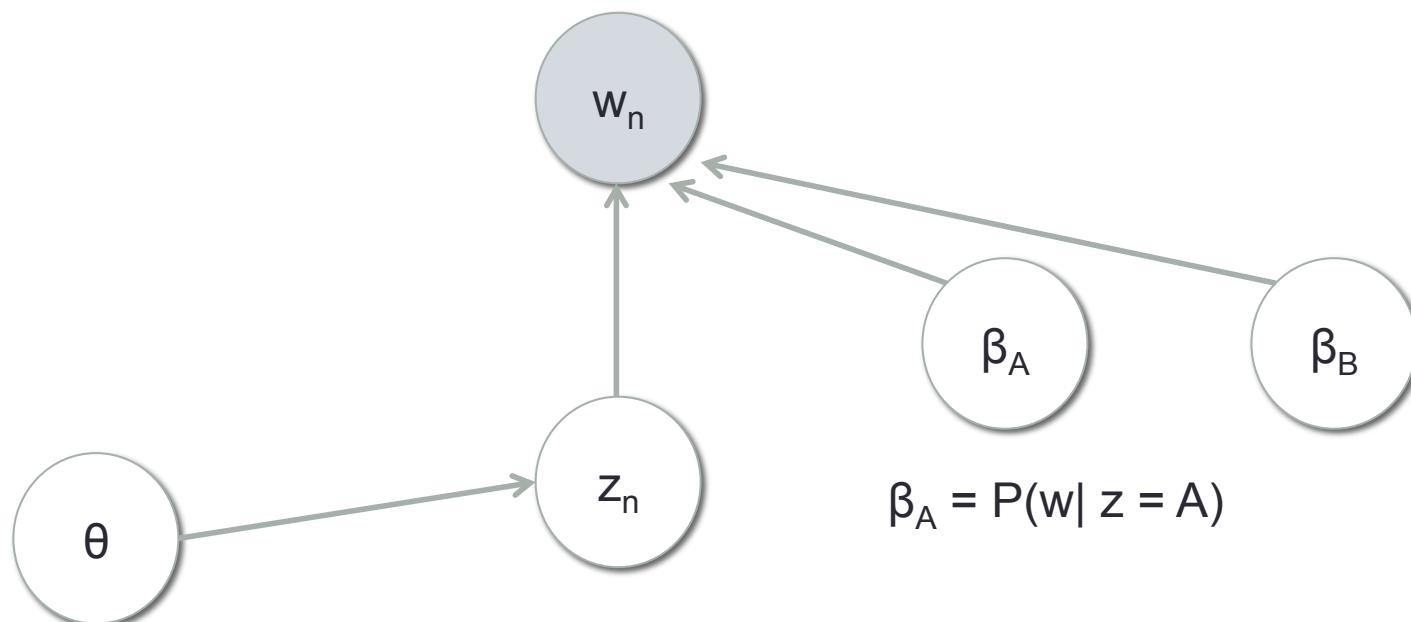
Topic modeling

- Most documents have multiple topics. Our model assumes 1 document 1 topic.
 - Let a document be a mixture of topics (language model interpolation). Each word has its own topic, z .
 - $P(w) = P(z = A) P(w | z = A) + P(z = B) P(w | z = B)$
 - $P(z = A) + P(z = B) = 1, \quad \theta = P(z = A) \quad \beta_A = P(w | z = A)$



Graphical model and generation

- You can generate a sample from a graphical model by following the arrows



Graphical model and generation

- Given θ, β_A, β_B



$A = 0.3$
 $B = 0.7$

Cat = 0.5
Dog = 0.5

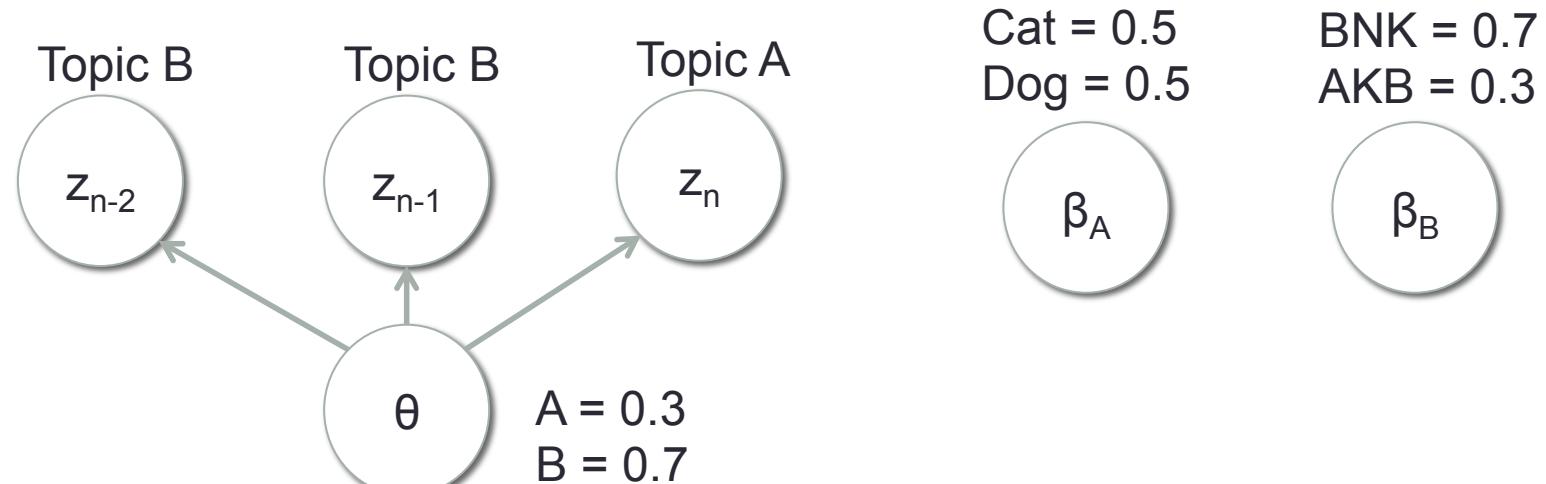


BNK = 0.7
AKB = 0.3



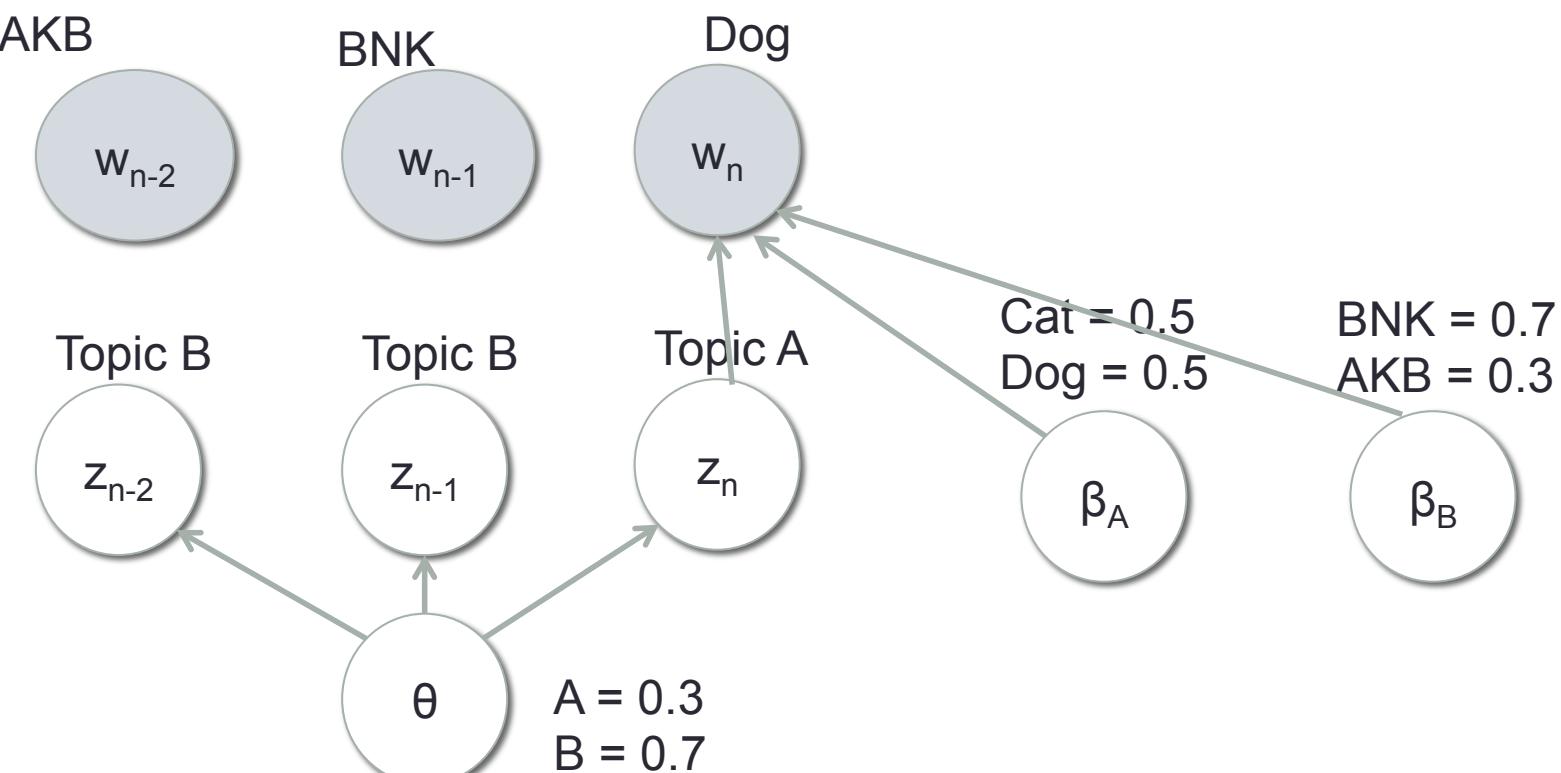
Graphical model and generation

- Given θ, β_A, β_B
- Generate (randomly create) z from θ



Graphical model and generation

- Given θ, β_A, β_B
- Generate (randomly create) z from θ
- Generate w_n from z_n, β_A, β_B

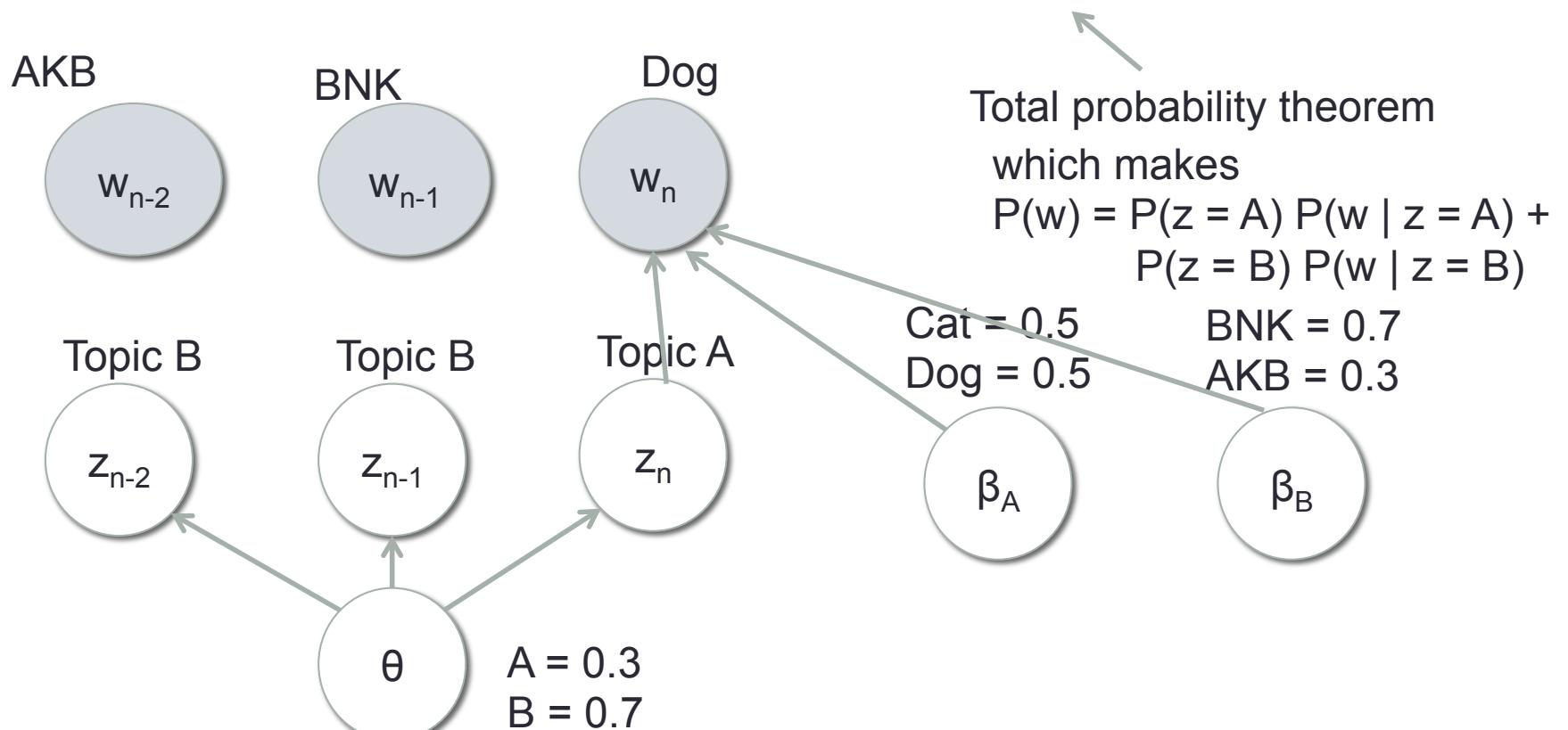


Graphical model and generation

How likely a sentence is likely to be generated follows this generation process
 $P(AKB, BNK, Dog, B, B, A) = P(B)P(B)P(A)P(AKB|B)P(BNK|B)P(Dog|A)$

Note

$$P(AKB, BNK, Dog) = P(AKB, BNK, Dog, A, A, A) + P(AKB, BNK, Dog, A, A, B) \\ P(AKB, BNK, Dog, A, B, A) + P(AKB, BNK, Dog, A, B, B) + \dots$$



Graphical models and plate notation

- Sometimes you have too many relationships.

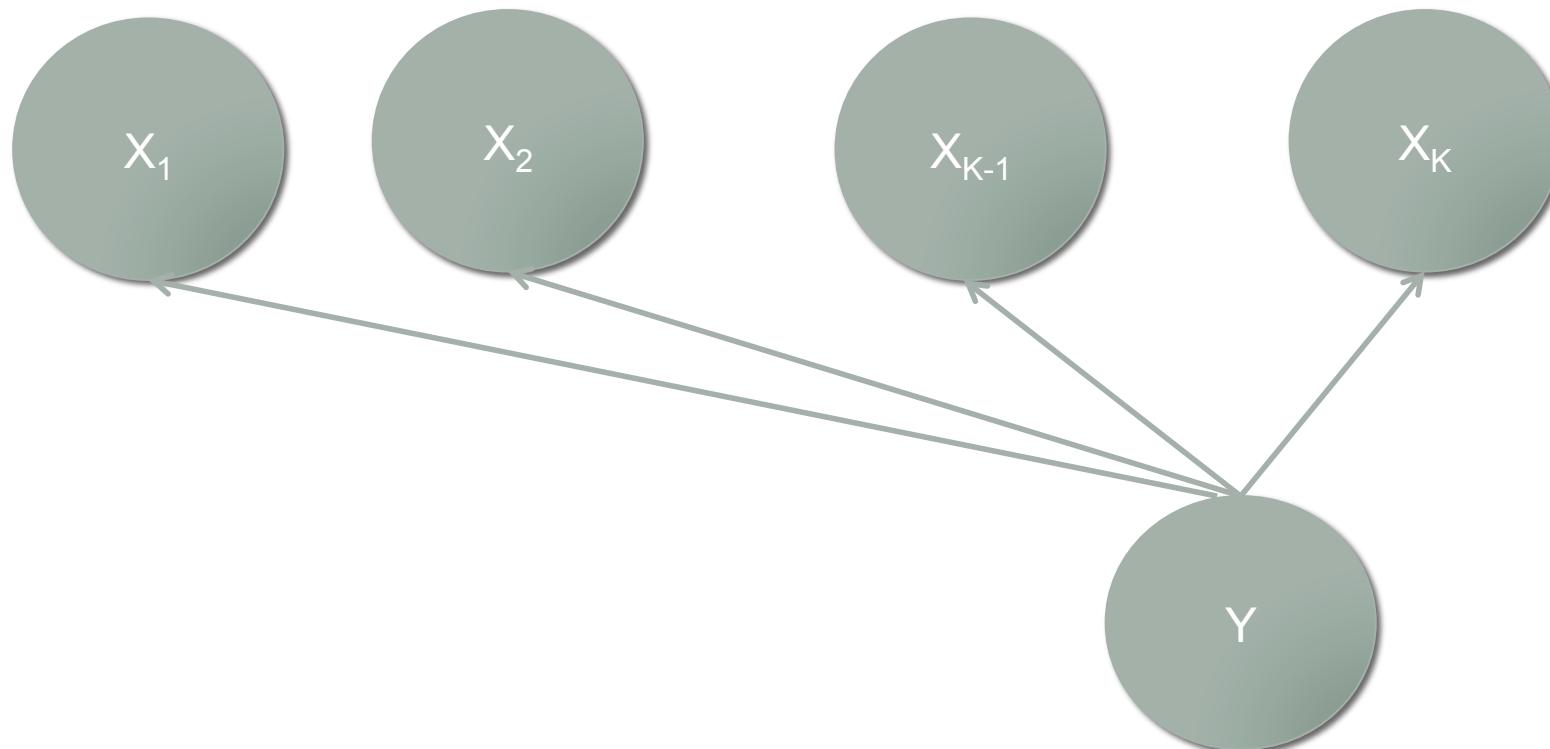
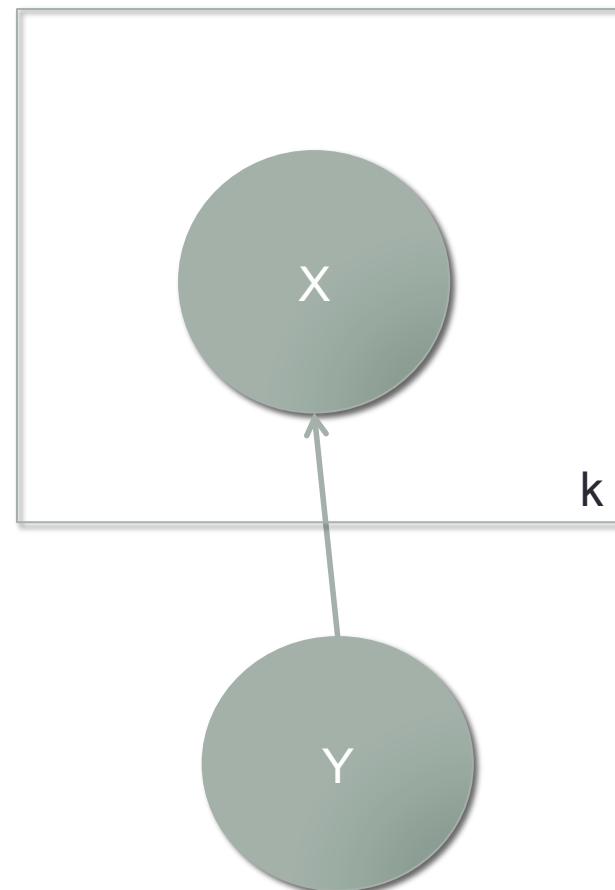
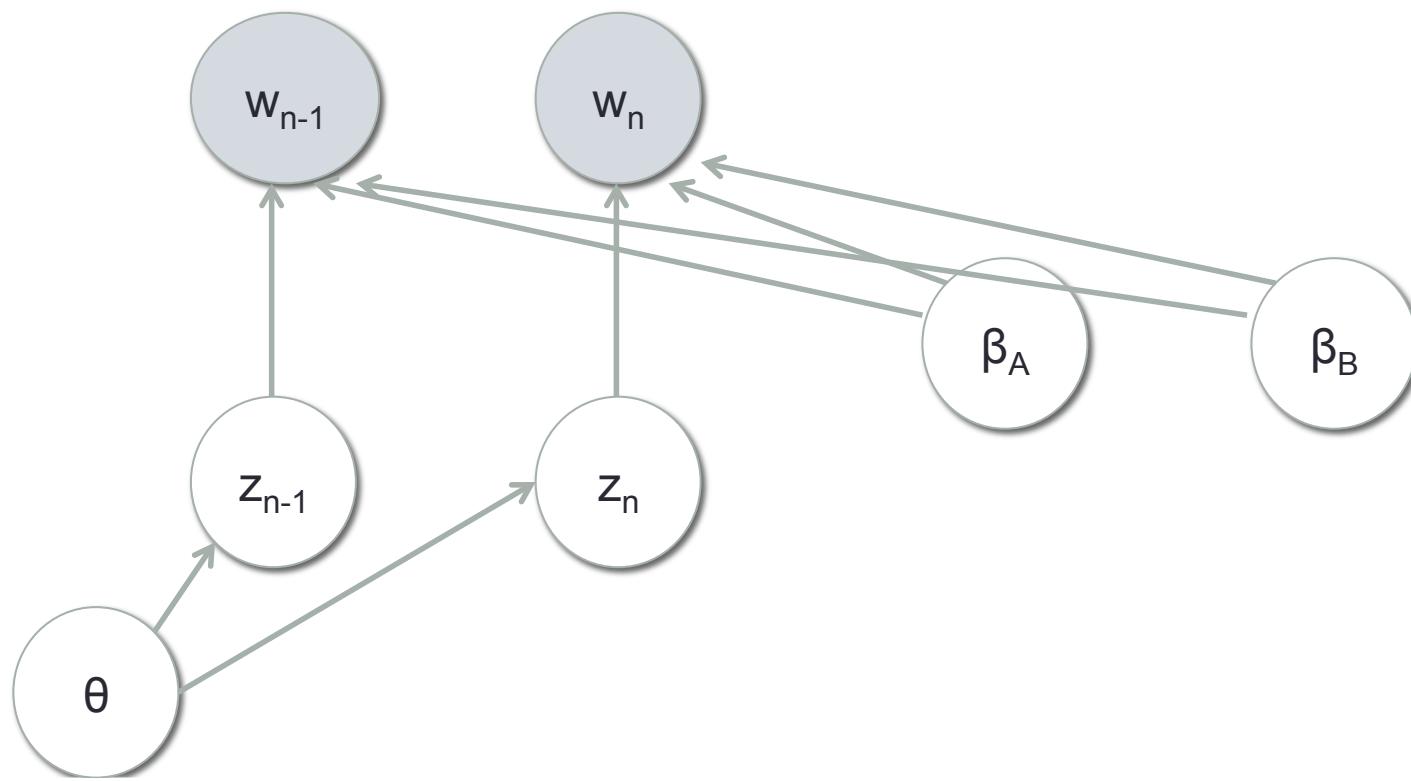


Plate notation

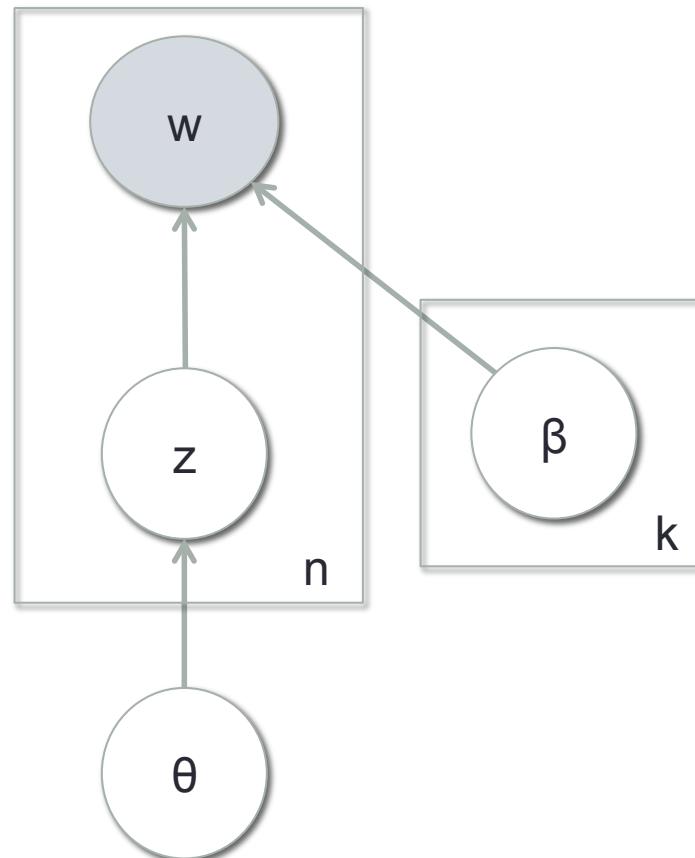
- Summarize by using a square box with number



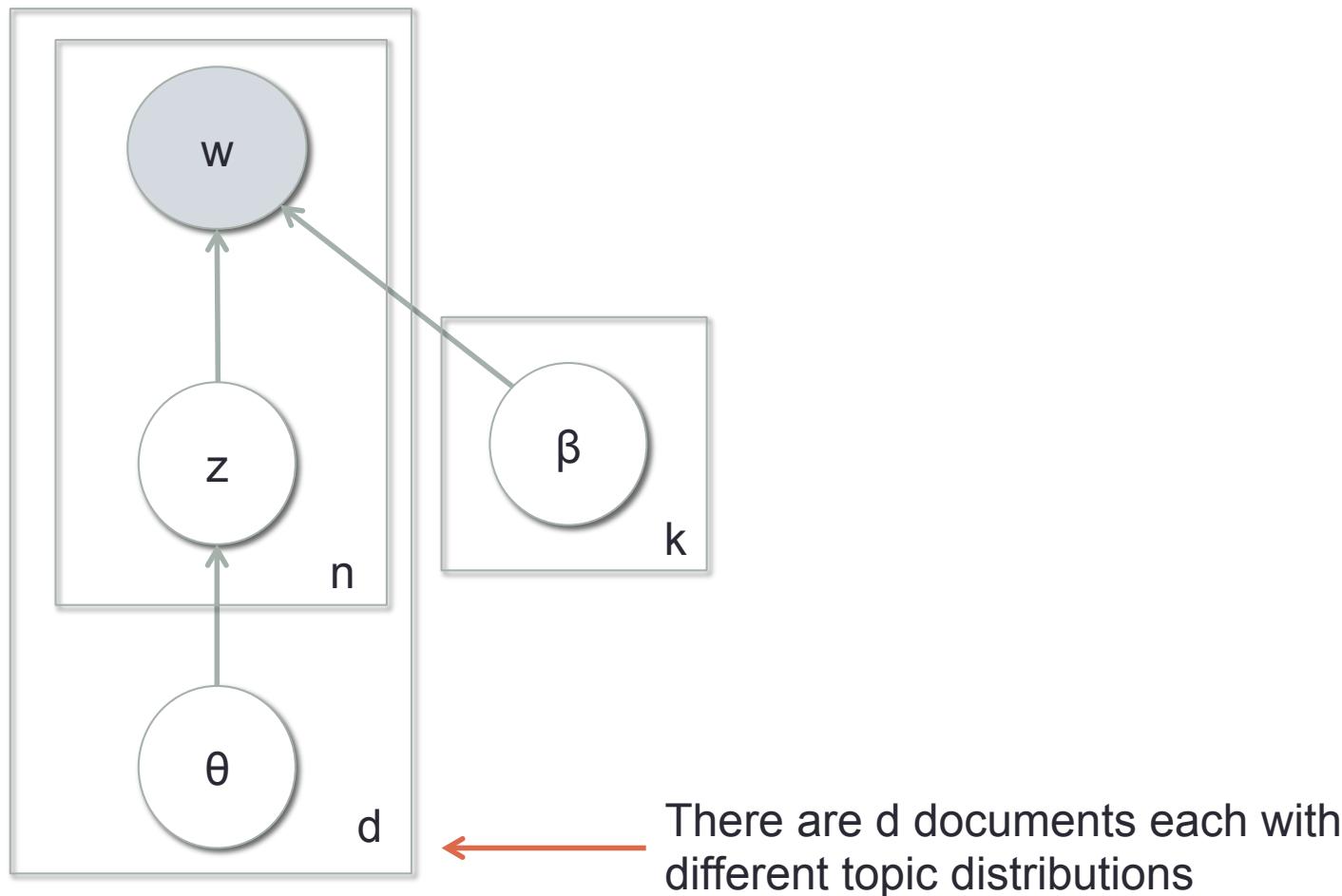
Topic modeling with plate notation



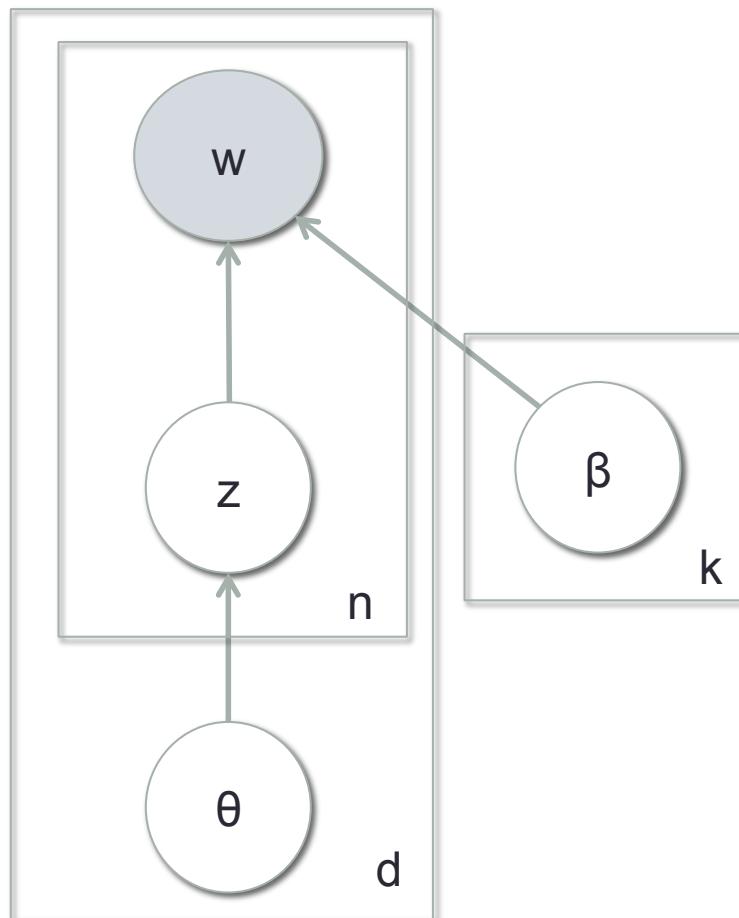
Topic modeling with plate notation



Topic modeling with plate notation

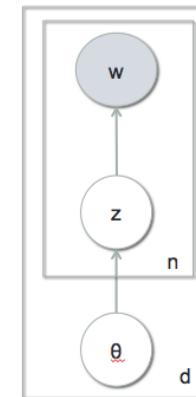


Topic modeling with plate notation



Called pLSA
probabilistic Latent Semantic Analysis

Note: if you look at other textbooks, you will see a slightly different picture



Learning topic latent model parameters

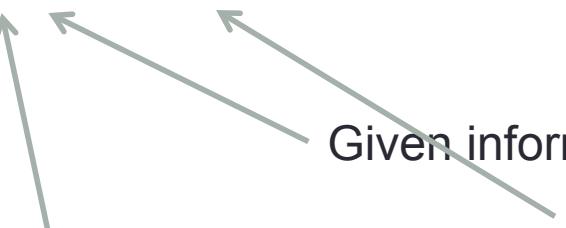
- How to find θ and β ?
 - “Cat,Dog,BNK”, “Cat, cat, cat, BNK”, “Dog, dog, BNK, dog, AKB”
- If we know, the latent topic z for each word we can use the counts
 - “Cat_A, Dog_A, BNK_B”
- $P(\text{Cat}|A) = \frac{\text{count}(\text{Cat_A})}{\text{count}(A)}$
- $P_1(A) = \frac{\text{count}(A)}{\text{count}(\text{all words in document 1})}$
- But we don't know the topic z for each word

Expectation maximization (EM)

- A method to iteratively maximize the likelihood of a model on training data
 - Expectation step (E step) – guess latent variables from model parameters (get soft counts)
 - Maximization step (M step) – re-estimate model parameters from latent variables (counts)

E-step

- Find an estimate for the latent variable given parameters θ, β
- $p(z_{di} | w_{di}, \theta, \beta)$


Topic of word i document d Given information of word i document d
and parameters θ, β

E-step

- Find an estimate for the latent variable given parameters θ, β

$$p(z_{di} | w_{di}, \theta, \beta) = \frac{p(z_{di}, w_{di}, \theta, \beta)}{\sum_{z'=1}^k p(z'_{di}, w_{di}, \theta, \beta)}$$

$p(w, \theta, \beta)$

E-step

- Find an estimate for the latent variable given parameters θ, β

$$p(z_{di} | w_{di}, \theta, \beta) = \frac{p(z_{di}, w_{di}, \theta, \beta)}{\sum_{z'=1}^k p(z'_{di}, w_{di}, \theta, \beta)}$$

$$= \frac{\theta_{z|d}\beta_{w|z}}{\sum_{z'=1}^k \theta_{z'|d}\beta_{w|z'}}$$

Index di for z and w dropped for clarity

E-step

- Find an estimate for the latent variable given parameters θ, β

$$p(z_{di} | w_{di}, \theta, \beta) = \frac{p(z_{di}, w_{di}, \theta, \beta)}{\sum_{z'=1}^k p(z'_{di}, w_{di}, \theta, \beta)}$$
$$= \frac{\theta_{z|d} \beta_{w|z}}{\sum_{z'=1}^k \theta_{z'|d} \beta_{w|z'}}$$


P(Word is from topic A | word is cat from document 1)

Probability that the word is from each topic. Use as counts

M-Step

- Instead of real counts use $P(z_{di})$ as the topic label
 - $P(\text{Cat}|A) = \frac{\text{count}(\text{Cat_A})}{\text{count}(A)}$
 - $P_1(A) = \frac{\text{count}(A)}{\text{count(all words in document 1)}}$

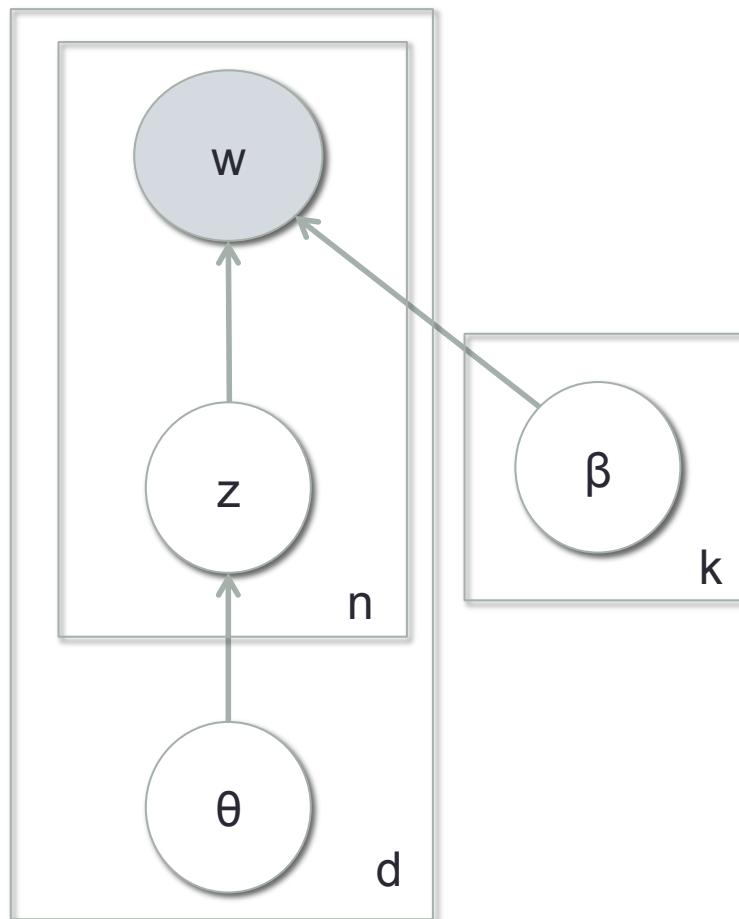
Expectation maximization (EM)

- Initialize θ, β
- Expectation step (E step) – guess latent variables from model parameters (get soft counts)
- Maximization step (M step) – re-estimate model parameters from latent variables (counts)
- Repeat E and M step until satisfied (Likelihood of the whole training set using the model does not change much)
- For more info on EM see pattern recognition lecture 6

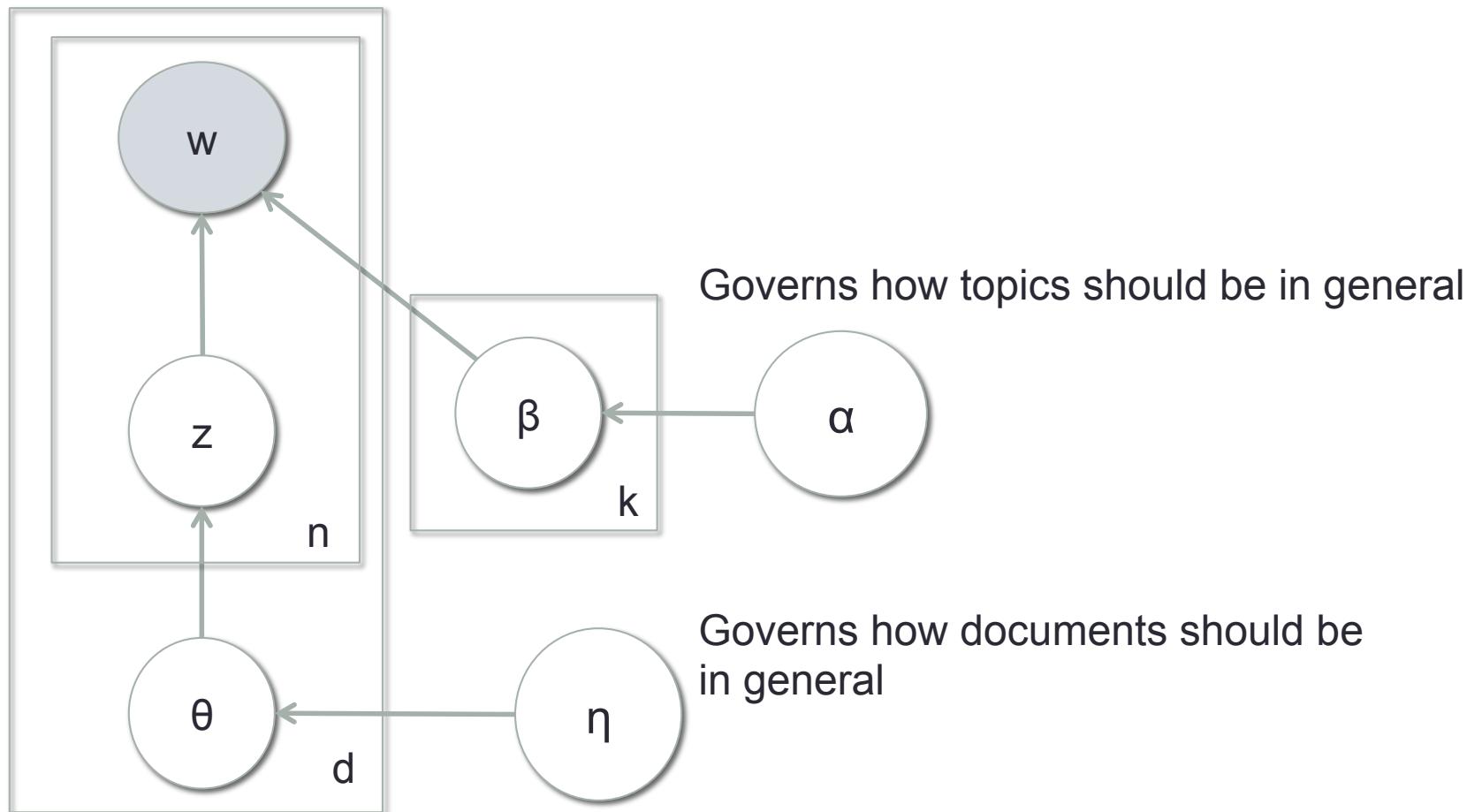
pLSA

- pLSA automatically clusters words into topic unigrams
 - Requires user to specify number of topics
- Automatically learn document representation based on the learned topics
 - $\text{DocA} = [0.7 \ 0.3]$ $\text{DocB} = [0.2 \ 0.8]$ $\text{DocC} = [0.5 \ 0.5]$
- Overfits easily to data outside of the training set
 - Nothing that ties all document together
 - A document from a document collection should be have topic distributions that are similar
- Solution: LDA (Latent Dirchlet Allocation)

pLSA



LDA



α and η

- α is a **Dirichlet distribution**
- Or a distribution of distributions...
- Example: rolling a die
- $P(x=1) = 0.3, P(x=2) = 0.1, \dots$
- This is a distribution. (**Multinomial distribution**)
- But what if I have a bag of dices, each with different distributions.
- α tells me the probability of what kind of die I will get

Dirichlet distribution

- pdf

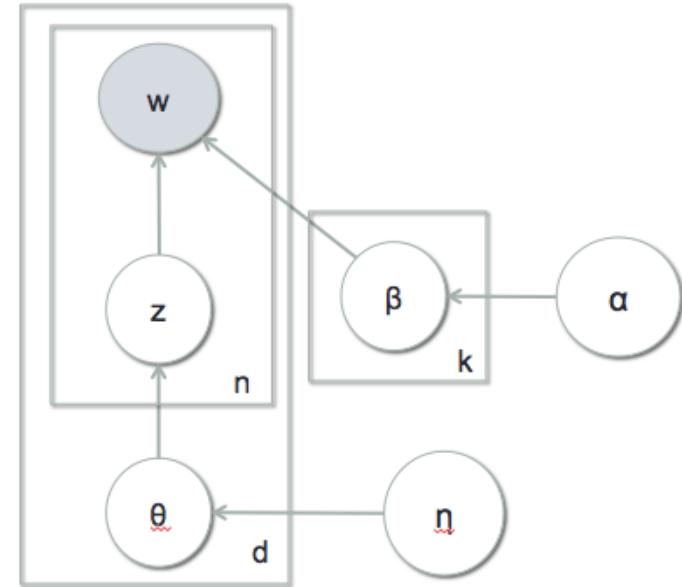
Parameters giving preference to each side of die/topic

$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^K x_i^{\alpha_i - 1}$$

Probability of each side of die, probability of each topic

Generative process

- For each topic $k \in \{1, \dots, K\}$:
 - Draw word distributions $\beta_k \sim \text{Dir}(\eta)$
- For each document $d \in \{1, \dots, D\}$:
 - Draw topic proportions $\theta_d \sim \text{Dir}(\alpha)$
 - For each word in a document $n \in \{1, \dots, N\}$:
 - Draw a topic index $z_{dn} \sim \text{Mult}(\theta)$
 - Generate word from chosen topic
 - $w_{dn} \sim \text{Mult}(\beta_{zdn})$



Topics

gene 0.04
dna 0.02
genetic 0.01
...

life 0.02
evolve 0.01
organism 0.01
...

brain 0.04
neuron 0.02
nerve 0.01
...

data 0.02
number 0.02
computer 0.01
...

Documents

Seeking Life's Bare (Genetic) Necessities

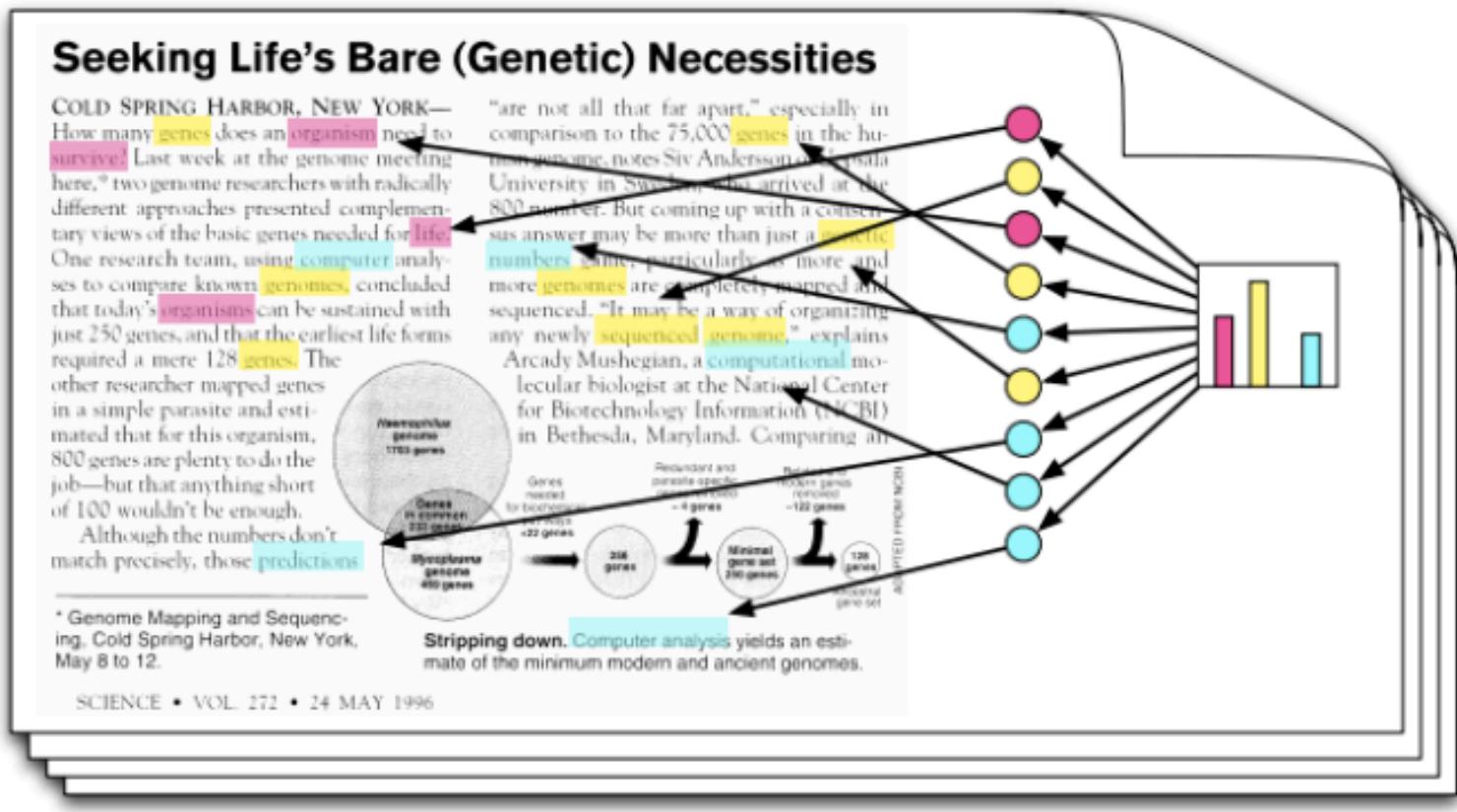
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



Introduction to Probabilistic Topic Models, Blei 2011

<http://menome.com/wp/wp-content/uploads/2014/12/Blei2011.pdf>

Parameter learning

- How do we actually learn these parameters and latent variables?
 - Gibbs sampling
 - Variational methods
- Totally beyond the scope of this class

LDA

- Automatically learns topics, and the word distribution of each topic
 - Just give a bunch of documents!
 - Each document is given a mixture of topics
 - Dirichlet prior prefers sparse topics – each document only have probability in few topics – easy for interpretability
- Requires user to pick number of topic
- Requires user to make sense of the learned topics

Example usage

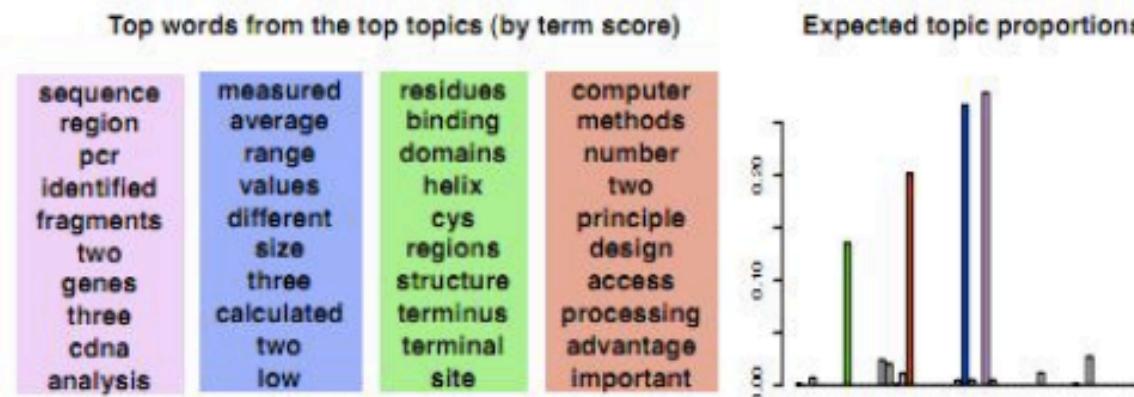
- Lots of lectures. Can we discover the topics? Can we classify the lectures? No topic labels given

Top 10 high probability nouns in word probabilities of Topics 3 ($\theta_{k=3}$), 15 ($\theta_{k=15}$), and 26 ($\theta_{k=26}$).

Topic 3 (~ classical mechanics)	Topic 15 (~astronomy)	Topic 26 (~ (time) unit)
m	Light	Percent
Energy	Degrees	Time
Force	Angle	Dollars
Mass	Frequency	Times
Point	Energy	Minutes
Velocity	Direction	Day
Direction	Waves	Bit
v	Sun	Year
Times	Star	Hour
Speed	Speed	Half

Chance and Statistical Significance in Protein and DNA Sequence Analysis

Samuel Karlin and Volker Brendel



Abstract with the most likely topic assignments

Statistical approaches help in the determination of significant configurations in protein and nucleic acid sequence data. Three recent statistical methods are discussed: (i) score-based sequence analysis that provides a means for characterizing anomalies in local sequence text and for evaluating sequence comparisons; (ii) quantile distributions of amino acid usage that reveal general compositional biases in proteins and evolutionary relations; and (iii) r-scan statistics that can be applied to the analysis of spacings of sequence markers.

Top Ten Similar Documents

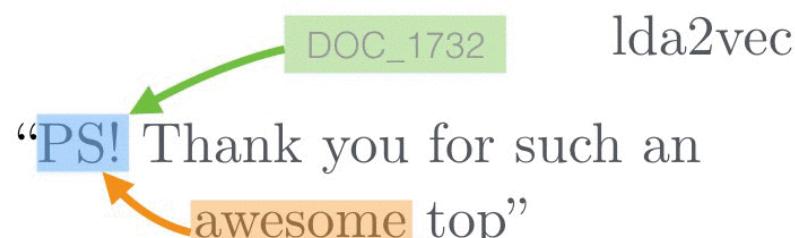
- Exhaustive Matching of the Entire Protein Sequence Database
- How Big Is the Universe of Exons?
- Counting and Discounting the Universe of Exons
- Detecting Subtle Sequence Signals: A Gibbs Sampling Strategy for Multiple Alignment
- Ancient Conserved Regions in New Gene Sequences and the Protein Databases
- A Method to Identify Protein Sequences that Fold into a Known Three-Dimensional Structure
- Testing the Exon Theory of Genes: The Evidence from Protein Structure
- Predicting Coiled Coils from Protein Sequences
- Genome Sequence of the Nematode *C. elegans*: A Platform for Investigating Biology

LDA with deep learning

- LDA was developed on discrete inputs (words)
 - Modified to work with dense representation (word vectors)
 - “Gaussian LDA for Topic Models with Word Embeddings”
 - <http://www.aclweb.org/anthology/P15-1077>
- Modified network structure and loss function to include LDA traits
 - Use like a neural network(just like how we use crf in neural networks)
 - LDA2vec

LDA2Vec

- Word2Vec estimates words around a pivot word (local information)
- LDA estimates words from a document topic (global information)
- LDA2Vec estimates words using both information



word2vec
“PS! Thank you for such an
awesome top”

LDA
“PS! Thank you for such an
awesome top”

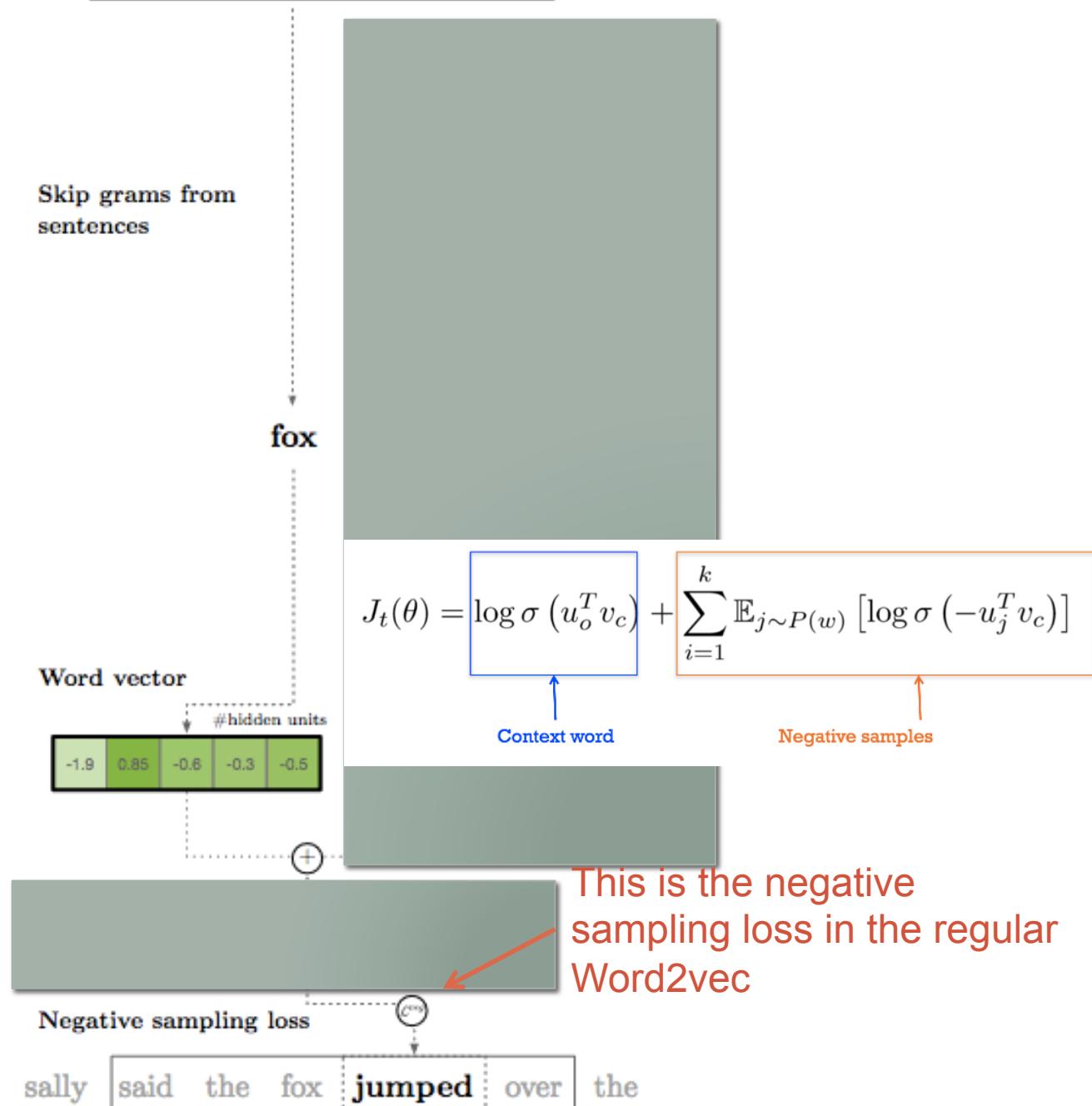
LDA2Vec

- Combines advantages from LDA and Word2Vec
 - Combine global info with local info
Word German + airline topic = Lufthansa
 - Dense word vector but sparse document vectors
Apple = [0.2 -1.2 -2.1 1.5]
Doc1 = [80% 20% 0% 0%]
 - Mixture of topics for each document vector give interpretability for humans

sally said the fox jumped over the

A typical skipgram
with negative sampling
A word is transformed to
a word vector

Maximize dot product
between pivot word and
context words
while reducing dot
product between pivot
word and negative
sampled words

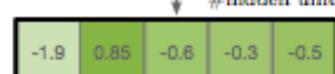


sally said the fox jumped over the

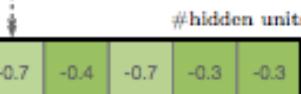
Skip grams from sentences

fox

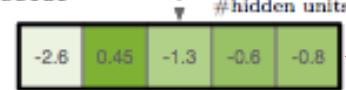
Word vector



Document vector



Context vector



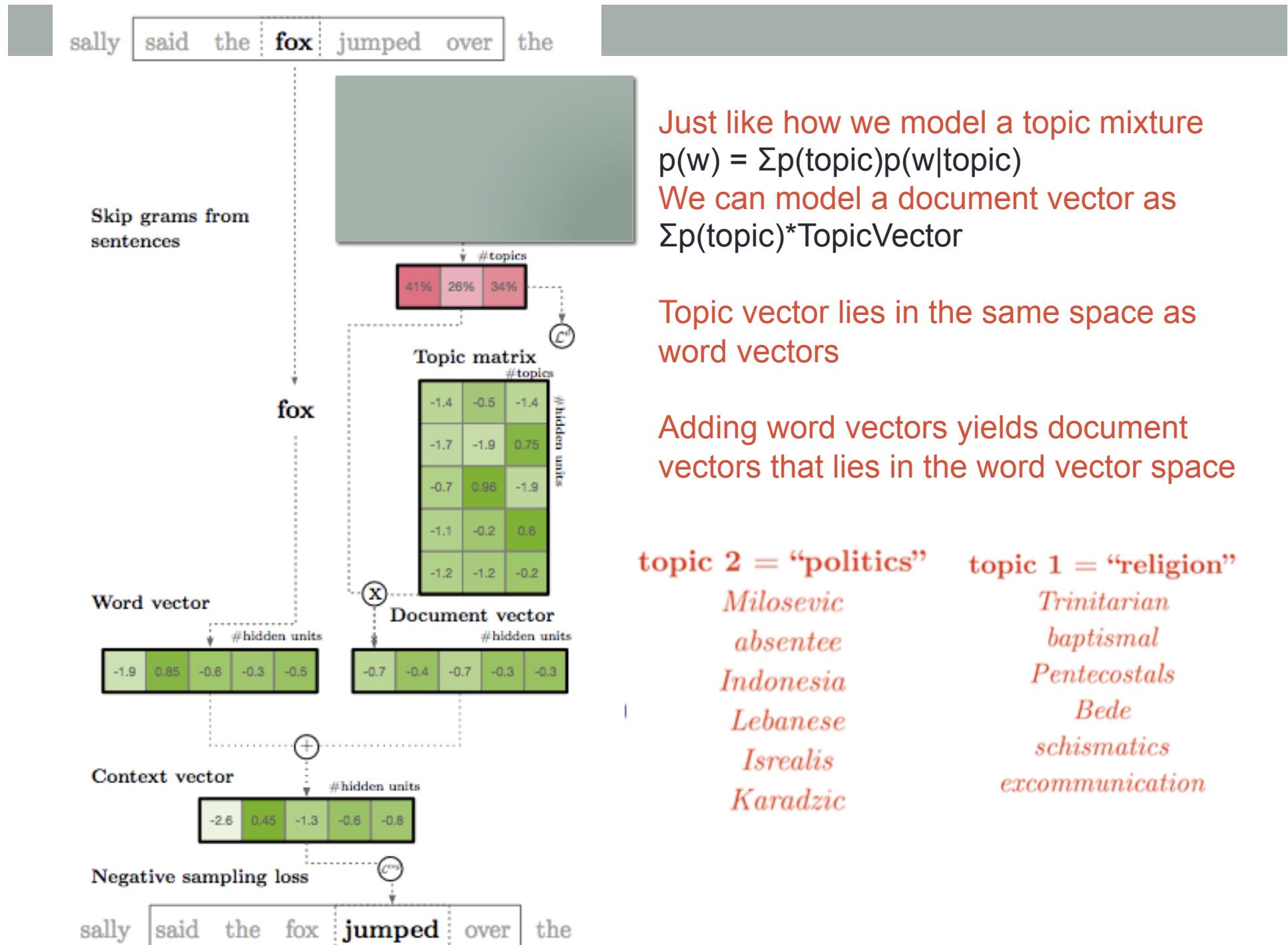
Negative sampling loss

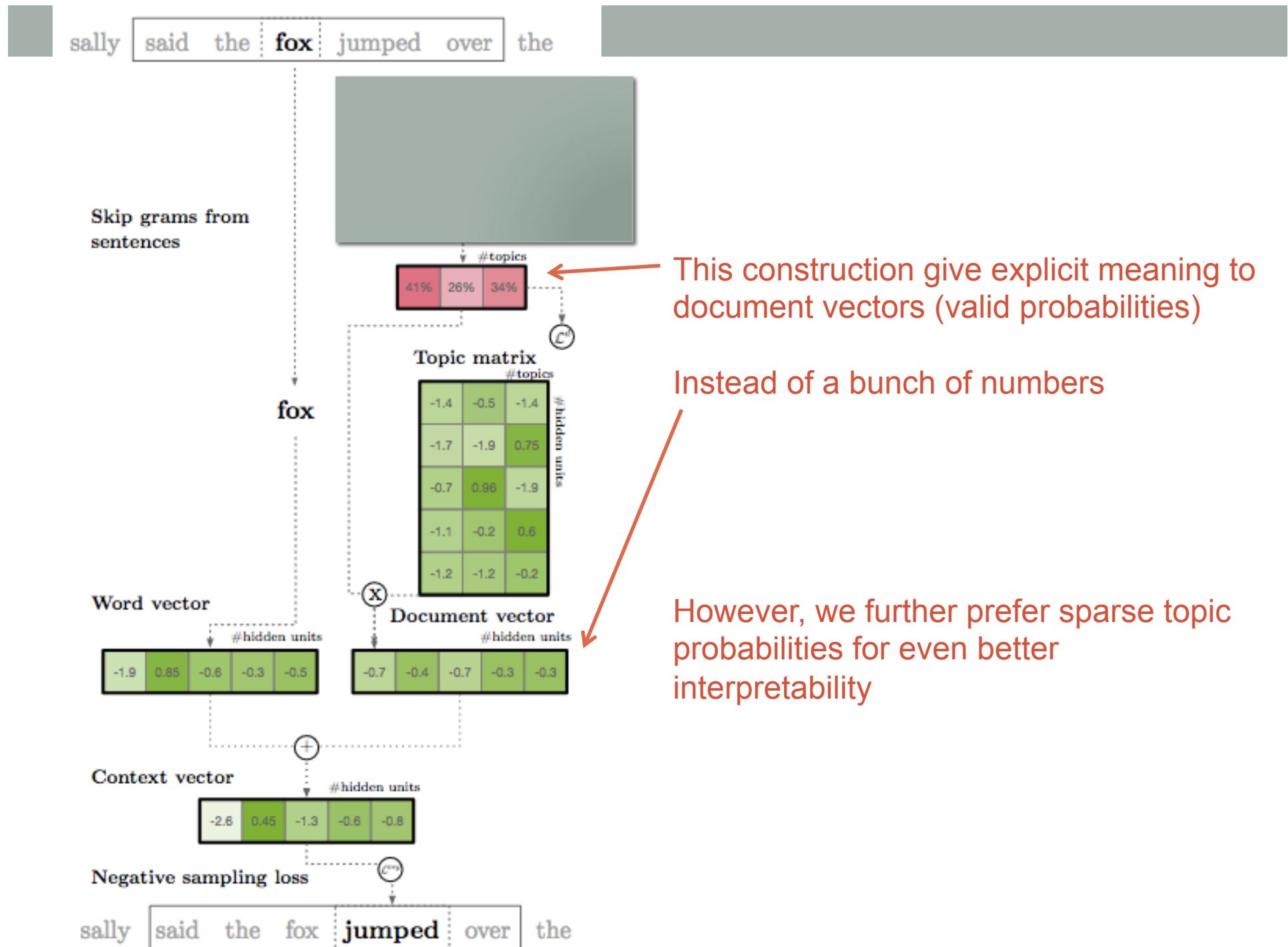
sally said the fox **jumped** over the

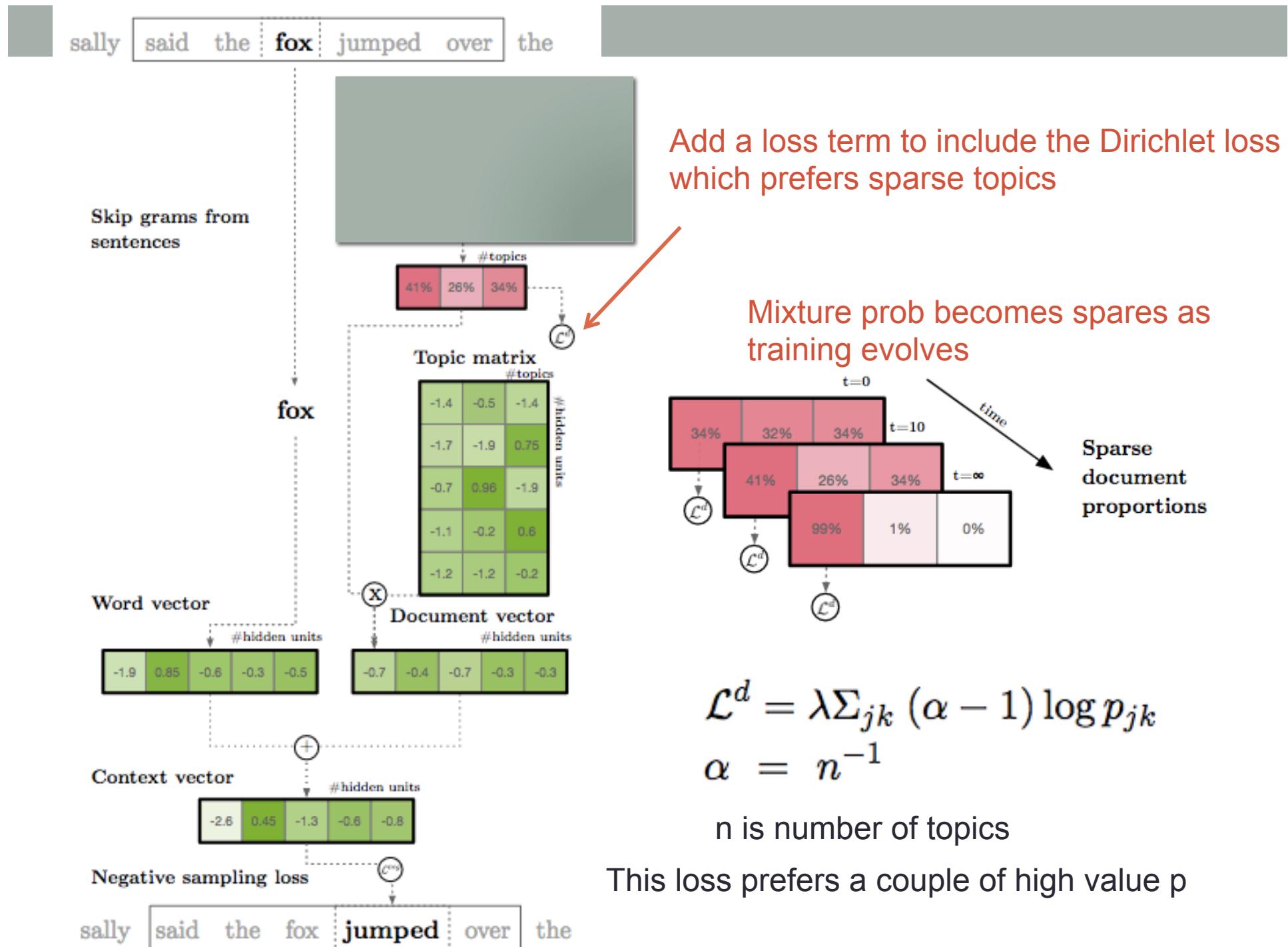
Add global info by having a document vector

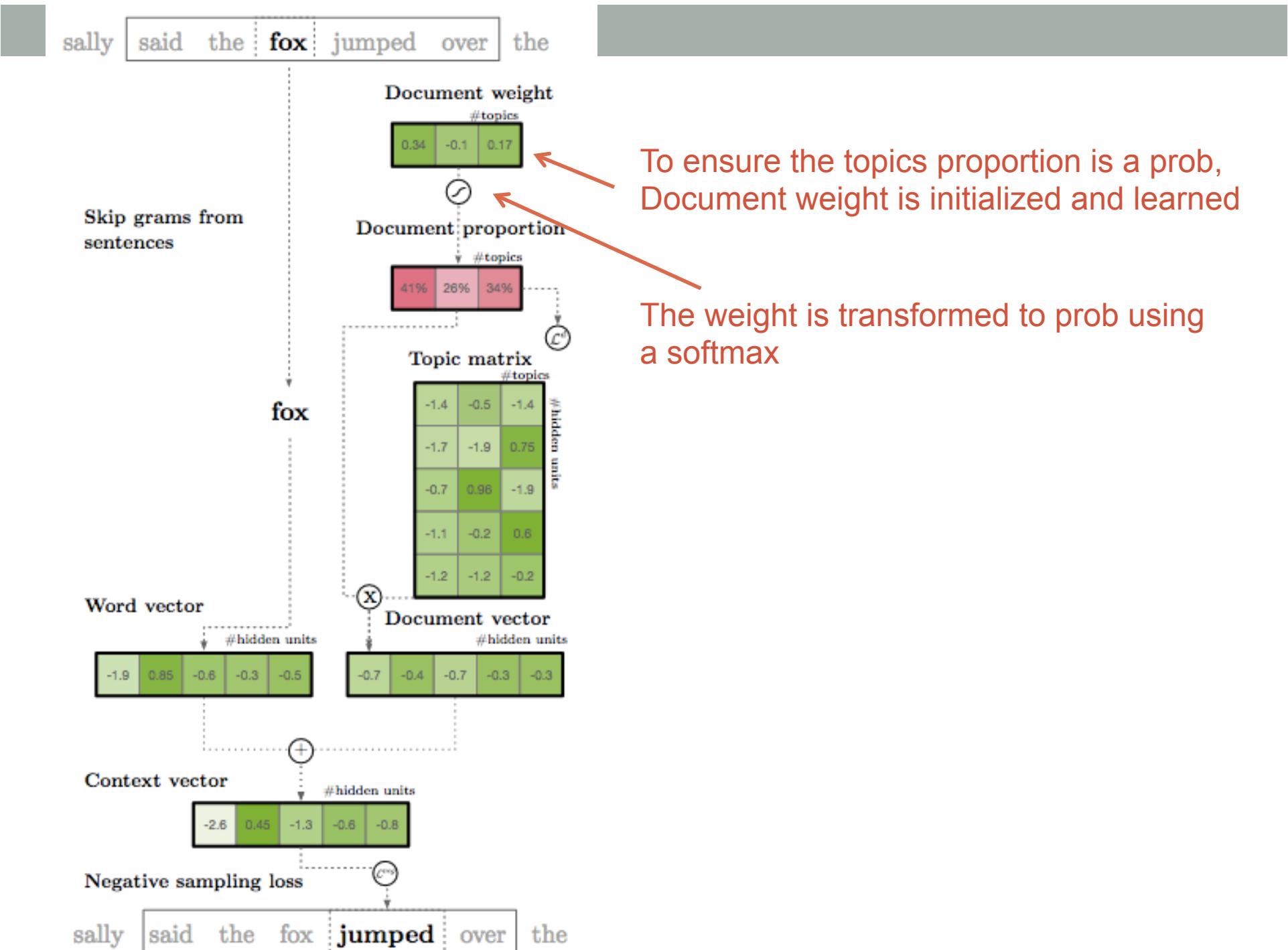
Document vector lies in the same space as word vector
can add together, just like word2vec's
king – queen + man

This vector (context vector) now has information
more than just a single word. More on context
vectors next lecture!









Overall loss function

- Loss from skipgram and Dirichlet

Dirichlet loss

$$\mathcal{L} = \sum_{word\ pairs\ (i,j)} [\sigma(\vec{c}_j \cdot \vec{w}_i) + \sigma(-\vec{c}_j \cdot \vec{w}_{negative})] + \sum_{documents\ k} \log q(d_k|\alpha)$$

Skipgram loss

$$\vec{c}_j = \vec{w}_j + \vec{d}_j$$

Context = word + document vector

$$\vec{d}_j = a_{j0} \cdot \vec{t}_0 + a_{j1} \cdot \vec{t}_1 + \dots$$

Document vector is mixture of topics

$$q(d_k|\alpha) = \sum_k (\alpha_k - 1) \log \alpha_k$$

Dirichlet likelihood term

LDA2Vec example

Selected Topic: 38

Previous Topic

Next Topic

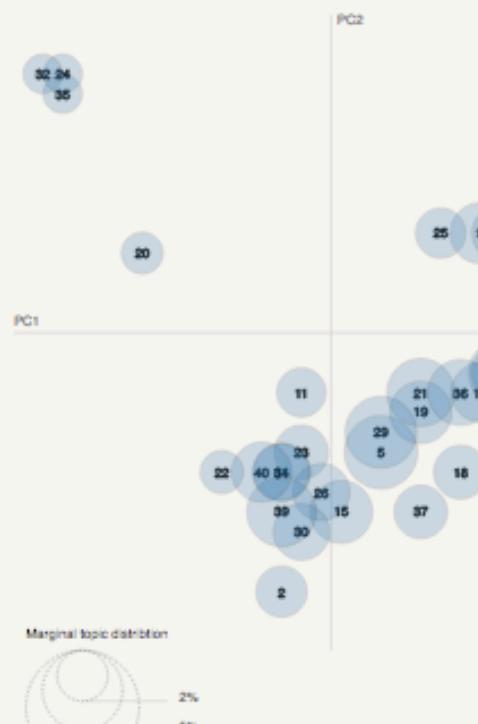
Clear Topic

Slide to adjust relevance metric:(2)

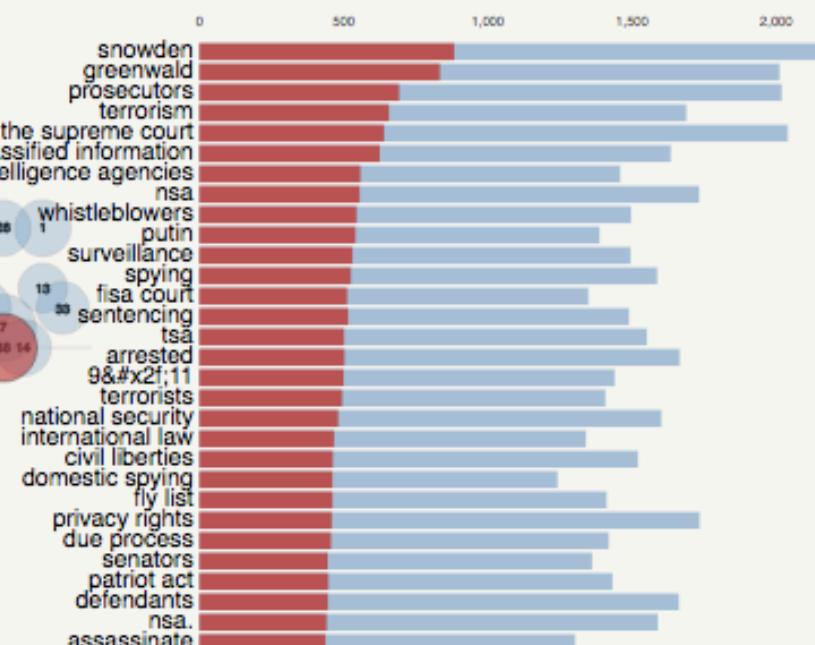
$\lambda = 1$



Intertopic Distance Map (via multidimensional scaling)



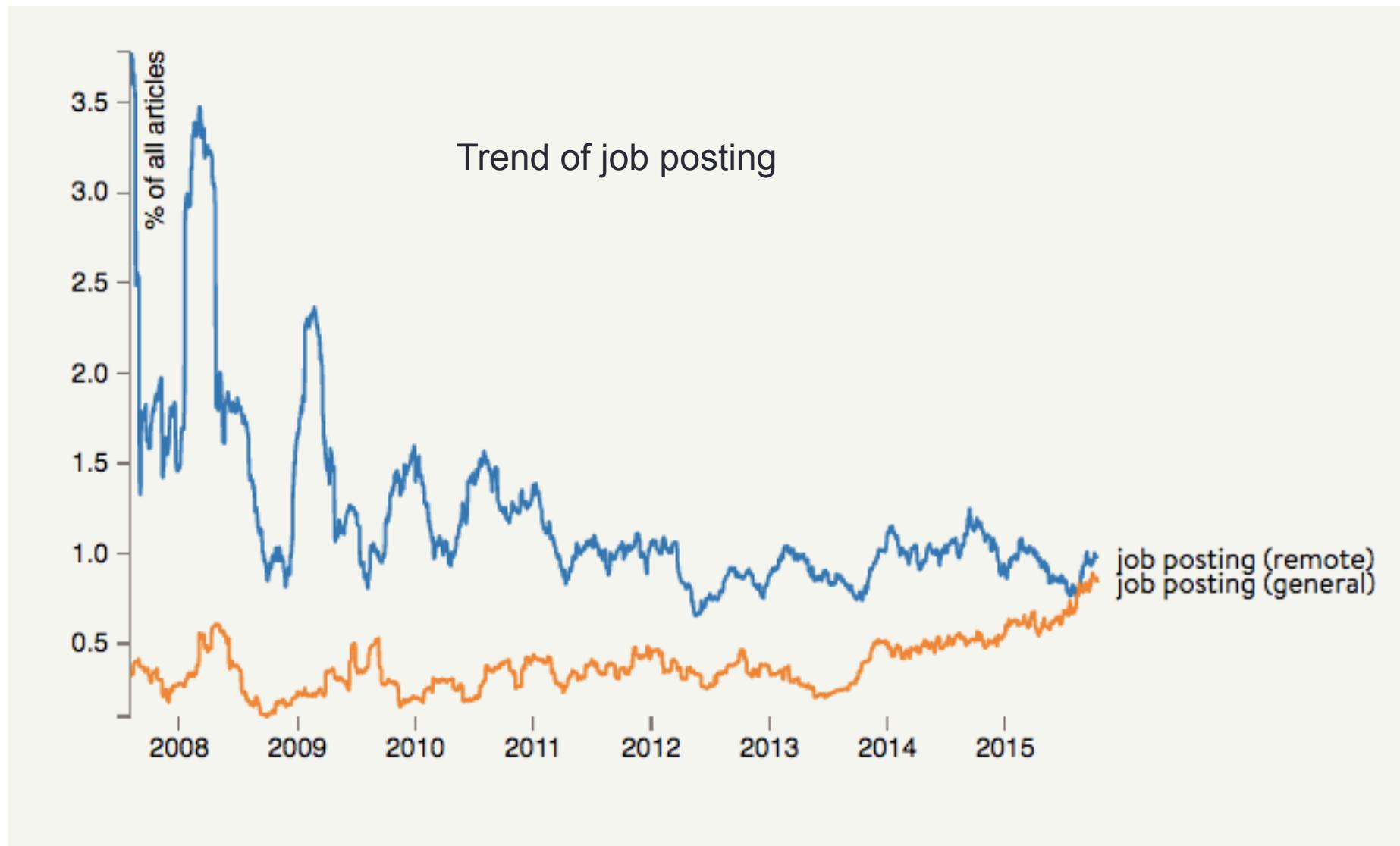
Top-30 Most Relevant Terms for Topic 38 (3.6% of tokens)



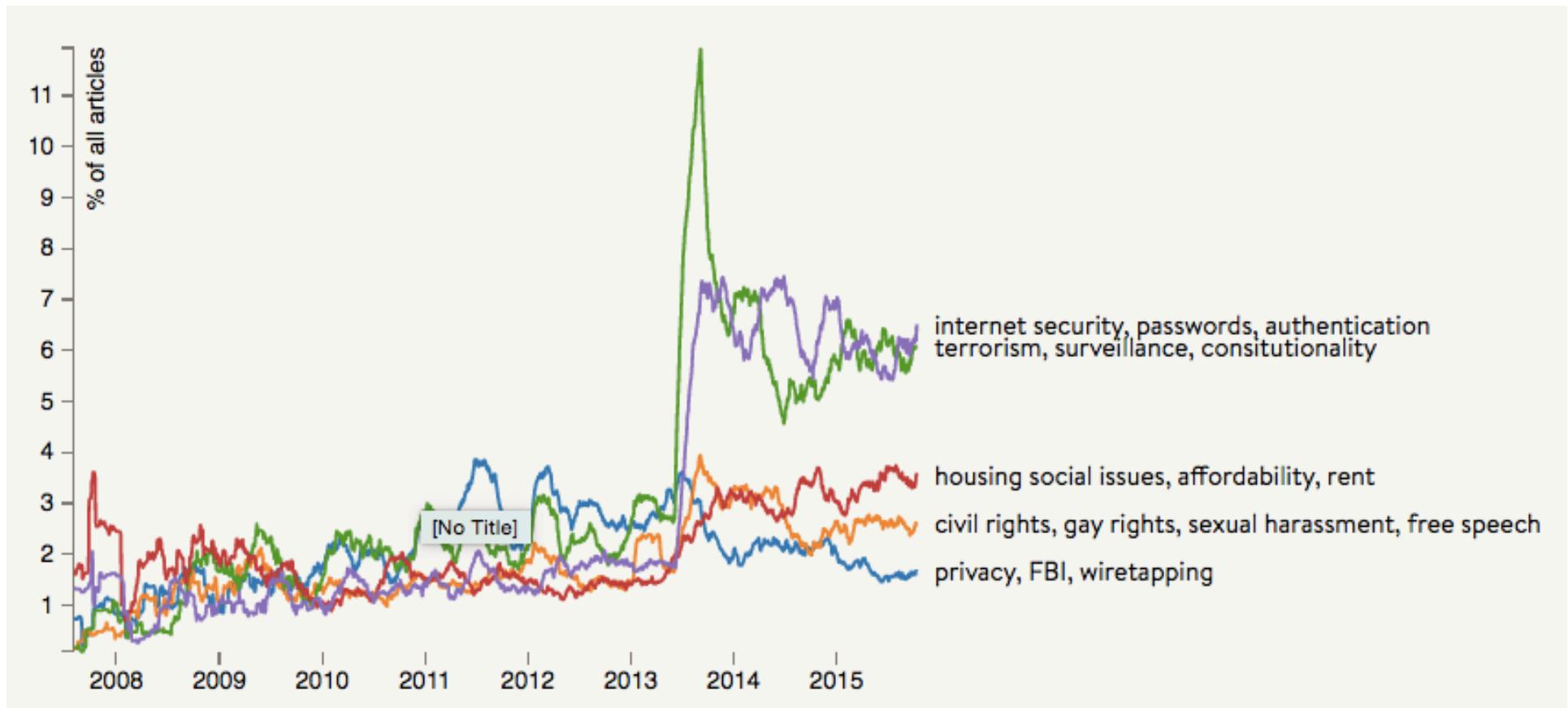
Overall term frequency
Estimated term frequency within the selected topic
1. saliency(term w) = frequency(w) * (sum_i p_i || w) * log(p_i || w) / p_i(||) for topics i; see Chuang et al. (2012)
(p(w); see Stevart & Shirley (2014))

<https://multithreaded.stitchfix.com/blog/2016/05/27/lda2vec/>

LDA2Vec example



Politics trend



LDA2Vec as a classification model

- Maybe not...
 - Model that offers human interpretability usually is worse
- But it's great for
 - Extract insights from trends by visualizations
 - Completely unsupervised model
 - Except for the need to specify number of topics and look through each topic words to determine the topic label
- This model gives a good example of
 - The design decisions for model building can be complex but logically motivated
 - Combining loss functions to get the desired behavior

Summary

- Text classification task
- Bag of words model
 - Naïve Bayes
 - Topic Models
 - Latent topic models (LDA)
 - LDA2vec
- Context vectors
 - More next time

