

# Deep latent models of word representation

Wilker Aziz  
University of Amsterdam

April 3, 2018

# Word representation

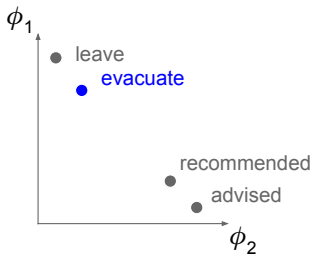
An abstraction that stands for the use of a word

*In the event of a chemical spill, most children know they should evacuate as advised by people in charge.*

# Word representation

An abstraction that stands for the use of a word

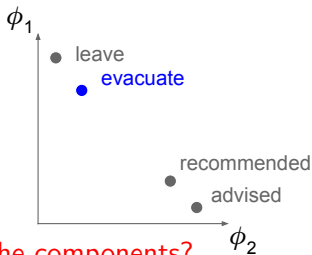
*In the event of a chemical spill, most children know they should  
**evacuate** as advised by people in charge.*



# Word representation

An abstraction that stands for the use of a word

*In the event of a chemical spill, most children know they should  
**evacuate** as advised by people in charge.*



How do we choose the components?

# Distributional Hypothesis

Context can represent the intended use of a word

*In the event of a chemical spill, most children know they should  
**evacuate** as advised by people in charge.*

- success hinges on **discriminative power** of available context

# Discriminative embedding models

*In the event of a chemical spill, most children know they should  
**evacuate** as advised by people in charge.*

Place words in  $\mathbb{R}^d$  as to answer questions like

*“Have I seen this word in this context?”*

# Discriminative embedding models

*In the event of a chemical spill, most children know they should  
evacuate as advised by people in charge.*

Place words in  $\mathbb{R}^d$  as to answer questions like

*"Have I seen this word in this context?"*

Fit a binary classifier

(Goldberg and Levy, 2014)

- positive examples
- negative examples

# Discriminative embedding models

*In the event of a chemical spill, most children know they should  
evacuate as advised by people in charge.*

Place words in  $\mathbb{R}^d$  as to answer questions like

*“Have I seen this word in this context?”*

Fit a binary classifier

(Goldberg and Levy, 2014)

- positive examples
- negative examples

“Far away” context expresses a form of negative correlation



# Discriminative embedding models

*In the event of a chemical spill, most children know they should  
evacuate as advised by people in charge.*

Place words in  $\mathbb{R}^d$  as to answer questions like

*“Have I seen this word in this context?”*

Fit a binary classifier

(Goldberg and Levy, 2014)

- positive examples
- negative examples

“Far away” context expresses a form of negative correlation

- But data only show positive examples

# Discriminative embedding models

*In the event of a chemical spill, most children know they should  
evacuate as advised by people in charge.*

*ambiguity*

Place words in  $\mathbb{R}^d$  as to answer questions like

*"Have I seen this word in this context?"*

Fit a binary classifier

(Goldberg and Levy, 2014)

- positive examples
- negative examples

"Far away" context expresses a form of negative correlation

- But data only show positive examples

# Limitations

Meaning representation is an **unsupervised** problem

- we cannot actually observe what we are trying to learn  
i.e. representations

# Limitations

Meaning representation is an **unsupervised** problem

- we cannot actually observe what we are trying to learn  
i.e. representations

Distributional hypothesis seems pretty strong

- but it fails when context is not sufficiently discriminative

# EMBEDALIGN

## Generative treatment

- model what we want to induce (i.e. representations)
- learn from **positive examples**
- learn from richer (less ambiguous) context

---

Rios et al. (2018)

# Outline

- 1 Embed-Align
- 2 Bayesian Skip-Gram
- 3 Practical

# Equivalence through translation

In the event of a chemical spill, most children know they should **evacuate** as advised by people in charge.

After ingestion, the substance speeds up the movement of your bowels encouraging you to **evacuate**.

# Equivalence through translation

In the event of a chemical spill, most children know they should **evacuate** as advised by people in charge.

After ingestion, the substance speeds up the movement of your bowels encouraging you to **evacuate**.

Em caso de vazamento químico, a maioria das crianças estão conscientes de que devem **deixar** o local, como sugerem autoridades.

Após ingerida, a substância acelera o movimento das paredes do intestino forçando o indivíduo a se **aliviar**.



# Equivalence through translation

In the event of a chemical spill, most children know they should **evacuate** as advised by people in charge.

Em caso de vazamento químico, a maioria das crianças estão conscientes de que devem **deixar** o local, como sugerem autoridades.

After ingestion, the substance speeds up the movement of your bowels encouraging you to **evacuate**.

Após ingerida, a substância acelera o movimento das paredes do intestino forçando o indivíduo a se **aliviar**.

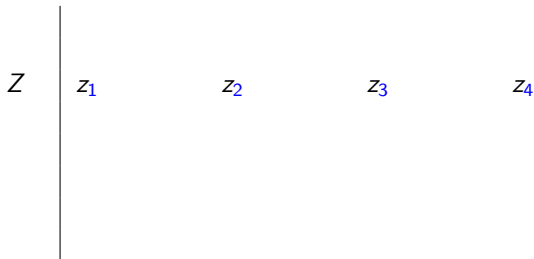
Observation from WSD community

- foreign text as **proxy to sense supervision**  
(Diab and Resnik, 2002)

quickly evacuate the area / deixe o local rapidamente

|

quickly evacuate the area / deixe o local rapidamente



quickly evacuate the area / deixe o local rapidamente

X	quickly <sub>1</sub>	evacuate <sub>2</sub>	the <sub>3</sub>	area <sub>4</sub>
	↑	↑	↑	↑
Z	z <sub>1</sub>	z <sub>2</sub>	z <sub>3</sub>	z <sub>4</sub>

quickly evacuate the area / deixe o local rapidamente

$X$	quickly <sub>1</sub>	evacuate <sub>2</sub>	the <sub>3</sub>	area <sub>4</sub>
	↑	↑	↑	↑
$Z$	$z_1$	$z_2$	$z_3$	$z_4$
$A$	$a_1 = 2$	$a_2 = 3$	$a_3 = 4$	$a_4 = 1$

quickly evacuate the area / deixe o local rapidamente

$X$	quickly <sub>1</sub>	evacuate <sub>2</sub>	the <sub>3</sub>	area <sub>4</sub>
	↑	↑	↑	↑
$Z$	$z_1$	$z_2$	$z_3$	$z_4$
$A$	$a_1 = 2$	$a_2 = 3$	$a_3 = 4$	$a_4 = 1$
$Z_a$	$z_2$	$z_3$	$z_4$	$z_1$

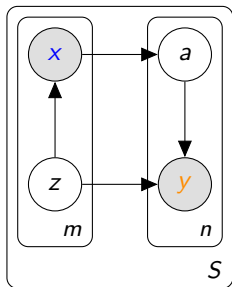
# Embed-Align

quickly evacuate the area / deixo o local rapidamente

X	quickly <sub>1</sub>	evacuate <sub>2</sub>	the <sub>3</sub>	area <sub>4</sub>
	↑	↑	↑	↑
Z	z <sub>1</sub>	z <sub>2</sub>	z <sub>3</sub>	z <sub>4</sub>
A	a <sub>1</sub> = 2	a <sub>2</sub> = 3	a <sub>3</sub> = 4	a <sub>4</sub> = 1
Z <sub>a</sub>	z <sub>2</sub>	z <sub>3</sub>	z <sub>4</sub>	z <sub>1</sub>
	↓	↓	↓	↓
Y	deixo <sub>1</sub>	o <sub>2</sub>	local <sub>3</sub>	rapidamente <sub>4</sub>

# Embed-Align

quickly evacuate the area / deixo o local rapidamente



$X$	quickly <sub>1</sub>	evacuate <sub>2</sub>	the <sub>3</sub>	area <sub>4</sub>
	↑	↑	↑	↑
$Z$	$z_1$	$z_2$	$z_3$	$z_4$
$A$	$a_1 = 2$	$a_2 = 3$	$a_3 = 4$	$a_4 = 1$
$Z_a$	$z_2$	$z_3$	$z_4$	$z_1$
	↓	↓	↓	↓
$Y$	deixo <sub>1</sub>	o <sub>2</sub>	local <sub>3</sub>	rapidamente <sub>4</sub>

Marginalising alignments collects additional training data for  $z$



# How does it disambiguate?

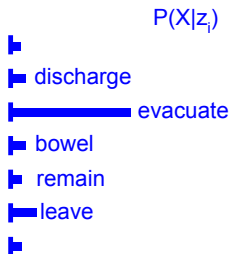
In the event of a chemical spill, most children know they should **evacuate**; as advised by people in charge.

# How does it disambiguate?

$z_i^{(1)}$

In the event of a chemical spill, most children know they should **evacuate** <sub>$i$</sub>  as advised by people in charge.

# How does it disambiguate?



$z_i^{(1)}$

In the event of a chemical spill, most children know they should **evacuate** <sub>$i$</sub>  as advised by people in charge.

# How does it disambiguate?

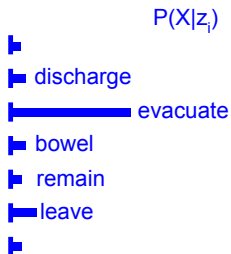


$z_i^{(1)}$

In the event of a chemical spill, most children know they should **evacuate**<sub>i</sub> as advised by people in charge.

After ingestion, the substance speeds up the movement of your bowels encouraging you to **evacuate**<sub>i</sub>.

# How does it disambiguate?



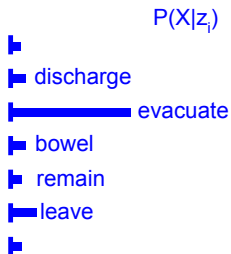
$z_i^{(1)}$

In the event of a chemical spill, most children know they should **evacuate**<sub>*i*</sub> as advised by people in charge.

$z_i^{(2)}$

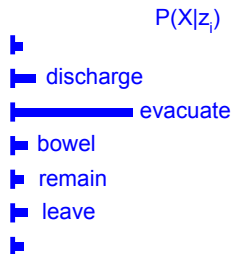
After ingestion, the substance speeds up the movement of your bowels encouraging you to **evacuate**<sub>*i*</sub>.

# How does it disambiguate?



$z_i^{(1)}$

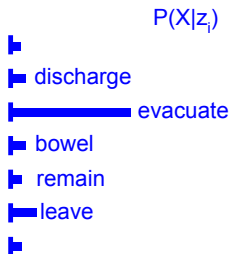
In the event of a chemical spill, most children know they should **evacuate**<sub>i</sub> as advised by people in charge.



$z_i^{(2)}$

After ingestion, the substance speeds up the movement of your bowels encouraging you to **evacuate**<sub>i</sub>.

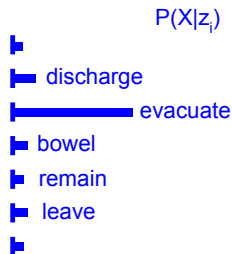
# How does it disambiguate?



$z_i^{(1)}$

In the event of a chemical spill, most children know they should **evacuate**<sub>i</sub> as advised by people in charge.

Em caso de vazamento químico, a maioria das crianças estão conscientes de que devem **deixar** o local, como sugerem autoridades.

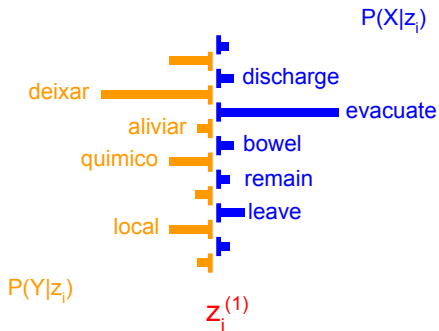


$z_i^{(2)}$

After ingestion, the substance speeds up the movement of your bowels encouraging you to **evacuate**<sub>i</sub>.

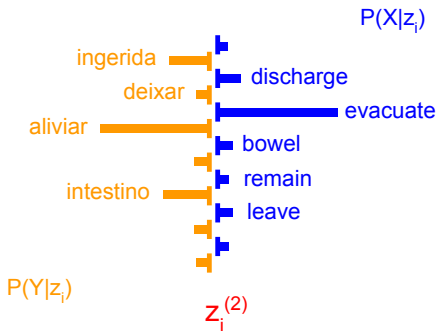
Após ingerida, a substância acelera o movimento das paredes do intestino forçando o indivíduo a se **aliviar**.

# How does it disambiguate?



In the event of a chemical spill, most children know they should **evacuate**<sub>i</sub> as advised by people in charge.

Em caso de vazamento químico, a maioria das crianças estão conscientes de que devem **deixar** o local, como sugerem autoridades.



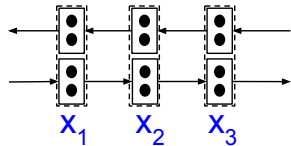
After ingestion, the substance speeds up the movement of your bowels encouraging you to **evacuate**<sub>i</sub>.

Após ingerida, a substância acelera o movimento das paredes do intestino forçando o indivíduo a se **aliviar**.



# Tractable inference

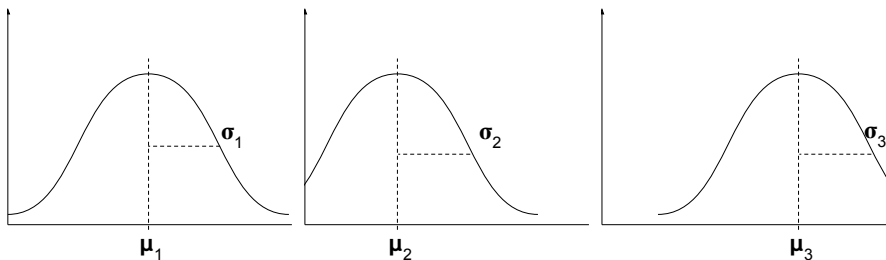
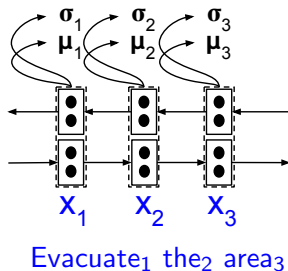
## ① Read sentence



Evacuate<sub>1</sub> the<sub>2</sub> area<sub>3</sub>

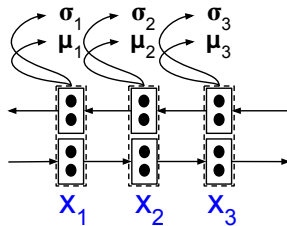
# Tractable inference

- 1 Read sentence
- 2 Predict posterior mean  $\mu_i$  and std  $\sigma_i$

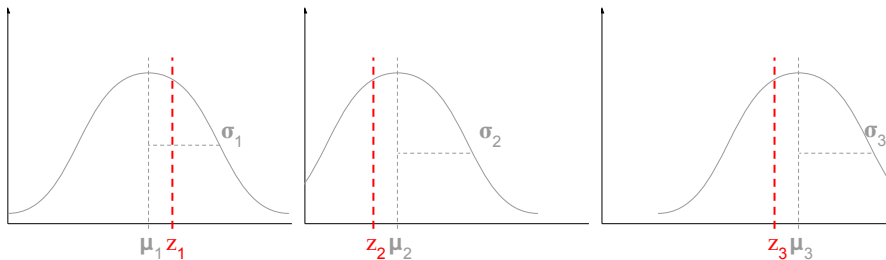


# Tractable inference

- 1 Read sentence
- 2 Predict posterior mean  $\mu_i$  and std  $\sigma_i$
- 3 Sample  $Z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$

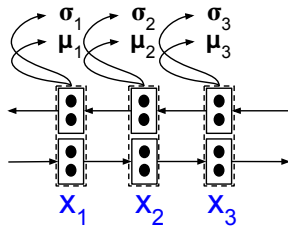


Evacuate<sub>1</sub> the<sub>2</sub> area<sub>3</sub>

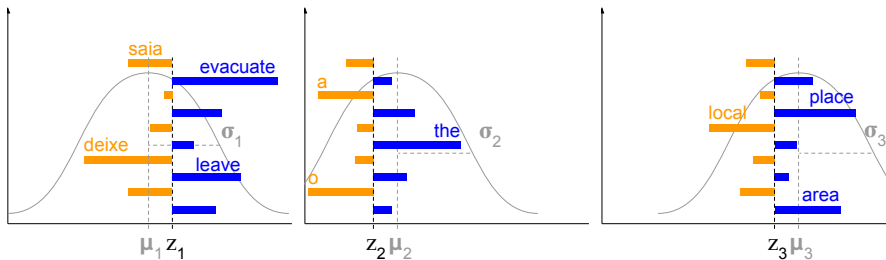


# Tractable inference

- 1 Read sentence
- 2 Predict posterior mean  $\mu_i$  and std  $\sigma_i$
- 3 Sample  $Z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$
- 4 Predict categorical distributions

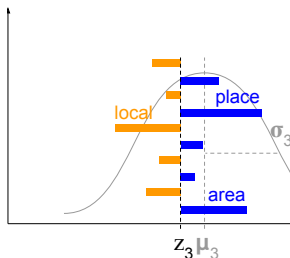
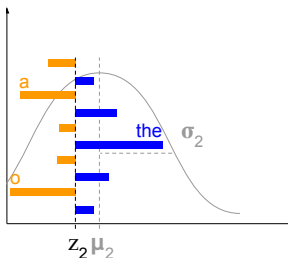
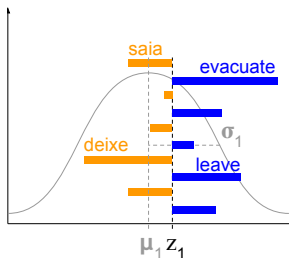
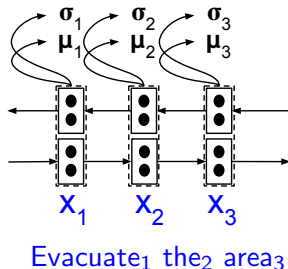


Evacuate<sub>1</sub> the<sub>2</sub> area<sub>3</sub>



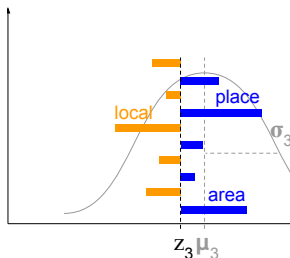
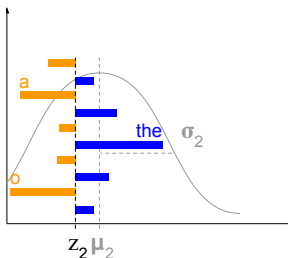
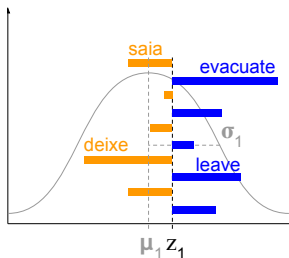
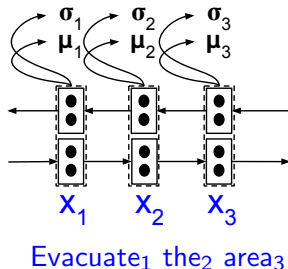
# Tractable inference

- 1 Read sentence
- 2 Predict posterior mean  $\mu_i$  and std  $\sigma_i$
- 3 Sample  $Z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$
- 4 Predict categorical distributions
- 5 Generate observations  
Evacuate<sub>1</sub> the<sub>2</sub> area<sub>3</sub> / Deix<sub>1</sub> o<sub>2</sub> local<sub>3</sub>



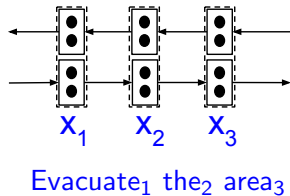
# Tractable inference

- 1 Read sentence
- 2 Predict posterior mean  $\mu_i$  and std  $\sigma_i$
- 3 Sample  $Z_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$
- 4 Predict categorical distributions
- 5 Generate observations  
Evacuate<sub>1</sub> the<sub>2</sub> area<sub>3</sub> / Deixe<sub>1</sub> o<sub>2</sub> local<sub>3</sub>
- 6 Maximise a lowerbound on likelihood  
(Kingma and Welling, 2014)



# What's special about it?

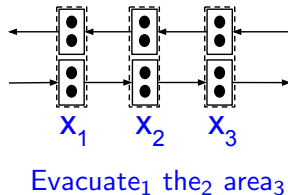
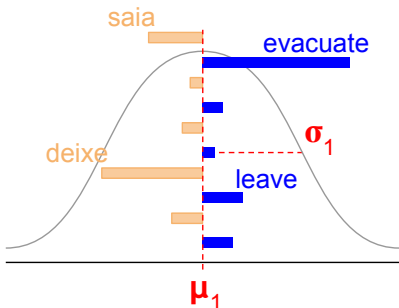
The model reads English text and



# What's special about it?

The model reads English text and

- predicts uncertainty
- describes “sense” using Portuguese words





# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$Z_i \sim \mathcal{N}(0, I)$$

# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$Z_i \sim \mathcal{N}(0, I)$$

$$X_i | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_{\theta}(z_i))$$

# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$Z_i \sim \mathcal{N}(0, I)$$

$$X_i | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_\theta(z_i))$$

$$A_j | m \sim \mathcal{U}(1/m)$$

# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$Z_i \sim \mathcal{N}(0, I)$$

$$X_i | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_\theta(z_i))$$

$$A_j | m \sim \mathcal{U}(1/m)$$

$$Y_j | z_1^m, a_j \sim \text{Cat}(\mathbf{g}_{a_j})$$

$$\mathbf{g}_{a_j} = \text{softmax}(\text{affine}_\theta(z_{a_j}))$$

Inference model: for  $i = 1, \dots, m$

# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$Z_i \sim \mathcal{N}(0, I)$$

$$X_i | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_\theta(z_i))$$

$$A_j | m \sim \mathcal{U}(1/m)$$

$$Y_j | z_1^m, a_j \sim \text{Cat}(\mathbf{g}_{a_j})$$

$$\mathbf{g}_{a_j} = \text{softmax}(\text{affine}_\theta(z_{a_j}))$$

Inference model: for  $i = 1, \dots, m$

$$Z_i | x_1^m \sim \mathcal{N}(\mathbf{u}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i))$$

# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$Z_i \sim \mathcal{N}(0, I)$$

$$X_i | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_\theta(z_i))$$

$$A_j | m \sim \mathcal{U}(1/m)$$

$$Y_j | z_1^m, a_j \sim \text{Cat}(\mathbf{g}_{a_j})$$

$$\mathbf{g}_{a_j} = \text{softmax}(\text{affine}_\theta(z_{a_j}))$$

Inference model: for  $i = 1, \dots, m$

$$Z_i | x_1^m \sim \mathcal{N}(\mathbf{u}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i))$$

$$\mathbf{h}_1^m = \text{enc}_\lambda(x_1^m)$$

# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$Z_i \sim \mathcal{N}(0, I)$$

$$X_i | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_\theta(z_i))$$

$$A_j | m \sim \mathcal{U}(1/m)$$

$$Y_j | z_1^m, a_j \sim \text{Cat}(\mathbf{g}_{a_j})$$

$$\mathbf{g}_{a_j} = \text{softmax}(\text{affine}_\theta(z_{a_j}))$$

Inference model: for  $i = 1, \dots, m$

$$Z_i | x_1^m \sim \mathcal{N}(\mathbf{u}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i))$$

$$\mathbf{h}_1^m = \text{enc}_\lambda(x_1^m)$$

$$\mathbf{u}_i = \text{affine}_\lambda(\text{relu}(\text{affine}_\lambda(\mathbf{h}_i)))$$



# Complete specification

Generative model: for  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$Z_i \sim \mathcal{N}(0, I)$$

$$X_i | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_\theta(z_i))$$

$$A_j | m \sim \mathcal{U}(1/m)$$

$$Y_j | z_1^m, a_j \sim \text{Cat}(\mathbf{g}_{a_j})$$

$$\mathbf{g}_{a_j} = \text{softmax}(\text{affine}_\theta(z_{a_j}))$$

Inference model: for  $i = 1, \dots, m$

$$Z_i | x_1^m \sim \mathcal{N}(\mathbf{u}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i))$$

$$\mathbf{h}_1^m = \text{enc}_\lambda(x_1^m)$$

$$\mathbf{u}_i = \text{affine}_\lambda(\text{relu}(\text{affine}_\lambda(\mathbf{h}_i)))$$

$$\mathbf{s}_i = \text{softplus}_\lambda(\text{relu}(\text{affine}_\lambda(\mathbf{h}_i)))$$

# ELBO

## ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z))$$

# ELBO

## ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z))$$

## KL term

$$\text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z)) = \underbrace{\sum_{i=1}^m \text{KL}(q_{\lambda}(z_i|x_1^m) \parallel p(z))}_{\text{mean field}}$$

# ELBO

## ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z))$$

## KL term

$$\begin{aligned} \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z)) &= \underbrace{\sum_{i=1}^m \text{KL}(q_{\lambda}(z_i|x_1^m) \parallel p(z))}_{\text{mean field}} \\ &= \sum_{i=1}^m \text{KL} \left( \underbrace{\mathcal{N}(\mathbf{u}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i))}_{\text{inference model}} \parallel \underbrace{\mathcal{N}(0, I)}_{\text{prior}} \right) \end{aligned}$$

# ELBO - L1 term

ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z))$$

# ELBO - L1 term

ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z))$$

Likelihood term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)]$$

# ELBO - L1 term

## ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z))$$

## Likelihood term

$$\begin{aligned} & \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] \\ &= \underbrace{\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m|z_1^m)] + \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)]}_{\text{conditional independence}} \end{aligned}$$

# ELBO - L1 term

## ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z))$$

## Likelihood term

$$\begin{aligned} & \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] \\ &= \underbrace{\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m|z_1^m)] + \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)]}_{\text{conditional independence}} \end{aligned}$$

## L1 term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m|z_1^m)] = \sum_{i=1}^m \mathbb{E}_{q_{\lambda}(z_i|x_1^m)} [\log P_{\theta}(x_i|z_i)]$$



# ELBO - L1 term

## ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p(z))$$

## Likelihood term

$$\begin{aligned} & \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m, y_1^n|z_1^m)] \\ &= \underbrace{\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m|z_1^m)] + \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)]}_{\text{conditional independence}} \end{aligned}$$

## L1 term

$$\begin{aligned} \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(x_1^m|z_1^m)] &= \sum_{i=1}^m \mathbb{E}_{q_{\lambda}(z_i|x_1^m)} [\log P_{\theta}(x_i|z_i)] \\ &= \sum_{i=1}^m \mathbb{E}_{q_{\lambda}(z_i|x_1^m)} [\log \text{Cat}(x_i|\mathbf{f}_i)] \end{aligned}$$

L2 term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)]$$

L2 term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] = \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{j=1}^n P_{\theta}(y_j|m, z_1^m) \right]$$

L2 term

$$\begin{aligned}\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{j=1}^n P_{\theta}(y_j|m, z_1^m) \right] \\ &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log P_{\theta}(y_j|m, z_1^m) \right]\end{aligned}$$

L2 term

$$\begin{aligned}\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{j=1}^n P_{\theta}(y_j|m, z_1^m) \right] \\ &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log P_{\theta}(y_j|m, z_1^m) \right] \\ &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P_{\theta}(y_j, a_j|m, z_1^m) \right]\end{aligned}$$

## L2 term

$$\begin{aligned}\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{j=1}^n P_{\theta}(y_j|m, z_1^m) \right] \\&= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log P_{\theta}(y_j|m, z_1^m) \right] \\&= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P_{\theta}(y_j, a_j|m, z_1^m) \right] \\&= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P(a_j|m) P_{\theta}(y_j|z_{a_j}) \right]\end{aligned}$$

L2 term

$$\begin{aligned}
\mathbb{E}_{q_{\lambda}(z_1^m | x_1^m)} [\log P_{\theta}(y_1^n | m, z_1^m)] &= \mathbb{E}_{q_{\lambda}(z_1^m | x_1^m)} \left[ \log \prod_{j=1}^n P_{\theta}(y_j | m, z_1^m) \right] \\
&= \mathbb{E}_{q_{\lambda}(z_1^m | x_1^m)} \left[ \sum_{j=1}^n \log P_{\theta}(y_j | m, z_1^m) \right] \\
&= \mathbb{E}_{q_{\lambda}(z_1^m | x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P_{\theta}(y_j, a_j | m, z_1^m) \right] \\
&= \mathbb{E}_{q_{\lambda}(z_1^m | x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P(a_j | m) P_{\theta}(y_j | z_{a_j}) \right] \\
&= \mathbb{E}_{q_{\lambda}(z_1^m | x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m \mathcal{U}(a_j | 1/m) \text{Cat}(y_j | \mathbf{g}_{a_j}) \right]
\end{aligned}$$

# ELBO - Lowerbound on L2 term

L2 term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] = \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P(a_j|m) P_{\theta}(y_j|z_{a_j}) \right]$$



# ELBO - Lowerbound on L2 term

L2 term

$$\begin{aligned}\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P(a_j|m) P_{\theta}(y_j|z_{a_j}) \right] \\ &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \mathbb{E}_{P(a_j|m)} [P_{\theta}(y_j|z_{a_j})] \right]\end{aligned}$$

# ELBO - Lowerbound on L2 term

L2 term

$$\begin{aligned}\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P(a_j|m) P_{\theta}(y_j|z_{a_j}) \right] \\ &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \mathbb{E}_{P(a_j|m)} [P_{\theta}(y_j|z_{a_j})] \right] \\ &\stackrel{\text{JL}}{\geq} \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \mathbb{E}_{P(a_j|m)} [\log P_{\theta}(y_j|z_{a_j})] \right]\end{aligned}$$

# ELBO - Lowerbound on L2 term

L2 term

$$\begin{aligned}\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P(a_j|m) P_{\theta}(y_j|z_{a_j}) \right] \\&= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \mathbb{E}_{P(a_j|m)} [P_{\theta}(y_j|z_{a_j})] \right] \\&\stackrel{\text{JL}}{\geq} \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \mathbb{E}_{P(a_j|m)} [\log P_{\theta}(y_j|z_{a_j})] \right] \\&= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \mathbb{E}_{\mathcal{U}(a_j|m)} [\log \text{Cat}(y_j|\mathbf{g}_{a_j})] \right]\end{aligned}$$

# ELBO - Lowerbound on L2 term

## L2 term

$$\begin{aligned}\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} [\log P_{\theta}(y_1^n|m, z_1^m)] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \sum_{a_j}^m P(a_j|m) P_{\theta}(y_j|z_{a_j}) \right] \\&= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \log \mathbb{E}_{P(a_j|m)} [P_{\theta}(y_j|z_{a_j})] \right] \\&\stackrel{\text{JL}}{\geq} \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \mathbb{E}_{P(a_j|m)} [\log P_{\theta}(y_j|z_{a_j})] \right] \\&= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{j=1}^n \mathbb{E}_{\mathcal{U}(a_j|m)} [\log \text{Cat}(y_j|\mathbf{g}_{a_j})] \right]\end{aligned}$$

Alignment model  $P(a_j|m) = \frac{1}{m}$  is not a function of  $\theta$

- we can make an MC estimate by sampling candidate alignments uniformly

# The softmax problem

Categorical parameters are expensive to compute

$$X \sim \text{Cat}(\mathbf{f})$$

$$\mathbf{f} = \text{softmax}(\hat{\mathbf{f}})$$

$$f_x = \frac{\exp(\hat{f}_x)}{\sum_{x' \in \mathcal{X}} \exp(\hat{f}_{x'})}$$

---

$v_1 = |\mathcal{X}|$  is the size of the vocabulary of  $L_1$ ;  $v_2 = |\mathcal{Y}|$  is the size of the vocabulary of  $L_2$

# The softmax problem

Categorical parameters are expensive to compute

$$X \sim \text{Cat}(\mathbf{f})$$

$$\mathbf{f} = \text{softmax}(\hat{\mathbf{f}})$$

$$f_x = \frac{\exp(\hat{f}_x)}{\sum_{x' \in \mathcal{X}} \exp(\hat{f}_{x'})}$$

- $\mathbf{f}_1^m$  requires normalising  $m$  distributions over the vocabulary of  $L_1$ , thus it takes time  $O(m \times v_1)$
- $\mathbf{g}_1^m$  requires normalising  $m$  distributions over the vocabulary of  $L_2$ , thus it takes time  $O(m \times v_2)$

---

$v_1 = |\mathcal{X}|$  is the size of the vocabulary of  $L_1$ ;  $v_2 = |\mathcal{Y}|$  is the size of the vocabulary of  $L_2$

# Efficient softmax

Logistic regression

$$P(X = x|z) = \frac{\exp(s(z, x))}{\sum_{x' \in \mathcal{X}} \exp(s(z, x'))}$$

# Efficient softmax

Logistic regression

$$P(X = x|z) = \frac{\exp(s(z, x))}{\sum_{x' \in \mathcal{X}} \exp(s(z, x'))}$$

Define

- a set  $\mathcal{C}(x)$  such that  $x \in \mathcal{C}$
- a set  $\mathcal{N}(x)$  such that  $\mathcal{C}(x) \cap \mathcal{N}(x) = \emptyset$



# Efficient softmax

Logistic regression

$$P(X = x|z) = \frac{\exp(s(z, x))}{\sum_{x' \in \mathcal{X}} \exp(s(z, x'))}$$

Define

- a set  $\mathcal{C}(x)$  such that  $x \in \mathcal{C}$
- a set  $\mathcal{N}(x)$  such that  $\mathcal{C}(x) \cap \mathcal{N}(x) = \emptyset$

Re-express normaliser for  $P(X = x|z)$

$$\sum_{x' \in \mathcal{X}} \exp(s(z, x')) = \sum_{x' \in \mathcal{C}(x)} \exp(s(z, x')) + \sum_{x' \in \mathcal{N}(x)} \kappa(x') \exp(s(z, x'))$$

- $\kappa(x') = \frac{1}{q(x')}$  and  $q(x')$  is an importance distribution

# Approximate $P_{x|z}$

Logistic regression

$$P_{\theta}(x|z) = \frac{\exp(s_{\theta}(z, x))}{\sum_{x' \in \mathcal{X}} \exp(s_{\theta}(z, x'))}$$

# Approximate $P_{x|z}$

Logistic regression

$$P_{\theta}(x|z) = \frac{\exp(s_{\theta}(z, x))}{\sum_{x' \in \mathcal{X}} \exp(s_{\theta}(z, x'))}$$

Build

- a set  $\mathcal{C}$  containing all  $L_1$  words in batch
- a set  $\mathcal{N}$  sampling uniformly without replacement from  $\mathcal{X} \setminus \mathcal{C}$

# Approximate $P_{x|z}$

Logistic regression

$$P_{\theta}(x|z) = \frac{\exp(s_{\theta}(z, x))}{\sum_{x' \in \mathcal{X}} \exp(s_{\theta}(z, x'))}$$

Build

- a set  $\mathcal{C}$  containing all  $L_1$  words in batch
- a set  $\mathcal{N}$  sampling uniformly without replacement from  $\mathcal{X} \setminus \mathcal{C}$

Approximate normaliser for  $P_{\theta}(x|z)$

$$\sum_{x' \in \mathcal{X}} \exp(s_{\theta}(z, x')) \approx \sum_{x' \in \mathcal{C}} \exp(s_{\theta}(z, x')) + \sum_{x' \in \mathcal{N}} \frac{|\mathcal{X} \setminus \mathcal{C}|}{|\mathcal{N}|} \exp(s_{\theta}(z, x'))$$

- $s_{\theta}(z, x) = z^{\top} \mathbf{c}_x + b_x$   
 $\mathbf{c}_x = \text{lookup}_{\theta}(x)$  is a deterministic embedding,  $b_x$  a bias term

# Approximate $P_{Y|Z}$

Logistic regression

$$P_{\theta}(y|z) = \frac{\exp(u_{\theta}(z, y))}{\sum_{y' \in \mathcal{Y}} \exp(u_{\theta}(z, y'))}$$

# Approximate $P_{Y|Z}$

Logistic regression

$$P_{\theta}(y|z) = \frac{\exp(u_{\theta}(z, y))}{\sum_{y' \in \mathcal{Y}} \exp(u_{\theta}(z, y'))}$$

Build

- a set  $\mathcal{C}$  containing all  $L_2$  words in batch
- a set  $\mathcal{N}$  sampling uniformly without replacement from  $\mathcal{Y} \setminus \mathcal{C}$

# Approximate $P_{Y|Z}$

Logistic regression

$$P_{\theta}(y|z) = \frac{\exp(u_{\theta}(z, y))}{\sum_{y' \in \mathcal{Y}} \exp(u_{\theta}(z, y'))}$$

Build

- a set  $\mathcal{C}$  containing all  $L_2$  words in batch
- a set  $\mathcal{N}$  sampling uniformly without replacement from  $\mathcal{Y} \setminus \mathcal{C}$

Approximate normaliser for  $P_{\theta}(y|z)$

$$\sum_{y' \in \mathcal{Y}} \exp(s_{\theta}(z, y')) \approx \sum_{y' \in \mathcal{C}} \exp(s_{\theta}(z, y')) + \sum_{y' \in \mathcal{N}} \frac{|\mathcal{Y} \setminus \mathcal{C}|}{|\mathcal{N}|} \exp(s_{\theta}(z, y'))$$

- $s_{\theta}(z, y) = z^{\top} \mathbf{c}_y + b_y$   
 $\mathbf{c}_y = \text{lookup}_{\theta}(y)$  is a deterministic embedding,  $b_y$  a bias term

## Complementary Sum Sampling (CSS) - Summary

The approximation effectively reduces the size of the support of the categorical variable

- in each batch, the support is made of the word types in the batch
- along with a random subset of “negative words”
- this is similar to “negative sampling” but improves on asymptotic behaviour
- it only affects the softmax: the model remains generative



# Outline

- 1 Embed-Align
- 2 Bayesian Skip-Gram
- 3 Practical

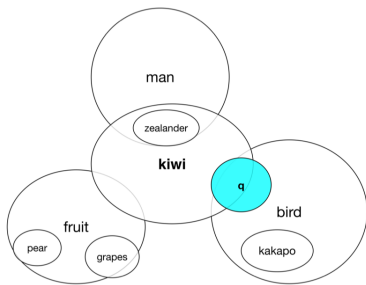


Figure 1: An idealized illustration of density embeddings. Unshaded ellipsoids encode prior densities of Gaussians. The shaded ellipsoid corresponds to the posterior for the word ‘kiwi’ when it appears in a context indicating that ‘kiwi’ refers to a bird.

*“Representing a word as a distribution provides many potential benefits. For example, such embeddings let us encode generality of terms (e.g., ‘kakapo’ is a type of ‘bird’), characterize uncertainty about semantic properties of the corresponding referent (e.g., a proper noun, such as ‘John’, encodes little about the person it refers to) or represent polysemy (e.g., ‘kiwi’ may refer to a fruit, a bird or a New Zealander).”*

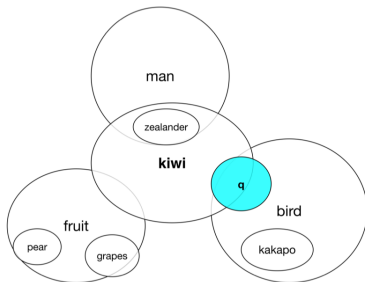


Figure 1: An idealized illustration of density embeddings. Unshaded ellipsoids encode prior densities of Gaussians. The shaded ellipsoid corresponds to the posterior for the word ‘kiwi’ when it appears in a context indicating that ‘kiwi’ refers to a bird.

*“In principle, using densities to represent words provides a natural way of encoding entailment: the decision regarding entailment relation can be made by testing the level sets of the distributions for ‘soft inclusion’. For example, in Figure 1, the ellipse for ‘kakapo’ lies within the ellipse for ‘bird’. ”*

# Complete specification

Generative model: for  $i = 1, \dots, m$

$$Z_i | x_i \sim \mathcal{N}(\mu_{x_i}, \text{diag}(\sigma_{x_i} \odot \sigma_{x_i}))$$

$$\mu_{x_i} = \text{lookup}_{\theta}(x_i)$$

$$\sigma_{x_i} = \text{softplus}(\text{lookup}_{\theta}(x_i))$$

# Complete specification

Generative model: for  $i = 1, \dots, m$

$$Z_i | x_i \sim \mathcal{N}(\mu_{x_i}, \text{diag}(\sigma_{x_i} \odot \sigma_{x_i}))$$

$$\mu_{x_i} = \text{lookup}_{\theta}(x_i)$$

$$\sigma_{x_i} = \text{softplus}(\text{lookup}_{\theta}(x_i))$$

for  $k \in \mathcal{K}_i = \underbrace{\{i - n, \dots, i - 1, i + 1, \dots, i + n\}}_{n \text{ words on each side of } x_i}$

# Complete specification

Generative model: for  $i = 1, \dots, m$

$$Z_i | x_i \sim \mathcal{N}(\mu_{x_i}, \text{diag}(\sigma_{x_i} \odot \sigma_{x_i}))$$

$$\mu_{x_i} = \text{lookup}_{\theta}(x_i)$$

$$\sigma_{x_i} = \text{softplus}(\text{lookup}_{\theta}(x_i))$$

$$\text{for } k \in \mathcal{K}_i = \underbrace{\{i - n, \dots, i - 1, i + 1, \dots, i + n\}}_{n \text{ words on each side of } x_i}$$

$$X_k | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_{\theta}(z_i))$$

# Complete specification

Generative model: for  $i = 1, \dots, m$

$$Z_i | x_i \sim \mathcal{N}(\mu_{x_i}, \text{diag}(\sigma_{x_i} \odot \sigma_{x_i}))$$

$$\mu_{x_i} = \text{lookup}_{\theta}(x_i)$$

$$\sigma_{x_i} = \text{softplus}(\text{lookup}_{\theta}(x_i))$$

$$\text{for } k \in \mathcal{K}_i = \underbrace{\{i - n, \dots, i - 1, i + 1, \dots, i + n\}}_{n \text{ words on each side of } x_i}$$

$$X_k | z_i \sim \text{Cat}(\mathbf{f}_i)$$

$$\mathbf{f}_i = \text{softmax}(\text{affine}_{\theta}(z_i))$$

Inference model

$$Z_i | x_{i-n}^{i+n} \sim \mathcal{N}(\mathbf{u}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i))$$

$$\mathbf{h}_i = \sum_{k \in \mathcal{K}_i} \text{relu}(\text{affine}_{\lambda}([x_k, x_i]))$$

$$\mathbf{u}_i = \text{affine}_{\lambda}(h_i)$$

$$\mathbf{s}_i = \text{softplus}(\text{affine}_{\lambda}(h_i))$$

## ELBO

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P(x_k|z_i) \right] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) || p_{\theta}(z_1^m|x_1^m))$$

- observation model is trained discriminatively  
latent variables generate overlapping subsets of observations



## ELBO - KL term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) || p_{\theta}(z_1^m|x_1^m))$$

KL term

$$\text{KL}(q_{\lambda}(z_1^m|x_1^m) || p_{\theta}(z_1^m|x_1^m))$$

---

Empirical Bayes: point estimate prior parameters

## ELBO - KL term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p_{\theta}(z_1^m|x_1^m))$$

KL term

$$\begin{aligned} & \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p_{\theta}(z_1^m|x_1^m)) \\ &= \sum_{i=1}^m \text{KL}\left(q_{\lambda}(z_i|x_{i-n}^{i+n}) \parallel p_{\theta}(z_i|x_i)\right) \end{aligned}$$

---

Empirical Bayes: point estimate prior parameters

## ELBO - KL term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p_{\theta}(z_1^m|x_1^m))$$

KL term

$$\begin{aligned} & \text{KL}(q_{\lambda}(z_1^m|x_1^m) \parallel p_{\theta}(z_1^m|x_1^m)) \\ &= \sum_{i=1}^m \text{KL}\left(q_{\lambda}(z_i|x_{i-n}^{i+n}) \parallel p_{\theta}(z_i|x_i)\right) \\ &= \sum_{i=1}^m \text{KL}\left(\underbrace{\mathcal{N}(\mathbf{u}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i))}_{\text{inference model}} \parallel \underbrace{\mathcal{N}(\boldsymbol{\mu}_{x_i}, \text{diag}(\boldsymbol{\sigma}_{x_i} \odot \boldsymbol{\sigma}_{x_i}))}_{\text{prior}}\right) \end{aligned}$$

---

Empirical Bayes: point estimate prior parameters

## ELBO - Likelihood term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) || p_{\theta}(z_1^m|x_1^m))$$

## ELBO - Likelihood term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) || p_{\theta}(z_1^m|x_1^m))$$

Likelihood term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right]$$

## ELBO - Likelihood term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] - \text{KL} (q_{\lambda}(z_1^m|x_1^m) || p_{\theta}(z_1^m|x_1^m))$$

Likelihood term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] = \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{i=1}^m \sum_{k \in \mathcal{K}_i} \log P_{\theta}(x_k|z_i) \right]$$

## ELBO - Likelihood term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) || p_{\theta}(z_1^m|x_1^m))$$

Likelihood term

$$\begin{aligned} \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{i=1}^m \sum_{k \in \mathcal{K}_i} \log P_{\theta}(x_k|z_i) \right] \\ &= \sum_{i=1}^m \sum_{k \in \mathcal{K}_i} \mathbb{E}_{q_{\lambda}(z_i|x_i^{i+n})} [\log P_{\theta}(x_k|z_i)] \end{aligned}$$

## ELBO - Likelihood term

$$\mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] - \text{KL}(q_{\lambda}(z_1^m|x_1^m) || p_{\theta}(z_1^m|x_1^m))$$

Likelihood term

$$\begin{aligned} \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \log \prod_{i=1}^m \prod_{k \in \mathcal{K}_i} P_{\theta}(x_k|z_i) \right] &= \mathbb{E}_{q_{\lambda}(z_1^m|x_1^m)} \left[ \sum_{i=1}^m \sum_{k \in \mathcal{K}_i} \log P_{\theta}(x_k|z_i) \right] \\ &= \sum_{i=1}^m \sum_{k \in \mathcal{K}_i} \mathbb{E}_{q_{\lambda}(z_i|x_{i-n}^{i+n})} [\log P_{\theta}(x_k|z_i)] \\ &= \sum_{i=1}^m \sum_{k \in \mathcal{K}_i} \mathbb{E}_{q_{\lambda}(z_i|x_{i-n}^{i+n})} [\log \text{Cat}(x_k|\mathbf{f}_i)] \end{aligned}$$



# The softmax problem

To circumvent an expensive softmax, change the likelihood term

$$P_{\theta}(x|z) = \frac{u_{\theta}(z, x)}{\sum_{x' \in \mathcal{X}} u_{\theta}(z, x')} \quad \text{with } u_{\theta}(\cdot, \cdot) > 0$$

# The softmax problem

To circumvent an expensive softmax, change the likelihood term

$$P_{\theta}(x|z) = \frac{u_{\theta}(z, x)}{\sum_{x' \in \mathcal{X}} u_{\theta}(z, x')} \quad \text{with } u_{\theta}(\cdot, \cdot) > 0$$

and re-write the likelihood term

$$\mathbb{E}_{q_{\lambda}(z)} [\log P_{\theta}(x|z)] = \mathbb{E}_{q_{\lambda}(z)} \left[ \log u_{\theta}(z, x) - \log \sum_{x' \in \mathcal{X}} u_{\theta}(z, x') \right]$$

# The softmax problem

To circumvent an expensive softmax, change the likelihood term

$$P_{\theta}(x|z) = \frac{u_{\theta}(z, x)}{\sum_{x' \in \mathcal{X}} u_{\theta}(z, x')} \quad \text{with } u_{\theta}(\cdot, \cdot) > 0$$

and re-write the likelihood term

$$\begin{aligned} \mathbb{E}_{q_{\lambda}(z)} [\log P_{\theta}(x|z)] &= \mathbb{E}_{q_{\lambda}(z)} \left[ \log u_{\theta}(z, x) - \log \sum_{x' \in \mathcal{X}} u_{\theta}(z, x') \right] \\ &= \mathbb{E}_{q_{\lambda}(z)} [\log u_{\theta}(z, x)] - \mathbb{E}_{q_{\lambda}(z)} \left[ \log \sum_{x' \in \mathcal{X}} u_{\theta}(z, x') \right] \end{aligned}$$

# Lowerbound on likelihood term

Then (by design) let

$$u_{\theta}(z, x) = \underbrace{P(x)}_{\text{fixed}} \mathcal{N}(z | \boldsymbol{\mu}_x, \text{diag}(\boldsymbol{\sigma}_x \odot \boldsymbol{\sigma}_x))$$

# Lowerbound on likelihood term

Then (by design) let

$$u_{\theta}(z, x) = \underbrace{P(x)}_{\text{fixed}} \mathcal{N}(z | \boldsymbol{\mu}_x, \text{diag}(\boldsymbol{\sigma}_x \odot \boldsymbol{\sigma}_x))$$

And bound  $\mathbb{E}_{q_{\lambda}(z)} \left[ \log \sum_{x' \in \mathcal{X}} u_{\theta}(z, x') \right]$

# Lowerbound on likelihood term

Then (by design) let

$$u_{\theta}(z, x) = \underbrace{P(x)}_{\text{fixed}} \mathcal{N}(z | \boldsymbol{\mu}_x, \text{diag}(\boldsymbol{\sigma}_x \odot \boldsymbol{\sigma}_x))$$

And bound  $\mathbb{E}_{q_{\lambda}(z)} [\log \sum_{x' \in \mathcal{X}} u_{\theta}(z, x')]$

$$= \mathbb{E}_{q_{\lambda}(z)} \left[ \log \sum_{x' \in \mathcal{X}} P(x) \mathcal{N}(z | \boldsymbol{\mu}_x, \text{diag}(\boldsymbol{\sigma}_x \odot \boldsymbol{\sigma}_x)) \right]$$

# Lowerbound on likelihood term

Then (by design) let

$$u_{\theta}(z, x) = \underbrace{P(x)}_{\text{fixed}} \mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x))$$

And bound  $\mathbb{E}_{q_{\lambda}(z)} [\log \sum_{x' \in \mathcal{X}} u_{\theta}(z, x')]$

$$\begin{aligned} &= \mathbb{E}_{q_{\lambda}(z)} \left[ \log \sum_{x' \in \mathcal{X}} P(x) \mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x)) \right] \\ &= \mathbb{E}_{q_{\lambda}(z)} [\log \mathbb{E}_{P(x)} [\mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x))]] \end{aligned}$$

# Lowerbound on likelihood term

Then (by design) let

$$u_{\theta}(z, x) = \underbrace{P(x)}_{\text{fixed}} \mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x))$$

And bound  $\mathbb{E}_{q_{\lambda}(z)} [\log \sum_{x' \in \mathcal{X}} u_{\theta}(z, x')]$

$$\begin{aligned} &= \mathbb{E}_{q_{\lambda}(z)} \left[ \log \sum_{x' \in \mathcal{X}} P(x) \mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x)) \right] \\ &= \mathbb{E}_{q_{\lambda}(z)} [\log \mathbb{E}_{P(x)} [\mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x))]] \\ &\stackrel{\text{J1}}{\geq} \mathbb{E}_{q_{\lambda}(z)} [\mathbb{E}_{P(x)} [\log \mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x))]] \end{aligned}$$



# Lowerbound on likelihood term

Then (by design) let

$$u_{\theta}(z, x) = \underbrace{P(x)}_{\text{fixed}} \mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x))$$

And bound  $\mathbb{E}_{q_{\lambda}(z)} [\log \sum_{x' \in \mathcal{X}} u_{\theta}(z, x')]$

$$\begin{aligned} &= \mathbb{E}_{q_{\lambda}(z)} \left[ \log \sum_{x' \in \mathcal{X}} P(x) \mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x)) \right] \\ &= \mathbb{E}_{q_{\lambda}(z)} [\log \mathbb{E}_{P(x)} [\mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x))]] \\ &\stackrel{\text{JL}}{\geq} \mathbb{E}_{q_{\lambda}(z)} [\mathbb{E}_{P(x)} [\log \mathcal{N}(z | \mu_x, \text{diag}(\sigma_x \odot \sigma_x))]] \end{aligned}$$

- $P(x)$  does not depend on  $\theta$ , we can compute an MC estimate e.g. empirical (unigram) distribution

# Outline

- 1 Embed-Align
- 2 Bayesian Skip-Gram
- 3 Practical**

# Practical

- Skip-gram (Mikolov et al., 2013)
- Bayesian skip-gram (Bražiņskas et al., 2017)
- Embed-Align (Rios et al., 2018)

# Comparison

	SkipGram	Bayesian SkipGram	EmbedAlign
LVM		✓	✓
Generative training			✓
Prior		type-specific Gaussian	$\mathcal{N}(0, I)$
Inference model		FFNN	BiLSTM
Softmax	negative sampling	JI	CSS

# Literature I

Aleksandar Botev, Bowen Zheng, and David Barber. Complementary sum sampling for likelihood approximation in large scale classification. In *Artificial Intelligence and Statistics*, pages 1030–1038, 2017.

Arthur Bražinskas, Serhii Havrylov, and Ivan Titov. Embedding words as distributions with a bayesian skip-gram model. *arXiv preprint arXiv:1711.11027*, 2017.

Yoav Goldberg and Omer Levy. word2vec explained: Deriving mikolov et al.’s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722*, 2014.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.

## Literature II

Miguel Rios, Wilker Aziz, and Khalil Sima'an. Deep generative model for joint alignment and word representation. *arXiv preprint arXiv:1802.05883*, 2018.