



# Discovering Geographical Topics in Twitter

Liangjie Hong, Lehigh University

**Amr Ahmed, Yahoo! Research**

Alexander J. Smola, Yahoo! Research

Siva Gurumurthy, Twitter

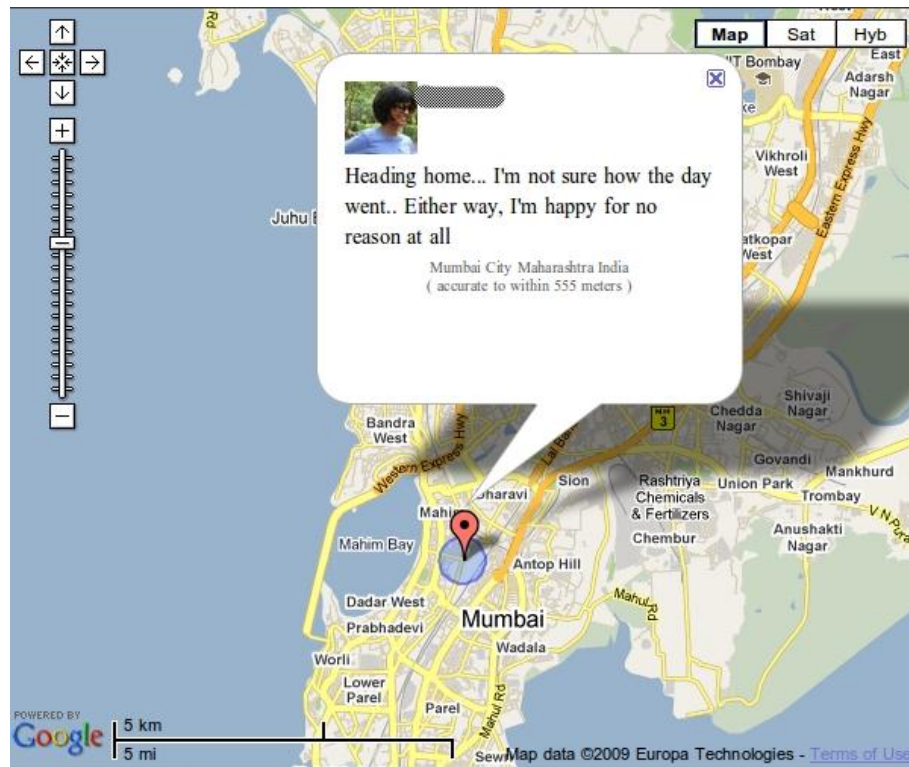
Kostas Tsioutsoulouklis, Twitter

# Overview

- Motivations
- Our Proposed Model
- Experiments
- Conclusions

# Motivations

## Twitter messages + Locations



# Motivations

## We want to know...

- How is information created and shared in different geographic locations? What is the inherent geographic variability of content?
- What are the spatial and linguistic characteristics of people? How does this vary across regions?
- Can we discover patterns in users' usage of micro-blogging services?
- **Can we predict user location from tweets?**

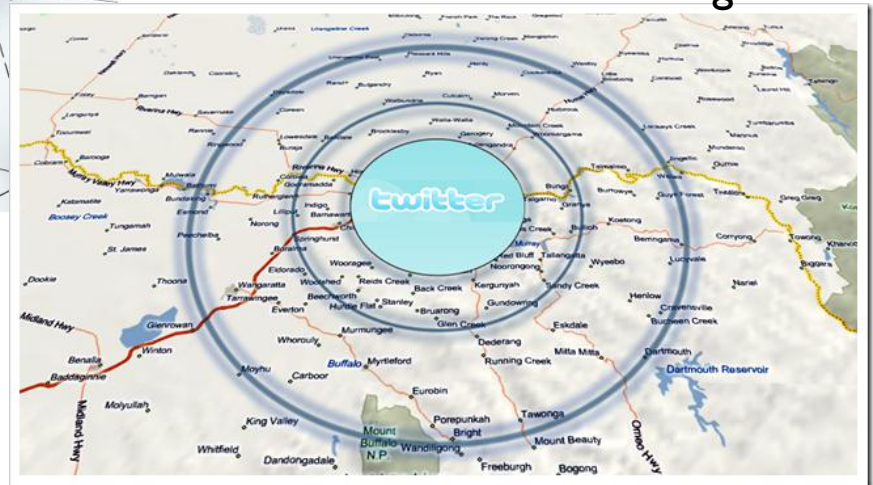
# Motivations

## Applications

Behavioral targeting and user modeling



Better local information filtering



# Motivations

## Challenges

- Tweets
  - noisy and short (140 characters)
- Only 1% of tweets geo-tagged
  - Can we predict locations for non-tagged tweets?
- Many intuitions to be combined
  - Background, regional language models, topics
  - Personal preferences, regional preferences...
- ...



Can we really infer locations for a tweet?

Yes via tweet decomposition



Just landed after a long flight. It is raining  
here at Lyon though!

What is the user's location?



Just landed after a long flight. It is raining  
here at Lyon though!

**background**

just  
after  
It  
be  
the  
can  
cant  
will

Just landed after a long flight. It is raining here at Lyon though!

**Travel**

landed  
flight  
delay  
TSE  
Gate  
terminal

**background**

just  
after  
It  
be  
the  
can  
cant  
will

Just landed after a long flight. It is raining here at Lyon though!

**Travel/airport**

landed  
flight  
delay  
TSE  
Gate  
terminal

**background**

just  
after  
It  
be  
the  
can  
cant  
will

**SE airport area**

Lyon  
Saint  
Exupery  
convention  
center  
raining

Semantic  
Topic

**Travel/airport**

landed  
flight  
delay  
TSE  
Gate  
terminal

Background  
Language  
Model

**background**

just  
after  
It  
be  
the  
can  
cant  
will

Regional  
Language  
Model

**SE airport area**

Lyon  
Saint  
Exupery  
convention  
center  
raining

Delayed again at the TSE check point and  
might miss my flight. way to go SF!

**Travel/airport**

landed  
flight  
delay  
TSE  
Gate  
terminal

**background**

just  
after  
It  
be  
the  
can  
cant  
will

**SFO**

SF  
SFO  
San  
Fransisco  
airport



Can we always do that?



Life is good! Feeling great today!

Life **is** good! Feeling great today!

**Daily life**

life  
feeling  
good  
today  
morning

**background**

just  
after  
It  
be  
the  
can  
cant  
will

?



Life is good! Feeling great today!

If we know something extra about the context and **user location preferences**, perhaps we can do better than random guessing!

# Motivations

## Previous work

- Simple regional language models
  - No factorization
- No personal preferences
- Complicated inference algorithms
  - Usually two step process
  - Fails to learn coherent regions

# Overview

- Motivations
- **Our Proposed Model**
- Experiments
- Conclusions

# Our Proposed Model

- A novel probabilistic model considers
  - Regional language models
  - Global topics
  - Personal preferences
- Sparse modeling + Bayesian treatment
- An efficient inference algorithm

# The Model

- Basic Intuition
  - Regions
  - Topics
  - Users
  - Tweets
- The generative process
  - Intuition
  - Glory details

# Basic Intuition: Region

- Must be coherent
  - There is **enough traffic** in it
  - Affects the way we write tweets
    - Has preference over **what topic discussed**
    - **Specific keywords**
  - Area over the map
  - Example
    - An airport
    - A park
    - A mall
    - A city

# Basic Intuition: Topic

- Classify the content of the tweet
- Might not tell us the location
- Puts a distribution over words
- Examples
  - Sports
  - Politics
  - Travel
  - Daily life, etc

# Basic Intuition: User

- Has preferences over locations
  - Where he usually spends his/her time
- Has preference over topics
  - What he tweets about



# Basic Intuition: Tweet

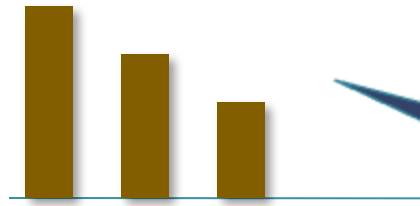
- Written by a **given user**
- At a specific **location** (region)
  - Depends on the user
- About a **specific topic**
  - Depends on
    - What the user talks about
    - What is being discussed at this location
- Composed of a **bag of words** from
  - Topic + location + background language models

# The Model

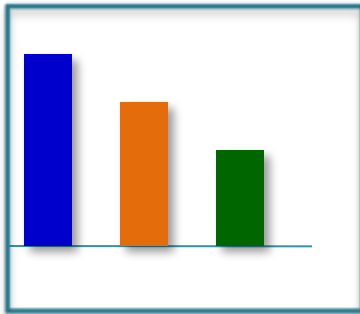
- Basic Intuition
  - Regions
  - Topics
  - Users
  - Tweets
- The generative process
  - Intuitive explanation
  - Glory details

# How a tweet is being generated?

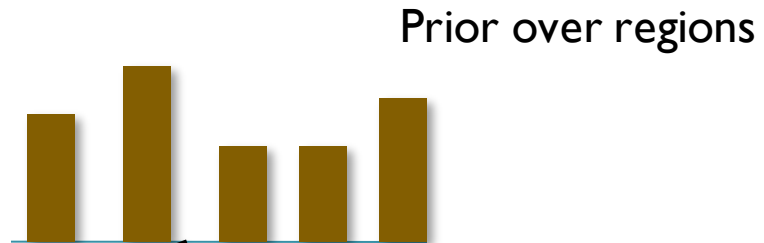
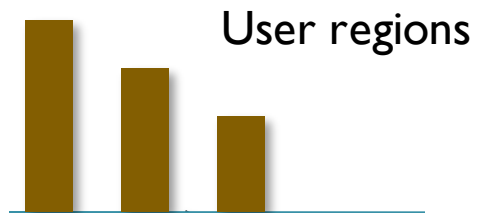
- Pick a location
- Pick a topic
- Generate the words



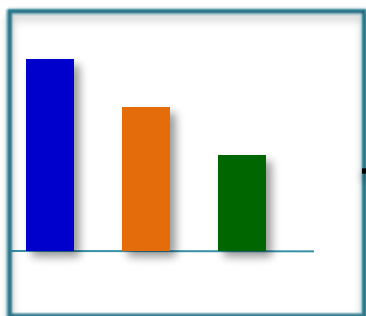
Preferences over regions  
Regions are unsupervised  
Just an area over the map



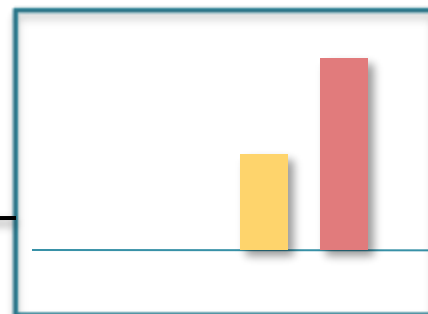
Preference over topics:  
What he likes to talk about



User topics



Region topics



Pick a region

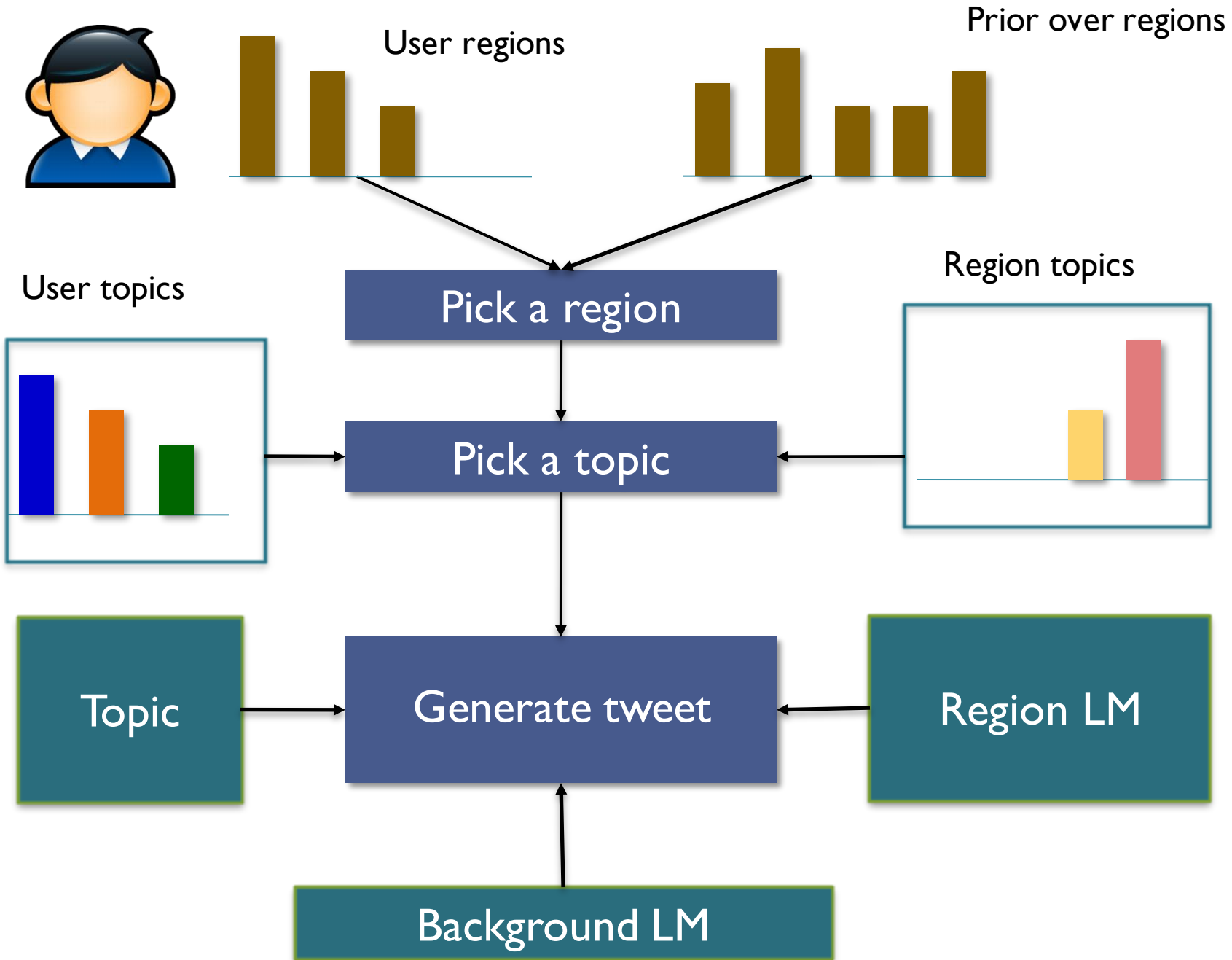
Pick a topic

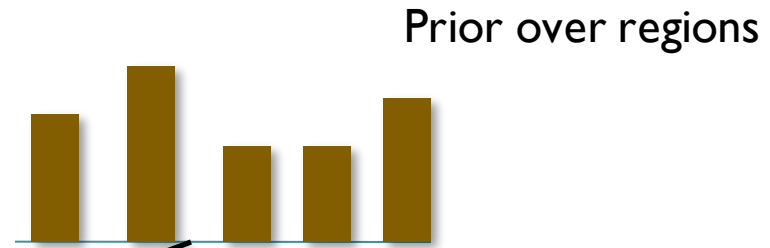
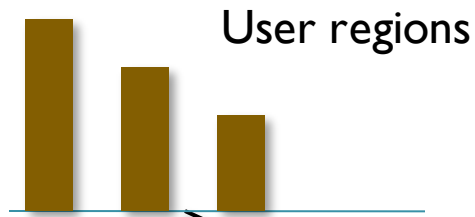
Topic

Generate tweet

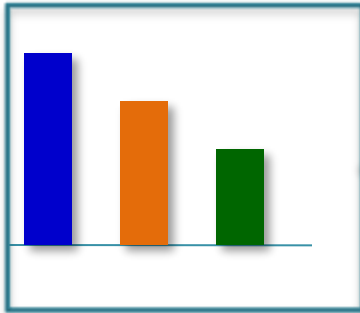
Region LM

Background LM

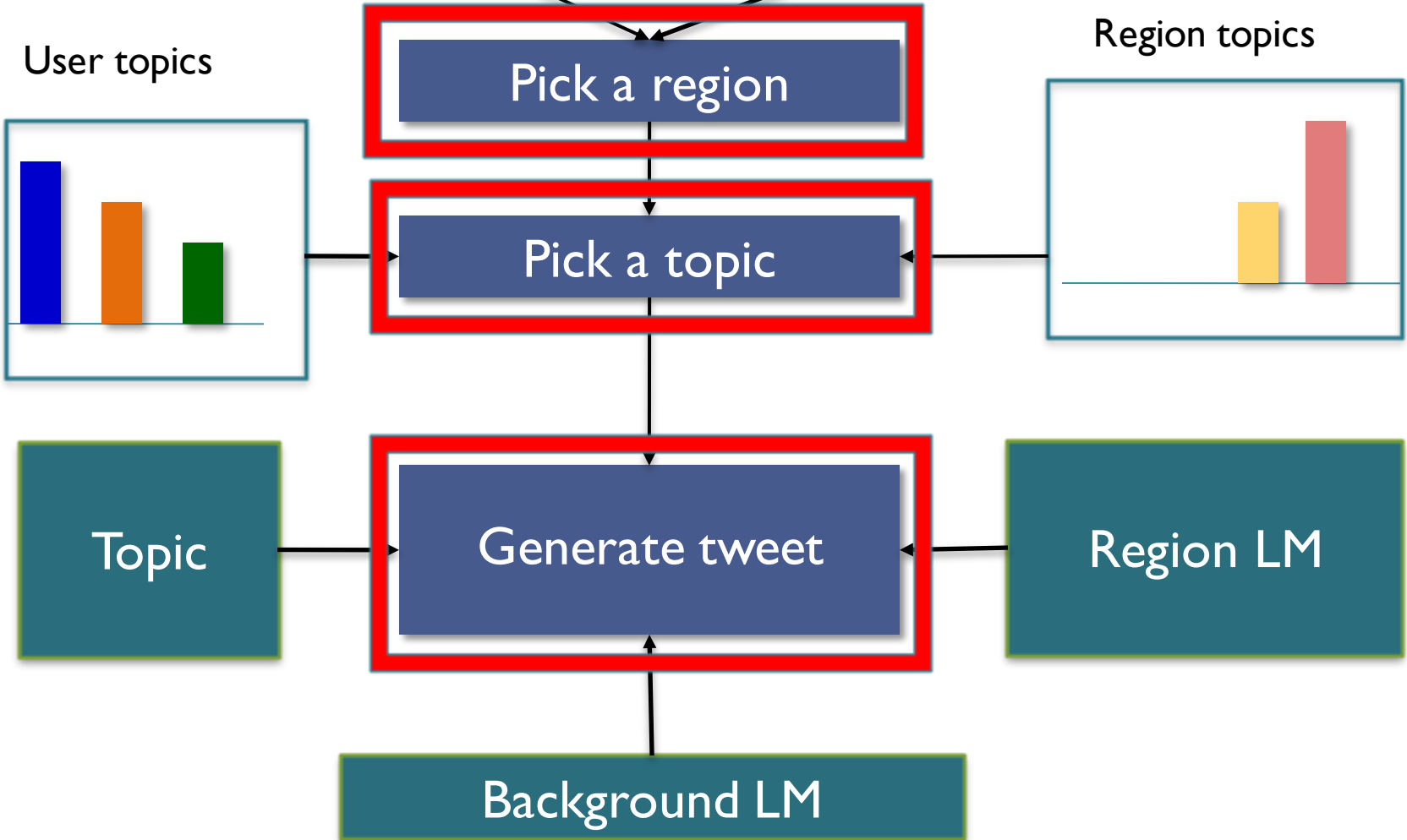
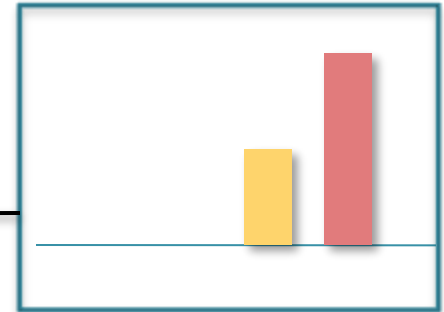




User topics



Region topics



# Discrete Additive Models

- Switch-based models
  - Normalized distributions
  - Pick one distribution
  - Sample from it
- **SAGE** (Eisenstein, Ahmed, Xing, 2011)
  - Un-normalized distribution
    - Log frequencies
  - Add them all together
  - Exponentiate and sample

# SAGE

## An Additive model for discrete distributions

- Discrete distribution via natural parameters

Example:

$$p(v|\phi) = \exp(\phi_v - g(\phi)) \quad \text{where } g(\phi) = \log \sum_v \exp(\phi_v)$$

- Log-frequency differences
- Addition of multiple models

Example:

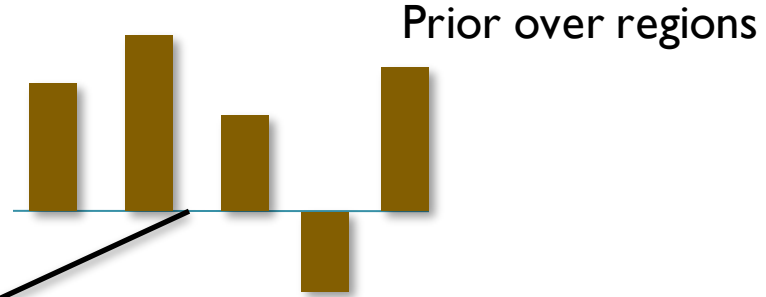
$$P(v|\phi_0, \phi_u, \phi_g) := p(v|\phi_0 + \phi_u + \phi_g)$$



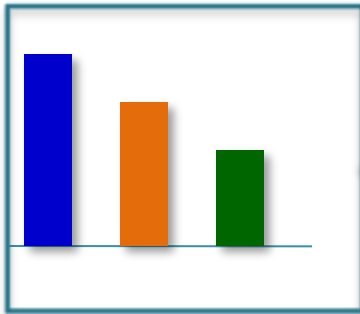
# SAGE

Use `SAGE` to replace “switch” variables to enable us incorporate multiple sources in different levels of our model easily

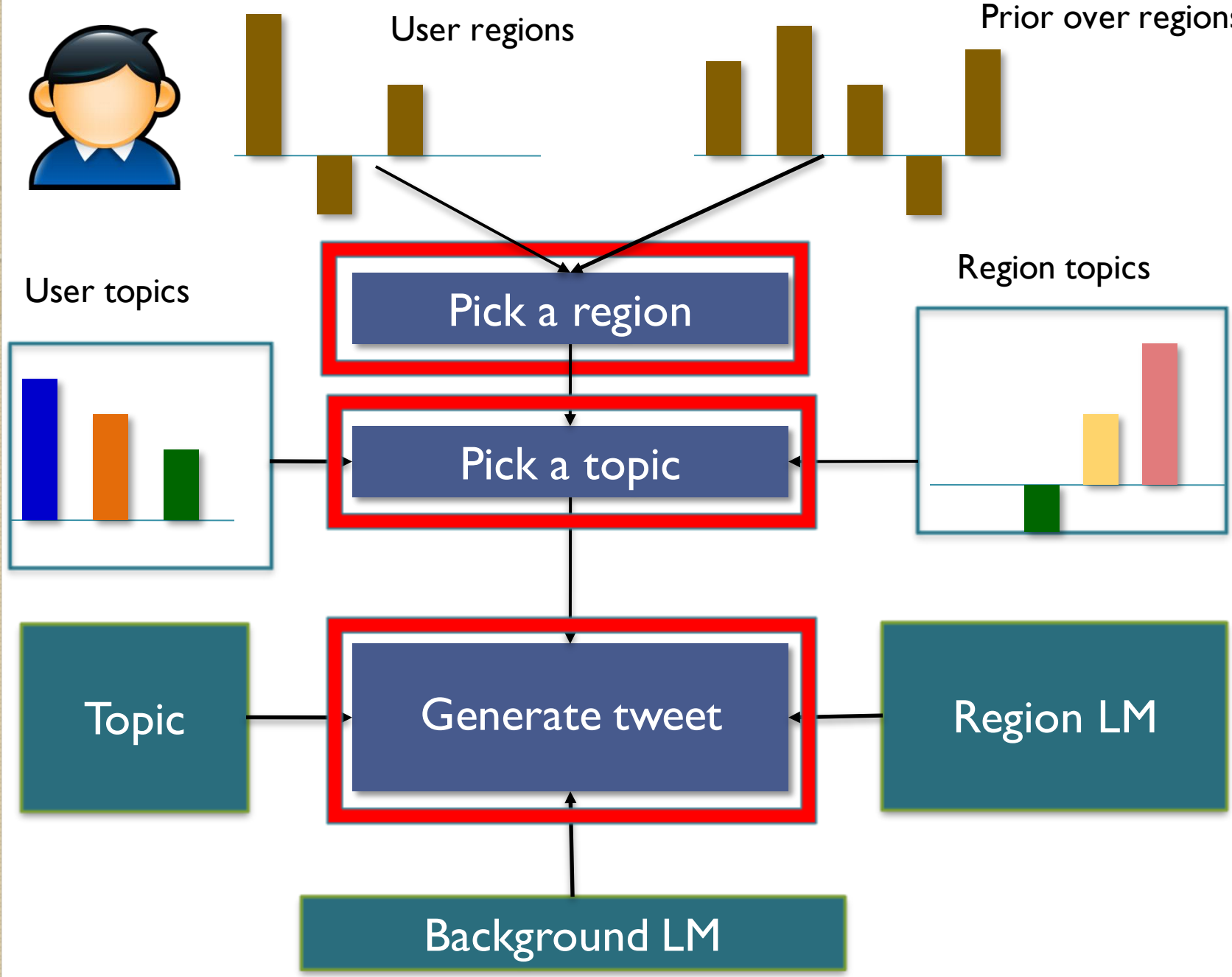
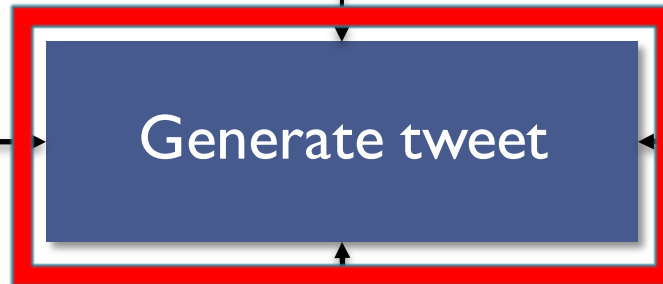
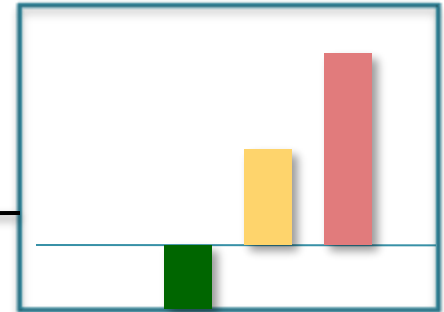
- Language models  
Example: background, regional, global...
- User preferences  
Example: global, regional, personal...
- ...



User topics



Region topics



# The Model

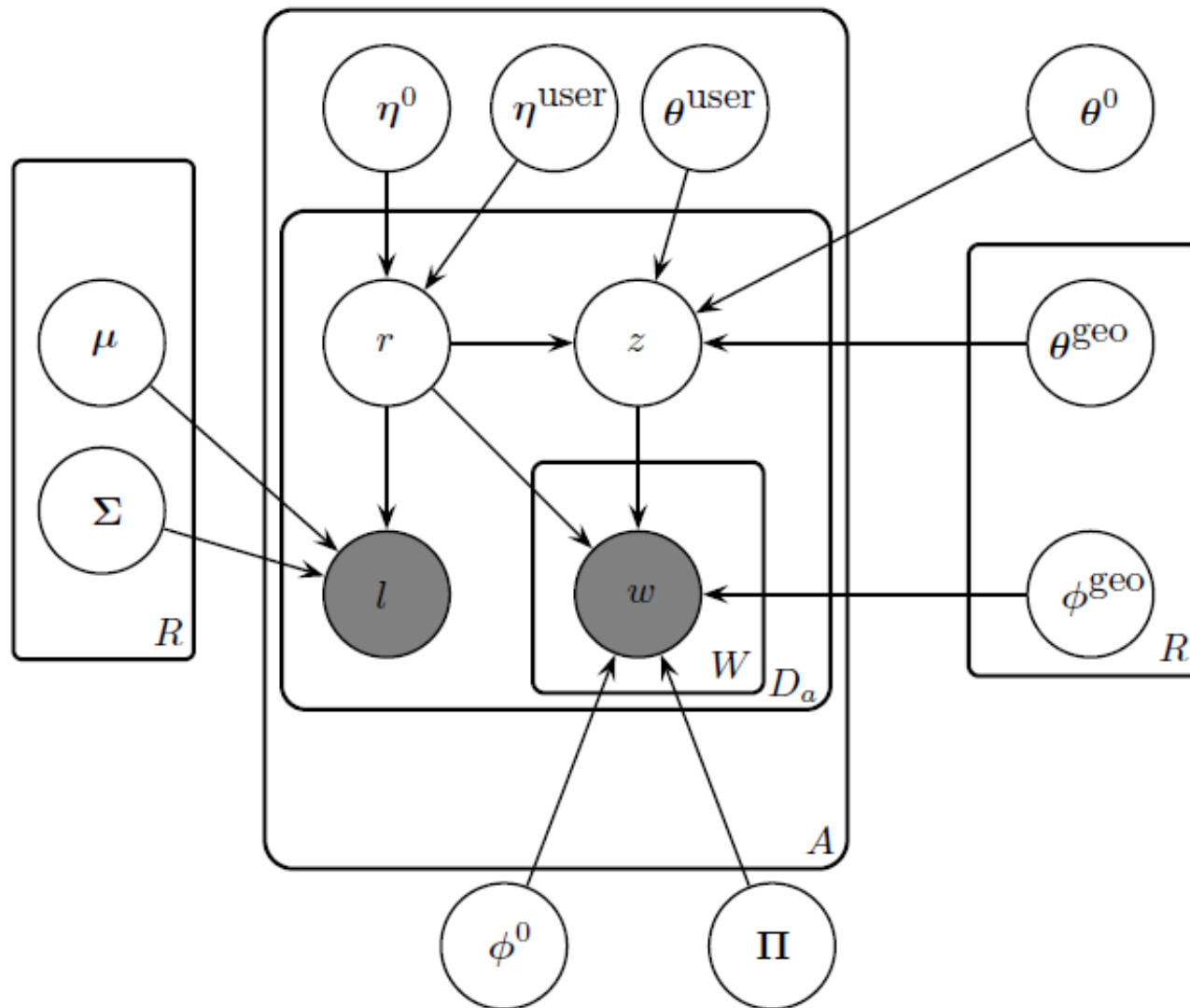
- Basic Intuition
  - Regions
  - Topics
  - Users
  - Tweets
- The generative process
  - Intuitive explanation
  - **Glory details**

# Generative Process

## Notations

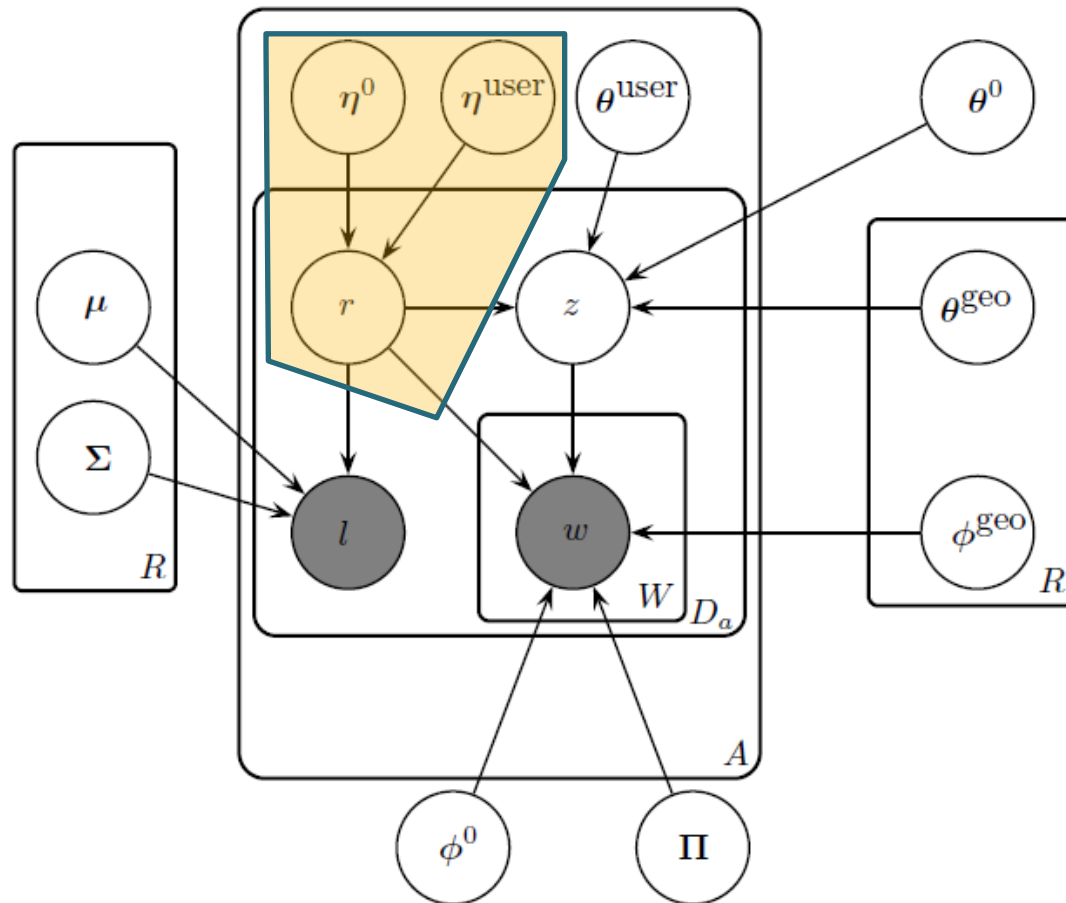
Symbol	Size	Usage
$\eta^0$	$1 \times \mathbb{R}$	global region distribution
$\eta^{\text{user}}$	$U \times \mathbb{R}$	user-dependent region distribution
$\theta^0$	$1 \times \mathbb{K}$	global topic distribution
$\theta^{\text{geo}}$	$\mathbb{R} \times \mathbb{K}$	region-dependent topic distribution
$\theta^{\text{user}}$	$U \times \mathbb{K}$	user-dependent topic distribution
$\phi^0$	$1 \times \mathbb{V}$	global term distribution
$\phi^{\text{geo}}$	$\mathbb{R} \times \mathbb{V}$	region-dependent term distribution
$\Pi$	$\mathbb{K} \times \mathbb{V}$	a global topic matrix
$\mu$	$\mathbb{R}^2$	mean location of a latent region
$\Sigma$	$\mathbb{R}^{2 \times 2}$	covariance matrix of a latent region

# The Graphical Model





# Region Selection



# Region Selection

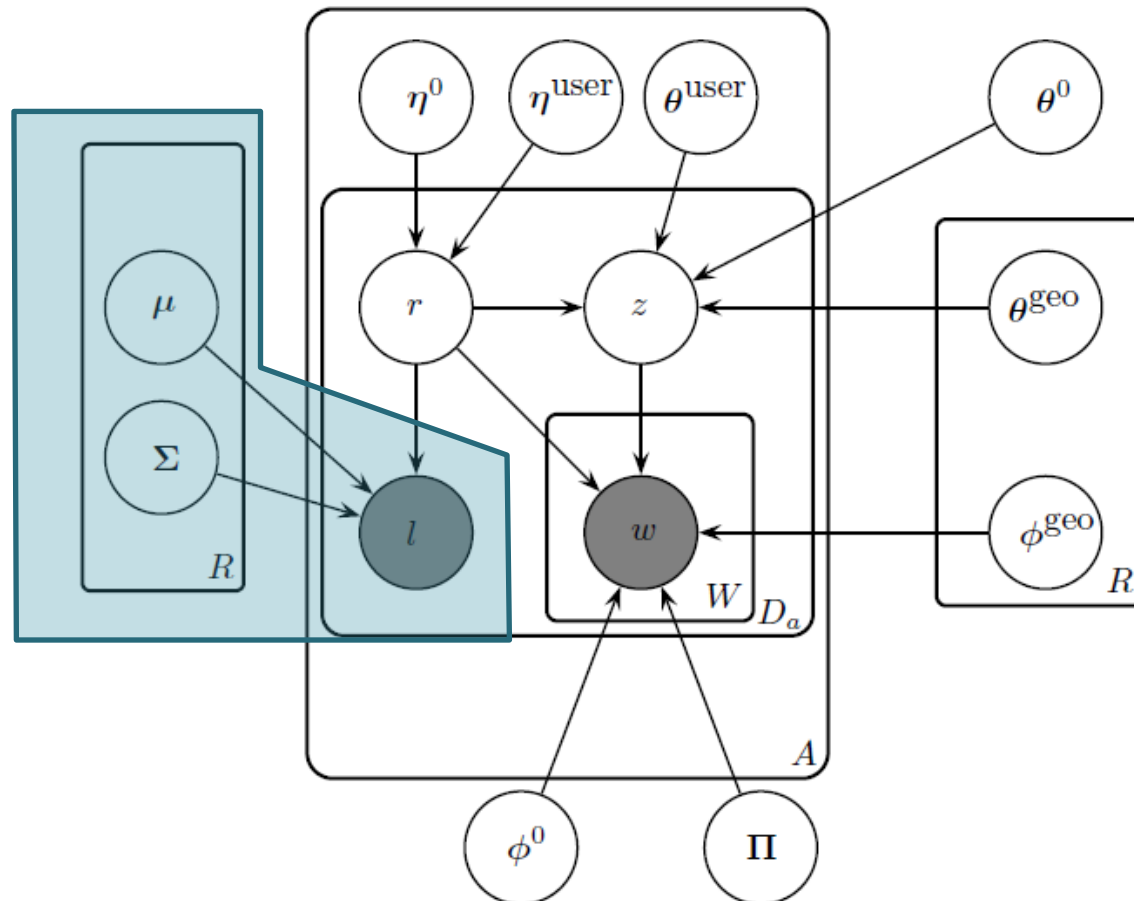
## Step-by-Step

- Users tend to appear in a handful geographical locations.

$$P(r|\eta^0, \eta_u^{\text{user}}) = p(r|\eta^0 + \eta_u^{\text{user}})$$



# Location Generation

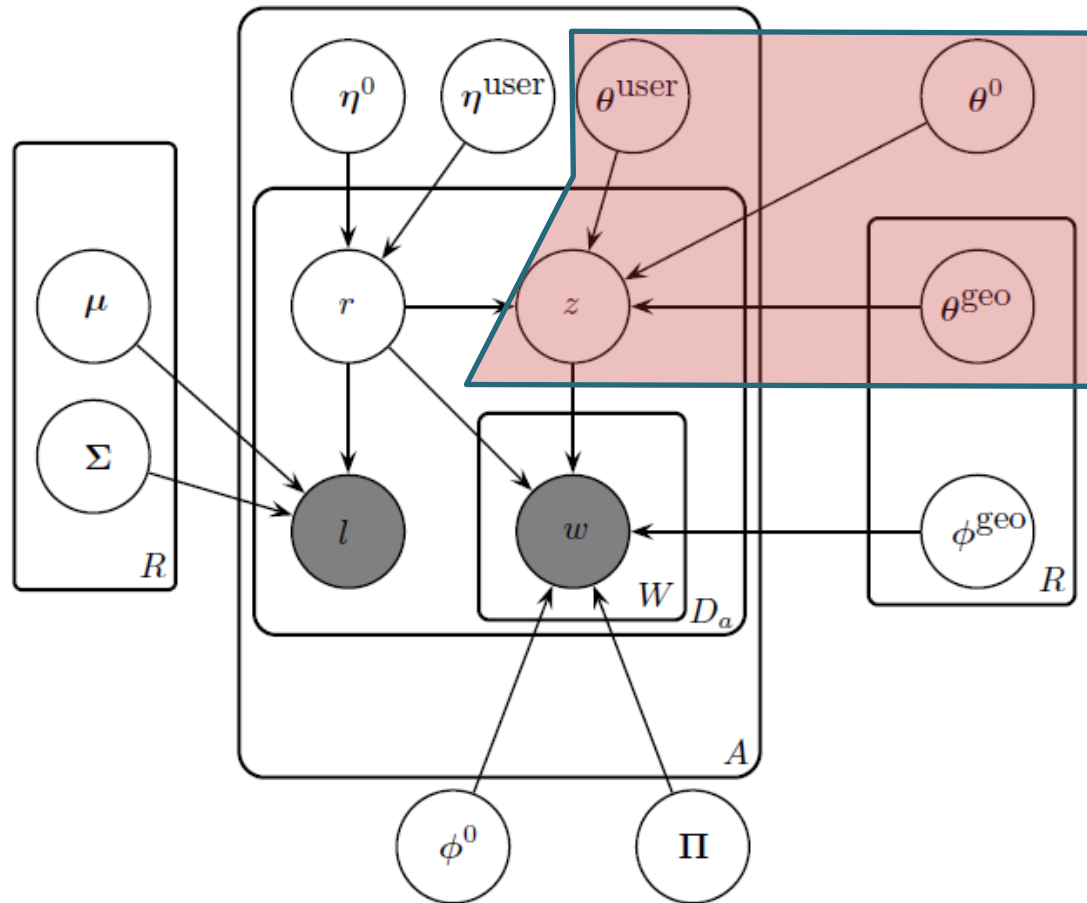


# Location Generation

- Once a region is selected, locations can be generated.

$$\mathbf{l}_d \sim \mathcal{N}(\boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r).$$

# Topic Selection

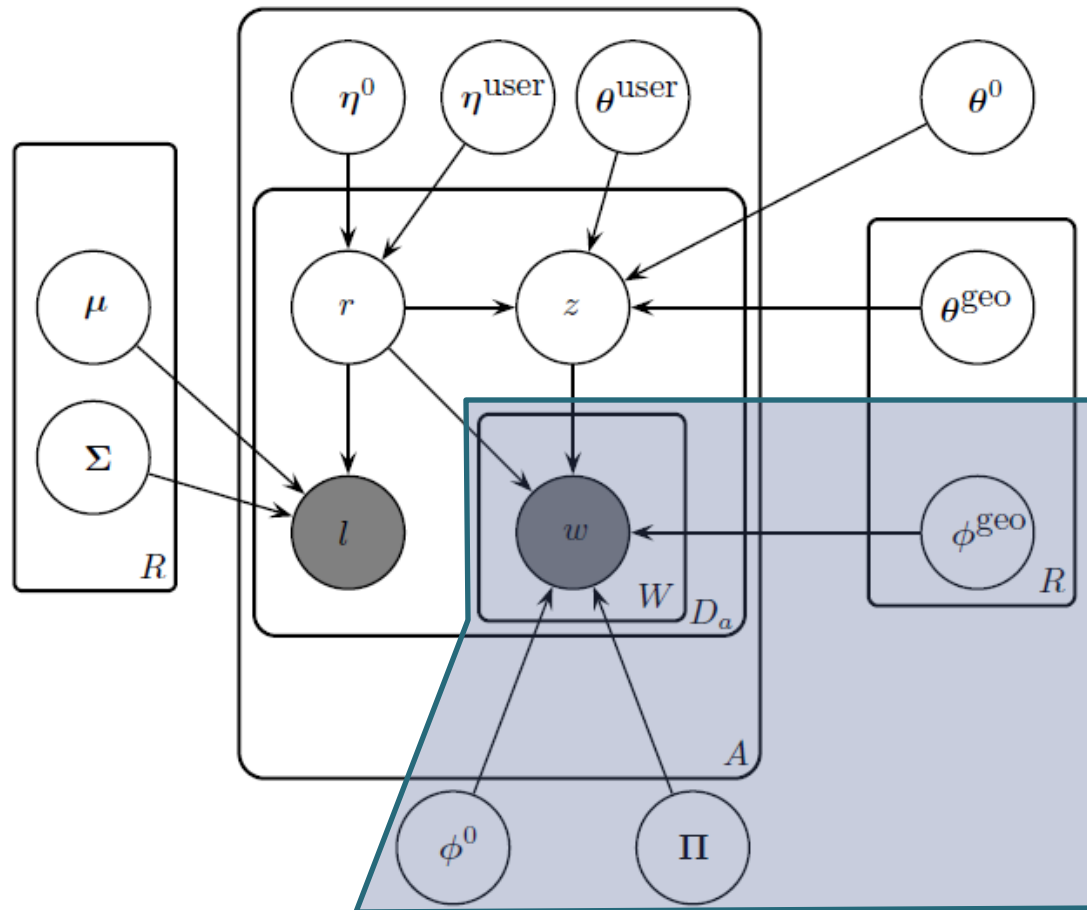


# Topic Selection

- Topics have different chances to be discussed in different regions by different users

$$P(z|\theta^0, \theta_u^{\text{user}}, \theta_r^{\text{geo}}) = p\left(z|\theta_j^0 + \theta_{u,j}^{\text{user}} + \theta_{r,j}^{\text{geo}}\right)$$

# Word Generation



# Word Generation

- Words used in a tweet depend on both the location and topic of the tweet.

$$P(w|z, \phi^0, \phi_r^{\text{geo}}, \mathbf{\Pi}_z) = p(w|\phi^0 + \phi_r^{\text{geo}} + \mathbf{\Pi}_{z_d})$$

# Sparse Modeling

- Laplace Priors

$$\begin{aligned}\eta_r^0 &\sim \mathcal{L}(0, \omega^0) & \eta_{u,r}^{\text{user}} &\sim \mathcal{L}(0, \omega_u) \\ \theta_z^{\text{geo}} &\sim \mathcal{L}(0, \lambda_l) & \theta_{u,z}^{\text{user}} &\sim \mathcal{L}(0, \lambda_u) & \theta_{r,z}^{\text{geo}} &\sim \mathcal{L}(0, \lambda_r) \\ \phi_v^0 &\sim \mathcal{L}(0, \psi^0) & \phi_{r,v}^{\text{geo}} &\sim \mathcal{L}(0, \psi_l) \\ \Pi_{z,v} &\sim \mathcal{L}(0, \psi_t)\end{aligned}$$

Sparsity results in  
predictive models

# Bayesian treatment

- Prior distributions over mean and covariance matrix
- Jeffery prior

$$\mu \sim \text{Unif.}$$

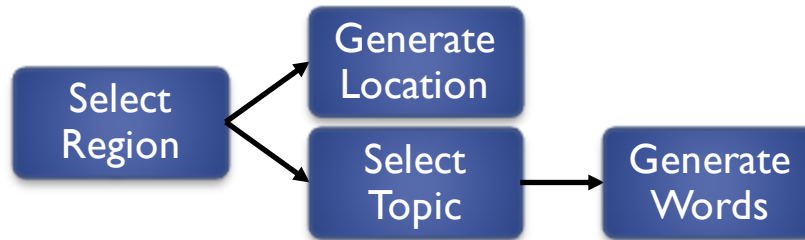
$$P(\Sigma) \propto |\Sigma|^{-(3/2)}.$$

- Penalize large regions
  - We want region to be predictive as much as the data supports



# Recap

- Generative Process



- Sparse Modeling
  - $L_1$  regularization (Laplace priors)
- Geographical Modeling
  - Bayesian treatment

# Inference Algorithm

- A variant of Monte Carlo EM
    - “E-Step”: Sample latent discrete variables
    - “M-step”: Update all model parameters
- 
- Sparse update of gradients
  - $L_1$  regularization: ISTA algorithm
  - Initialize regions with K-means algorithm

# Overview

- Motivations
- Our Proposed Model
- **Experiments**
- Conclusions

# Experiments

## Dataset

- Twitter data
  - Randomly sample 1,000 users
  - All tweets from Jan 2011 to May 2011
  - 573,203 distinct tweets
- Twitter geographical data
  - Locations + Twitter Places

# Experiments

## Location Prediction

- Metric
  - average error distance
  - Kilometers



# Experiments

## Location Prediction

- Baselines

- Yin et al. WWW 2011 paper
  - PLSA formalism
  - No personalization
- Our model without  $\phi^{\text{geo}}$ ,  $\eta^{\text{user}}$  and  $\theta^{\text{user}}$ 
  - Similar to Yin et al.'s formalism but SAGE model
- Our model without  $\eta^{\text{user}}$  and  $\theta^{\text{user}}$

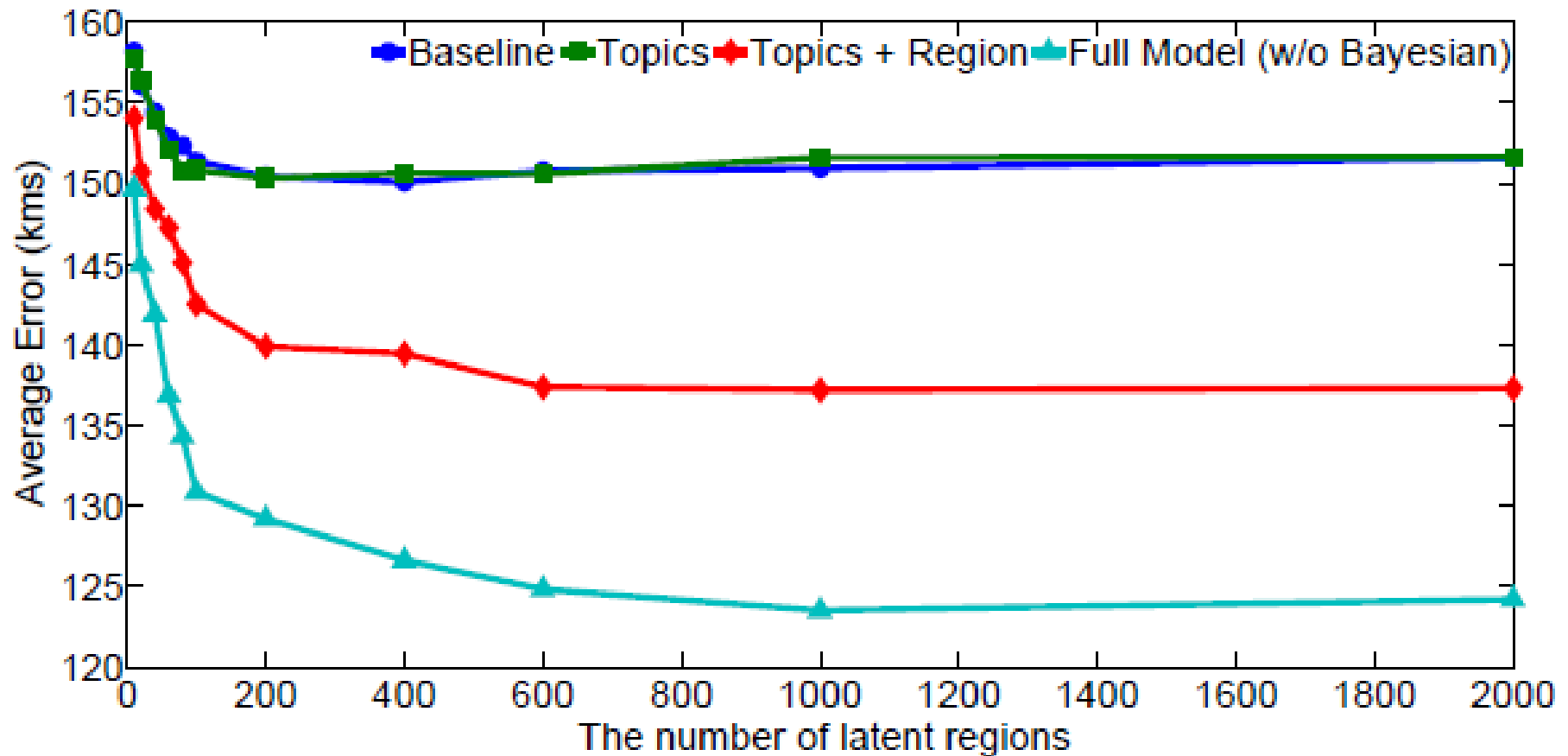
# Experiments

## Location Prediction

- Baselines

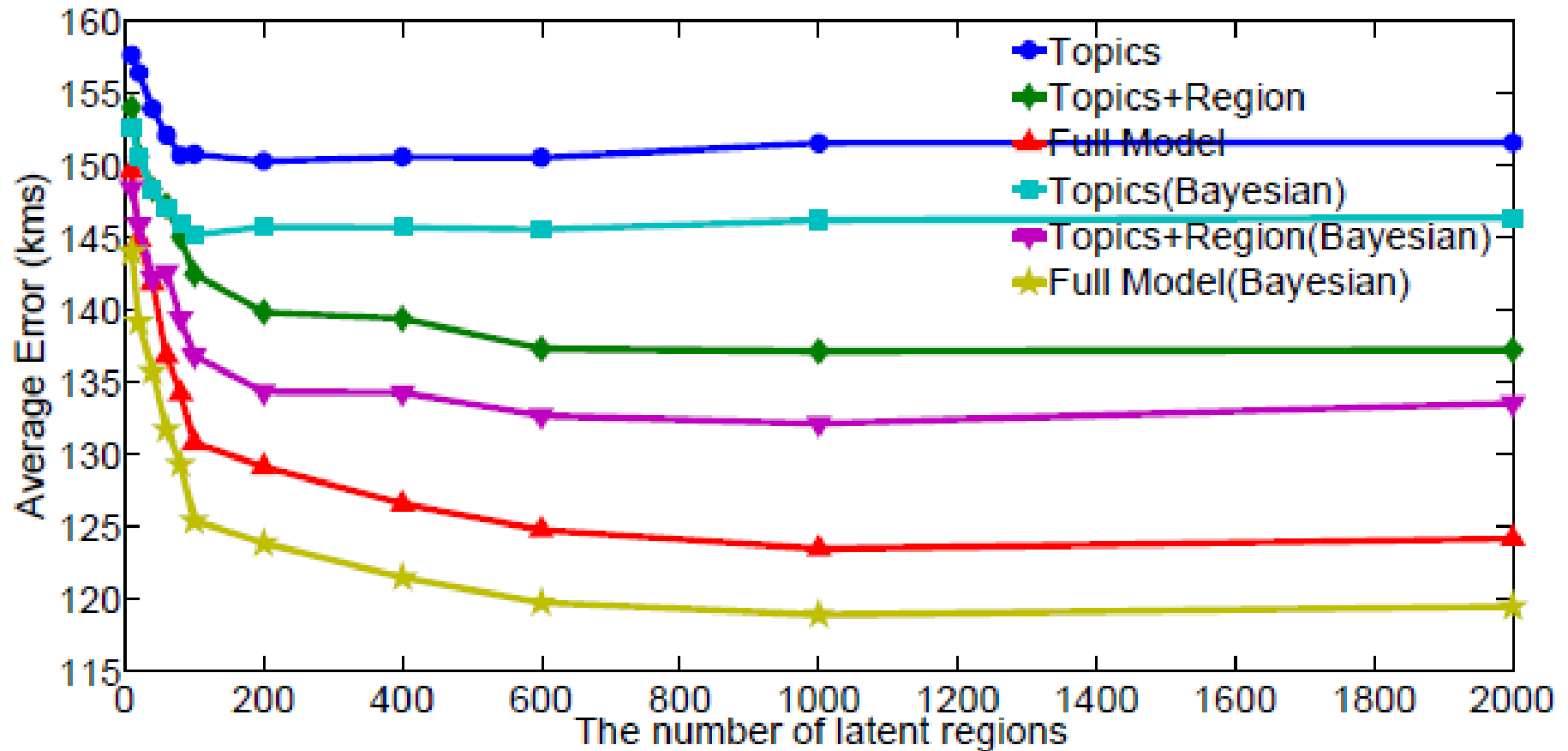
- Yin et al. WWW 2011 paper
  - PLSA formalism
  - No personalization
- Our model without  $\phi^{\text{geo}}$ ,  $\eta^{\text{user}}$  and  $\theta^{\text{user}}$ 
  - Similar to Yin et al.'s formalism but SAGE model
- Our model without  $\eta^{\text{user}}$  and  $\theta^{\text{user}}$

# Location Prediction

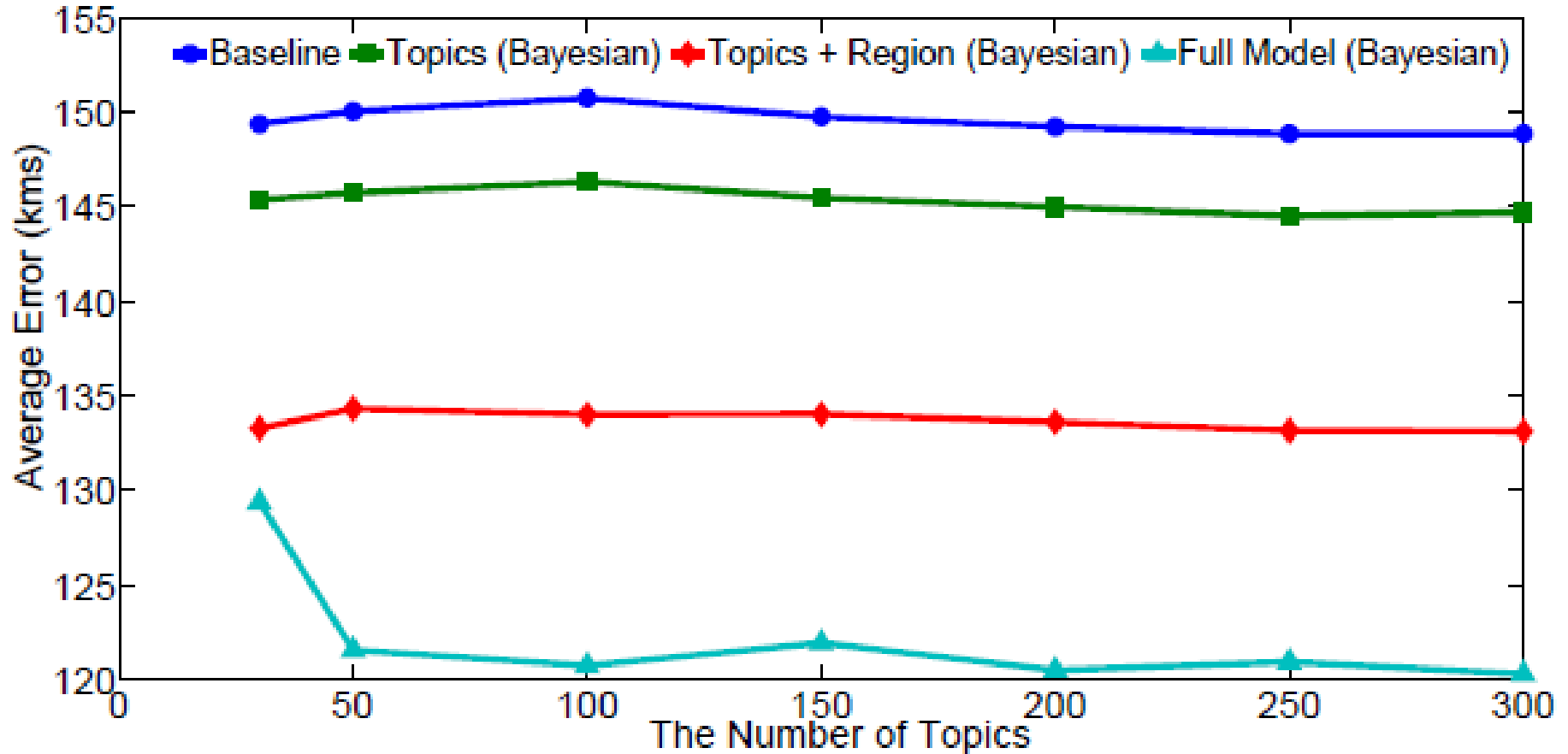




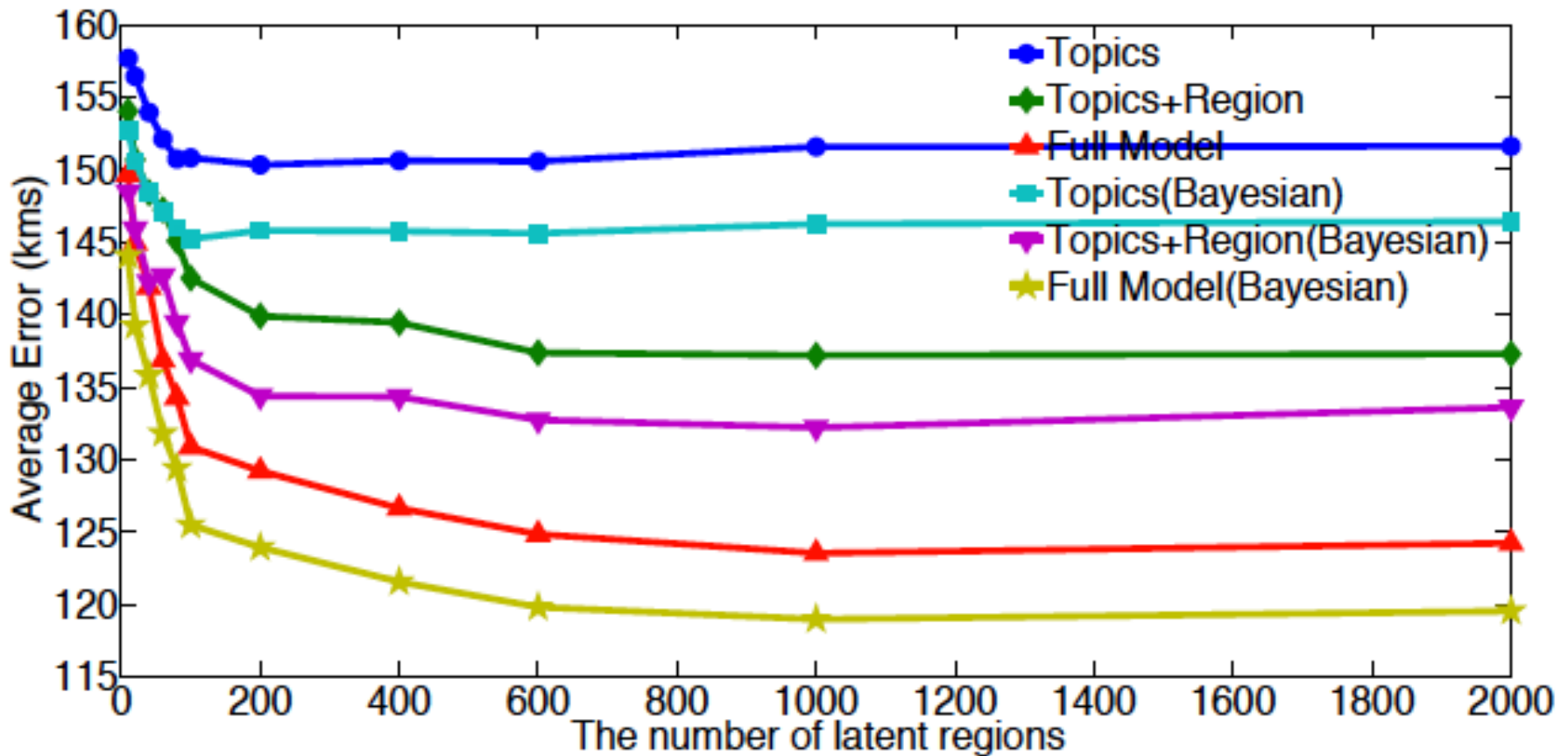
# Bayesian Treatment



# Number of Topics



# Number of Regions



# Experiments (Public Data)

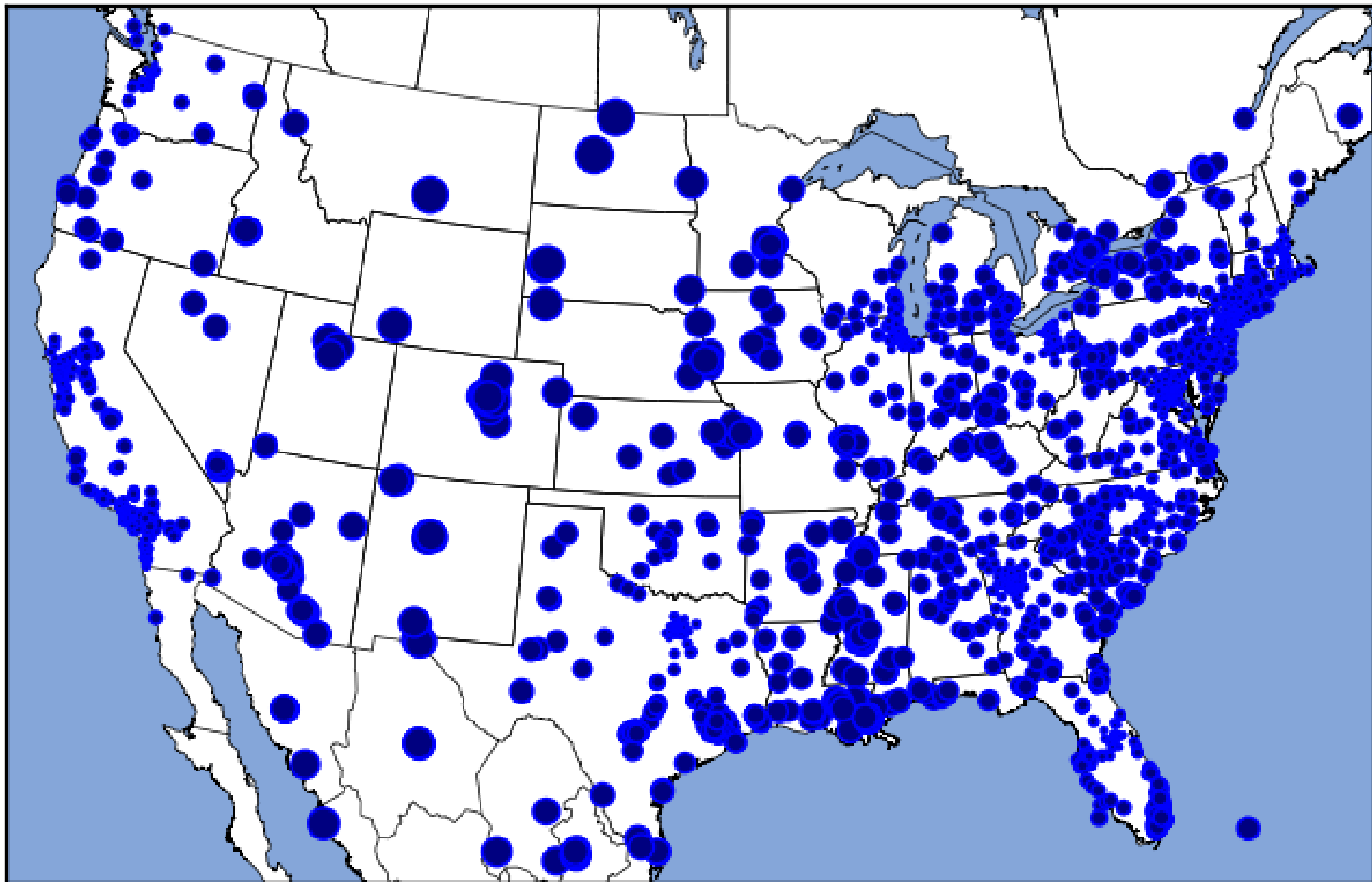
# of regions	[3]	[2]	[1]	Topics	Topics + Region	Full Model
10	494	479	501	540.60	481.58	449.45
20	494	479	501	522.18	446.03	420.83
40	494	479	501	513.06	414.95	395.13
60	494	479	501	507.37	410.09	380.04
80	494	479	501	499.42	408.38	374.01
100	494	479	501	498.94	407.78	<b>372.99</b>

[1] Eisenstein et al. EMNLP 2010.

[2] Wing and J. Baldrige. ACL 2011.

[3] Eisenstein, Ahmed, Xing ICML 2011.

# Error Analysis



# Global and local topics

## Entertainments

lady beiber album music beats artist video listen  
itunes apple produced movies #bieber lol new songs

## Sports

yankees match nba football giants wow win winner game  
weekend horse #nba

## Politics

obama election middle east china uprising egypt russian  
tunisia #egypt afghanistan people eu

## Location with Top Ranked Terms

### United States->New York->Brooklyn

brooklyn ave flatbush avenue mta prospect 5th #brooklyn spotlight carroll bushwick museum broadway madison  
vanderbilt coney slope eastern subway new york pkwy #viernesnayobon #mets otsego greenwich starbucks

### United States->California->San Francisco

sfo francisco san airport international millbrae terminal flight burlingame bart mateo boarding bayshore telecommute  
landed heading bay airlines united bound flying #sfo caminogroupon caltrain moon tsa baggage california engineer valley

### United States->Pennsylvania->Philadelphia

philadelphia #philadelphia phl #jobs market others #job street philly walnut septa chestnut the cherry  
sansom arch spruce citizens locust btw temple pennsylvania rittenhouse passyunk bitlyetq7a6 bookrenters pike international

### United Kingdom->England->London

winds lhr hounslow terminal the cloudy mph ickenham bath heathrow temperature airport car only airways uxbridge sun  
splendid fair london british lounge tothers harmondsworth speedbird whens for stars day flight dominos navigation brunel

### Australia->New South Wales->Sydney

sydney #sydney bondi george street mascot domestic syd surry station cnr platforms harbour darlinghurst qantas hoteloxford  
eddy haymarket terminal wales australia chalmers uts pitt #marketing junction darling centre #citijobs citigroup druit

# Conclusions

- Probabilistic model for geographical information
  - Regional variations
  - Personal preferences
- Effective inference algorithm
- Best location prediction
- Discriminatively learned language models
- Future work
  - Hierarchical model
  - Hash tags
  - Temporal location model