

Intro. to Topic Modeling (cont'd) + Factor Analysis

Kayhan Batmanghelich

Topic Modeling

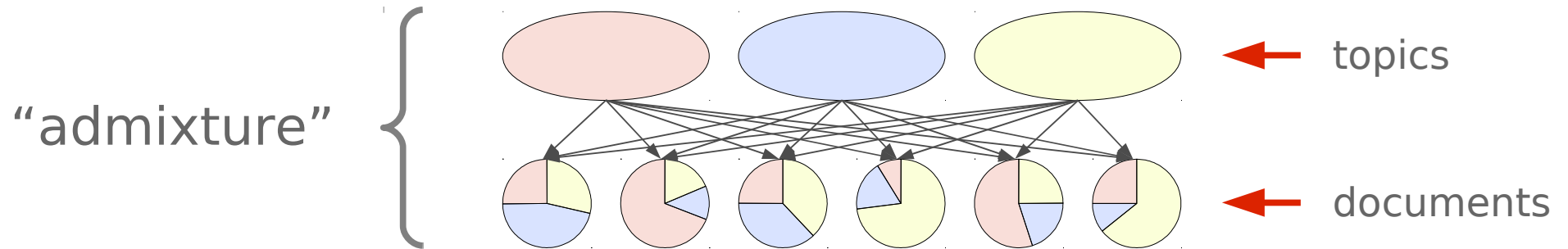
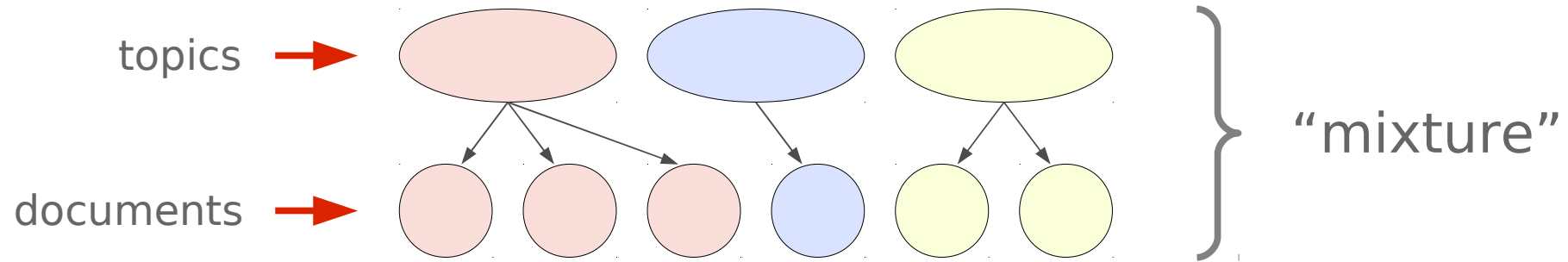
Motivation:

Suppose you're given a massive corpora and asked to carry out the following tasks

- **Organize** the documents into **thematic categories**
- **Describe** the evolution of those categories **over time**
- Enable a domain expert to **analyze and understand** the content
- Find **relationships** between the categories
- Understand how **authorship** influences the content

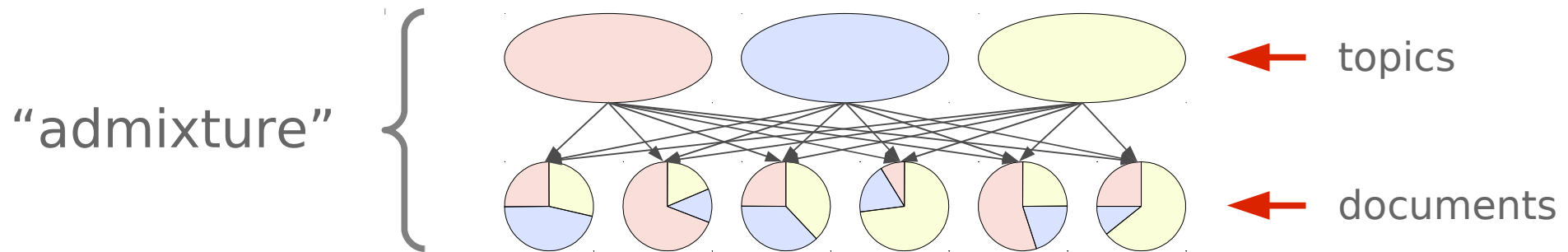


Mixture vs. Admixture (LDA)



Latent Dirichlet Allocation

- Generative Process



- Example corpus

the	he	is
x_{11}	x_{12}	x_{13}

Document 1

the	and	the
x_{21}	x_{22}	x_{23}

Document 2

she	she	is	is
x_{31}	x_{32}	x_{33}	x_{34}

Document 3

Latent Dirichlet Allocation

- Generative Process

For each topic $k \in \{1, \dots, K\}$:
 $\phi_k \sim \text{Dir}(\beta)$ [draw distribution over words]
For each document $m \in \{1, \dots, M\}$:
 $\theta_m \sim \text{Dir}(\alpha)$ [draw distribution over topics]
 For each word $n \in \{1, \dots, N_m\}$:
 $z_{mn} \sim \text{Mult}(1, \theta_m)$ [draw topic assignment]
 $x_{mn} \sim \phi_{z_{mn}}$ [draw word]

- Example corpus

the	he	is
x_{11}	x_{12}	x_{13}

Document 1

the	and	the
x_{21}	x_{22}	x_{23}

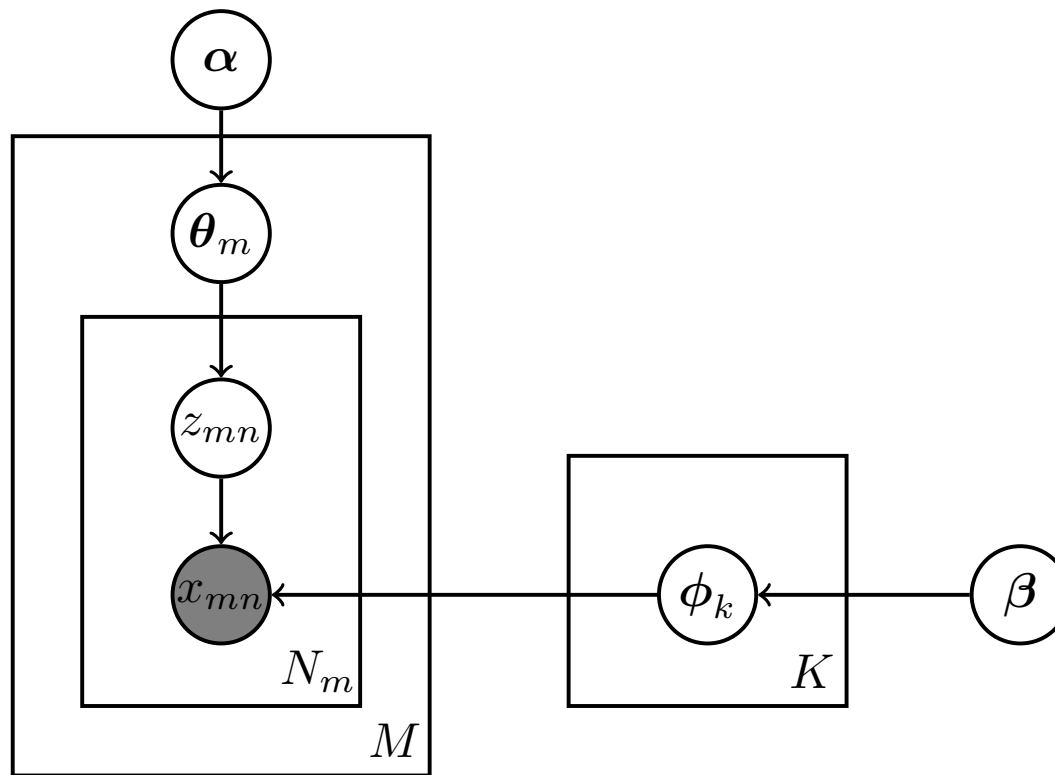
Document 2

she	she	is	is
x_{31}	x_{32}	x_{33}	x_{34}

Document 3

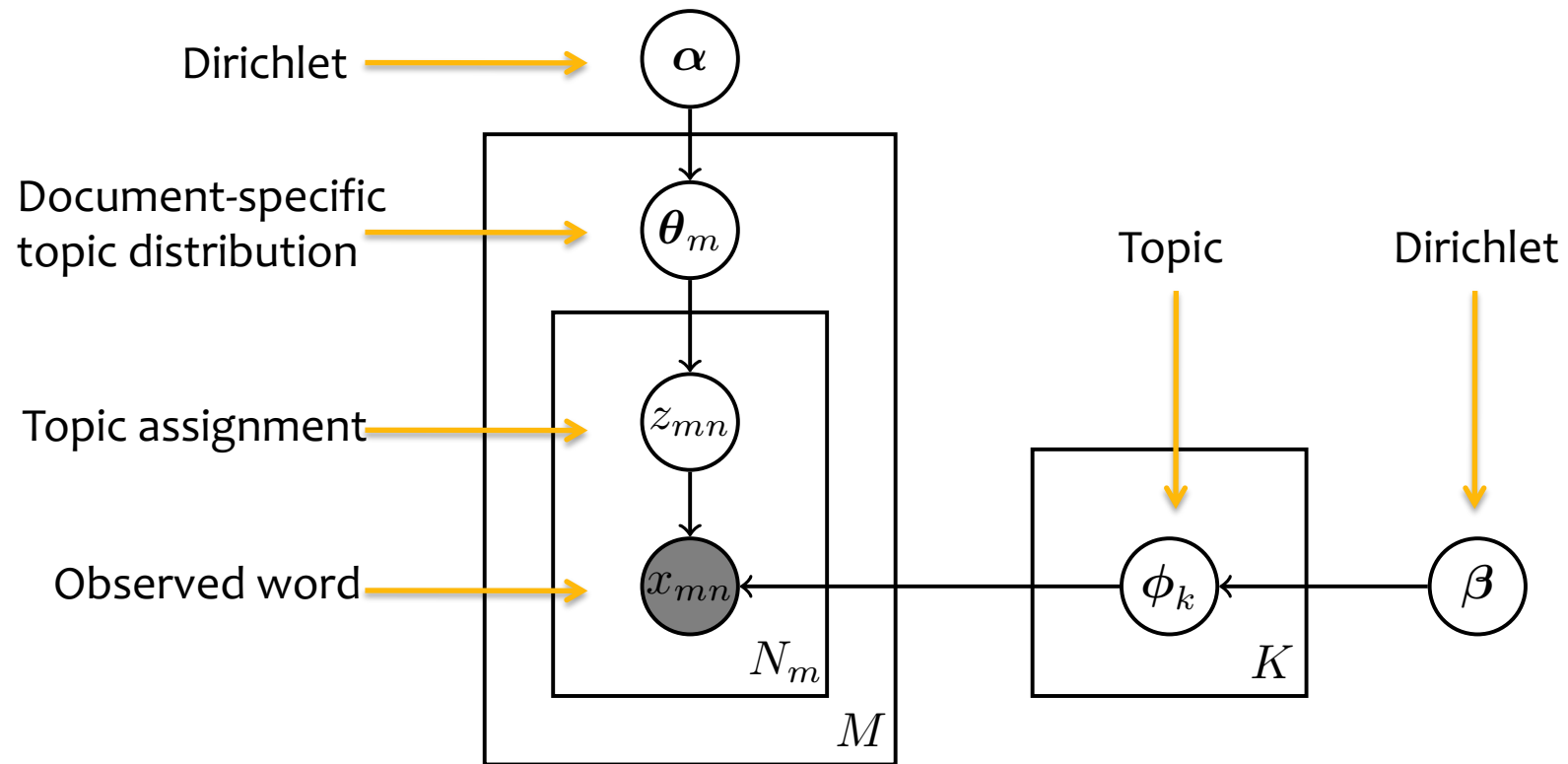
Latent Dirichlet Allocation

- Plate Diagram

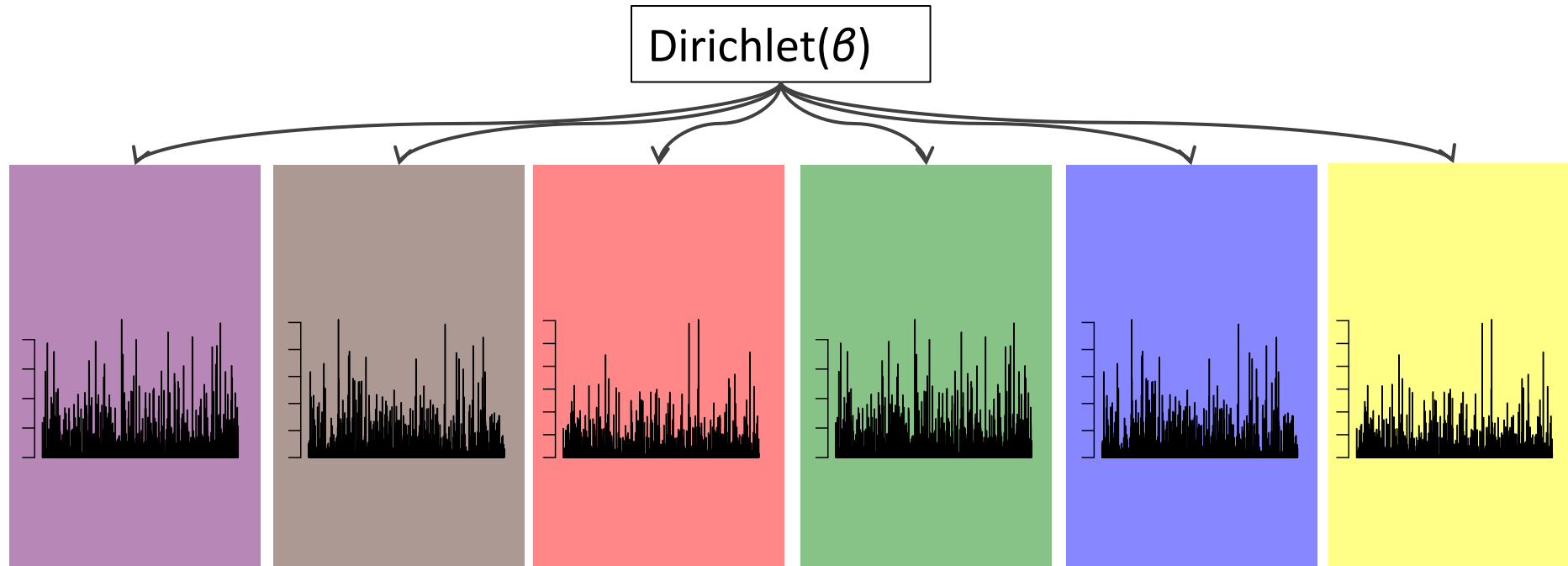


Latent Dirichlet Allocation

- Plate Diagram

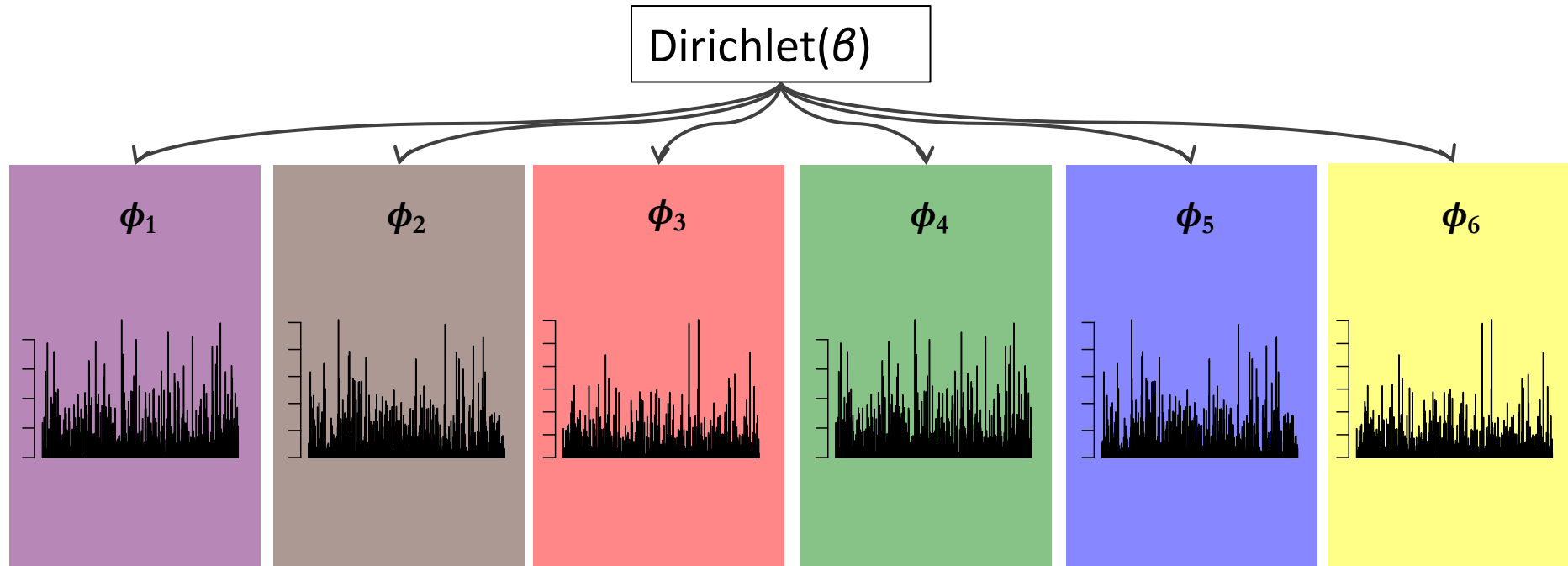


LDA for Topic Modeling



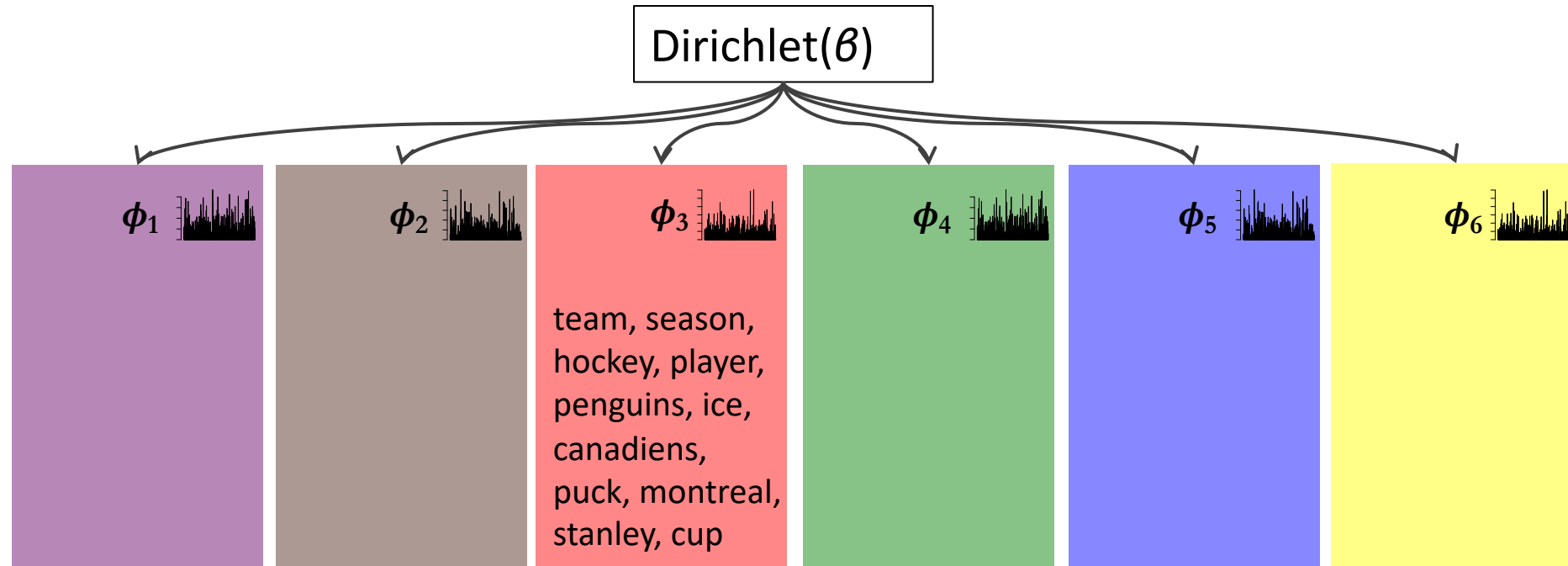
- The **generative story** begins with only a **Dirichlet prior** over the topics.
- Each **topic** is defined as a **Multinomial distribution** over the vocabulary, parameterized by ϕ_k

LDA for Topic Modeling



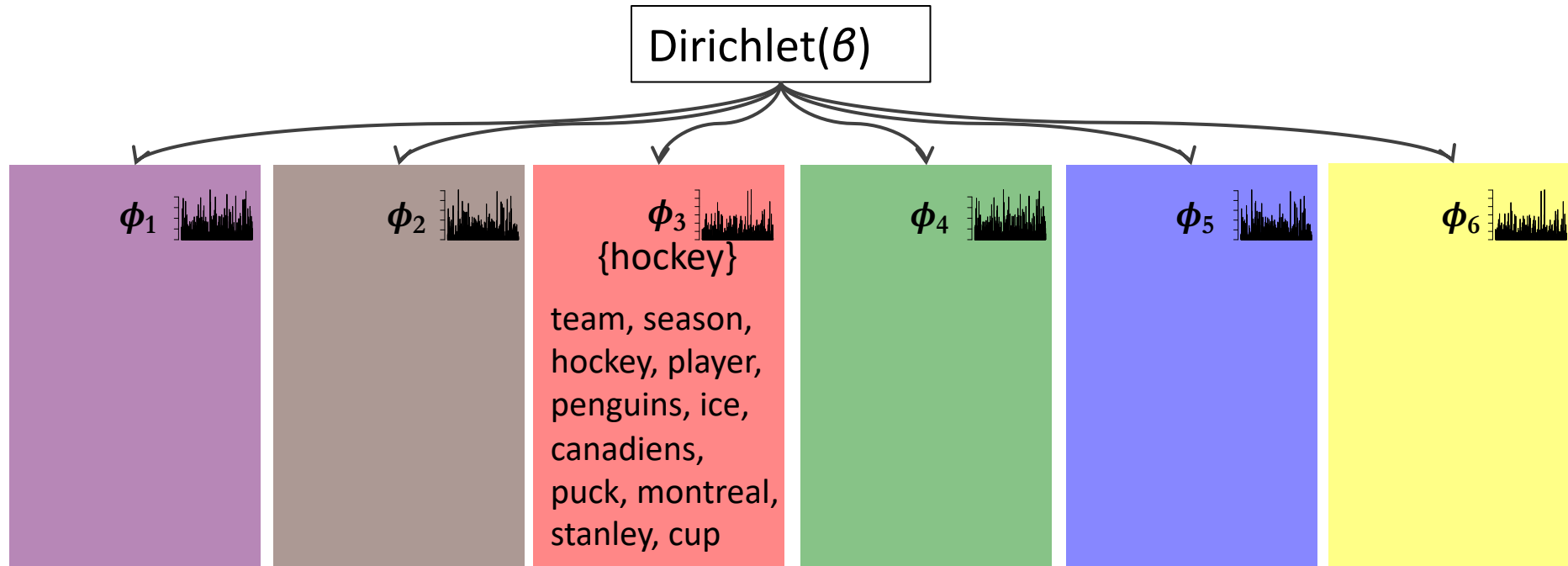
- The **generative story** begins with only a **Dirichlet prior** over the topics.
- Each **topic** is defined as a **Multinomial distribution** over the vocabulary, parameterized by ϕ_k

LDA for Topic Modeling



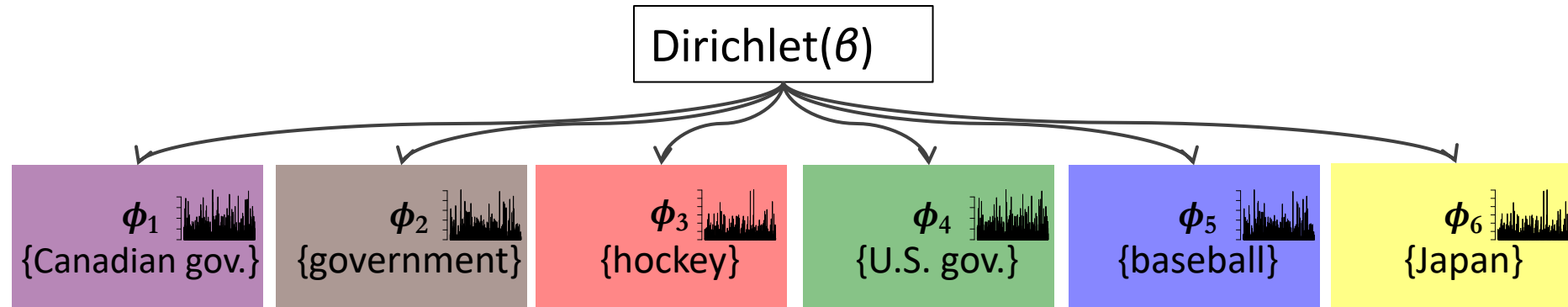
- A topic is visualized as its **high probability words**.

LDA for Topic Modeling



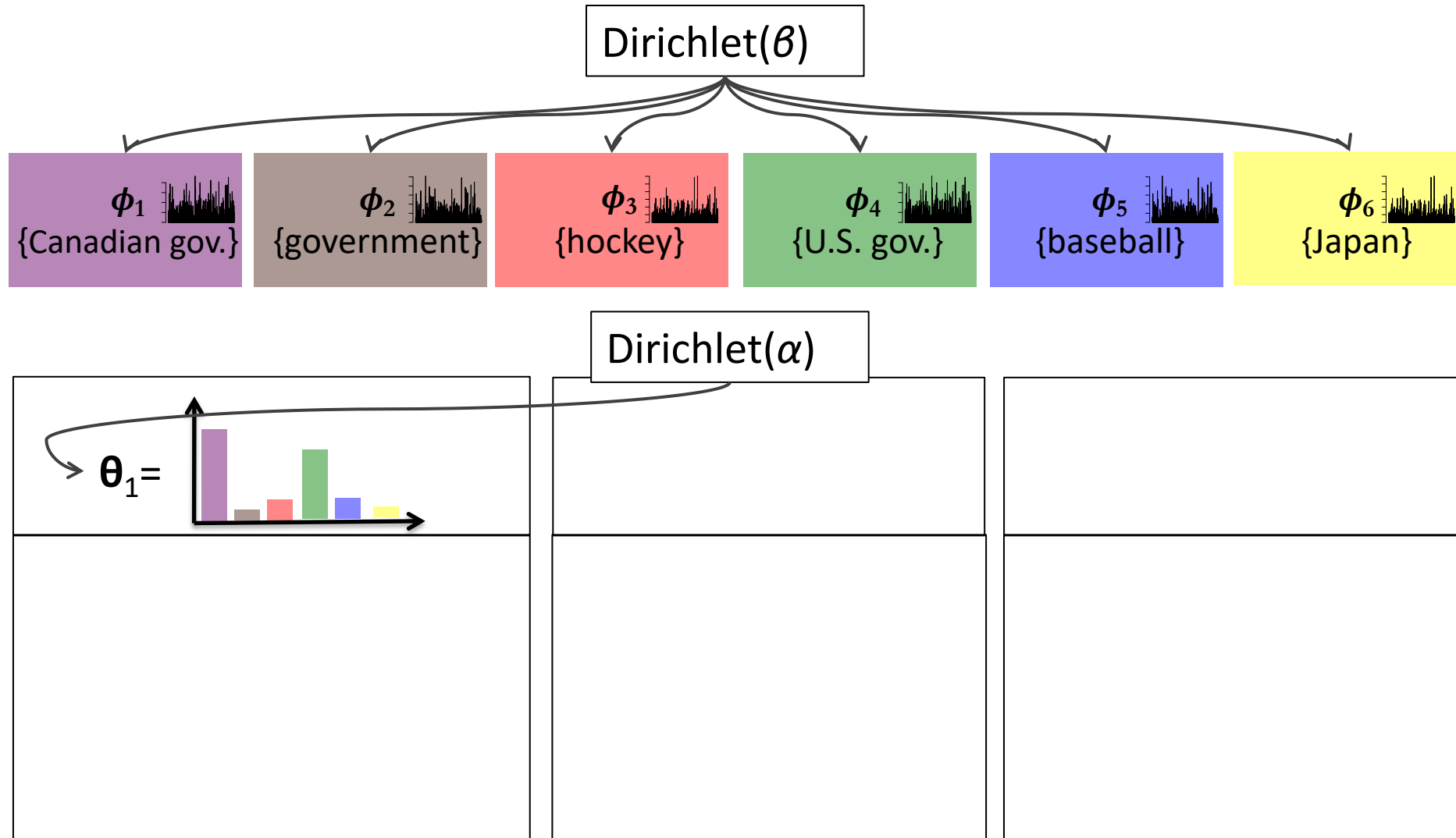
- A topic is visualized as its **high probability words**.
- A pedagogical **label** is used to identify the topic.

LDA for Topic Modeling

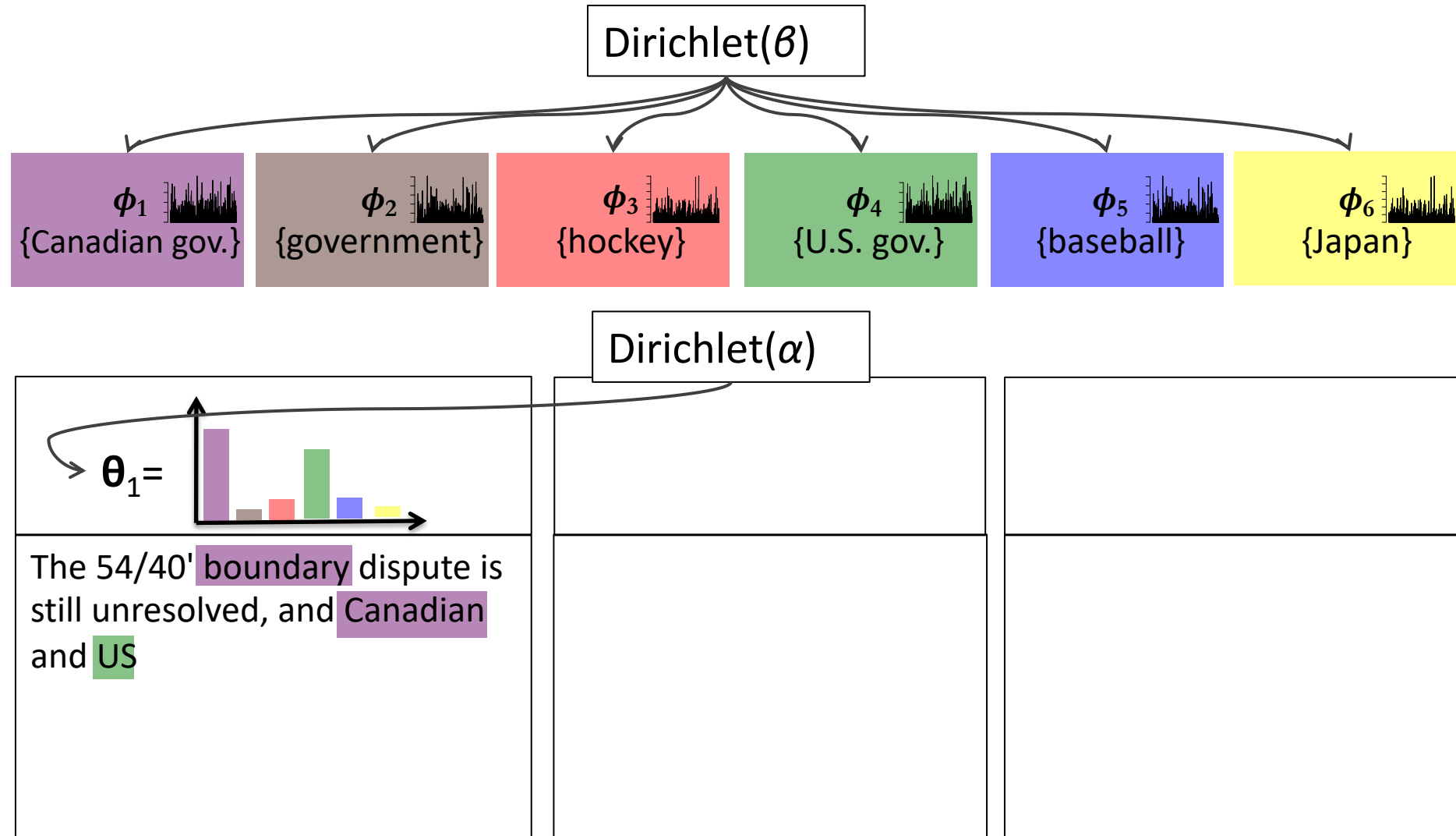


- A topic is visualized as its high probability words.
- A pedagogical **label** is used to identify the topic.

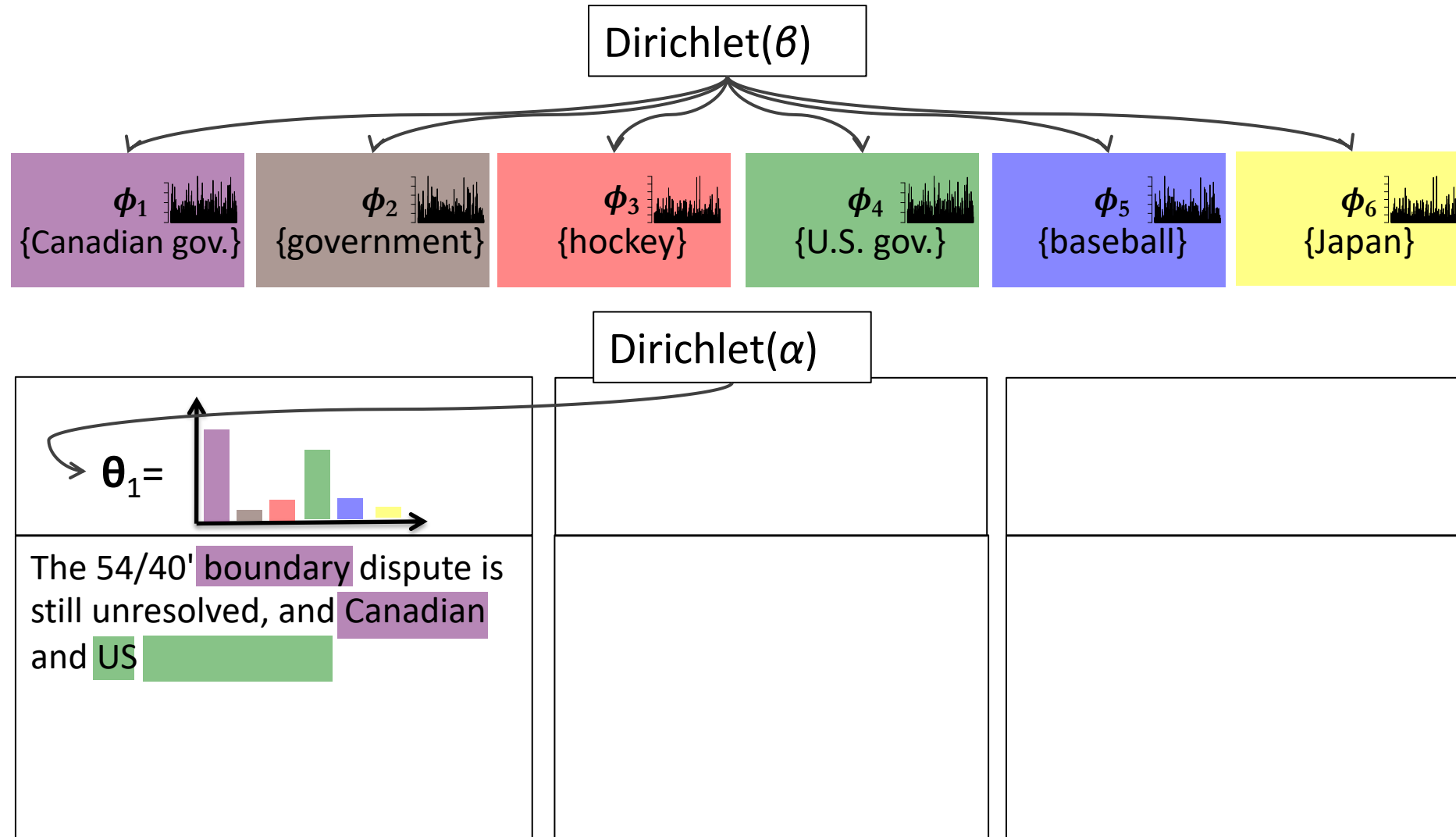
LDA for Topic Modeling



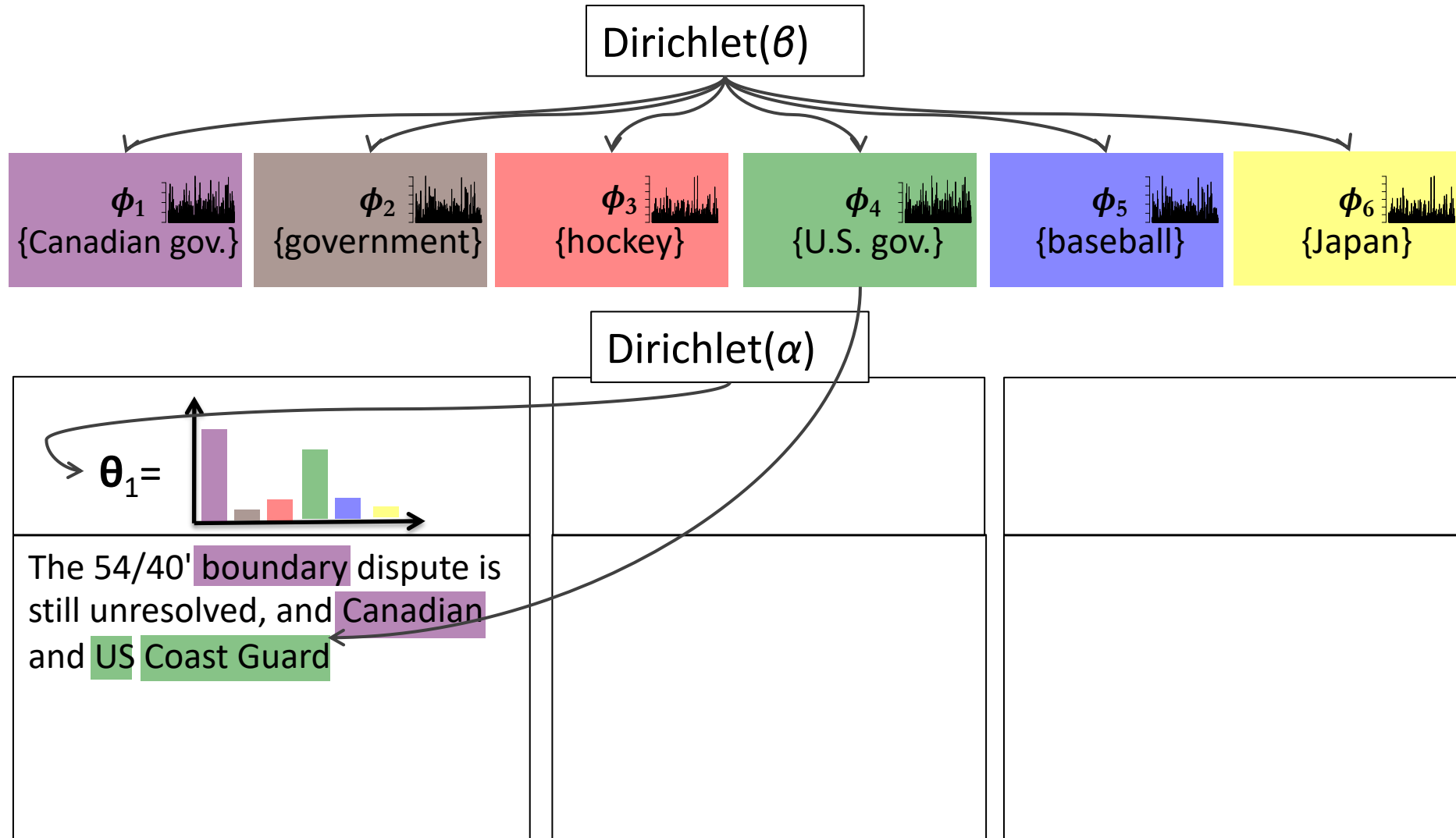
LDA for Topic Modeling



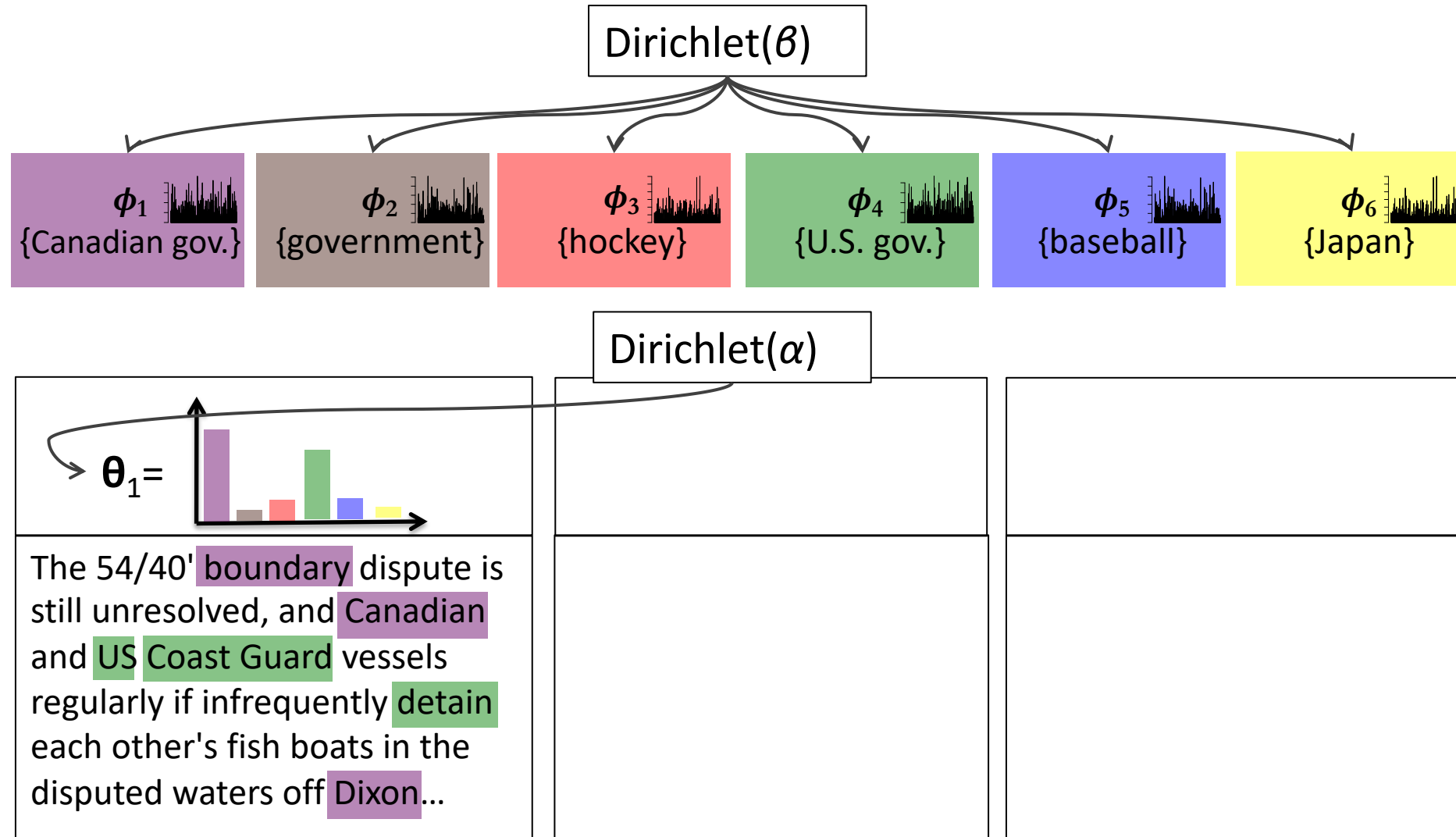
LDA for Topic Modeling



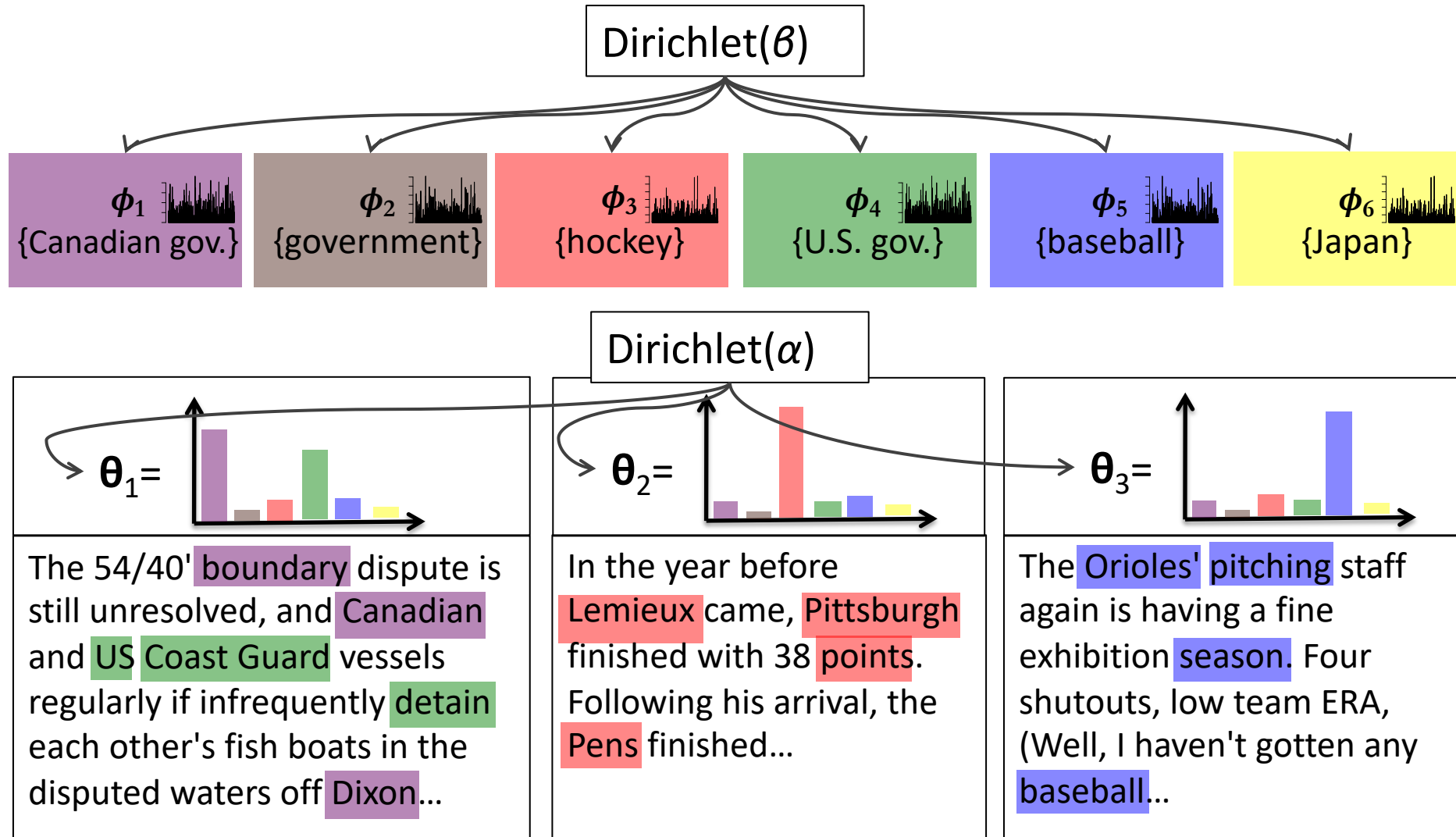
LDA for Topic Modeling



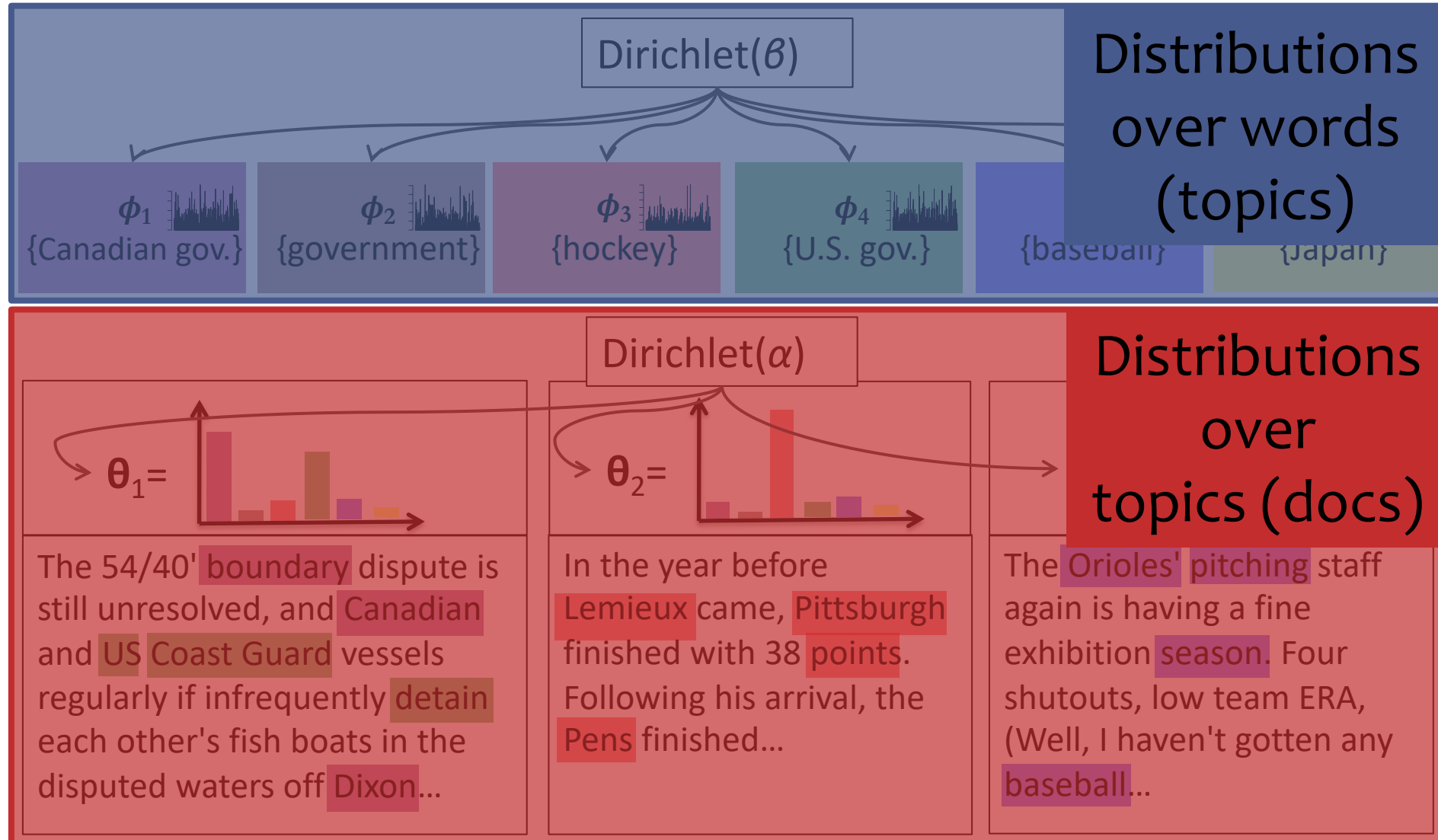
LDA for Topic Modeling



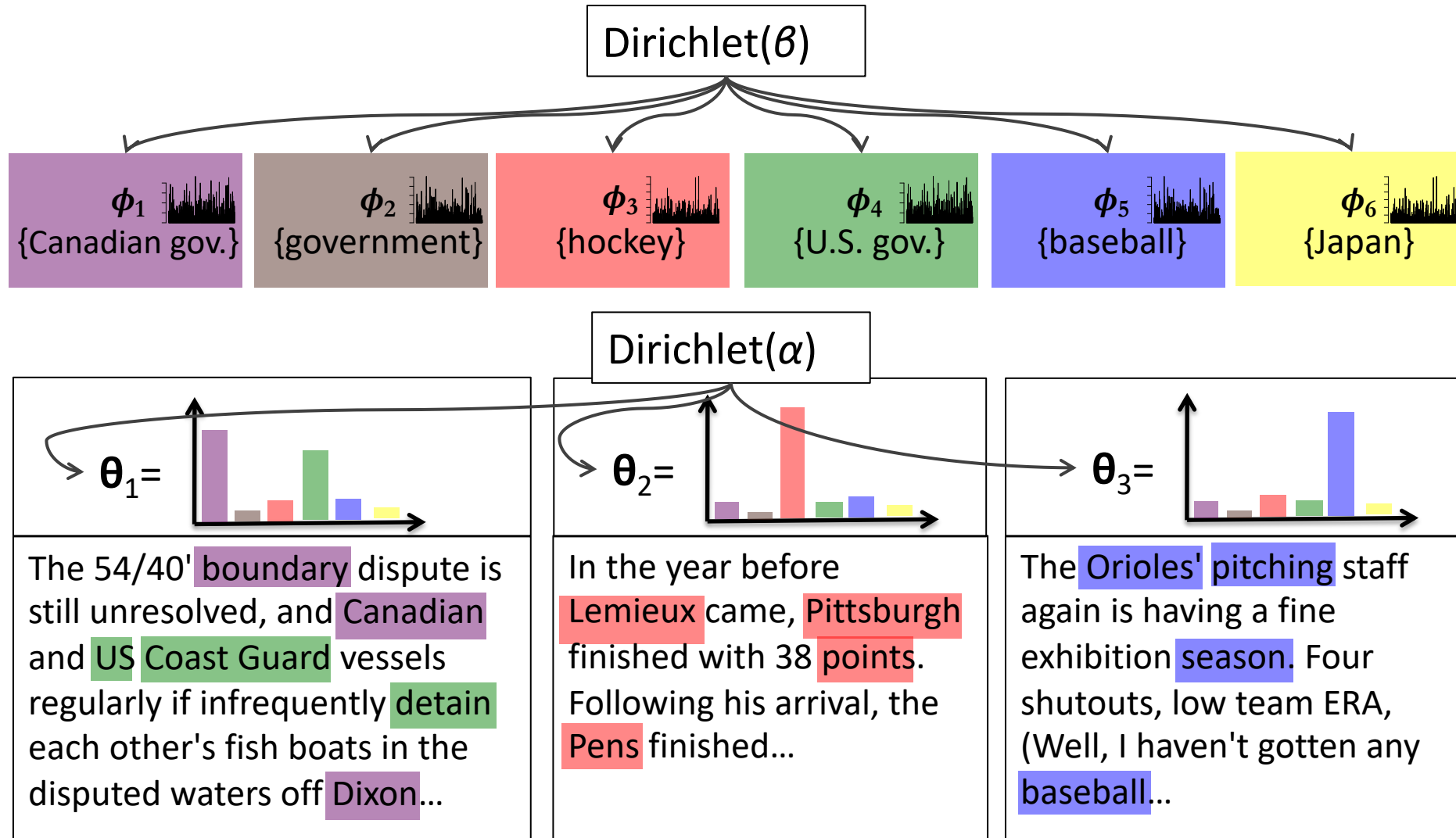
LDA for Topic Modeling



LDA for Topic Modeling

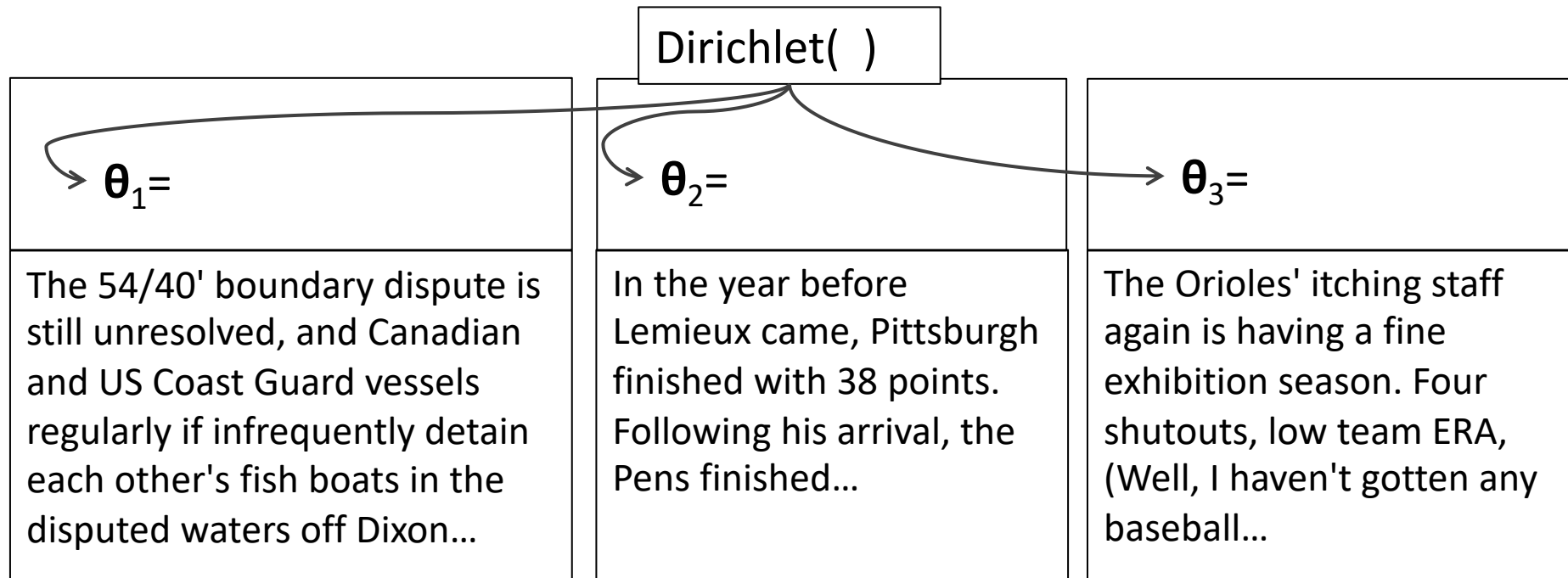
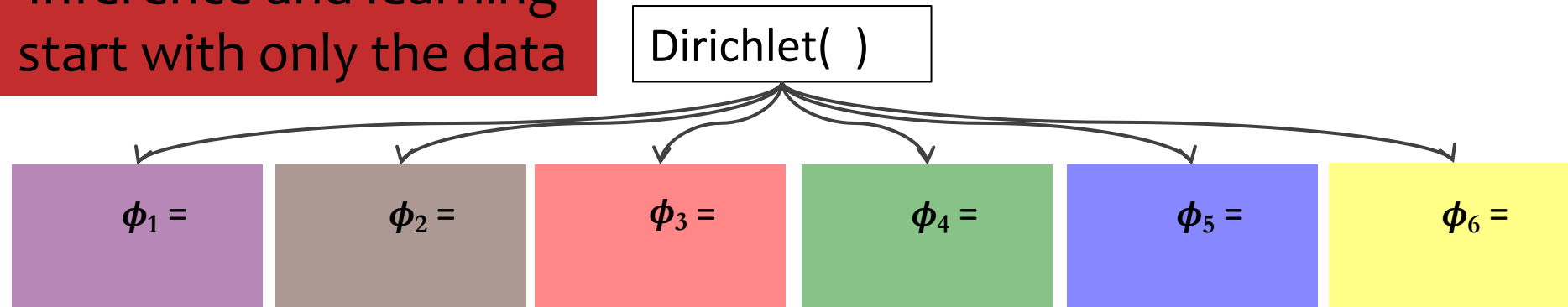


LDA for Topic Modeling



LDA for Topic Modeling

Inference and learning start with only the data



Latent Dirichlet Allocation

Questions:

- Is this a believable story for the generation of a corpus of documents?
- Why might it work well anyway?

Latent Dirichlet Allocation

Why does LDA “work”?

- LDA trades off two goals.
 - ① For each document, allocate its words to as few topics as possible.
 - ② For each topic, assign high probability to as few terms as possible.
- These goals are at odds.
 - Putting a document in a single topic makes #2 hard:
All of its words must have probability under that topic.
 - Putting very few words in each topic makes #1 hard:
To cover a document's words, it must assign many topics to it.
- Trading off these goals finds groups of tightly co-occurring words.

Latent Dirichlet Allocation

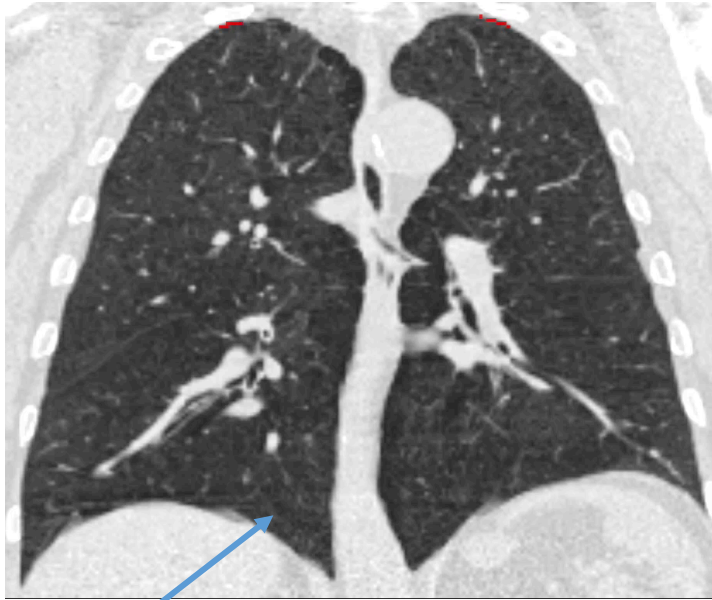
How does this relate to my other favorite model for capturing low-dimensional representations of a corpus?

- Builds on latent semantic analysis (Deerwester et al., 1990; Hofmann, 1999)
- It is a mixed-membership model (Erosheva, 2004).
- It relates to PCA and matrix factorization (Jakulin and Buntine, 2002)
- Was independently invented for genetics (Pritchard et al., 2000)

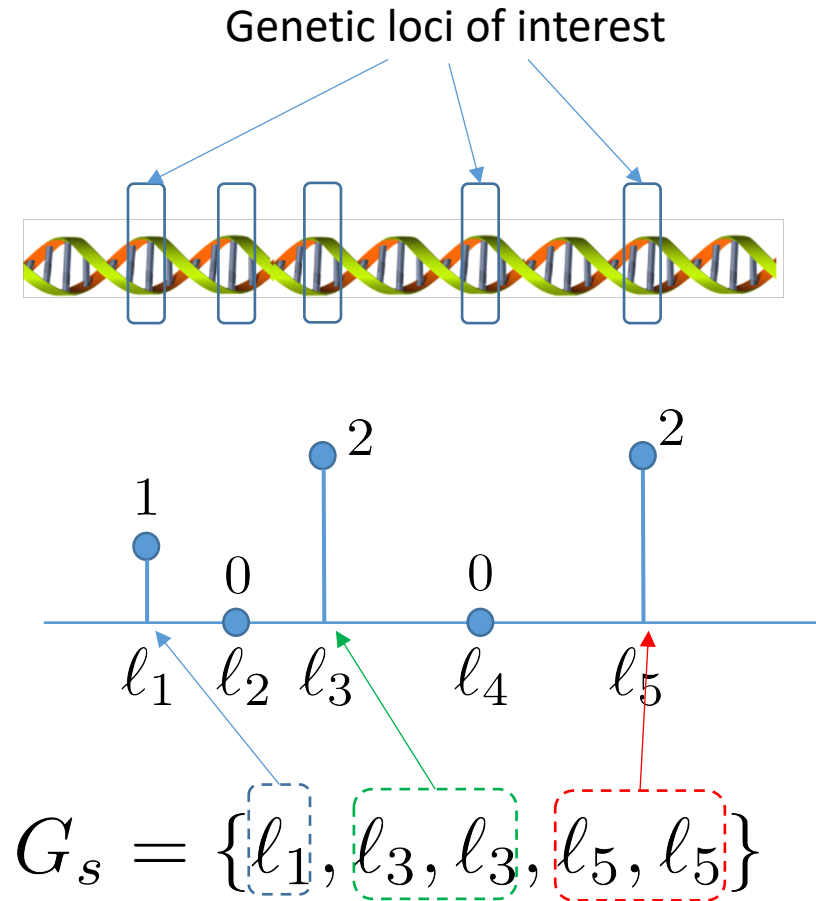
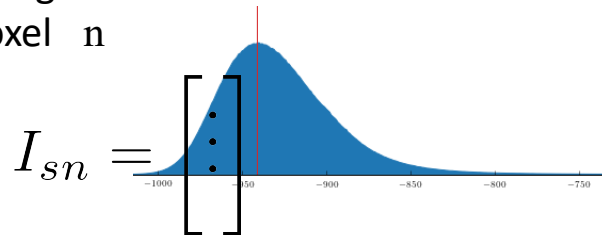
Case Study:
Modeling Joint Imaging and Genetic data

Imaging and Genetic Data

Subject s

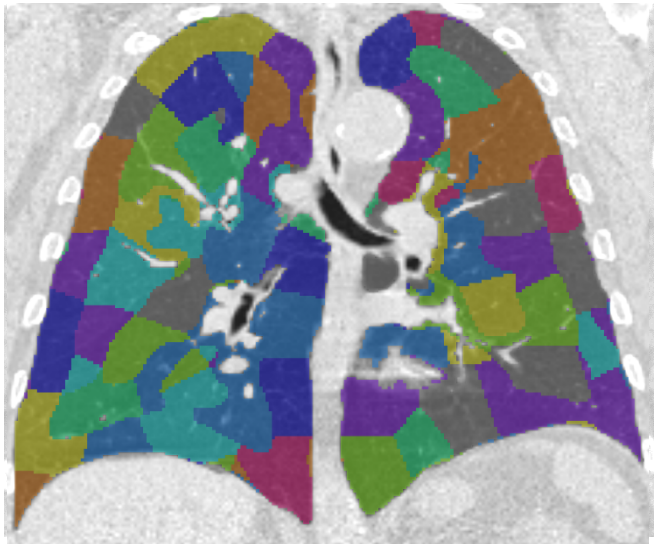


Imaging signature of
Supervoxel n

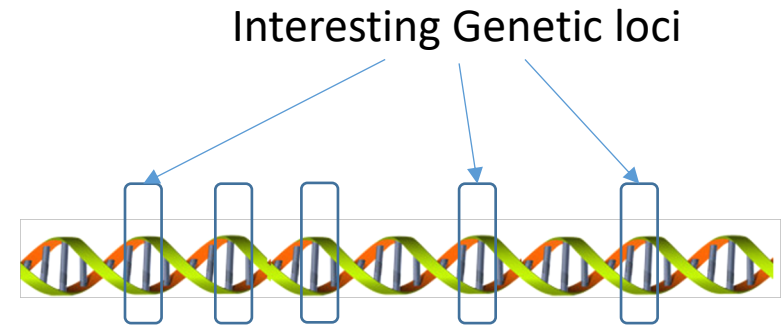


Bag of Words Model

Subject s



Visual Words (I_{sn})



$$G_s = \{\ell_1, \ell_3, \ell_3, \ell_5, \ell_5\}$$

Genetic Words
(Genetic variants)

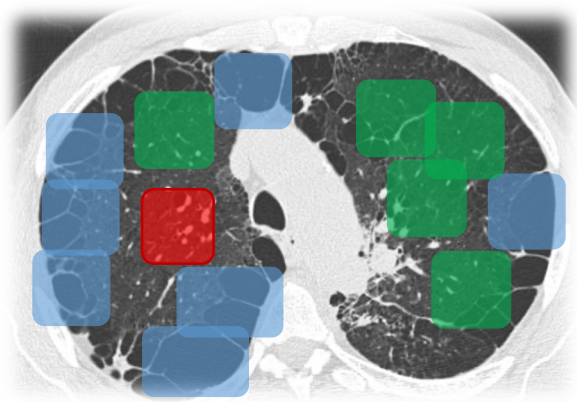
Subject



Document

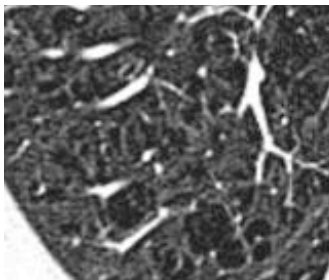
Analogy: Subject as a Document

Pattern 1
Pattern 2
Pattern 3

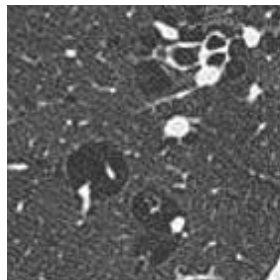


Topics (Image Patterns):

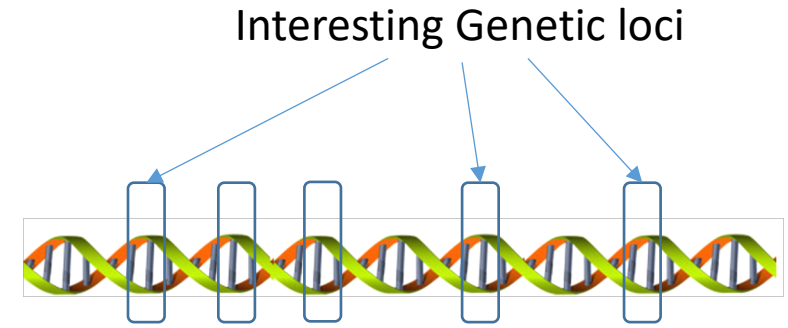
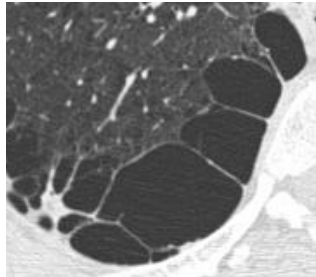
Pattern 1



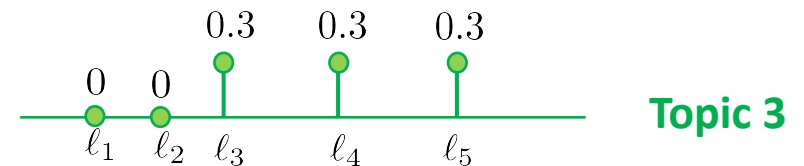
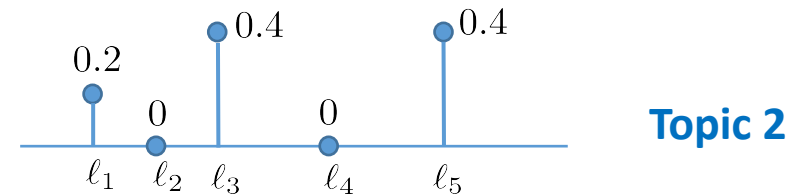
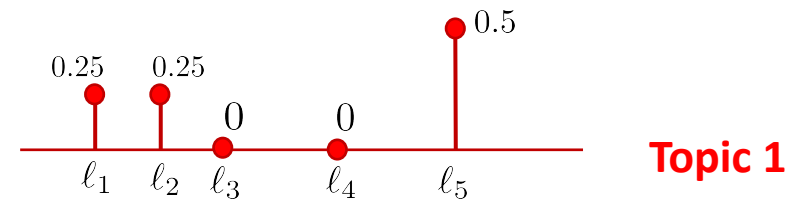
Pattern 2



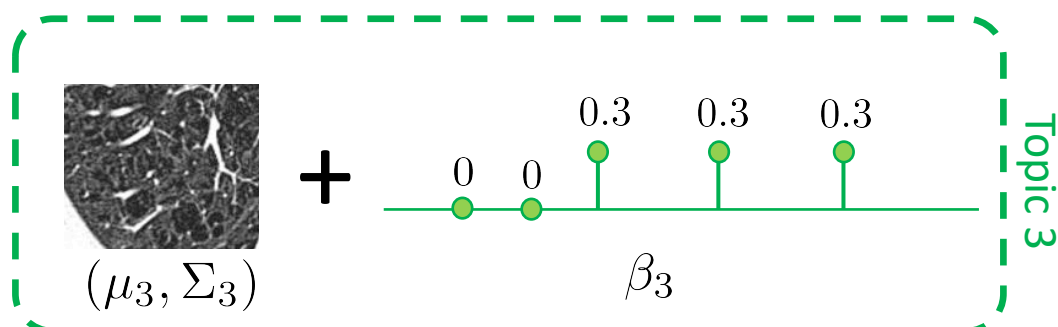
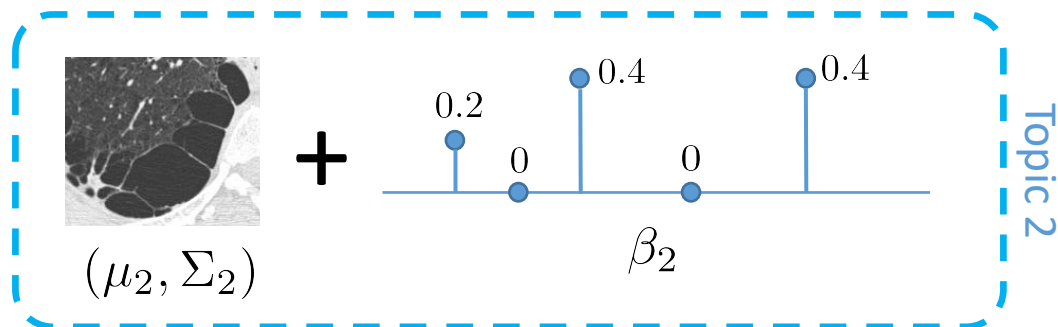
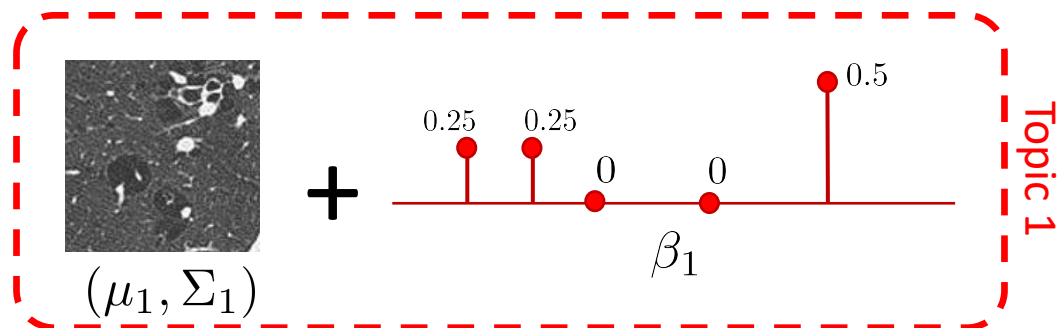
Pattern 3



Distribution of genetic variants



Imaging – Genetic Pair Topics Signatures

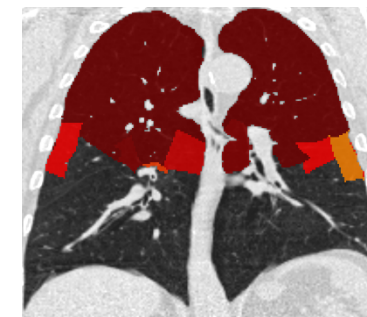


Subject s

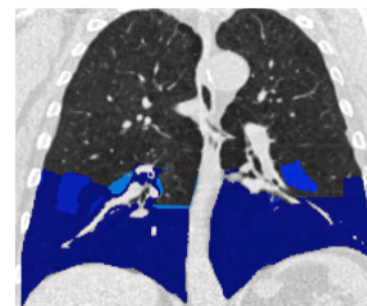
Subject
Proportion

Supervoxel membership

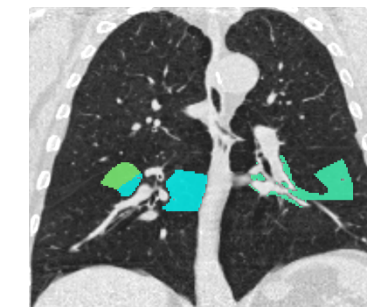
40%



40%



20%



Probabilistic Model

$$(\mu_1, \Sigma_1), \beta_1$$

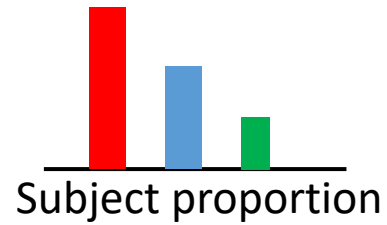
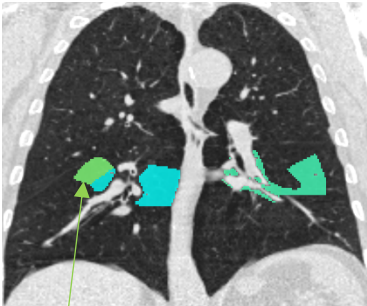
$$(\mu_2, \Sigma_2), \beta_2$$

$$(\mu_3, \Sigma_3), \beta_3$$

• • •

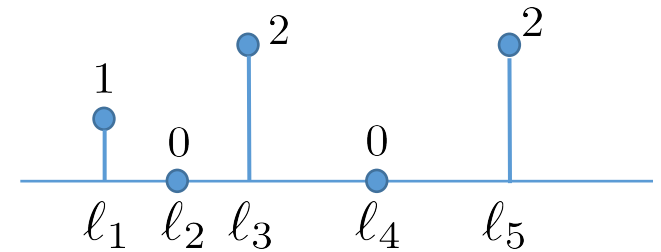
$$(\mu_K, \Sigma_K), \beta_K$$

Subject s



$$z_{sm}^I \sim \text{Cat}(\text{---})$$

$$I_{sm} \sim \mathcal{N}(\cdot; \mu_3, \Sigma_3)$$

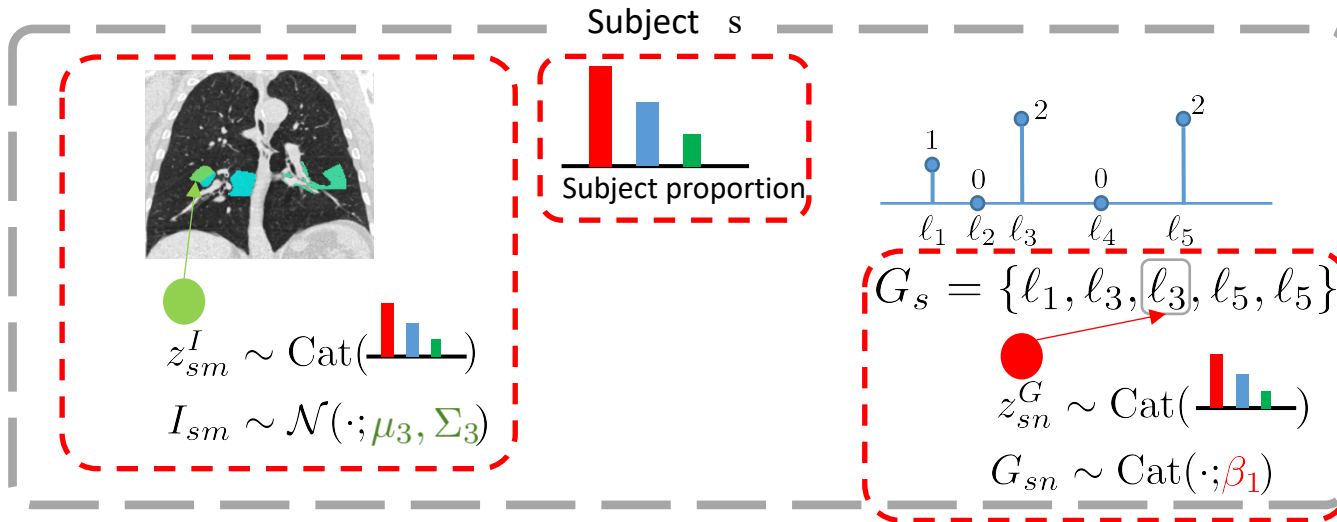
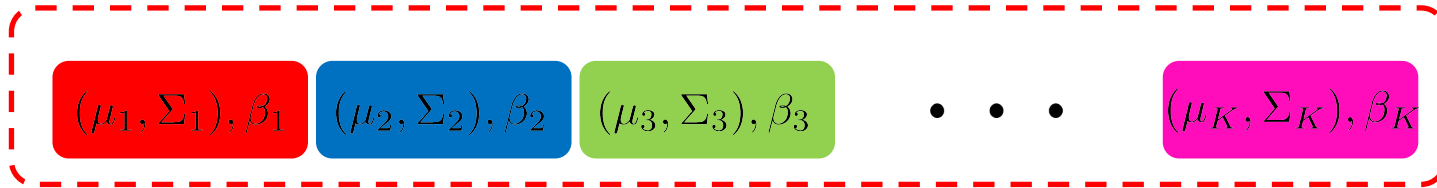


$$G_s = \{l_1, l_3, \boxed{l_3}, l_5, l_5\}$$

$$z_{sn}^G \sim \text{Cat}(\text{---})$$

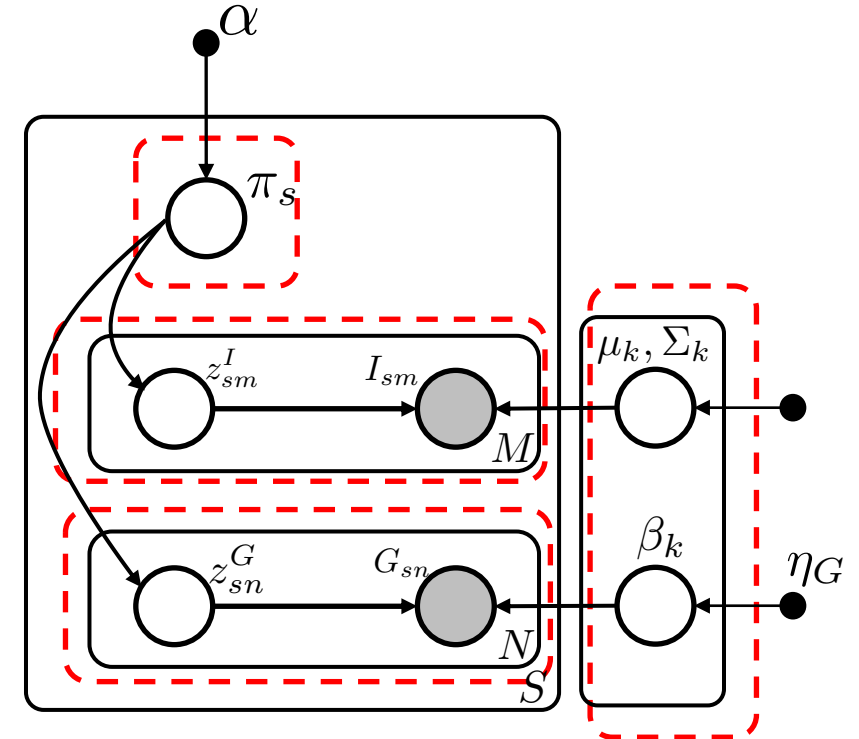
$$G_{sn} \sim \text{Cat}(\cdot; \beta_1)$$

Graphical Model

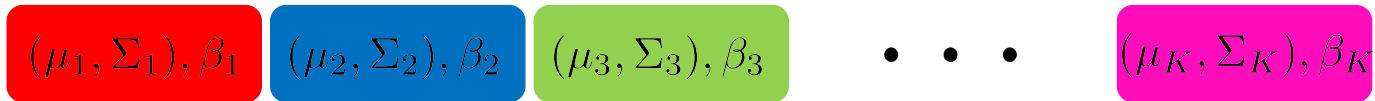


$$(\mu_k, \Sigma_k) \sim \text{NIW}(\eta^I)$$

$$\beta_k \sim \text{Dir}(\eta^G)$$

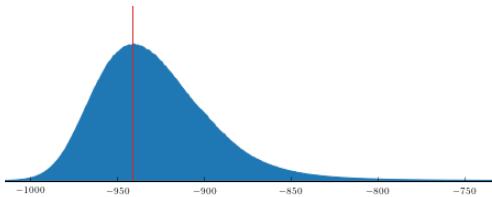


Inference



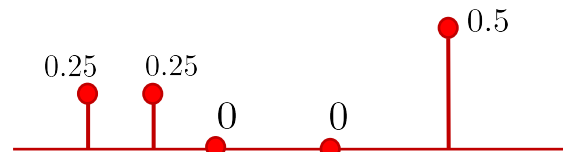
Topic pairs

$$p(\mu_k | \{I_{sm}\}, \{G_{sn}\}; \pi)$$



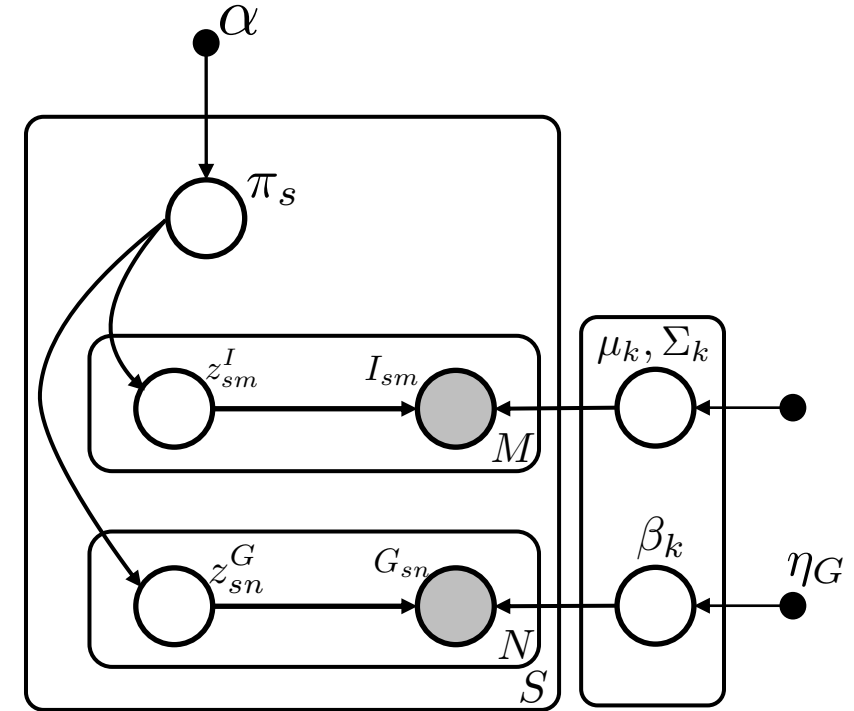
(Imaging signature)

$$p(\beta_k | \{I_{sm}\}, \{G_{sn}\}; \pi)$$



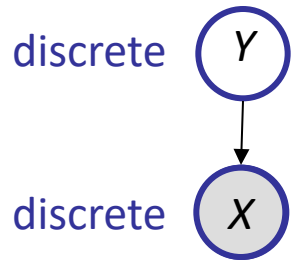
(Different Ranking SNPs)

$$\pi = \{\alpha, \omega, \eta^I, \eta^G\} \text{ (hyper-parameters)}$$

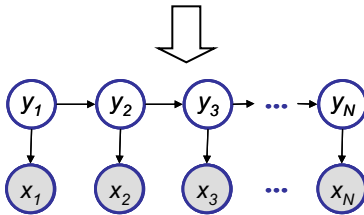


Factor Analysis

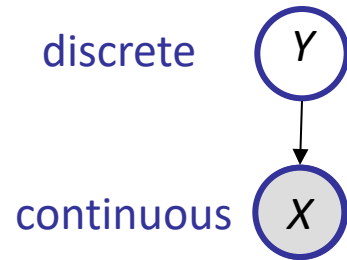
A road map to more complex dynamic models



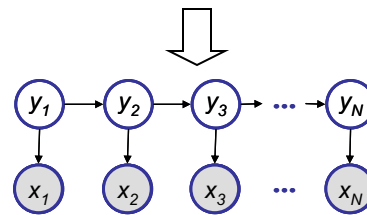
Mixture model
e.g., mixture of multinomials



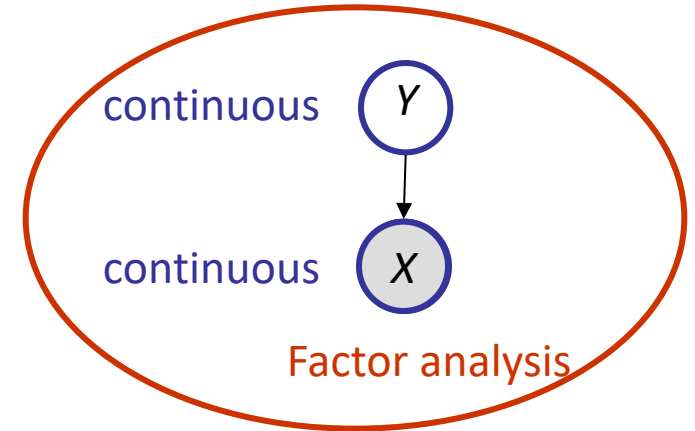
HMM
(for discrete sequential data, e.g., text)



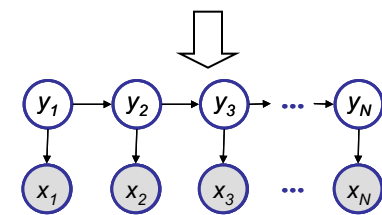
Mixture model
e.g., mixture of Gaussians



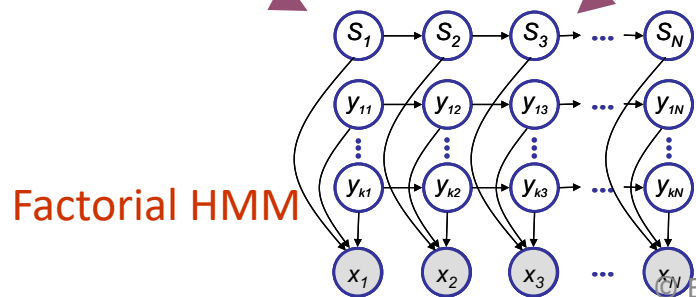
HMM
(for continuous sequential data, e.g., speech signal)



Factor analysis

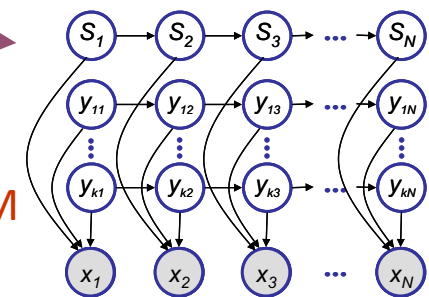


State space model



Factorial HMM

Switching SSM



Recall multivariate Gaussian

- Multivariate Gaussian density:

$$p(\mathbf{x} | \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right\}$$

- A joint Gaussian:

$$p\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \middle| \mu, \Sigma\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \middle| \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

- How to write down $p(\mathbf{x}_1)$, $p(\mathbf{x}_1 | \mathbf{x}_2)$ or $p(\mathbf{x}_2 | \mathbf{x}_1)$ using the block elements in μ and Σ ?
 - Formulas to remember:

$$p(\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_2 | \mathbf{m}_2^m, \mathbf{V}_2^m)$$

$$\mathbf{m}_2^m = \mu_2$$

$$\mathbf{V}_2^m = \Sigma_{22}$$

$$p(\mathbf{x}_1 | \mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 | \mathbf{m}_{12}, \mathbf{V}_{12})$$

$$\mathbf{m}_{12} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{x}_2 - \mu_2)$$

$$\mathbf{V}_{12} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Review: The matrix inverse lemma

- Consider a block-partitioned matrix: $M = \begin{bmatrix} E & F \\ G & H \end{bmatrix}$

- First we diagonalize M

$$\begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} E & F \\ G & H \end{bmatrix} \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} = \begin{bmatrix} E-FH^{-1}G & 0 \\ 0 & H \end{bmatrix}$$

- Schur complement: $M/H = E-FH^{-1}G$
- Then we inverse, using this formula: $XYZ = W \Rightarrow Y^{-1} = ZW^{-1}X$

$$\begin{aligned} M^{-1} &= \begin{bmatrix} E & F \\ G & H \end{bmatrix}^{-1} = \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} \begin{bmatrix} (M/H)^{-1} & 0 \\ 0 & H^{-1} \end{bmatrix} \begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix} \\ &= \begin{bmatrix} (M/H)^{-1} & -(M/H)^{-1}FH^{-1} \\ -H^{-1}G(M/H)^{-1} & H^{-1} + H^{-1}G(M/H)^{-1}FH^{-1} \end{bmatrix} = \begin{bmatrix} E^{-1} + E^{-1}F(M/E)^{-1}GE^{-1} & -E^{-1}F(M/E)^{-1} \\ -(M/E)^{-1}GE^{-1} & (M/E)^{-1} \end{bmatrix} \end{aligned}$$

- Matrix inverse lemma

$$(E-FH^{-1}G)^{-1} = E^{-1} + E^{-1}F(H-GE^{-1}F)^{-1}GE^{-1}$$

Review: Some matrix algebra

- Trace and derivatives

$$\text{tr}[A] \stackrel{\text{def}}{=} \sum_i a_{ii}$$

- Cyclical permutations

$$\text{tr}[ABC] = \text{tr}[CAB] = \text{tr}[BCA]$$

- Derivatives

$$\frac{\partial}{\partial A} \text{tr}[BA] = B^T$$

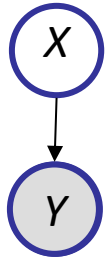
$$\frac{\partial}{\partial A} \text{tr}[x^T A x] = \frac{\partial}{\partial A} \text{tr}[x x^T A] = x x^T$$

- Determinants and derivatives

$$\frac{\partial}{\partial A} \log|A| = A^{-1}$$

Factor analysis

- An **unsupervised linear regression** model

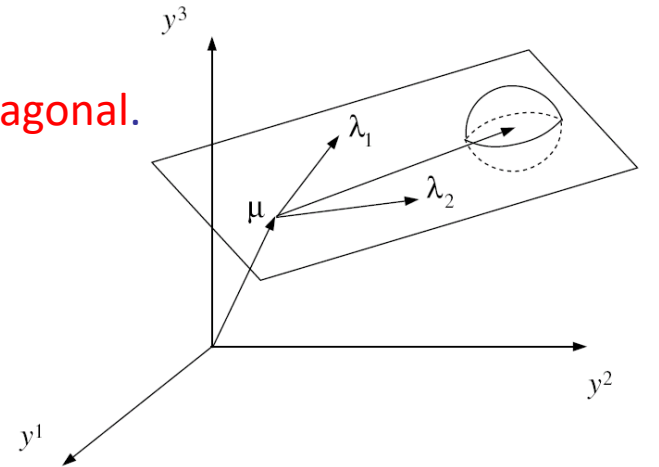


$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, I)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mu + \Lambda \mathbf{x}, \Psi)$$

where Λ is called a factor **loading matrix**, and Ψ is **diagonal**.

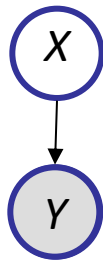
- Geometric interpretation



- To generate data, first generate a point within the manifold then add noise. Coordinates of point are components of latent variable.

Marginal data distribution

- A marginal Gaussian (e.g., $p(\mathbf{x})$) times a conditional Gaussian (e.g., $p(\mathbf{y} | \mathbf{x})$) is a **joint Gaussian**
- Any marginal (e.g., $p(\mathbf{y})$ of a **joint Gaussian** (e.g., $p(\mathbf{x}, \mathbf{y})$) is also a Gaussian
 - Since the marginal is Gaussian, we can determine it by just computing its mean and variance. (Assume noise uncorrelated with data.)



$$\begin{aligned} E[\mathbf{Y}] &= E[\mu + \Lambda \mathbf{X} + \mathbf{W}] && \text{where } \mathbf{W} \sim \mathcal{N}(\mathbf{0}, \Psi) \\ &= \mu + \Lambda E[\mathbf{X}] + E[\mathbf{W}] \\ &= \mu + \mathbf{0} + \mathbf{0} = \mu \end{aligned}$$

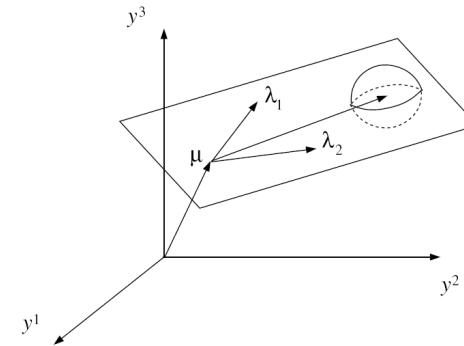
FA = Constrained-Covariance Gaussian

- Marginal density for factor analysis (\mathbf{y} is p -dim, \mathbf{x} is k -dim):

$$p(\mathbf{y} | \theta) = \mathcal{N}(\mathbf{y}; \mu, \Lambda \Lambda^T + \Psi)$$

- So the effective covariance is the low-rank outer product of two long skinny matrices plus a diagonal matrix:

$$\text{Cov}[\mathbf{y}] = \Lambda \Lambda^T + \Psi$$



- In other words, **factor analysis is just a constrained Gaussian model (number of free params of the covariance is limited)**. (If Ψ were not diagonal then we could model any Gaussian and it would be pointless.)

FA **joint** distribution

- Model

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, I)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mu + \Lambda\mathbf{x}, \Psi)$$

- Covariance between \mathbf{x} and \mathbf{y}

$$\begin{aligned} \text{Cov}[\mathbf{X}, \mathbf{Y}] &= E[(\mathbf{X} - \mathbf{0})(\mathbf{Y} - \mu)^T] = E[\mathbf{X}(\mu + \Lambda\mathbf{X} + \mathbf{W} - \mu)^T] \\ &= E[\mathbf{X}\mathbf{X}^T \Lambda^T + \mathbf{X}\mathbf{W}^T] \\ &= \Lambda^T \end{aligned}$$

- Hence the joint distribution of \mathbf{x} and \mathbf{y} :

$$p\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda\Lambda^T + \Psi \end{bmatrix}\right)$$

- Assume noise is uncorrelated with data or latent variables.

$$\begin{bmatrix} (M/H)^{-1} & -(M/H)^{-1}FH^{-1} \\ -H^{-1}G(M/H)^{-1} & H^{-1} + H^{-1}G(M/H)^{-1}FH^{-1} \end{bmatrix} = \begin{bmatrix} E^{-1} + E^{-1}F(M/E)^{-1}GE^{-1} & -E^{-1}F(M/E)^{-1} \\ -(M/E)^{-1}GE^{-1} & (M/E)^{-1} \end{bmatrix}$$

Inference in Factor Analysis

- Apply the Gaussian conditioning formulas to the joint distribution we derived above, where

$$\Sigma_{11} = I$$

$$\Sigma_{12} = \Sigma_{12}^T = \Lambda^T$$

$$\Sigma_{22} = (\Lambda\Lambda^T + \Psi)$$

we can now derive the **posterior** of the latent variable \mathbf{x} given observation \mathbf{y} , $p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x} | \mathbf{m}_{1|2}, \mathbf{V}_{1|2})$, where

$$\mathbf{m}_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{y} - \mu_2)$$

$$\mathbf{V}_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

$$= \Lambda^T(\Lambda\Lambda^T + \Psi)^{-1}(\mathbf{y} - \mu)$$

$$= I - \Lambda^T(\Lambda\Lambda^T + \Psi)^{-1}\Lambda$$

Applying the matrix inversion lemma

$$(E - FH^{-1}G)^{-1} = E^{-1} + E^{-1}F(H - GE^{-1}F)^{-1}GE^{-1}$$

$$\Rightarrow \mathbf{V}_{1|2} = (I + \Lambda^T\Psi^{-1}\Lambda)^{-1} \quad \mathbf{m}_{1|2} = \mathbf{V}_{1|2}\Lambda^T\Psi^{-1}(\mathbf{y} - \mu)$$

- Here we only need to invert a matrix of size $|\mathbf{x}| \times |\mathbf{x}|$, instead of $|\mathbf{y}| \times |\mathbf{y}|$.

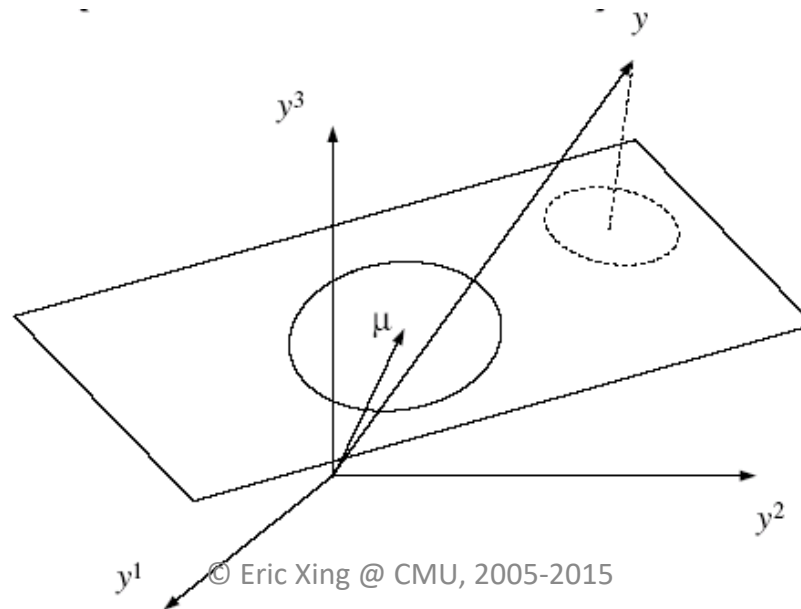
Geometric interpretation: inference is linear projection

- The posterior is:

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}; \mathbf{m}_{1|2}, \mathbf{V}_{1|2})$$

$$\mathbf{m}_{1|2} = \mathbf{V}_{1|2} \Lambda^T \Psi^{-1} (\mathbf{y} - \boldsymbol{\mu}) \quad \mathbf{V}_{1|2} = (\mathbf{I} + \Lambda^T \Psi^{-1} \Lambda)^{-1}$$

- **Posterior covariance** does not depend on observed data \mathbf{y} !
- Computing the **posterior mean** is just a **linear operation**:



Learning FA

- Now, assume that we are given $\{y_n\}$ (the observation on high-dimensional data) only
- We have derived how to estimate x_n from $P(X|Y)$
- How can we learning the model?
 - Loading matrix Λ
 - Manifold center μ
 - Variance Ψ

EM for Factor Analysis

- Incomplete data log likelihood function (marginal density of \mathbf{y})

$$\begin{aligned}\ell(\theta, \mathcal{D}) &= -\frac{N}{2} \log |\Lambda \Lambda^T + \Psi| - \frac{1}{2} \sum_n (\mathbf{y}_n - \boldsymbol{\mu})^T (\Lambda \Lambda^T + \Psi)^{-1} (\mathbf{y}_n - \boldsymbol{\mu}) \\ &= -\frac{N}{2} \log |\Lambda \Lambda^T + \Psi| - \frac{1}{2} \text{tr}[(\Lambda \Lambda^T + \Psi)^{-1} \mathbf{S}], \quad \text{where } \mathbf{S} = \sum_n (\mathbf{y}_n - \boldsymbol{\mu})(\mathbf{y}_n - \boldsymbol{\mu})^T\end{aligned}$$

- Estimating $\boldsymbol{\mu}$ is trivial: $\hat{\boldsymbol{\mu}}^{ML} = \frac{1}{N} \sum_n \mathbf{y}_n$
- Parameters Λ and Ψ are coupled nonlinearly in log-likelihood
- Complete log likelihood

$$\begin{aligned}\ell_c(\theta, \mathcal{D}) &= \sum_n \log p(\mathbf{x}_n, \mathbf{y}_n) = \sum_n \log p(\mathbf{x}_n) + \log p(\mathbf{y}_n | \mathbf{x}_n) \\ &= -\frac{N}{2} \log |I| - \frac{1}{2} \sum_n \mathbf{x}_n^T \mathbf{x}_n - \frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n (\mathbf{y}_n - \Lambda \mathbf{x}_n)^T \Psi^{-1} (\mathbf{y}_n - \Lambda \mathbf{x}_n) \\ &= -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr}[\mathbf{x}_n \mathbf{x}_n^T] - \frac{N}{2} \text{tr}[\mathbf{S} \Psi^{-1}], \quad \text{where } \mathbf{S} = \frac{1}{N} \sum_n (\mathbf{y}_n - \Lambda \mathbf{x}_n)(\mathbf{y}_n - \Lambda \mathbf{x}_n)^T\end{aligned}$$

E-step for Factor Analysis

- Compute $\langle \ell_{\epsilon}(\theta, \mathcal{D}) \rangle_{p(\mathbf{x}|\mathbf{y})}$

$$\langle \ell_{\epsilon}(\theta, \mathcal{D}) \rangle = -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr}[\langle \mathbf{X}_n \mathbf{X}_n^T \rangle] - \frac{N}{2} \text{tr}[\langle \mathbf{S} \rangle \Psi^{-1}]$$

$$\langle \mathbf{S} \rangle = \frac{1}{N} \sum_n (\mathbf{y}_n \mathbf{y}_n^T - \mathbf{y}_n \langle \mathbf{X}_n^T \rangle \Lambda^T - \Lambda \langle \mathbf{X}_n^T \rangle \mathbf{y}_n^T + \Lambda \langle \mathbf{X}_n \mathbf{X}_n^T \rangle \Lambda^T)$$

$$\langle \mathbf{X}_n \rangle = E[\mathbf{X}_n | \mathbf{y}_n]$$

$$\langle \mathbf{X}_n \mathbf{X}_n^T \rangle = \text{Var}[\mathbf{X}_n | \mathbf{y}_n] + E[\mathbf{X}_n | \mathbf{y}_n] E[\mathbf{X}_n | \mathbf{y}_n]^T$$

- Recall that we have derived:

$$\mathbf{V}_{1|2} = (I + \Lambda^T \Psi^{-1} \Lambda)^{-1} \quad \mathbf{m}_{1|2} = \mathbf{V}_{1|2} \Lambda^T \Psi^{-1} (\mathbf{y} - \mu)$$

$$\Rightarrow \langle \mathbf{X}_n \rangle = \mathbf{m}_{\mathbf{x}_n|\mathbf{y}_n} = \mathbf{V}_{1|2} \Lambda^T \Psi^{-1} (\mathbf{y}_n - \mu) \quad \text{and} \quad \langle \mathbf{X}_n \mathbf{X}_n^T \rangle = \mathbf{V}_{1|2} + \mathbf{m}_{\mathbf{x}_n|\mathbf{y}_n} \mathbf{m}_{\mathbf{x}_n|\mathbf{y}_n}^T$$

M-step for Factor Analysis

- Take the derivatives of the expected complete log likelihood wrt. parameters.
 - Using the trace and determinant derivative rules:

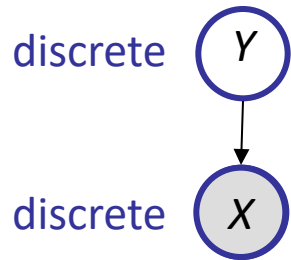
$$\begin{aligned}\frac{\partial}{\partial \Psi^{-1}} \langle \ell_c \rangle &= \frac{\partial}{\partial \Psi^{-1}} \left(-\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr}[\langle \mathbf{x}_n \mathbf{x}_n^T \rangle] - \frac{N}{2} \text{tr}[\langle \mathbf{S} \rangle \Psi^{-1}] \right) \\ &= \frac{N}{2} \Psi - \frac{N}{2} \langle \mathbf{S} \rangle \quad \Rightarrow \quad \Psi^{t+1} = \langle \mathbf{S} \rangle\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \Lambda} \langle \ell_c \rangle &= \frac{\partial}{\partial \Lambda} \left(-\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr}[\langle \mathbf{x}_n \mathbf{x}_n^T \rangle] - \frac{N}{2} \text{tr}[\langle \mathbf{S} \rangle \Psi^{-1}] \right) = -\frac{N}{2} \Psi^{-1} \frac{\partial}{\partial \Lambda} \langle \mathbf{S} \rangle \\ &= -\frac{N}{2} \Psi^{-1} \frac{\partial}{\partial \Lambda} \left(\frac{1}{N} \sum_n (\mathbf{y}_n \mathbf{y}_n^T - \mathbf{y}_n \langle \mathbf{x}_n^T \rangle \Lambda^T - \Lambda \langle \mathbf{x}_n^T \rangle \mathbf{y}_n^T + \Lambda \langle \mathbf{x}_n \mathbf{x}_n^T \rangle \Lambda^T) \right) \\ &= \Psi^{-1} \sum_n \mathbf{y}_n \langle \mathbf{x}_n^T \rangle - \Psi^{-1} \Lambda \sum_n \langle \mathbf{x}_n \mathbf{x}_n^T \rangle \quad \Rightarrow \quad \Lambda^{t+1} = \left(\sum_n \mathbf{y}_n \langle \mathbf{x}_n^T \rangle \right) \left(\sum_n \langle \mathbf{x}_n \mathbf{x}_n^T \rangle \right)^{-1}\end{aligned}$$

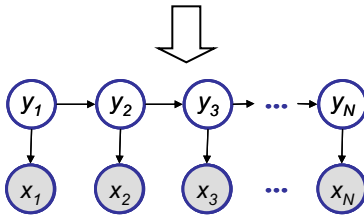
Model Invariance and Identifiability

- There is *degeneracy* in the FA model.
- Since Λ only appears as outer product $\Lambda\Lambda^T$, the model is invariant to rotation and axis flips of the latent space.
- We can replace Λ with ΛQ for any orthonormal matrix Q and the model remains the same: $(\Lambda Q)(\Lambda Q)^T = \Lambda(QQ^T)\Lambda^T = \Lambda\Lambda^T$.
- This means that there is no “one best” setting of the parameters. An infinite number of parameters all give the ML score!
- Such models are called *un-identifiable* since two people both fitting ML parameters to the identical data will not be guaranteed to identify the same parameters.

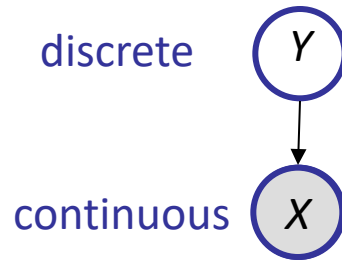
A road map to more complex dynamic models



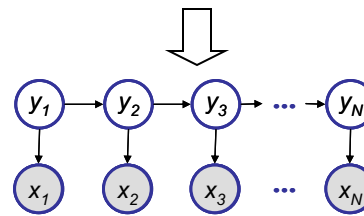
Mixture model
e.g., mixture of multinomials



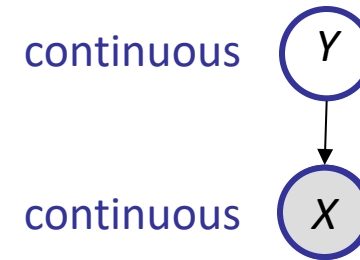
HMM
(for discrete sequential data, e.g., text)



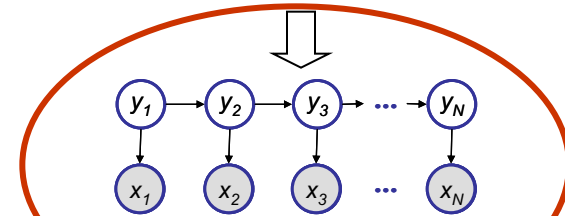
Mixture model
e.g., mixture of Gaussians



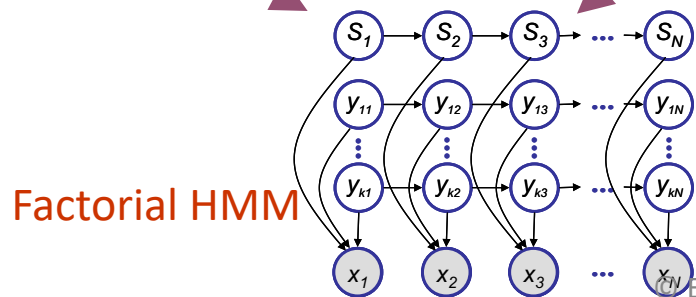
HMM
(for continuous sequential data, e.g., speech signal)



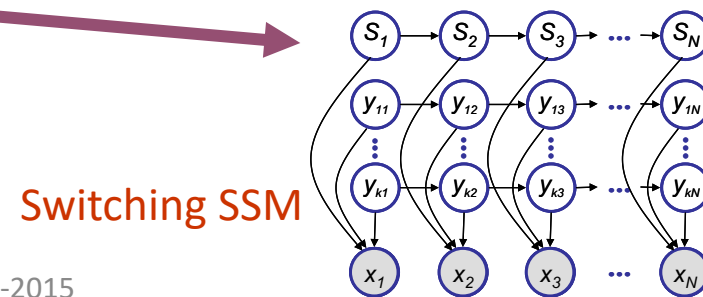
Factor analysis



State space model



Factorial HMM



Switching SSM