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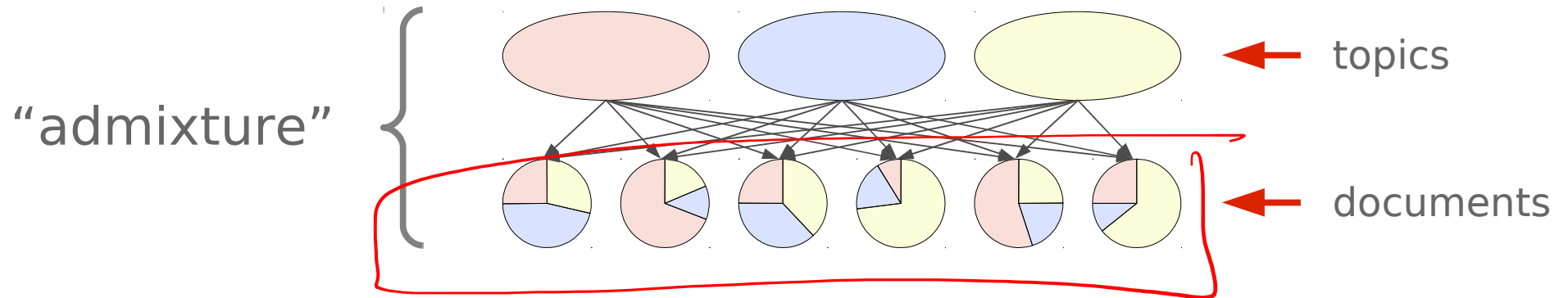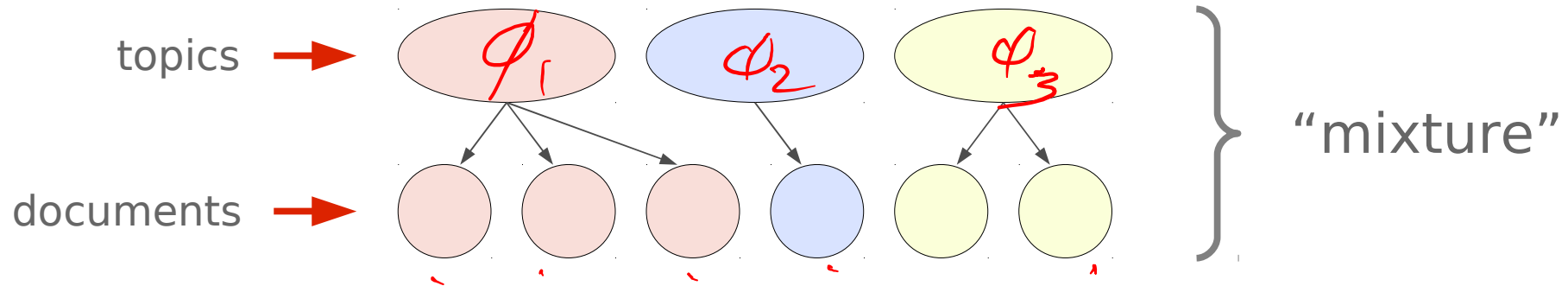# Factor Analysis

**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
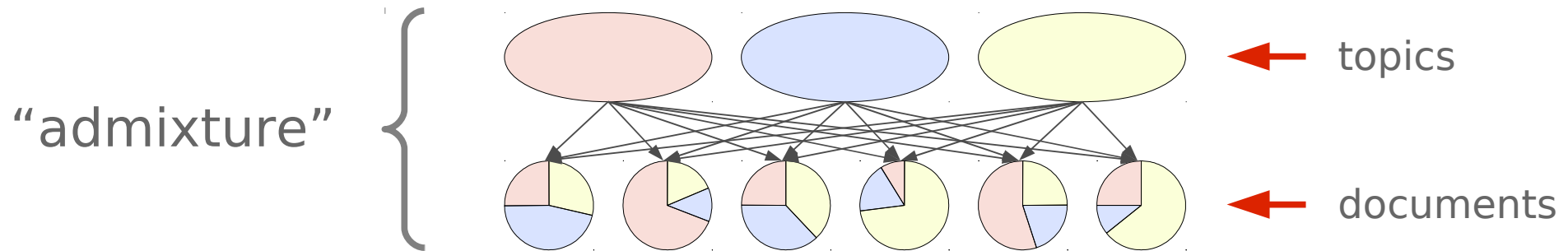- Understand how **authorship** influences the content

# Mixture vs. Admixture (LDA)



topics →  $\phi_1$   $\phi_2$   $\phi_3$   } "mixture"

documents →

"admixture" {   ← topics

← documents

Diagrams from Wallach, JHU 2011, slides

# Latent Dirichlet Allocation

- Generative Process



- Example corpus

# Latent Dirichlet Allocation

- Generative Process

For each topic $k \in \{1, \ldots, K\}$:
$\phi_k \sim \text{Dir}(\boldsymbol{\beta})$       *Dir*         [*draw distribution over words*]
For each document $m \in \{1, \ldots, M\}$
$\boldsymbol{\theta}_m \sim \text{Dir}(\boldsymbol{\alpha})$    $\Theta_k$   # topic [*draw distribution over topics*]
    For each word $n \in \{1, \ldots, N_m\}$
      $z_{mn} \sim \text{Mult}(1, \boldsymbol{\theta}_m)$        [*draw topic assignment*]
      $x_{mn} \sim \boldsymbol{\phi}_{z_{mi}}$            [*draw word*]

*(handwritten annotations in red:)* Docum Specif Dist of topic

- 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



K  #topic

dic (word)

$\theta_m \in \mathbb{R}^K$

$Dir(\alpha)$

$Dir(\beta)$

$\phi_k \in \mathbb{R}^D$
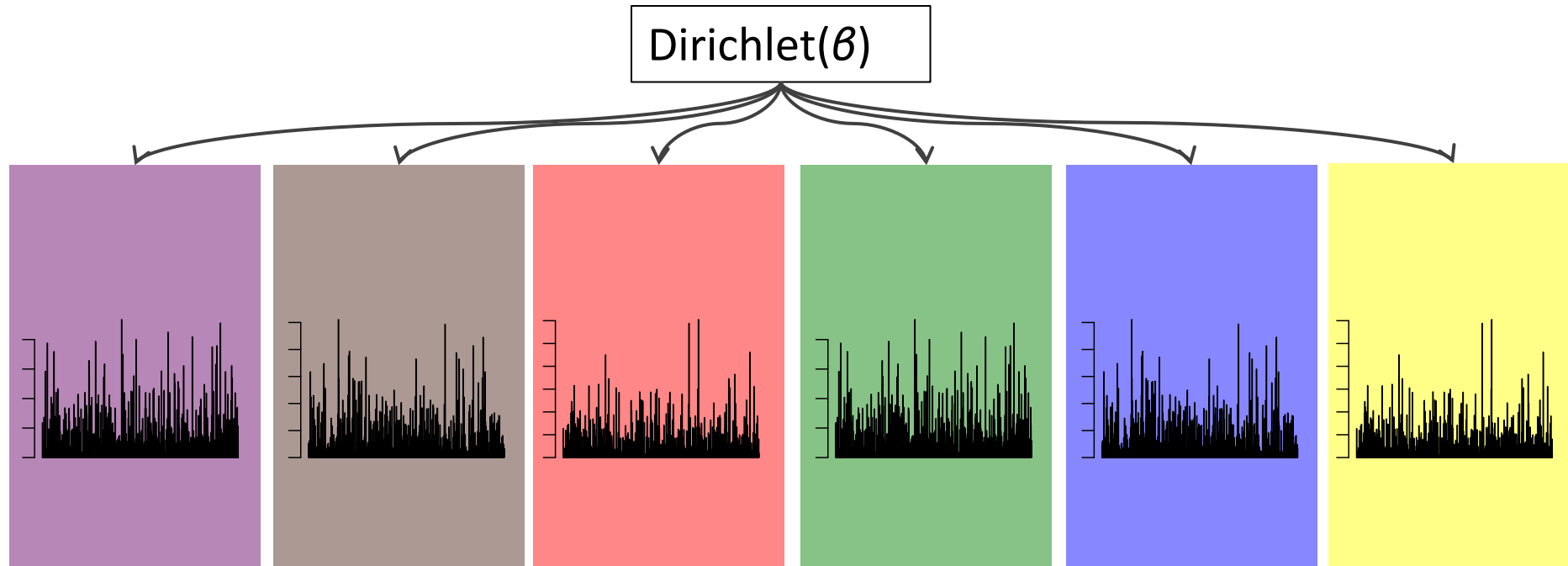
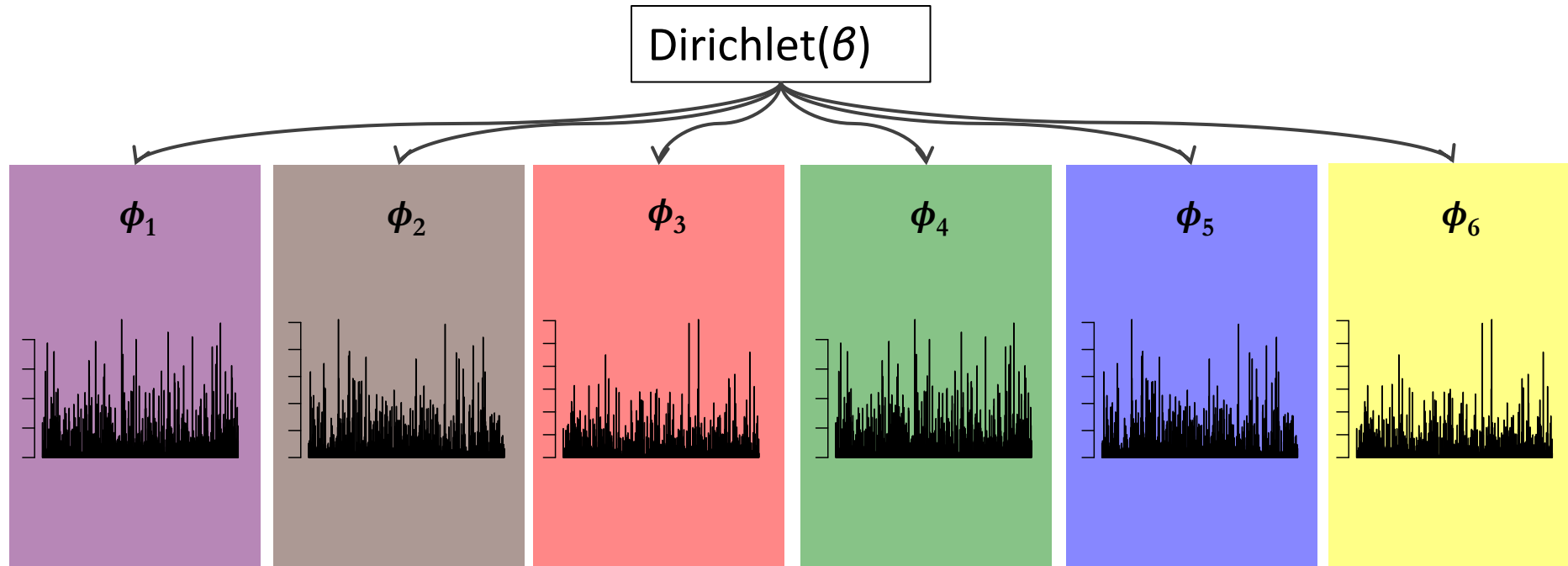# Latent Dirichlet Allocation

- Plate Diagram
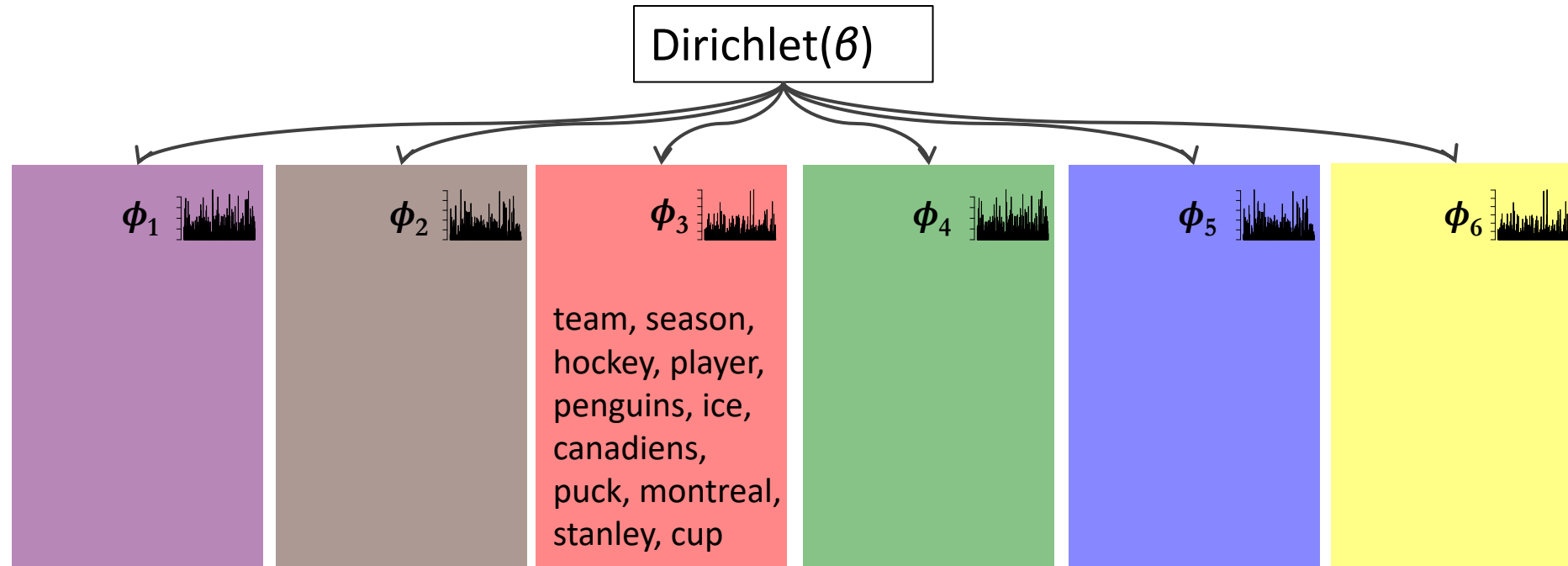
# LDA for Topic Modeling



- 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 $\phi_k$
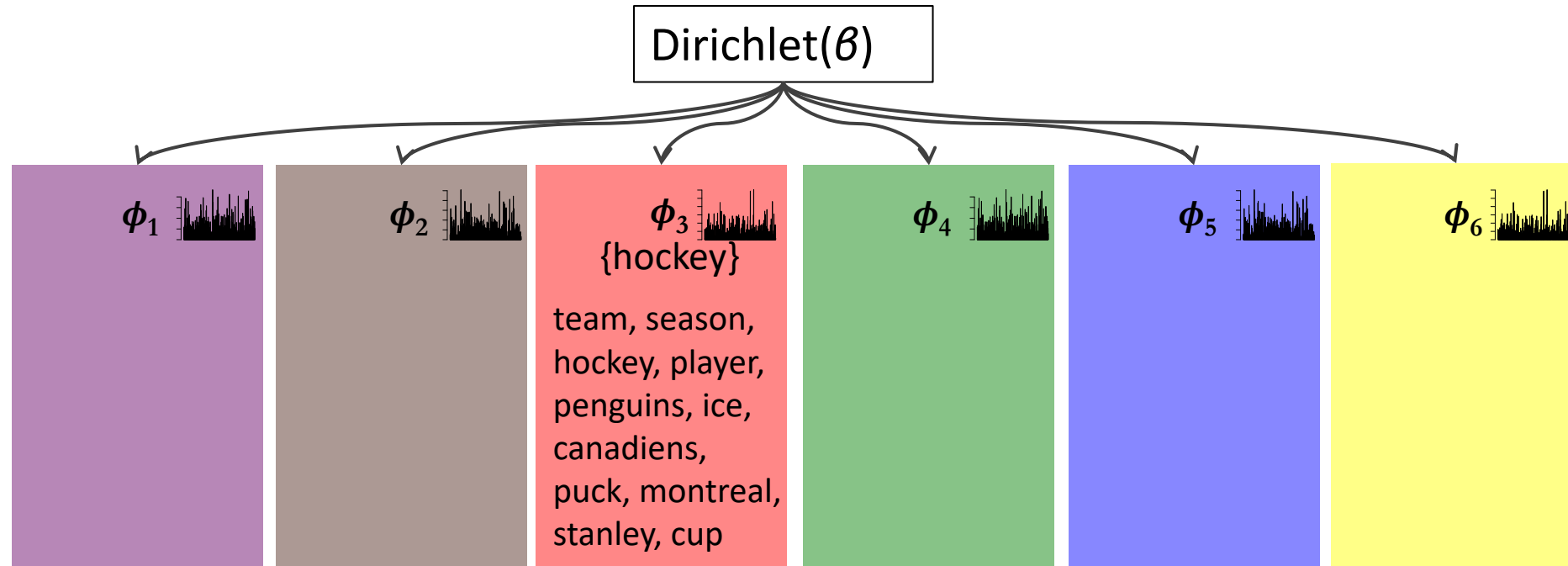
# LDA for Topic Modeling



- 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 $\phi_k$
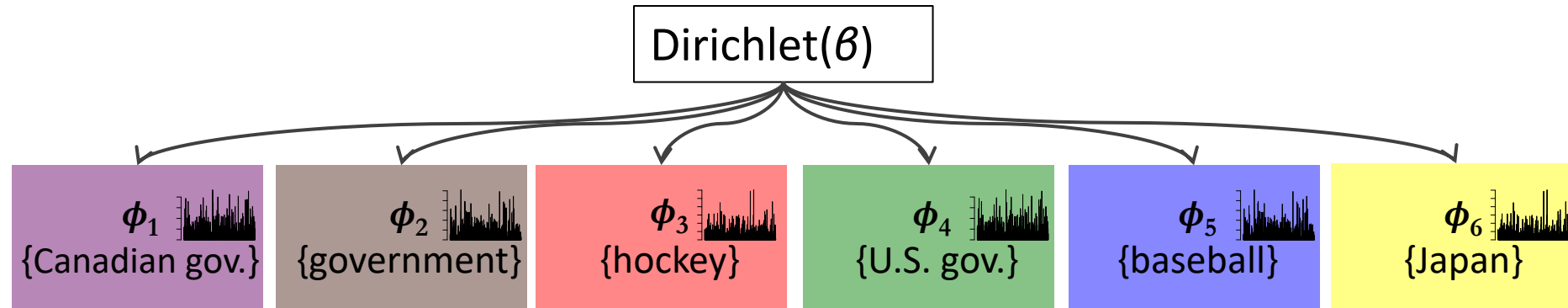
# LDA for Topic Modeling

Dirichlet($\beta$)

$\phi_1$　$\phi_2$　$\phi_3$　$\phi_4$　$\phi_5$　$\phi_6$

team, season, hockey, player, penguins, ice, canadiens, puck, montreal, stanley, cup

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

# LDA for Topic Modeling



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

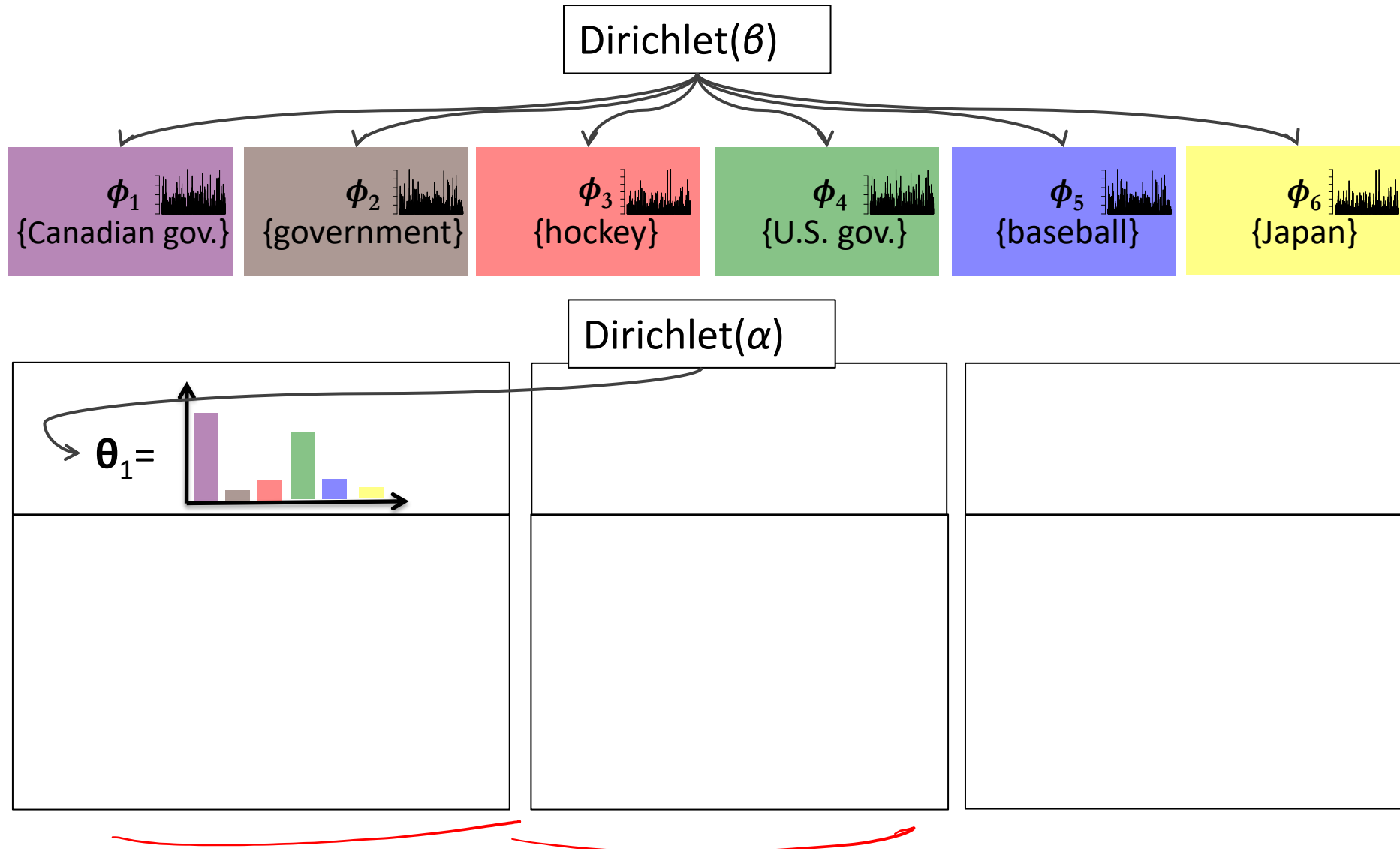- A pedagogical **label** is used to identify the topic.
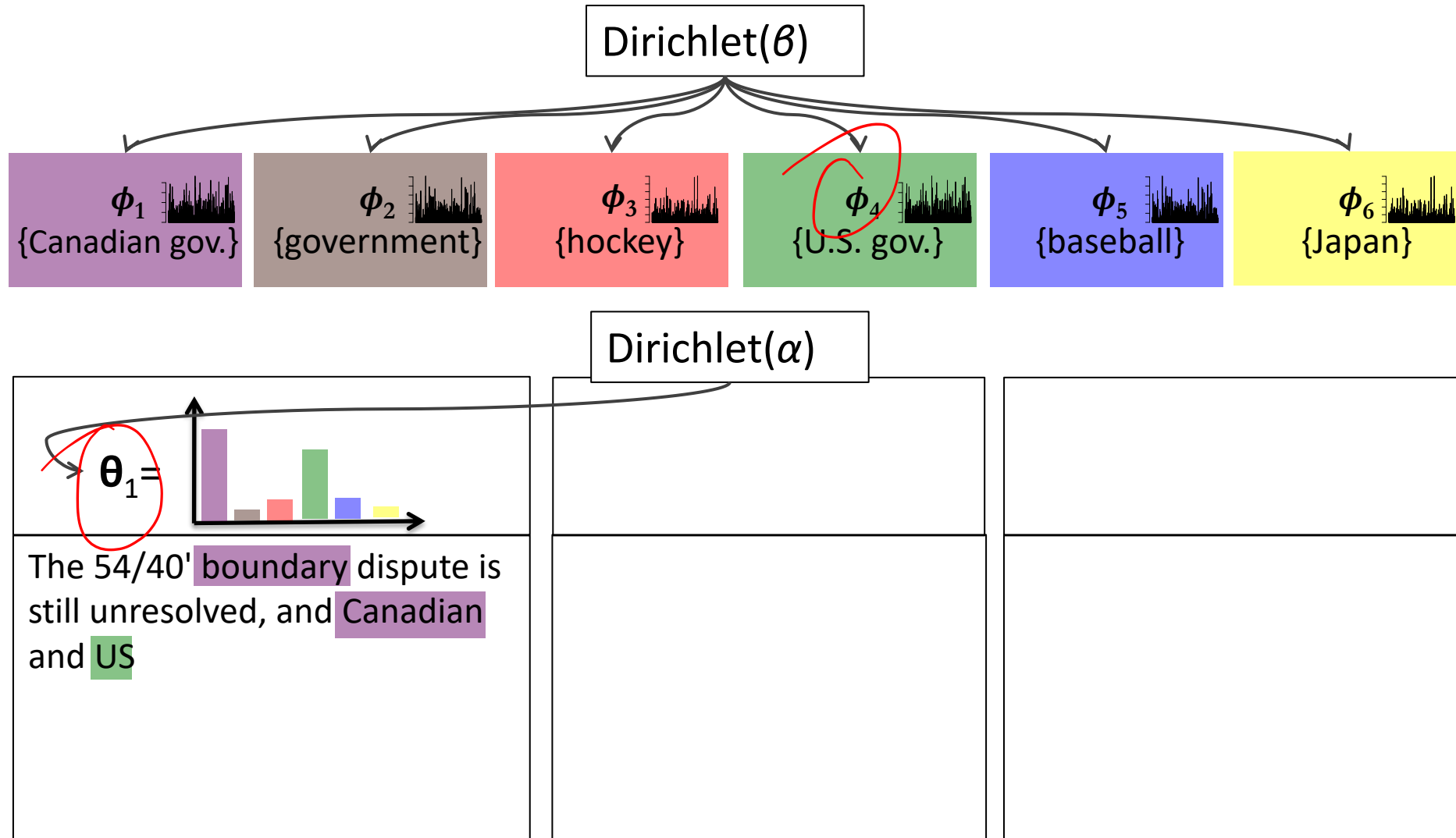
# LDA for Topic Modeling

Dirichlet($\beta$)

$\phi_1$ {Canadian gov.} | $\phi_2$ {government} | $\phi_3$ {hockey} | $\phi_4$ {U.S. gov.} | $\phi_5$ {baseball} | $\phi_6$ {Japan}

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

26

# LDA for Topic Modeling

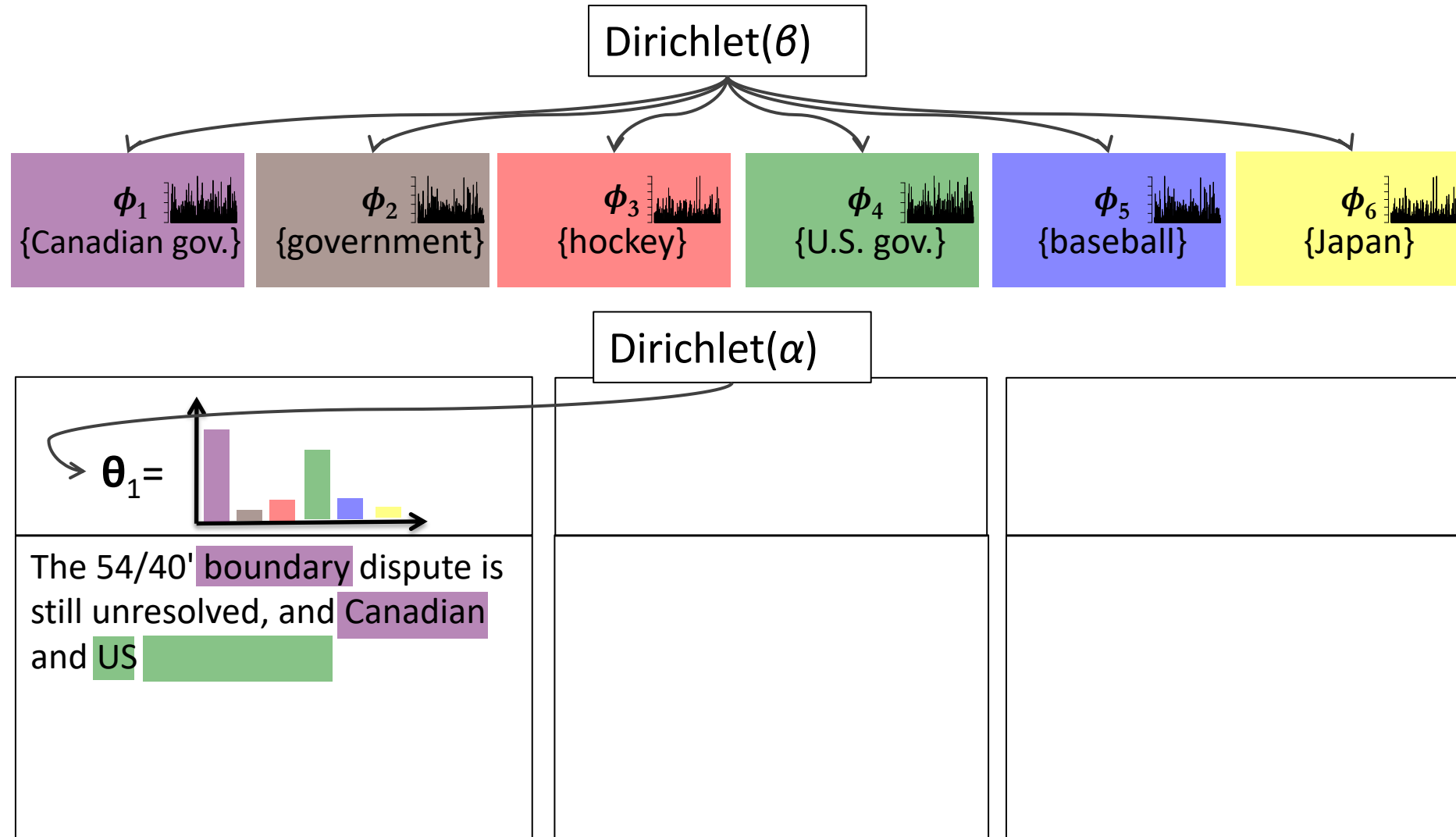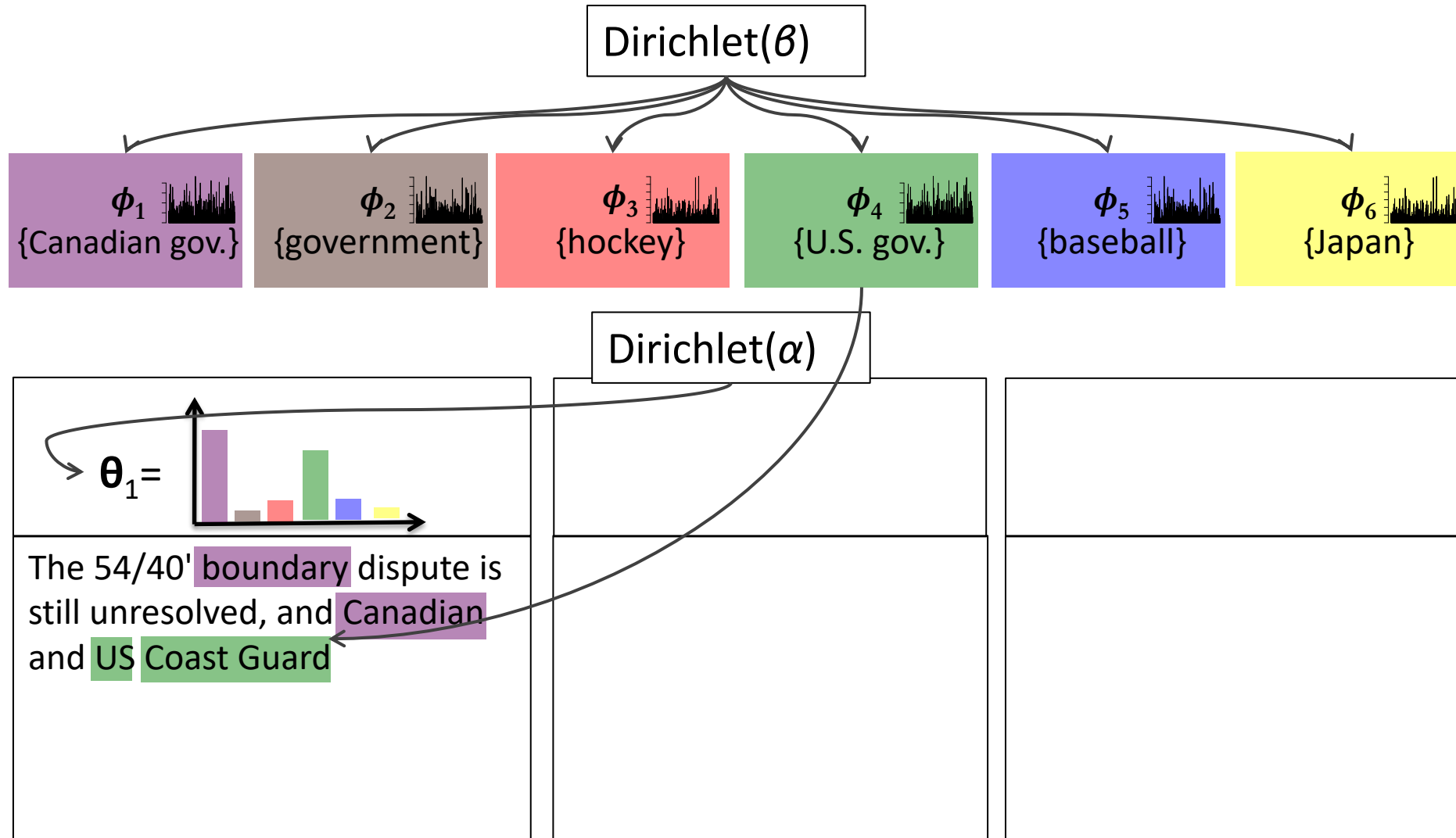# LDA for Topic Modeling

(Blei, Ng, & Jordan, 2003)

# LDA for Topic Modeling

Dirichlet($\beta$)

$\phi_1$ {Canadian gov.}   $\phi_2$ {government}   $\phi_3$ {hockey}   $\phi_4$ {U.S. gov.}   $\phi_5$ {baseball}   $\phi_6$ {Japan}

Dirichlet($\alpha$)

$\theta_1 =$

The 54/40' boundary dispute is still unresolved, and Canadian and US

# LDA for Topic Modeling

Dirichlet($\beta$)

$\phi_1$ {Canadian gov.}

$\phi_2$ {government}

$\phi_3$ {hockey}

$\phi_4$ {U.S. gov.}

$\phi_5$ {baseball}

$\phi_6$ {Japan}

Dirichlet($\alpha$)

$\theta_1 =$

The 54/40' boundary dispute is still unresolved, and Canadian and US Coast Guard

# LDA for Topic Modeling

Dirichlet($\beta$)

$\phi_1$ {Canadian gov.}

$\phi_2$ {government}

$\phi_3$ {hockey}

$\phi_4$ {U.S. gov.}

$\phi_5$ {baseball}

$\phi_6$ {Japan}

Dirichlet($\alpha$)

$\theta_1 =$

The 54/40' boundary dispute is still unresolved, and Canadian and US Coast Guard vessels regularly if infrequently detain each other's fish boats in the disputed waters off Dixon…

# LDA for Topic Modeling

# LDA for Topic Modeling

(Blei, Ng, & Jordan, 2003)

**Distributions over words (topics)**

Dirichlet($\beta$)

$\phi_1$ {Canadian gov.}  $\phi_2$ {government}  $\phi_3$ {hockey}  $\phi_4$ {U.S. gov.}  {baseball}  {Japan}

**Distributions over topics (docs)**

Dirichlet($\alpha$)

$\theta_1 =$

The 54/40' boundary dispute is still unresolved, and Canadian and US Coast Guard vessels regularly if infrequently detain each other's fish boats in the disputed waters off Dixon...

$\theta_2 =$

In the year before Lemieux came, Pittsburgh finished with 38 points. Following his arrival, the Pens finished...

The Orioles' pitching staff again is having a fine exhibition season. Four shutouts, low team ERA, (Well, I haven't gotten any baseball...

33

# LDA for Topic Modeling

Dirichlet($\beta$)

$\phi_1$ {Canadian gov.}   $\phi_2$ {government}   $\phi_3$ {hockey}   $\phi_4$ {U.S. gov.}   $\phi_5$ {baseball}   $\phi_6$ {Japan}

Dirichlet($\alpha$)

$\theta_1 =$

The 54/40' boundary dispute is still unresolved, and Canadian and US Coast Guard vessels regularly if infrequently detain each other's fish boats in the disputed waters off Dixon…

$\theta_2 =$

In the year before Lemieux came, Pittsburgh finished with 38 points. Following his arrival, the Pens finished…

$\theta_3 =$

The Orioles' pitching staff again is having a fine exhibition season. Four shutouts, low team ERA, (Well, I haven't gotten any baseball…

# 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 Join Imaging and Genetic data

# Imaging and Genetic Data

## Subject $s$



Genetic loci of interest

Imaging signature of
Supervoxel  n

$$I_{sn} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

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

# Bag of Words Model



Subject $s$

Interesting Genetic loci

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
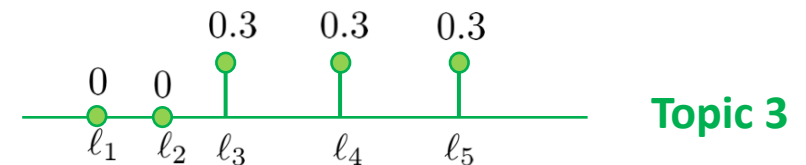Pattern 2
Pattern 3

**Topics (Image Patterns):**

Pattern 1    Pattern 2    Pattern 3

Interesting Genetic loci

Distribution of genetic variants

0.25  0.25  0  0  0.5
$\ell_1$  $\ell_2$  $\ell_3$  $\ell_4$  $\ell_5$    Topic 1

0.2  0  0.4  0  0.4
$\ell_1$  $\ell_2$  $\ell_3$  $\ell_4$  $\ell_5$    Topic 2

0  0  0.3  0.3  0.3
$\ell_1$  $\ell_2$  $\ell_3$  $\ell_4$  $\ell_5$    Topic 3

42

Imaging – Genetic Pair Topics Signatures

Topic 1: $(\mu_1, \Sigma_1)$ + $\beta_1$ with values 0.25, 0.25, 0, 0, 0.5

Topic 2: $(\mu_2, \Sigma_2)$ + $\beta_2$ with values 0.2, 0, 0.4, 0, 0.4

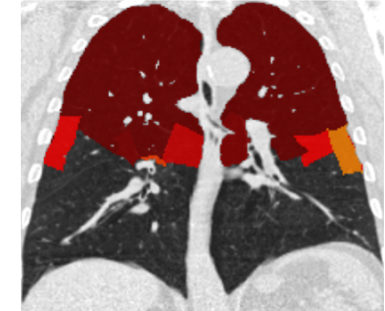Topic 3: $(\mu_3, \Sigma_3)$ + $\beta_3$ with values 0, 0, 0.3, 0.3, 0.3
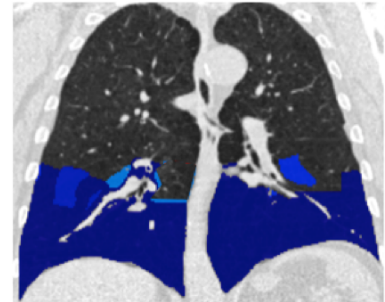
Subject $s$

Subject Proportion — Supervoxel membership

40%

40%

20%

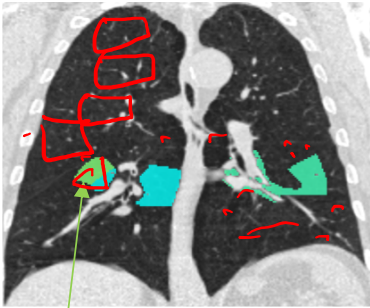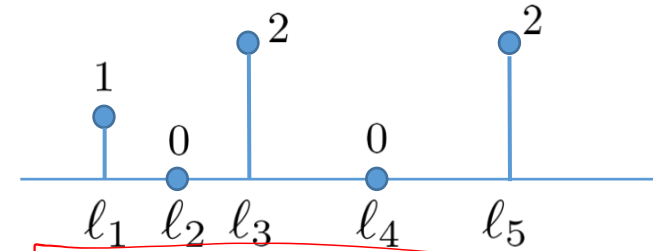# Probabilistic Model

$(\mu_1, \Sigma_1), \beta_1$    $(\mu_2, \Sigma_2), \beta_2$    $(\mu_3, \Sigma_3), \beta_3$    $\bullet \ \bullet \ \bullet$    $(\mu_K, \Sigma_K), \beta_K$

## Subject  s



Subject proportion

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

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

$1 \qquad 2 \qquad\qquad 2$
$0 \qquad\qquad 0$
$\ell_1 \quad \ell_2 \ \ell_3 \qquad \ell_4 \qquad \ell_5$

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

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

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

# Graphical Model



$$(\mu_1, \Sigma_1), \beta_1 \quad (\mu_2, \Sigma_2), \beta_2 \quad (\mu_3, \Sigma_3), \beta_3 \quad \cdots \quad (\mu_K, \Sigma_K), \beta_K$$

Subject $s$

Subject proportion

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

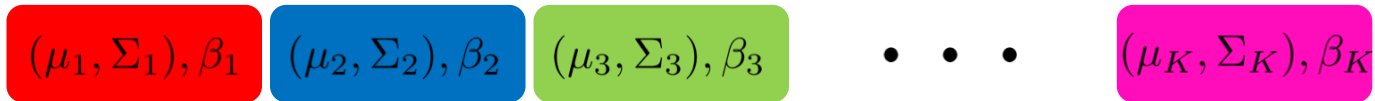$$z_{sm}^I \sim \mathrm{Cat}(\quad)$$

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

$$z_{sn}^G \sim \mathrm{Cat}(\quad)$$

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

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

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

Patient specif distri

label image Patch

label genetic

45

# Inference

$(\mu_1, \Sigma_1), \beta_1$  $(\mu_2, \Sigma_2), \beta_2$  $(\mu_3, \Sigma_3), \beta_3$  $\cdots$  $(\mu_K, \Sigma_K), \beta_K$
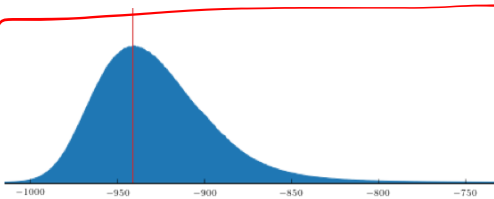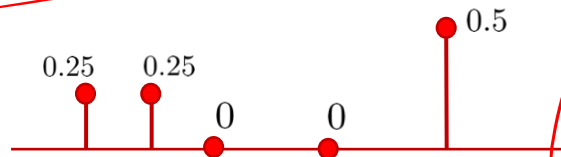
## Topic pairs

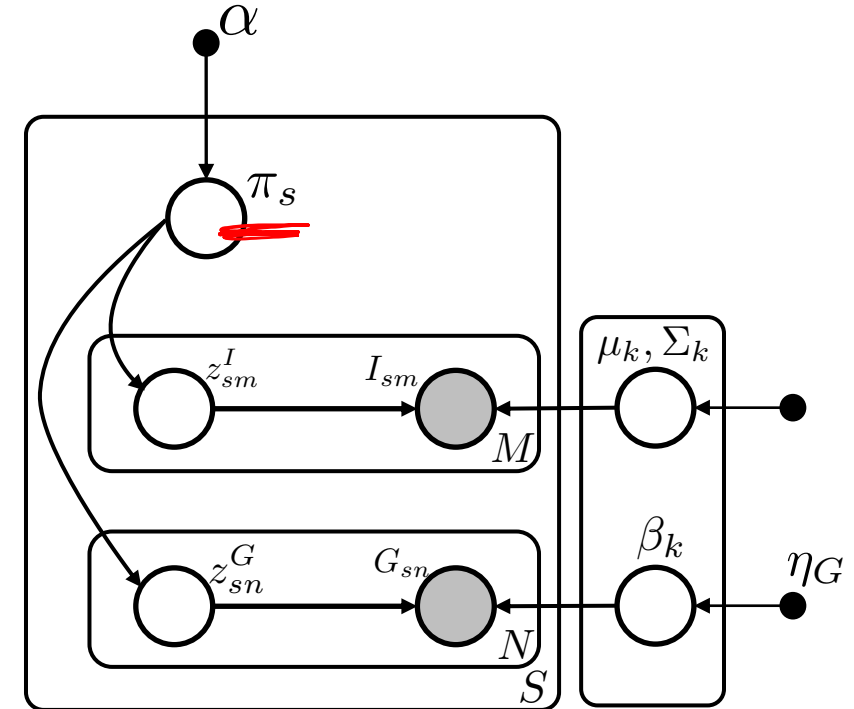$$p(\mu_k | \{I_{sm}\}, \{G_{sn}\}; \pi)$$

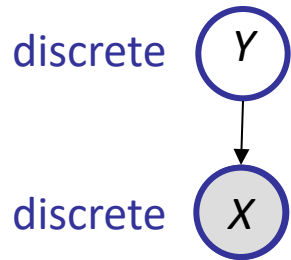$$p(\beta_k | \{I_{sm}\}, \{G_{sn}\}; \pi)$$

(Imaging signature)

(Different Ranking SNPs)

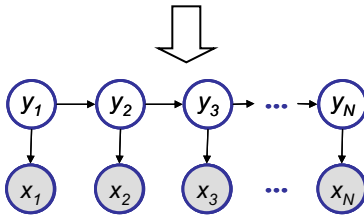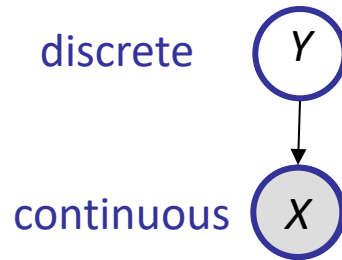$$\pi = \{\alpha, \omega, \eta^I, \eta^G\} \text{ (hyper-parameters)}$$

# Factor Analysis

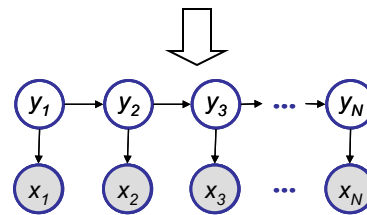# A road map to more complex dynamic models

# Recall multivariate Gaussian

- Multivariate Gaussian density:

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

- A joint Gaussian:

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

- 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 $\mu$ and $\Sigma$?

  - Formulas to remember:

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

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

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

$$p(\mathbf{x}_1|\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 \mid \mathbf{m}_{1|2}, \mathbf{V}_{1|2})$$

$$\mathbf{m}_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x}_2 - \mu_2)$$

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

# 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: $XYZ = W \implies Y^{-1} = ZW^{-1}X$

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

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

- Matrix inverse lemma

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

# Review: Some matrix algebra

$f(x) = ax$

$\frac{\partial f}{\partial x} = a$

- Trace and derivatives

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

  - Cyclical permutations

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

  - Derivatives

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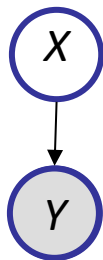
$$\frac{\partial}{\partial A} \mathrm{tr}[BA] = B^T$$

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

- Determinants and derivatives

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

# Factor analysis

*P(y)*
*P(y|x)*

$$y = \mu + \Lambda x + w$$
$$w \sim N(0, \Psi)$$

- An unsupervised linear regression model

X

Y

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; 0, I)$$
$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mu + \Lambda \mathbf{x}, \Psi)$$

where $\Lambda$ is called a factor loading matrix, and $\Psi$ is diagonal.

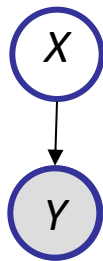- Geometric interpretation

$y^3$

$\lambda_1$

$\lambda_2$

$\mu$

$y^2$

$y^1$

- 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.)

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

# FA = Constrained-Covariance Gaussian

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

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

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



- In other words, factor analysis is just a constrained Gaussian model (number of free params of the covariance is limited). (If $\Psi$ 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 **x** and **y**

$$Cov[\mathbf{X}, \mathbf{Y}] = E\left[(\mathbf{X} - \mathbf{0})(\mathbf{Y} - \mu)^T\right] = E\left[\mathbf{X}(\mu + \Lambda\mathbf{X} + \mathbf{W} - \mu)^T\right]$$

$$= E\left[\mathbf{X}\mathbf{X}^T \Lambda^T + \mathbf{X}\mathbf{W}^T\right]$$

$$= \Lambda^T$$

- Hence the joint distribution of **x** and **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.

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

$-H^{-1}G(M/H)^{-1}$  $H^{-1}+H^{-1}G(M/H)^{-}FH^{-1}$

# Inference in Factor Analysis

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

$$\Sigma_{11} = I$$
$$\Sigma_{12} = \Sigma_{12}{}^{T} = \Lambda^{T}$$
$$\Sigma_{22} = \left(\Lambda\Lambda^{T} + \Psi\right)$$

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

$$\mathbf{m}_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{y}-\mu_2)$$
$$\mathbf{V}_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

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

Applying the matrix inversion lemma

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

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

- Here we only need to invert a matrix of size $|\mathbf{x}|\times|\mathbf{x}|$, instead of $|\mathbf{y}|\times|\mathbf{y}|$.
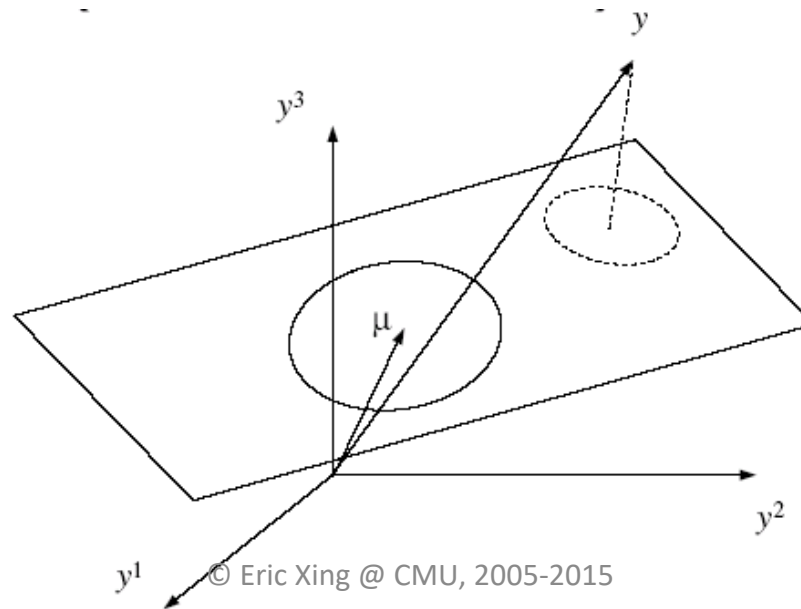
# Geometric interpretation: inference is linear projection

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

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

- 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 $\Lambda$
  - Manifold center $\mu$
  - Variance $\Psi$

# EM for Factor Analysis

$P(x;\theta) = \sum \pi_k N(x_i, \mu_k, \Sigma_k)$

$\theta = \{\mu_k, \Sigma_k, \pi_k\}$

- Incomplete data log likelihood function (marginal density of y)

$$\ell(\theta, D) = -\frac{N}{2}\log|\Lambda\Lambda^T + \Psi| - \frac{1}{2}\sum_n (y_n - \mu)^T(\Lambda\Lambda^T + \Psi)^{-1}(y_n - \mu)$$

$\rightarrow trace(\Lambda\Lambda^T + \Psi)^{-1}$

$$= -\frac{N}{2}\log|\Lambda\Lambda^T + \Psi| - \frac{1}{2}\mathrm{tr}\left[(\Lambda\Lambda^T + \Psi)^{-1}S\right], \qquad \text{where } S = \sum_n (y_n - \mu)(y_n - \mu)^T$$

$P(x, y)$

$P(y)$

- Estimating $\mu$ is trivial: $\hat{\mu}^{ML} = \frac{1}{N}\sum_n y_n$
- Parameters $\Lambda$ and $\Psi$ are coupled nonlinearly in log-likelihood

- Complete log likelihood

$$\ell_c(\theta, D) = \sum_n \log p(x_n, y_n) = \sum_n \log p(x_n) + \log p(y_n | x_n)$$

$\langle \rangle$

$P(X|D)$

$$= -\frac{N}{2}\log|I| - \frac{1}{2}\sum_n x_n^T x_n - \frac{N}{2}\log|\Psi| - \frac{1}{2}\sum_n (y_n - \Lambda x_n)^T \Psi^{-1}(y_n - \Lambda x_n)$$

$$= -\frac{N}{2}\log|\Psi| - \frac{1}{2}\sum_n \mathrm{tr}\left[x_n x_n^T\right] - \frac{N}{2}\mathrm{tr}\left[S\Psi^{-1}\right], \qquad \text{where } S = \frac{1}{N}\sum_n (y_n - \Lambda x_n)(y_n - \Lambda x_n)^T$$

unobs

Date

# E-step for Factor Analysis

- Compute $\left\langle \ell_c (\theta, D) \right\rangle_{p(x|y)}$

$$\left\langle \ell_c (\theta, D) \right\rangle = -\frac{N}{2}\log|\Psi| - \frac{1}{2}\sum_n \text{tr}\left[\left\langle X_n X_n^T \right\rangle\right] - \frac{N}{2}\text{tr}\left[\left\langle \mathbf{S} \right\rangle \Psi^{-1}\right]$$

$$\left\langle \mathbf{S} \right\rangle = \frac{1}{N}\sum_n (y_n y_n^T - y_n \left\langle X_n^T \right\rangle \Lambda^T - \Lambda \left\langle X_n^T \right\rangle y_n^T + \Lambda \left\langle X_n X_n^T \right\rangle \Lambda^T)$$

$$\left\langle X_n \right\rangle = E\left[X_n \mid y_n\right]$$

$$\left\langle X_n X_n^T \right\rangle = Var\left[X_n \mid y_n\right] + E\left[X_n \mid y_n\right]E\left[X_n \mid y_n\right]^T$$

- Recall that we have derived:

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

$$\Longrightarrow \quad \left\langle X_n \right\rangle = \mathbf{m}_{x_n|y_n} = \mathbf{V}_{1|2}\Lambda^T \Psi^{-1}(y_n - \mu) \quad \text{and} \quad \left\langle X_n X_n^T \right\rangle = \mathbf{V}_{1|2} + \mathbf{m}_{x_n|y_n}\mathbf{m}_{x_n|y_n}^T$$

# M-step for Factor Analysis

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

$$\frac{\partial}{\partial \Psi^{-1}} \langle \ell_c \rangle = \frac{\partial}{\partial \Psi^{-1}} \left( -\frac{N}{2} \log|\Psi| - \frac{1}{2}\sum_n \text{tr}\left[\langle X_n X_n^T \rangle\right] - \frac{N}{2} \text{tr}\left[\langle \mathbf{S} \rangle \Psi^{-1}\right] \right)$$

$$= \frac{N}{2} \Psi - \frac{N}{2} \langle \mathbf{S} \rangle \qquad \Rightarrow \qquad \boxed{\Psi^{t+1} = \langle \mathbf{S} \rangle}$$

$$\frac{\partial}{\partial \Lambda} \langle \ell_c \rangle = \frac{\partial}{\partial \Lambda} \left( -\frac{N}{2} \log|\Psi| - \frac{1}{2}\sum_n \text{tr}\left[\langle X_n X_n^T \rangle\right] - \frac{N}{2} \text{tr}\left[\langle \mathbf{S} \rangle \Psi^{-1}\right] \right) = -\frac{N}{2} \Psi^{-1} \frac{\partial}{\partial \Lambda} \langle \mathbf{S} \rangle$$

$$= -\frac{N}{2} \Psi^{-1} \frac{\partial}{\partial \Lambda} \left( \frac{1}{N}\sum_n (y_n y_n^T - y_n \langle X_n^T \rangle \Lambda^T - \Lambda \langle X_n^T \rangle y_n^T + \Lambda \langle X_n X_n^T \rangle \Lambda^T) \right)$$

$$= \Psi^{-1} \sum_n y_n \langle X_n^T \rangle - \Psi^{-1} \Lambda \sum_n \langle X_n X_n^T \rangle \qquad \Rightarrow \qquad \Lambda^{t+1} = \left( \sum_n y_n \langle X_n^T \rangle \right) \left( \sum_n \langle X_n X_n^T \rangle \right)^{-1}$$
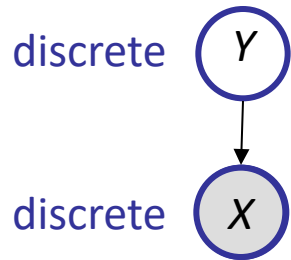
$y|x$　　　　　$y = \Delta x + \mu$
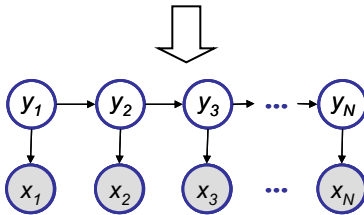
# Model Invariance and Identifiability

- There is *degeneracy* in the FA model.

- Since $\Lambda$ only appears as outer product $\Lambda\Lambda^{\mathrm{T}}$, the model is invariant to rotation and axis flips of the latent space.

- We can replace $\Lambda$ with $\Lambda Q$ for any orthonormal matrix Q and the model remains the same: $(\Lambda Q)(\Lambda Q)^{\mathrm{T}} = \Lambda(QQ^{\mathrm{T}})\Lambda^{\mathrm{T}} = \Lambda\Lambda^{\mathrm{T}}$.

- 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.
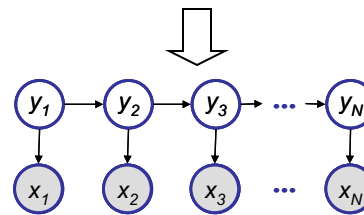
# A road map to more complex dynamic models