Department of Computer Science
UNIVERSITY OF COLORADO **BOULDER**

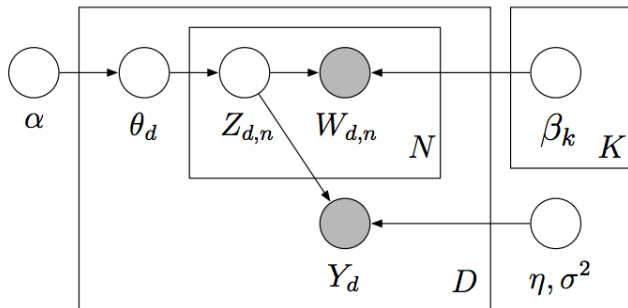# **Supervised Topic Models**

Advanced Machine Learning for NLP
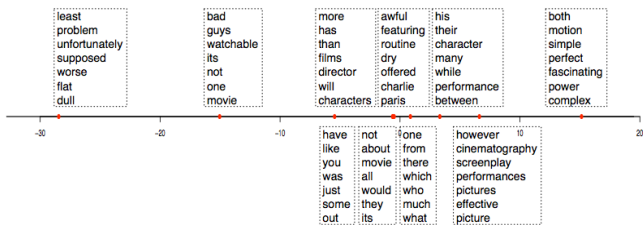Jordan Boyd-Graber
MULTILINGUAL

- Normal LDA generative story
- Document also has label $y_d$

$$y_d \sim \mathcal{N}\left(y_d \,|\, y_d, \eta^\top \mathbb{E}_\theta\left[\bar{Z}\right]\right) \qquad (1)$$

# How does this change topics?

**How does this change topics?**

Recall the joint likelihood:

$$p(\boldsymbol{z} \,|\, \alpha, \lambda, \boldsymbol{w}, \boldsymbol{\eta}) \propto \tag{2}$$

$$\prod_d \prod_d \frac{\prod_d \Gamma\big(n_{d,k} + \alpha_{d,k}\big)}{\Gamma\big(\sum_d n_{d,k} + \alpha_{d,k}\big)} \prod_k \frac{\prod_k \Gamma\big(t_{k,v} + \lambda_{k,v}\big)}{\Gamma\big(\sum_k t_{k,v} + \lambda_{k,v}\big)} \tag{3}$$

$$\prod_d \exp\big\{-(y_d - \eta^\top \bar{z})\big\} \tag{4}$$

Apply gibbs sampling equations:

$$p(z_{d,n} = k \,|\, \ldots) \propto (n_{d,k}^{-d,n} + \alpha_k) \frac{t_{k,w_{d,n}}^{-d,n} + \lambda_{w_{d,n}}}{t^{-d,n} + V\lambda} \exp\big\{-(y_d - \eta^\top \bar{z})^2\big\} \tag{5}$$

Recall the joint likelihood:

$$p(\boldsymbol{z} \mid \alpha, \lambda, \boldsymbol{w}, \boldsymbol{\eta}) \propto \tag{2}$$

$$\prod_d \prod_d \frac{\prod_d \Gamma(n_{d,k} + \alpha_{d,k})}{\Gamma(\sum_d n_{d,k} + \alpha_{d,k})} \prod_k \frac{\prod_k \Gamma(t_{k,v} + \lambda_{k,v})}{\Gamma(\sum_k t_{k,v} + \lambda_{k,v})} \tag{3}$$

$$\prod_d \exp\left\{-(y_d - \eta^\top \bar{z})\right\} \tag{4}$$

Apply gibbs sampling equations:

$$p(z_{d,n} = k \mid \ldots) \propto (n_{d,k}^{-d,n} + \alpha_k) \frac{t_{k,w_{d,n}}^{-d,n} + \lambda_{w_{d,n}}}{t^{-d,n} + V\lambda} \textcolor{red}{\exp\left\{-(y_d - \eta^\top \bar{z})^2\right\}} \tag{5}$$

Let's expand last term

$$\exp\left\{-(y_d - \boldsymbol{\eta}^\top \bar{z})^2\right\} \tag{6}$$

$$\exp\left\{-(y_d - \boldsymbol{\eta}^\top \bar{z})^2\right\} = \exp\left\{-y_d^2 + 2\boldsymbol{\eta}^\top \bar{z}_d\, y_d - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{6}$$

$$\tag{7}$$

$$\exp\left\{-(y_d - \boldsymbol{\eta}^\top \bar{z})^2\right\} = \exp\left\{-y_d^2 + 2\boldsymbol{\eta}^\top \bar{z}_d\, y_d - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{6}$$

$$\propto \exp\left\{2\sum_j \eta_j z_{d,j}^{-d,n}\, y_d + 2\frac{y_d}{N_d}\eta_k - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{7}$$

$$\tag{8}$$

Expand product

$$\exp\left\{-(y_d - \boldsymbol{\eta}^\top \bar{z})^2\right\} = \exp\left\{-y_d^2 + 2\boldsymbol{\eta}^\top \bar{z}_d\, y_d - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{6}$$

$$\propto \exp\left\{2\sum_j \eta_j z_{d,j}^{-d,n}\, y_d + 2\frac{y_d}{N_d}\eta_k - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{7}$$

$$\propto \exp\left\{2\frac{y_d}{N_d}\eta_k - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{8}$$

$$\tag{9}$$

Remove constant term, explicitly write dot product

**How does this change topics?**

$$\exp\left\{-(y_d - \boldsymbol{\eta}^\top \bar{z})^2\right\} = \exp\left\{-y_d^2 + 2\boldsymbol{\eta}^\top \bar{z}_d y_d - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{6}$$

$$\propto \exp\left\{2\sum_j \eta_j z_{d,j}^{-d,n} y_d + 2\frac{y_d}{N_d}\eta_k - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{7}$$

$$\propto \exp\left\{2\frac{y_d}{N_d}\eta_k - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{8}$$

$$\propto \exp\left\{2\frac{y_d}{N_d}\eta_k - \left(\frac{\eta_k}{N_d} + \sum_j \eta_j \bar{z}_{d,j}^{-d,n}\right)^2\right\} \tag{9}$$

$$\tag{10}$$

Break dot product into $k$ and non-$k$ terms

**How does this change topics?**

$$\exp\left\{-(y_d - \boldsymbol{\eta}^\top \bar{z})^2\right\} \propto \exp\left\{2\sum_j \eta_j z_{d,j}^{-d,n} y_d + 2\frac{y_d}{N_d}\eta_k - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \quad (6)$$

$$\propto \exp\left\{2\frac{y_d}{N_d}\eta_k - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \quad (7)$$

$$\propto \exp\left\{2\frac{y_d}{N_d}\eta_k - \left(\frac{\eta_k}{N_d} + \sum_j \eta_j \bar{z}_{d,j}^{-d,n}\right)^2\right\} \quad (8)$$

$$\propto \exp\left\{2\frac{y_d}{N_d}\eta_k - \left(\frac{\eta_k}{N_d}\right)^2 - \frac{2\eta_k}{N_d}\boldsymbol{\eta}^\top \bar{z}^{-d,n}\right\} \quad (9)$$

$$(10)$$

Expand product, drop constant terms

$$\exp\left\{-(y_d - \boldsymbol{\eta}^\top \bar{z})^2\right\} \propto \exp\left\{2\frac{y_d}{N_d}\eta_k - (\boldsymbol{\eta}^\top \bar{z}_d)^2\right\} \tag{6}$$

$$\propto \exp\left\{2\frac{y_d}{N_d}\eta_k - \left(\frac{\eta_k}{N_d} + \sum_j \eta_j \bar{z}_{d,j}^{-d,n}\right)^2\right\} \tag{7}$$

$$\propto \exp\left\{2\frac{y_d}{N_d}\eta_k - \left(\frac{\eta_k}{N_d}\right)^2 - \frac{2\eta_k}{N_d}\boldsymbol{\eta}^\top \bar{z}^{-d,n}\right\} \tag{8}$$

$$\propto \exp\left\{2\frac{\eta_k}{N_d}\left(y_d - \boldsymbol{\eta}^\top \bar{z}^{-d,n}\right) - \left(\frac{\eta_k}{N_d}\right)^2\right\} \tag{9}$$
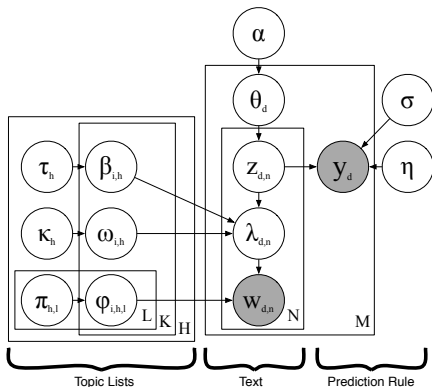
$$\tag{10}$$

Factor terms

**Let's go a step further**

- Latent space is really useful
- Let's make it coherent across languages
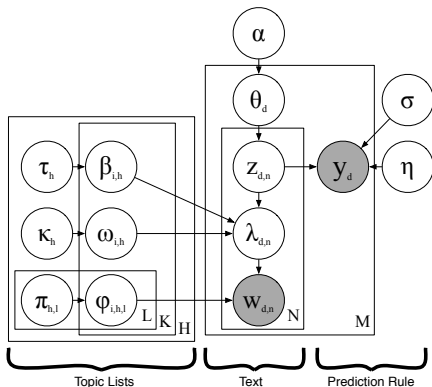- Requires a glue across languages

1. For each topic $k = 1 \ldots K$, draw correlated multilingual word distribution $\{\boldsymbol{\beta}_k, \boldsymbol{\omega}_k, \boldsymbol{\phi}_k\}$

2. For each document $d$, $\theta_d \sim \text{Dir}(\alpha)$
   1. $z_{d,n} \sim \text{Discrete}(\theta_d)$
   2. Draw path $\lambda_{d,n}$ through multilingual tree $z_{d,n}$, emit $w_{d,n}$
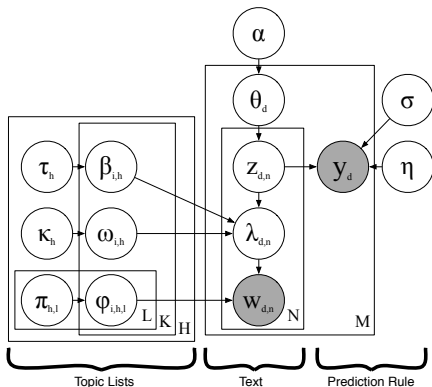
3. $y_d \sim \text{Norm}(\eta^\top \bar{z}, \sigma^2)$



Topic Lists      Text      Prediction Rule

1. For each topic $k = 1 \ldots K$, draw correlated multilingual word distribution $\{\boldsymbol{\beta}_k, \boldsymbol{\omega}_k, \boldsymbol{\phi}_k\}$

2. For each document $d$, $\theta_d \sim \text{Dir}(\alpha)$
   1. $z_{d,n} \sim \text{Discrete}(\theta_d)$
   2. Draw path $\lambda_{d,n}$ through multilingual tree $z_{d,n}$, emit $w_{d,n}$

3. $y_d \sim \text{Norm}(\eta^\top \bar{z}, \sigma^2)$

① For each topic $k = 1 \ldots K$, draw correlated multilingual word distribution $\{\boldsymbol{\beta}_k, \boldsymbol{\omega}_k, \boldsymbol{\phi}_k\}$

② For each document $d$, $\theta_d \sim \text{Dir}(\alpha)$

  ① $z_{d,n} \sim \text{Discrete}(\theta_d)$

  ② Draw path $\lambda_{d,n}$ through multilingual tree $z_{d,n}$, emit $w_{d,n}$

③ $y_d \sim \text{Norm}(\eta^\top \bar{z}, \sigma^2)$

- Statistical NLP typically uses Dirichlet distributions because of conjugacy
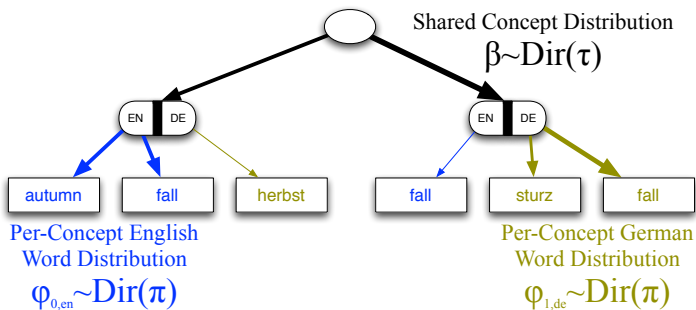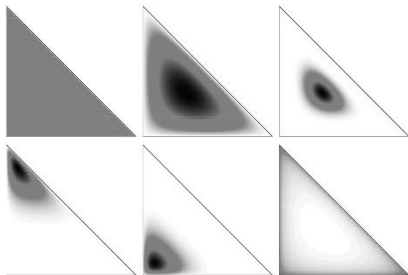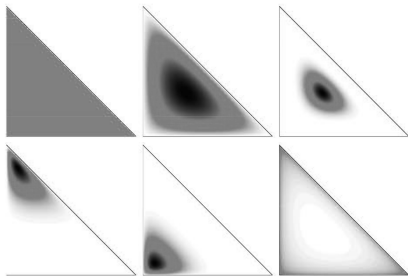- Parameter of Dirichlet encode mean and variance

## Encoding Correlations

- Statistical NLP typically uses Dirichlet distributions because of conjugacy
- Parameter of Dirichlet encode mean and variance
- But we want correlations!

- Statistical NLP typically uses Dirichlet distributions because of conjugacy



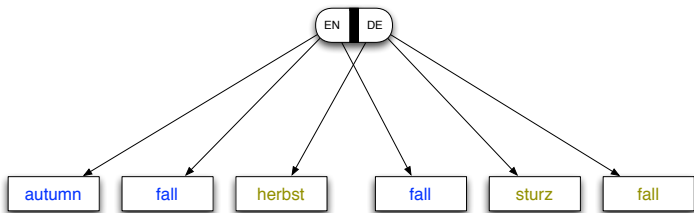- Parameter of Dirichlet encode mean and variance

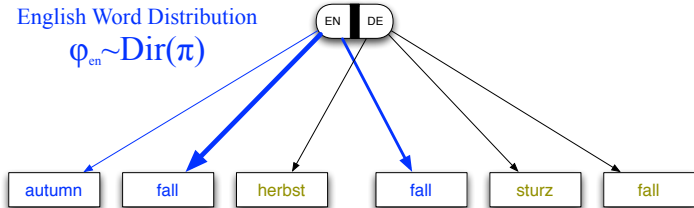- Statistical NLP typically uses Dirichlet distributions because of conjugacy



- Parameter of Dirichlet encode mean and variance
- But we want correlations!

gut   hǎo   good

English Word Distribution
$$\varphi_{en} \sim Dir(\pi)$$

English Word Distribution
$\varphi_{en} \sim \text{Dir}(\pi)$

German Word Distribution
$\varphi_{de} \sim \text{Dir}(\pi)$

EN | DE

autumn | fall | herbst | fall | sturz | fall

Shared Concept Distribution
$\beta \sim \mathrm{Dir}(\tau)$

Shared Concept Distribution
$$\beta \sim \mathrm{Dir}(\tau)$$

Per-Concept German
Word Distribution
$$\varphi_{1,de} \sim \mathrm{Dir}(\pi)$$

## Encoding Correlations

- CEDICT (Chinese/English)
- HanDeDict (Chinese/German)
- Ding (German/English)

GermaNet

- Jointly sample $z$ and path $\lambda$ through multilingual tree

$$p(z_n = k, \lambda_n = r | \boldsymbol{z}_{-n}, \boldsymbol{\lambda}_{-n}, w_n, \eta, \sigma, \Theta) =$$
$$p(y_d | \boldsymbol{z}, \eta, \sigma) p(\lambda_n = r | z_n = k, \boldsymbol{\lambda}_{-n}, w_n, \boldsymbol{\tau}, \boldsymbol{\kappa}, \boldsymbol{\pi})$$
$$p(z_n = k | \boldsymbol{z}_{-n}, \alpha).$$

- Collapse out multinomial distributions in tree
- Slice sample hyperparameters
- After pass of $z$, update $\eta$

# Multilingual Supervised LDA

# Evaluation: Learned Topics (Chinese - German)



(harry) harry 哈利 (harry)
(volume) band 带 (belt)
(sky) himmel 天 (sky)
(universe) all 都 (both)
(vampire) vampir 部 (section)
(last) letzt 吸血鬼 (vampire)
(part) teil 强烈 (strong)
最后 (last)

(god) gott 上帝 (god)
(lord) herr 先生 (lord)
(religion) religion 都 (both)
(universe) all 宗教 (religion)
(world) welt 科学 (science)
(science) wissenschaft 社会 (community)
(medicine) medizin
(society) gesellschaft

-1.6    -0.8    0.0    0.8
  -1.2    -0.4    0.4

书 (book)
恐怕 ([I'm afraid that...])
(book) buch 亲身 (myself)
(itself) sich 都 (both)
(that) dass 大部分 (mostly)
(much) viel 书本 (book)
(no) kein 并非 ([really isn't])
(good) gut 脱 (discard)
(when) wenn

(woman) frau
(point) punkt
(man) mann
(equal) gleich
(fast) schnell
(female) weiblich
(soon) bald

快 (quick)
点 (a little)
女人 (woman)
男人 (man)
女 (female)
男 (male)
女性 (female)
都 (both)

(good) gut
(sentence) satz
(two) zwei
(story) story
(treasure) schatz
(attractive) attraktiv
(elegant) elegant
(gem) juwel

好 (good)
套 (set)
宝 (treasure)
帅 (handsome)
两 (both)
故事 (story)
小 (small)

- Take large corpus (6000) of English movie reviews rated from 0–100
- Combine them with smaller German corpus (300) rated using same system
- Compute mean squared error (lower is better) on held out data

| Train | Test | GermaNet | Dictionary | Flat |
|-------|------|----------|------------|------|
| DE | DE | 73.8 | 24.8 | 92.2 |
| EN | DE | 7.44 | 2.68 | 18.3 |
| EN + DE | DE | **1.17** | **1.46** | **1.39** |

Moral: More data, even in another language, helps