# Topic Models

**Review** The <u>generative</u> <u>model</u> for Naïve Bayes

$$X = \{x_1, \dots x_D\} \qquad K \text{ labels/classes}$$
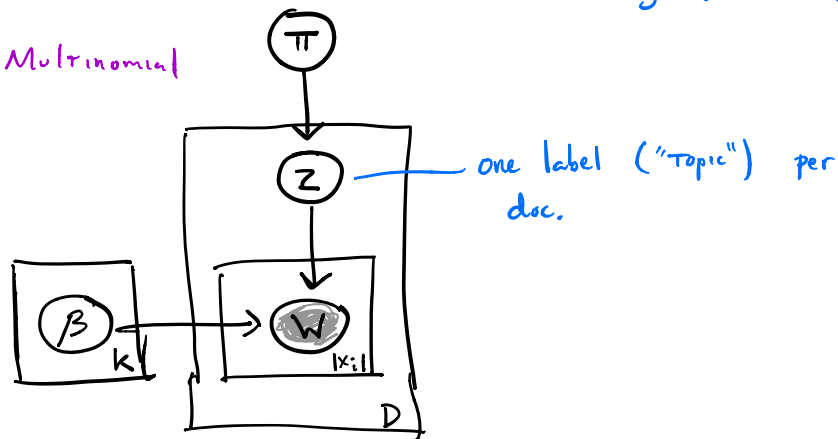
$$Z_d \sim \text{Categorical}(\pi)$$

$$X_{dn} \mid Z_d \sim \text{Categorical}(\beta_{Z_d})$$

Probability over words $w$ given class $Z_d$.

★ Terminology note:
"Categorical" = Multinomial w/ 1 trial.



one label ("Topic") per doc.

Assumes <u>one</u> topic $Z_d$ per document or instance.

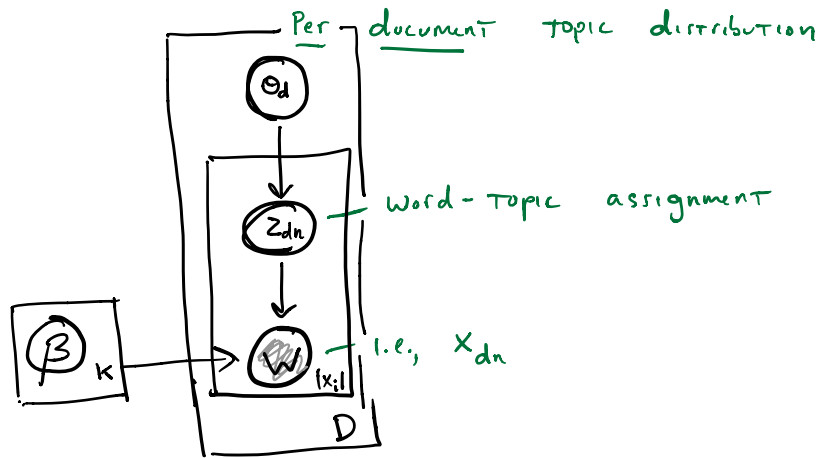Most documents will span <u>multiple</u> topics.

To respect this, we need a new model.

$$Z_{dn} \sim \text{Categorical}(\theta_d)$$

$$X_{dn} | Z_{dn} \sim \text{Categorical}(\beta_{Z_{dn}})$$

distribution over words for topic $Z_{dn}$.

Graphically

Per-document topic distribution

$\theta_d$

word-topic assignment $Z_{dn}$

$\beta$ $K$

$W$ i.e., $X_{dn}$

$[x_i]$

$D$

# EM for Topic Models: PLSA
### (Probabistic Latent Semantic Analysis)

$$Z_{dn} \sim \text{Categorical}(\Theta_d)$$

$$X_{dn} \mid Z_{dn} \sim \text{Categorical}(\beta_{Z_{dn}})$$

distribution over words for Topic $Z_{dn}$.

## E-step

Update soft assignments

$$\phi_{dnk} = \frac{\Theta_{dk}\, \beta_{kv}}{\sum_{\ell} \Theta_{d\ell}\, \beta_{\ell v}} \qquad v \overset{def}{=} x_{dn}$$

P that word $n$ in document $d$ was drawn from Topic $k$

## M-step

$$\beta_{kv} = \frac{N_{kv}}{\sum_{w} N_{kw}} \qquad N_{kv} \overset{def}{=} \sum_{d,n} \phi_{dnk}\, x_{nv}$$

$$= \text{expected count of } v \text{ in Topic } k$$

$$\Theta_{dk} = \frac{N_{dk}}{\sum_{\ell} N_{d\ell}} \qquad N_{\ell} \overset{def}{=} \sum_{n} \phi_{dn\ell}$$

$$= \text{Expected number of words from Topic } \ell \text{ in } d$$