# Discriminative Topic Modeling with Logistic LDA

Iryna Korshunova[1*]    Hanchen Xiong[2]    Mateusz Fedoryszak[2]    Lucas Theis[2]

[1]Ghent University    [2]Twitter
[*]Work done at Twitter

## Overview

**Logistic LDA** is a novel discriminative variant of latent Dirichlet allocation (LDA) which is easy to apply to arbitrary inputs, such as images or text embeddings.

**Logistic LDA** preserves LDA's extensibility and interpretability. In particular, it explicitly models item topics and group-level topic distributions, while integrating deep neural networks in a principled manner.

Among other desirable properties, **logistic LDA**:

✓ can be supervised, semi-supervised or unsupervised
✓ is scalable to large datasets
✓ can benefit from the vast literature on LDA
✓ applicable to a wide range of problems with group structure present in the data

## Latent Dirichlet Allocation

$D$ - number of documents in a corpus
$N_d$ - number of words in a document d
$K$ - number of topics
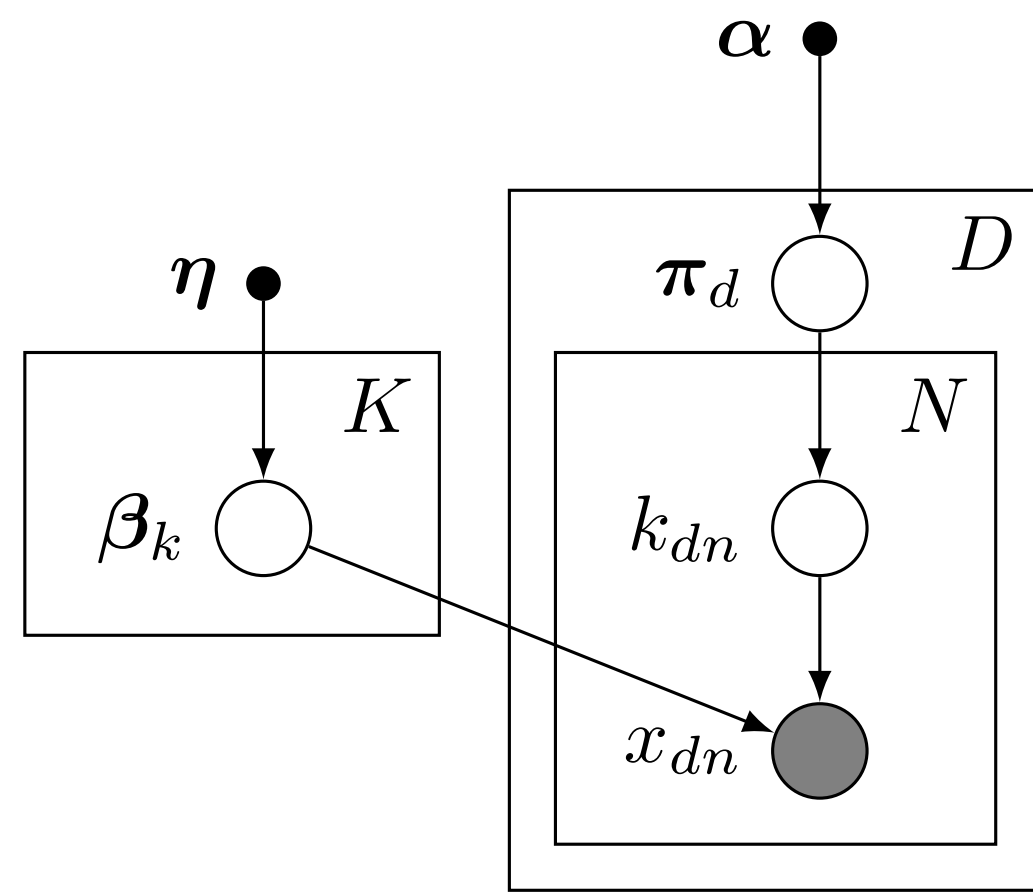$V$ - number of words in the vocabulary

$x_{dn}$ - n-th observed word in d-th document
$k_{dn}$ - latent topic of word $x_{dn}$
$\pi_d$ - distribution over topics
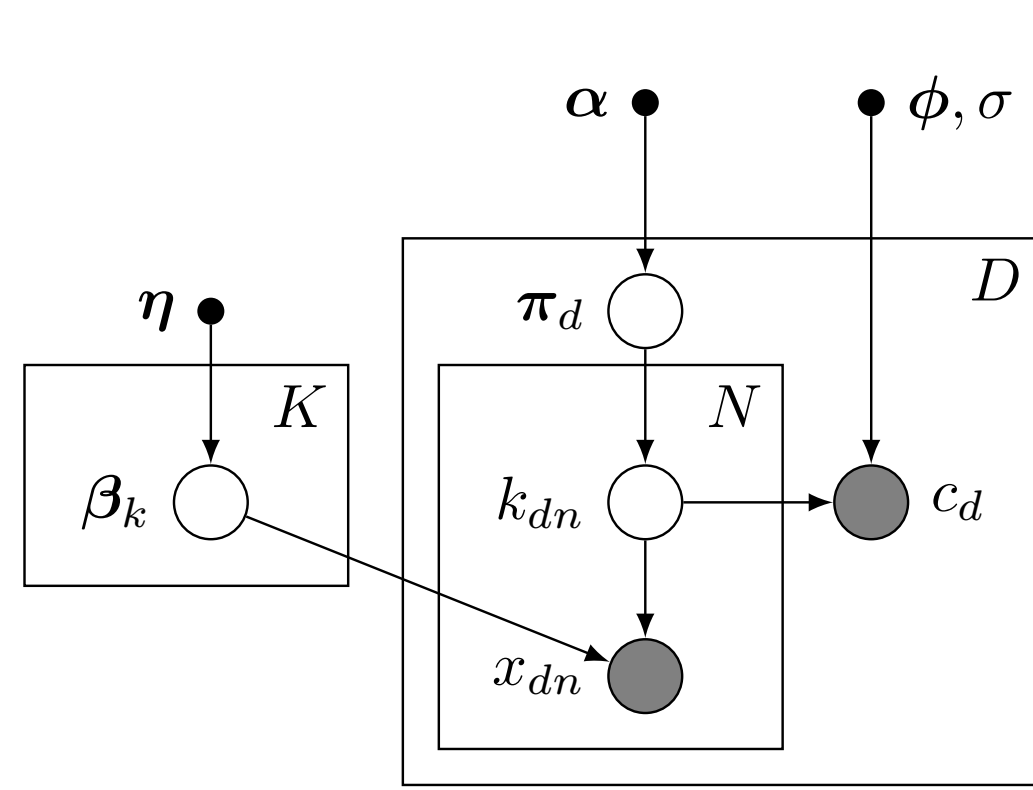$\beta$ - $K \times V$ matrix of topic-word distributions

**LDA:**



**Generative process**:

1. Draw topic proportions $\pi_d \sim \text{Dir}(\alpha)$
2. Draw topic-word distributions $\beta_k \sim \text{Dir}(\eta)$
3. For each word $x_{dn}$:
   3.1 Draw a topic assignment $k_{dn} \sim \text{Cat}(\pi_d)$
   3.2 Draw a word $x_{dn} \sim \text{Cat}(\beta^\top k_{dn})$
4. In supervised LDA, draw a response variable:
   $c_d \sim \mathcal{N}(\phi^\top(\frac{1}{N_d}\sum_n k_{dn}), \sigma^2)$

**sLDA [1]:'**



## Generative or Discriminative? [5]

**Generative**

$$p(c, x, \theta) = p(x, c \mid \theta)p(\theta)$$
$$= p(c \mid \pi)p(x \mid c, \lambda)p(\theta), \text{ with } \theta = \{\pi, \lambda\}$$

e.g. LDA, naive Bayes classifier, linear discriminant analysis, GMM

**Discriminative**

$$p(c, x, \theta) = p(c \mid x, \theta)p(\theta)p(x)$$

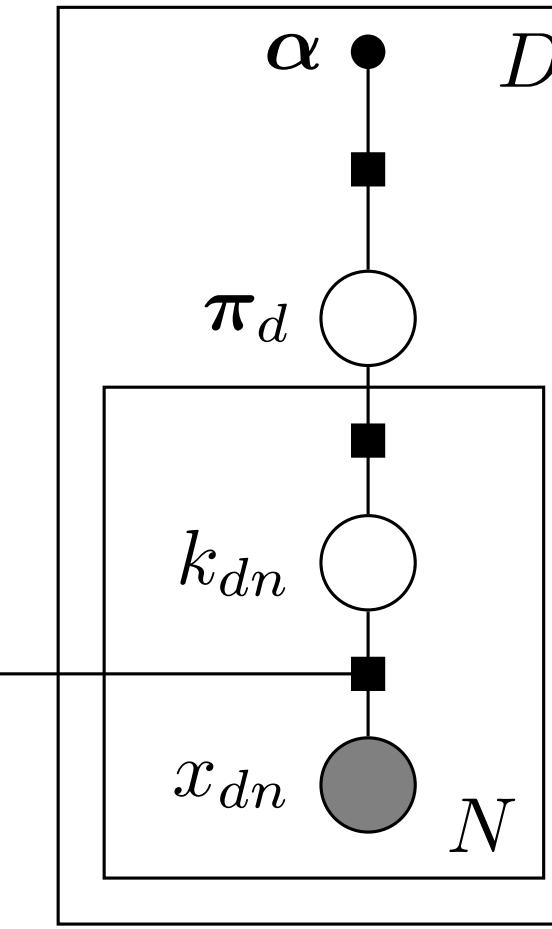e.g. logistic regression, SVM, CRF

**Logisitc LDA**

$$p(c, x, \theta) = p(c, \theta \mid x)p(x)$$

## Alternative View of LDA

▶ We specify a set of full conditional probabilities:

$$p(\pi_d \mid k_d) = \text{Dir}\left(\pi_d; \alpha + \sum_n k_{dn}\right)$$

$$p(k_{dn} \mid x_{dn}, \pi_d, \theta) = k_{dn}^\top \text{softmax}(g(x_{dn}, \theta) + \ln \pi_d)$$

$$p(\theta \mid x, k) \propto \exp\left(r(\theta) + \sum_{dn} k_{dn}^\top g(x_{dn}, \theta)\right)$$



▶ That result in a valid joint distribution:

$$p(\pi, k, \theta \mid x) \propto \exp\left((\alpha - 1)^\top \sum_d \ln \pi_d + \sum_{dn} k_{dn}^\top \ln \pi_d + \sum_{dn} k_{dn}^\top g(x_{dn}, \theta) + r(\theta)\right)$$

▶ Where LDA is a special case:

$$g(x_{dn}, \beta) = \ln \beta \, x_{dn} \qquad r(\beta) = (\eta - 1)^\top \sum_k \ln \beta_k \qquad \sum_j \beta_{kj} = 1$$

## Logistic LDA

$$g(x_{dn}, \theta) = \ln \text{softmax} \, f(x_{dn}, \theta) \qquad r(\theta, x) = \gamma \cdot \mathbf{1}^\top \ln \sum_{dn} \exp g(x_{dn}, \theta)$$

**Supervised Logistic LDA**

$$p(\pi_d \mid k_d, c_d) = \text{Dir}\left(\pi_d; \alpha + \sum_n k_{dn} + \lambda c_d\right)$$

$$p(k_{dn} \mid x_{dn}, \pi_d, \theta) = k_{dn}^\top \text{softmax}(g(x_{dn}, \theta) + \ln \pi_d)$$

$$p(\theta \mid x, k) \propto \exp\left(r(\theta) + \sum_{dn} k_{dn}^\top g(x_{dn}, \theta)\right)$$

$$p(c_d \mid \pi_d) = \text{softmax}(\lambda c_d^\top \ln \pi_d)$$



## Training and Inference

$$\min D_{\text{KL}}\left[q(\theta)\left(\prod_d q(c_d)\right)\left(\prod_d q(\pi_d)\right)\left(\prod_{dn} q(k_{dn})\right) \mid\mid p(\pi, k, c, \theta \mid x)\right]$$

▶ Coordinate descent updates for variational parameters when $\theta$ is fixed:

$$q(c_d) = c_d^\top \hat{p}_d \qquad\qquad \hat{p}_d = \text{softmax}\left(\lambda \psi(\hat{\alpha}_d)\right)$$
$$q(\pi_d) = \text{Dir}(\pi_d; \hat{\alpha}_d) \qquad \hat{\alpha}_d = \alpha + \sum_n \hat{p}_{dn} + \lambda \hat{p}_d$$
$$q(k_{dn}) = k_{dn}^\top \hat{p}_{dn} \qquad\qquad \hat{p}_{dn} = \text{softmax}\left(f(x_{dn}, \hat{\theta}) + \psi(\hat{\alpha}_d)\right)$$

▶ VI loss wrt. $\theta$:

$$\ell(\hat{\theta}) \approx -\sum_{dn}(\hat{p}_{dn} + \gamma \cdot \hat{r}_{dn})^\top g(x_{dn}, \hat{\theta})$$

▶ Alternative empirical loss when $c_d$ is observed:
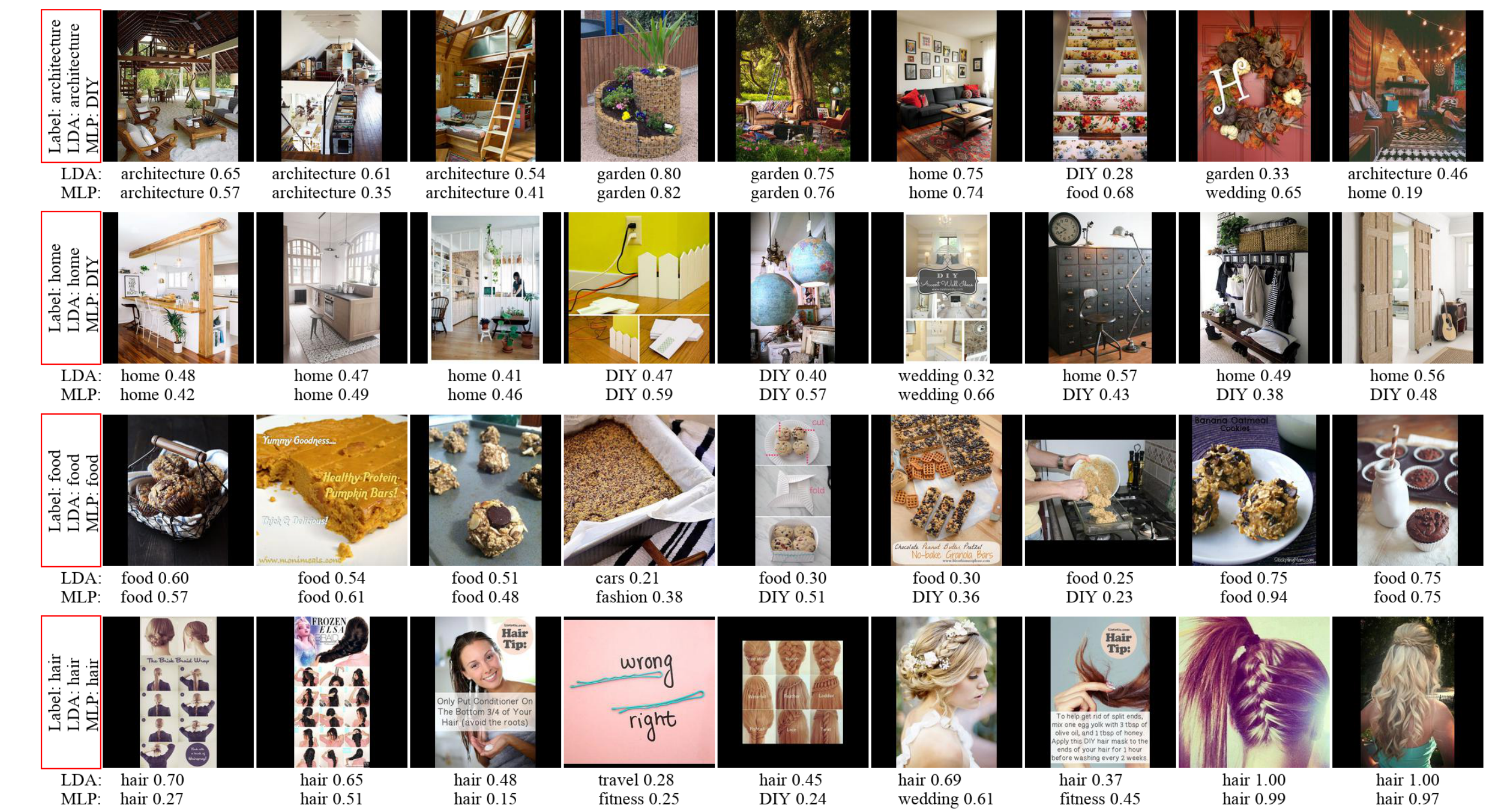
$$\ell(\hat{\theta}) = -\sum_d c_d^\top \ln \hat{p}_d$$

## Experiments: Twitter

▶ Dataset of ∼4M tweets from ∼100K authors where some tweets and authors were annotated with one of 300 topics.
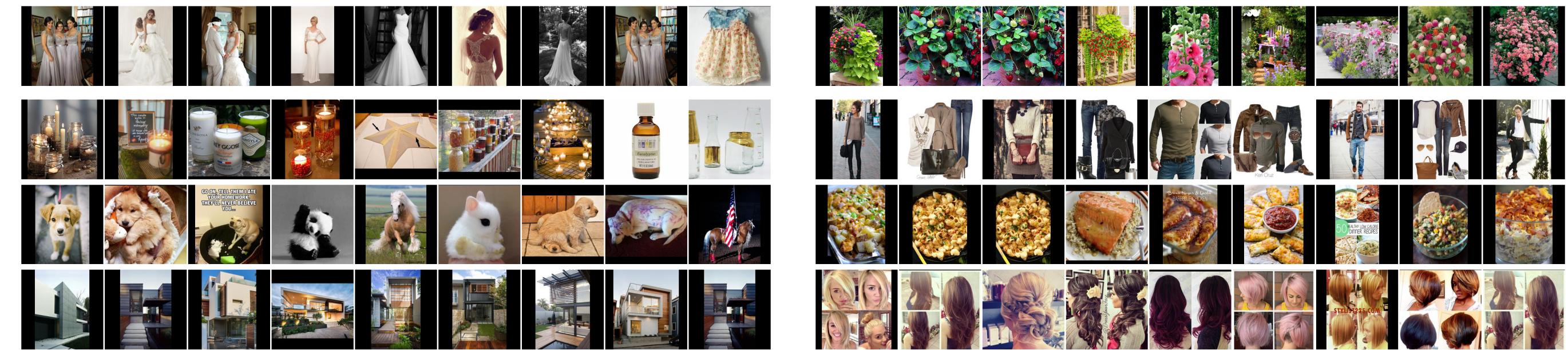
▶ In production: timeline filtering according to topics.

| Model | Author | Tweet |
|---|---|---|
| MLP (individual) | 26.6% | 32.4% |
| MLP (majority) | 35.0% | n/a |
| LDA | 33.1% | 25.4% |
| **Logistic LDA** | **38.7%** | **35.6%** |

## Experiments: Pinterest

**Predictions for Pinterest boards and pins**



**Unsupervised topics**



## Experiments: 20-Newsgroups

**Document classification accuracy**

| Logistic LDA | SVM [4] | LSTM [2] | SA-LSTM [2] | oh-2LSTMp [3] |
|---|---|---|---|---|
| 84.4% | 82.9% | 82.0% | 84.4% | 86.5% |

**Unsupervised topics**

1 bmw, motor, car, honda, motorcycle, auto, mg, engine, ford, bike
2 christianity, prophet, atheist, religion, holy, scripture, biblical, catholic, religious
3 spacecraft, orbit, probe, ship, satellite, rocket, surface, shipping, moon, launch
4 user, computer, microsoft, monitor, programmer, electronic, processing, data, app, systems
5 congress, administration, economic, accord, trade, criminal, seriously, fight, responsible, future

## Bibliography

[1] D. Blei and J. McAuliffe. *Supervised Topic Models.* NIPS, 2008.
[2] A. Dai, Q. Le. *Semi-supervised Sequence Learning.* NIPS, 2015.
[3] R. Johnson, T.Zhang. *Supervised and Semi-supervised Text Categorization Using LSTM for Region Embeddings.* ICML, 2016.
[4] A. Cardoso-Cachopo. *Improving Methods for Single-label Text Categorization.* PhD thesis, Universidade Tecnica de Lisboa, 2007.
[5] J.Lasserre, C. Bishop *Generative or Discriminative? Getting the Best of Both Worlds.* Bayesian Statistics 8, 2007