
Prediction-Constrained Topic Models for Antidepressant Recommendation

Michael C. Hughes¹, Gabriel Hope², Leah Weiner³,
Thomas H. McCoy, M.D.⁴, Roy H. Perlis, M.D.⁴, Erik B. Sudderth², and Finale Doshi-Velez¹

¹School of Engineering and Applied Sciences, Harvard University
²School of Information & Computer Sciences, Univ. of California, Irvine
³Dept. of Computer Science, Brown University
⁴Massachusetts General Hospital

Abstract

Supervisory signals can help topic models discover low-dimensional data representations that are more interpretable for clinical tasks. We propose a framework for training supervised latent Dirichlet allocation that balances two goals: faithful generative explanations of high-dimensional data and accurate prediction of associated class labels. Existing approaches fail to balance these goals by not properly handling a fundamental asymmetry: the intended task is always predicting labels from data, not data from labels. Our new prediction-constrained objective trains models that predict labels from heldout data well while also producing good generative likelihoods and interpretable topic-word parameters. In a case study on predicting depression medications from electronic health records, we demonstrate improved recommendations compared to previous supervised topic models and high-dimensional logistic regression from words alone.

1 Introduction

Patient history in electronic health records (EHR) can be represented as counts of predefined concepts like procedures, labs, and medications. For such datasets, topic models such as *latent Dirichlet allocation* (LDA, [2]) are popular for extracting insightful low-dimensional structure for clinicians [15; 5]. A natural goal is to use such low-dimensional representations as features for a specific *supervised prediction*, such as recommending drugs to a patient with depression. Many general-purpose efforts have attempted to train *supervised topic models* [22; 11; 3], including the well-known supervised Latent Dirichlet Allocation (sLDA, [14]).

However, a recent survey of healthcare prediction tasks [8] finds that many of these approaches have little benefit, if any, over standard unsupervised LDA for heldout predictions. In this work, we expose and correct several deficiencies in these previous formulations of supervised topic models. We introduce a learning objective that directly enforces the intuitive goal of representing the data in a way that enables accurate downstream predictions. Our objective acknowledges the inherent asymmetry of prediction tasks: clinicians want to predict medication outcomes given medical records, not medical records given outcomes. Approaches like sLDA that optimize the *joint* likelihood of labels and words ignore this crucial asymmetry. Our new *prediction-constrained* (PC) objective for training latent variable models allows practitioners to effectively balance explaining abundant count data while ensuring high-quality predictions of labels from this data. We hope achieving strong gains in this predictive framework will pave the way for causal latent variable models for drug recommendation.

This short paper was accepted at the NIPS Machine Learning for Health workshop (NIPS ML4H 2017). A longer tech report with further experiments [9] is available online: <https://arxiv.org/abs/1707.07341>.

2 Limitations of Existing Topic Models

Supervised LDA. The sLDA model learns from a collection of D documents. Each document d is represented by counts of V discrete words, $x_d \in \mathbb{Z}_+^V$. In the EHR context, these are often ICD-9 or ICD-10 codes, such as “F33.0: Major depressive disorder, recurrent, mild”. For our supervised case, each document (patient) d also has a binary label $y_d \in \{0, 1\}$, indicating whether a medication was successful. Both words x_d and label y_d are generated by a document-specific mixture of K topics:

$$x_d | \pi_d, \phi \sim \text{Mult}(x_d | \sum_{k=1}^K \pi_{dk} \phi_k, N_d), \quad y_d | \pi_d, \eta \sim \text{Bern}(y_d | \sigma(\sum_{k=1}^K \pi_{dk} \eta_k)). \quad (1)$$

The key latent variable is π_d , the document-topic probability vector, with prior $\pi_d \sim \text{Dir}(\alpha)$. The trainable parameters are topic-word probabilities ϕ_k and regression weights η_k (we fix α for simplicity). Let $\sigma(z) = (1 + e^{-z})^{-1}$ be the logit function, and N_d the observed size of document d .

There are a host of objectives and inference methods for supervised LDA, including [14; 19; 21; 22]. A key contribution of this work is identifying a fundamental shortcoming of all these objectives: they do not train topic models that are also effective at label prediction. This myriad of methods, and their shortcomings, arise because of model misspecification. If our count data truly came from a topic model, and those topics truly led to good label predictions, then even *unsupervised* topic models would do well. Trouble arises when we desire the dimensionality reduction provided by a topic model for interpretability or efficiency, but the data were not produced by the LDA generative process.

Limitations of Joint Bayesian or Maximum-Likelihood Training of the sLDA Model. Supervised LDA [14] and related work [19; 20; 16] assumes a graphical model in which the target label y_d can be viewed as yet another output of document-topic probabilities π_d . When the number of counts in x_d is significantly larger than the cardinality of y_d , as is typical in practice, the likelihood associated with x_d will be much larger in magnitude than the likelihood associated with y_d . That is, the *correct* application of Bayesian inference within this model will essentially *ignore* the task of predicting the target y_d . Thus, [8] finds that for large K , sLDA is no better than LDA.

Limitations of Label Replication. The Power-sLDA approach of Zhang and Kjellström [21] suggests improving sLDA’s predictions by artificially replicating the label y_d multiple times. Standard Bayesian methods use both data x_d and replicated labels y_d to infer the document-topic probabilities π_d while training. However, the predictive posterior $p(\pi_d | x_d)$ may be very different from the training posterior $p(\pi_d | x_d, y_d)$. Put another way, label replication strengthens the connection between π_d and y_d , but it does not strengthen the task we care about: prediction of y_d from x_d alone. Figure B.1 in the appendix demonstrates this issue: regardless of the replication level, when the model is misspecified Power-sLDA fails to find topics that are good for predicting y_d from x_d .

Other popular sLDA objectives reduce to label replication. Posterior regularization (PR)-based methods [4; 6] enforce explicit performance constraints on the posterior. The MedLDA approach of Zhu et al. [22, 23, 24] is instead derived from a maximum entropy discrimination framework, and uses a hinge loss to penalize errors in the prediction of y_d . In extended derivations in our longer tech report [9], we show that both MedLDA and PR training objectives can be written as instances of label replication, and thus inherit Power-sLDA’s failure to generalize well.

Limitations of Fully Discriminative Learning. Unlike the above approaches, backpropagation supervised LDA (BP-sLDA, [3]) focuses *entirely* on the prediction of y_d . Chen et al. [3] do handle the direct prediction of y_d from x_d , but no term in their objective forces topics to accurately model the data x_d at all. Our objective can be seen as a principled generalization that balances the explanation of data x_d (which [3] ignores) and prediction of targets y_d . Our improved generative modeling leads to more interpretable topic-word distributions.

3 Prediction-Constrained sLDA

We propose a novel, *prediction-constrained* (PC) objective that explicitly encodes the asymmetry of the discriminative label prediction task. In particular, we ensure that topics learned during joint training can also be used to make accurate predictions about y given x , by solving:

$$\min_{\phi, \eta} - \left[\sum_{d=1}^D \log p(x_d | \phi, \alpha) \right] \text{ subject to } - \sum_{d=1}^D \log p(y_d | x_d, \phi, \eta, \alpha) \leq \epsilon. \quad (2)$$

The scalar ϵ is the highest aggregate loss we are willing to tolerate. There are many variations on this theme; for example, one could instead use a hinge loss as in Zhu et al. [22]. The structure of Eq. (2) matches the goals of a domain expert who wishes to explain as much of the data x as possible, while

still making sufficiently accurate predictions. We recommend adding standard regularization terms (log priors for ϕ and η), though we leave these out to keep notation focused on our contributions.

Applying the Karush-Kuhn-Tucker conditions, Eq. (2) becomes an equivalent unconstrained problem:

$$\min_{\phi, \eta} - \sum_{d=1}^D [\log p(x_d | \phi, \alpha) + \lambda_\epsilon \log p(y_d | x_d, \phi, \eta, \alpha)] \quad (3)$$

For any prediction tolerance ϵ , there exists a scalar multiplier $\lambda_\epsilon > 0$ such that the optimum of Eq. (2) is a minimizer of Eq. (3). The relationship between λ_ϵ and ϵ is monotonic but has no analytic form. We must search over one-dimensional penalties λ_ϵ for an appropriate value.

While our PC objective is superficially similar to Power-sLDA [21] and MedLDA [22], it is distinct: the multiplier λ_ϵ rescales the log-posterior $\log p(y_d | x_d)$, while label-replication rescales the log-likelihood $\log p(y_d | \pi_d)$. By “replicating” the *entire* posterior, rather than just the link between latent and target variables, our PC objective achieves the asymmetric goal of predicting y_d from x_d alone.

Computing $p(x_d | \phi)$ and $p(y_d | x_d, \phi, \eta)$ requires marginalizing π_d over the simplex. However, these integrals are intractable. To gain traction, we first contemplate *instantiating* π_d :

$$\min_{\pi, \phi, \eta} - \sum_{d=1}^D [\log p(\pi_d | \alpha) + \log p(x_d | \pi_d, \phi) + \lambda_\epsilon \log p(y_d | \pi_d, \eta)] \quad (4)$$

As discussed above, solutions to this objective would lead to weighted *joint* training and its symmetry problems. Since we wish to train under the same asymmetric conditions needed at test time, where we have x_d but not y_d , we instead *fix* π_d to a deterministic mapping of the words x_d to the topic simplex. Specifically, we fix to the *maximum a posteriori* (MAP) solution $\pi_d = \operatorname{argmax}_{\pi_d \in \Delta^K} \log p(\pi_d | x_d, \phi, \alpha)$, which we write as an embedding: $\pi_d \leftarrow \operatorname{MAP}_{\phi, \alpha}(x_d)$.

Our chosen embedding can be seen as a feasible approximation to the posterior $p(\pi_d | x_d, \phi, \alpha)$. This choice respects the need to use the same embedding of observed words x_d into low-dimensional π_d in both training and test scenarios. We can now write our tractable training objective for PC-sLDA:

$$- \sum_{d=1}^D [\log p(\operatorname{MAP}_{\phi, \alpha}(x_d) | \alpha) + \log p(x_d | \operatorname{MAP}_{\phi, \alpha}(x_d), \phi) + \lambda_\epsilon \log p(y_d | \operatorname{MAP}_{\phi, \alpha}(x_d), \eta)] \quad (5)$$

While this objective is similar to BP-sLDA [3], the key difference is that our method *balances* the generative and discriminative terms via the multiplier λ_ϵ . In contrast, Chen et al. [3] consider only fully unsupervised (labels y are ignored) or fully supervised (the distribution of x is ignored) cases.

MAP via Exponentiated Gradient. The document-topic MAP problem for unsupervised LDA is $\max_{\pi_d \in \Delta^K} \log p(\pi_d | x_d, \phi, \alpha)$ [17]. It is convex for $\alpha \geq 1$ and non-convex otherwise. For the convex case, we start from uniform probabilities and iteratively do *exponentiated gradient* updates [10]:

$$\text{init: } \pi_d^0 \leftarrow \left[\frac{1}{K} \dots \frac{1}{K} \right], \quad \text{repeat: } \pi_{dk}^t \leftarrow \frac{p_{dk}^t}{\sum_{j=1}^K p_{dj}^t}, \quad p_{dk}^t = \pi_{dk}^{t-1} \circ e^{\nu \nabla \log p(\pi_d^{t-1} | x_d)}. \quad (6)$$

With small enough steps $\nu > 0$, exponentiated gradient converges to the MAP solution. We define our embedding $\operatorname{MAP}_{\phi, \alpha}(x_d)$ to be the deterministic outcome of T iterations of Eq. (6). $T \approx 100$ and $\nu \approx 0.005$ work well. The non-convex case can be solved similarly after reparameterization [18; 12].

Learning via gradient descent. Our entire objective function, including the MAP estimation procedure, is fully *differentiable* with respect to the parameters ϕ, η . Thus, modern gradient descent methods like Adam may be applied to estimate ϕ, η from observed data. We have developed Python implementations using both Autograd [13] and Tensorflow [1], which we will release to the public.

Hyperparameter selection. The key hyperparameter for our PC-sLDA algorithm is the multiplier λ_ϵ . For topic models, λ_ϵ typically needs to scale like the number of tokens in the average document, though it may need to be larger depending on tension between the unsupervised and supervised terms of the objective. In our experiments, we try logarithmically spaced values $\lambda_\epsilon \in \{10, 100, 1000, \dots\}$ and select the best using validation data, although this requires training multiple models. This cost can be somewhat mitigated by using the final parameters at one λ_ϵ value to initialize the next λ_ϵ , although this may not escape to new preferred basins of attraction in the overall non-convex objective.

4 Antidepressant Case Study

We consider predicting which subset of 11 common antidepressants will be successful for a patient with major depressive disorder given a bag-of-words representation x_d of the patient’s electronic health record (EHR). These are real deidentified data from tertiary care hospitals, split into

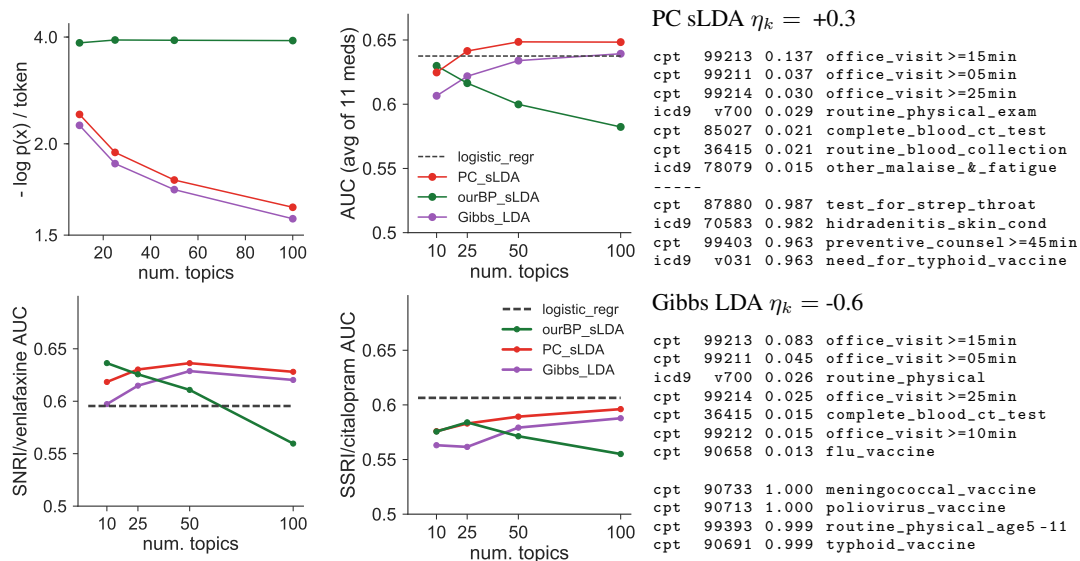


Figure 1: Antidepressant prediction. *Left*: Performance as number of topics increases. We show heldout negative log likelihood (generative, lower is better) and heldout AUC (discriminative, higher is better) for the average of all 11 drugs as well as two specific drugs (one SNRI and one SSRI). We use our own implementation of BP-sLDA for this multiple binary label prediction task. Both PC-sLDA and BP-sLDA results use runs initialized from Gibbs. While BP-sLDA exhibits severe overfitting, our PC-sLDA improves on the baseline Gibbs predictions reliably. *Right*: Interpretation of topic #11 of $K = 25$ for both Gibbs-LDA and our PC-sLDA initialized from Gibbs. We show the regression coefficient η_k for this topic when predicting patient success with citalopram. The top list is ranked by $p(\text{word}|\text{topic})$; the bottom list by $p(\text{topic}|\text{word})$, indicating potential *anchor words*. The original Gibbs topic is mostly about routine preventative care and vaccination. PC sLDA training evolves the topic to emphasize longer duration encounters focused on counseling, with a few unfocused terms.

29774/3721/3722 documents (one per patient) with $V = 5126$ EHR codewords (diagnoses/procedures/medicines). The appendix gives details on data preprocessing. Our results are:

PC-sLDA has better label prediction. Overall, antidepressant recommendation is challenging even for nonlinear classifiers, so we do not expect AUC scores to be very high. However, our PC-sLDA is competitive, beating Gibbs LDA and logistic regression at average prediction across 11 drugs in Fig. 1 when given enough topics. BP-sLDA does well with few topics, but overfits with too many.

PC-sLDA recovers better heldout data likelihoods than BP-sLDA. Fig. 1 shows trends in negative log likelihood on heldout data (lower is better). As expected, unsupervised Gibbs-LDA consistently achieves the best scores, because explaining data is its sole objective. BP-sLDA is consistently poor, having per-token likelihoods about >1.0 nats higher than others. These results show that the solely discriminative approach of BP-sLDA cannot explain the data well. In contrast, our PC-sLDA can capture essential data properties while still predicting labels accurately.

PC-sLDA’s learned topic-word probabilities ϕ are interpretable for the prediction task. We emphasize that our PC training estimates topic-word parameters ϕ that are distinct from unsupervised training and more appropriate for the label prediction task. Fig. 1 shows that PC-sLDA initialized from Gibbs indeed causes an original Gibbs topic to significantly evolve its regression weight η_k and associated top words. The original Gibbs topic covers routine outpatient preventative care and vaccination. The evolved PC-sLDA topic prefers long-duration primary care encounters focused on behavior change (“counseling”). With clinical collaborators, we hypothesize that this more focused topic leads to a positive η_k value because the drug in question (citalopram/Celexa) is often a treatment of choice for patients with uncomplicated MDD diagnosed and treated in primary care.

5 Conclusion

We have presented a new training objective for topic models that can effectively incorporate supervised labels to improve parameter training, even when the model is misspecified. Future work can explore this same PC objective with improved models for observational health records that account for important factors such as patient demographics, temporal evolution, or causality.

References

- [1] M. Abadi, A. Agarwal, P. Barham, et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [2] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [3] J. Chen, J. He, Y. Shen, L. Xiao, X. He, J. Gao, X. Song, and L. Deng. End-to-end learning of LDA by mirror-descent back propagation over a deep architecture. In *Neural Information Processing Systems*, 2015.
- [4] K. Ganchev, J. Graça, J. Gillenwater, and B. Taskar. Posterior regularization for structured latent variable models. *Journal of Machine Learning Research*, 11:2001–2049, Aug. 2010.
- [5] M. Ghassemi, T. Naumann, F. Doshi-Velez, N. Brimmer, R. Joshi, A. Rumshisky, and P. Szolovits. Unfolding physiological state: mortality modelling in intensive care units. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2014.
- [6] J. Graça, K. Ganchev, and B. Taskar. Expectation maximization and posterior constraints. In *Neural Information Processing Systems*, 2008.
- [7] T. L. Griffiths and M. Steyvers. Finding scientific topics. *Proceedings of the National Academy of Sciences*, 2004.
- [8] Y. Halpern, S. Horng, L. A. Nathanson, N. I. Shapiro, and D. Sontag. A comparison of dimensionality reduction techniques for unstructured clinical text. In *ICML workshop on clinical data analysis*, 2012.
- [9] M. C. Hughes, L. Weiner, G. Hope, T. H. McCoy, R. H. Perlis, E. B. Sudderth, and F. Doshi-Velez. Prediction-constrained training for semi-supervised mixture and topic models. *arXiv preprint 1707.07341*, 2017. URL <https://arxiv.org/pdf/1707.07341.pdf>.
- [10] J. Kivinen and M. K. Warmuth. Exponentiated gradient versus gradient descent for linear predictors. *Information and Computation*, 132(1):1–63, 1997.
- [11] S. Lacoste-Julien, F. Sha, and M. I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, 2009.
- [12] D. J. C. MacKay. Choice of basis for Laplace approximation. *Machine Learning*, 33(1), 1998.
- [13] D. Maclaurin, D. Duvenaud, M. Johnson, and R. Adams. Autograd: Reverse-mode differentiation of native python. <http://github.com/HIPS/autograd>, 2015.
- [14] J. D. McAuliffe and D. M. Blei. Supervised topic models. In *Neural Information Processing Systems*, pages 121–128, 2008.
- [15] M. J. Paul and M. Dredze. Discovering health topics in social media using topic models. *PLoS One*, 9(8): e103408, 2014.
- [16] Y. Ren, Y. Wang, and J. Zhu. Spectral learning for supervised topic models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [17] D. Sontag and D. Roy. Complexity of inference in latent dirichlet allocation. In *Neural Information Processing Systems*, 2011.
- [18] M. Taddy. On estimation and selection for topic models. In *Artificial Intelligence and Statistics*, 2012.
- [19] C. Wang, D. Blei, and F.-F. Li. Simultaneous image classification and annotation. In *IEEE Conf. on Computer Vision and Pattern Recognition*, 2009.
- [20] Y. Wang and J. Zhu. Spectral methods for supervised topic models. In *Advances in Neural Information Processing Systems*, pages 1511–1519, 2014.
- [21] C. Zhang and H. Kjellström. How to supervise topic models. In *ECCV Workshop on Graphical Models in Computer Vision*, 2014.
- [22] J. Zhu, A. Ahmed, and E. P. Xing. MedLDA: maximum margin supervised topic models. *The Journal of Machine Learning Research*, 13(1):2237–2278, 2012.
- [23] J. Zhu, N. Chen, H. Perkins, and B. Zhang. Gibbs max-margin topic models with fast sampling algorithms. In *International Conference on Machine Learning*, 2013.
- [24] J. Zhu, N. Chen, and E. P. Xing. Bayesian inference with posterior regularization and applications to infinite latent svms. *Journal of Machine Learning Research*, 15(1):1799–1847, 2014.

A EHR dataset description

We study a cohort of hundreds of thousands patients drawn from two large academic medical centers and their affiliated outpatient networks over a period of several years. Each patient has at least one ICD9 diagnostic code for major depressive disorder (ICD9s 296.2x or 3x or 311, or ICD10 equivalent). Each included patient had an identified successful treatment which included one of 25 possible common anti-depressants marked as “primary” treatments for major depressive disorder by clinical collaborators. We labeled an interval of a patient’s record “successful” if all prescription events in the interval used the same subset of primary drugs, the interval lasted at least 90 days, and encounters occurred at least every 13 months. Applying this criteria, we identified 64431 patients who met our definition of success. For each patient, we extracted a bag-of-codewords x_d of 5126 possible codewords (representing medical history before any successful treatment) and binary label vector y_d , marking which of 11 prevalent anti-depressants (if any) were used in known successful treatment.

Extracting data x_d . For each patient with known successful treatment, we build a data vector x_d to summarize all facts known about the patient in the EHR before any successful treatment was given. Thus, we must confine our records to the interval from the patient’s first encounter to the last encounter before any of the drugs on his or her successful list were first prescribed. To summarize this patient’s interval of “pre-successful treatment”, we built a sparse count vector of all procedures, diagnoses, labs, and medications from the EHR which fit within the interval (22,000 possible codewords). By definition, none of the anti-depressant medications on the patient’s eventual success list appear in x_d . To simplify, we reduced this to a final vocabulary of 5126 codewords that occurred in at least 1000 distinct patients. We discard any patients with fewer than 2 tokens in x_d (little to no history).

Extracting labels y_d . Among the 25 primary drugs, we identified a smaller set of 11 anti-depressants which were used in “successful treatment” for at least 1000 patients. The remaining 15 primary drugs did not occur commonly enough that we could accurately assess prediction quality (build large enough heldout sets). Our chosen list of drugs to predict are: nortriptyline, amitriptyline, bupropion, fluoxetine, sertraline, paroxetine, venlafaxine, mirtazapine, citalopram, escitalopram, and duloxetine. Because these drugs can be given in combination, this is a multiple binary label problem. Future work could look into structured prediction tasks.

B Toy Bars Case Study

To study tradeoffs between models of $p(x)$ and $p(y|x)$, we built a toy dataset that is deliberately *misspecified*: neither the unsupervised LDA maximum likelihood solution nor the standard sLDA joint likelihood solution performs much better than chance at label prediction. We look at 500 training documents, each with $V = 9$ possible vocabulary words that can be arranged in a 3-by-3 grid to indicate some bar-like co-occurrence structure. Each binary label y_d is unrelated to the bar structure, but is unambiguously indicated by the top-left word. We visualize some documents and their associated labels in the top row of Fig. B.1.

We compare our proposed PC sLDA training procedure with several competitors (MED sLDA [22], Gibbs unsupervised LDA [7], BP sLDA [3], etc.). We also include a method that maximizes the joint likelihood $\log p(x, y)$ with different amounts of label replication. We call this method ML-sLDA with replication value λ . The special case of $\lambda = 1$ is standard sLDA, while the case of $\lambda \gg 1$ is known as Power sLDA [21].

Each algorithm is run to convergence on training data, and its best parameters – topic-word probabilities ϕ and regression weights η – are chosen to minimize method-specific training loss. Each method’s best solution is then located on a 2-dimensional fitness landscape: the x-axis is negative log likelihood of data x averaged per token (lower is better) and y-axis is the negative log likelihood of labels y averaged per document (lower is better). These averages are computed on the *training* set. We show these fitness scores under two possible modes for estimating each document-topic vector π_d . *Train mode* computes the joint likelihood MAP estimate $\max_{\pi_d} \log p(\pi_d|x_d, y_d, \phi, \eta, \alpha)$. *Predict mode* computes the data-only MAP estimate $\max_{\pi_d} \log p(\pi_d|x_d, \phi, \alpha)$. This distinction highlights the key difference between PC-sLDA with high λ , which deliberately trains parameters ϕ, η to be good at prediction, and alternatives like maximum likelihood with label replication (ML with $\lambda > 1$),

which trains models that do well in training mode but fail miserably in a predictive setting (even on the training set). We further see that methods that purely optimize label prediction such as BP-sLDA achieve reasonable prediction scores but *terrible* data likelihood scores.

The visualized parameters show an important trend: Our PC-sLDA with $\lambda \geq 10$ is the only method to use just one topic to explain the signal word. Thus, it is the only method to reach the sweet spot of good $y|x$ predictions and good x explanations. Gibbs sensibly finds 4 bars and places the signal word slightly in each one, as does MED-sLDA and Power sLDA.

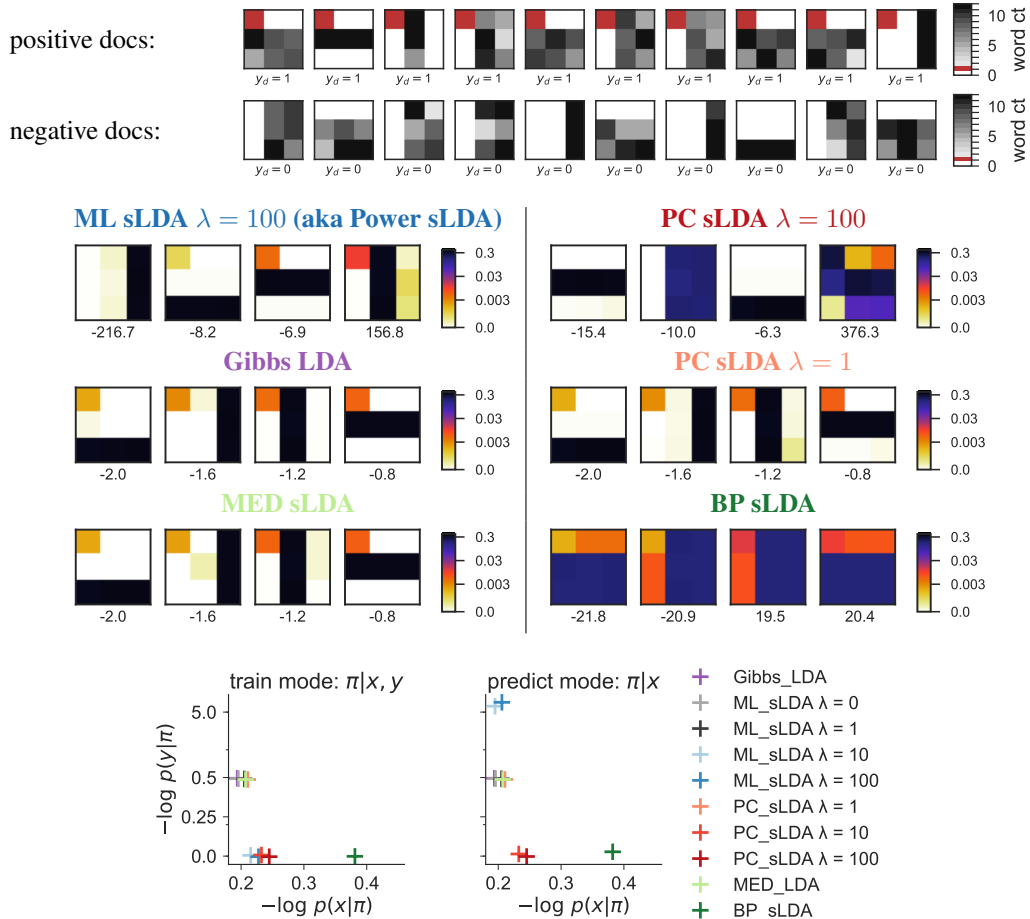


Figure B.1: 3×3 bars task: advantages of PC training under misspecification. We compare several training procedures to see which can simultaneously model the bar-like co-occurrence structure while making accurate binary label predictions. *Top rows*: Example labeled training documents: the 3×3 heatmap shows the word count vector x_d , and the caption indicates the label y_d . Colormap chosen to highlight the top-left-corner symbol that, when it appears just once, perfectly signals the document belongs to the positive class ($y_d = 1$). Remaining vocabulary symbols in each document are drawn from one or two of 4 possible horizontal and vertical “bar” topics. These symbols, when non-zero, have much higher counts than the top-left signal word. *Middle rows*: Visualization of the topic-word probabilities for the best $K = 4$ topic model trained by each method. Colormap has a *logarithmic scale* to show how the rare signal word is explained. *Bottom row*: Location of each method’s estimated parameters on the fitness landscape where x-axis is generative model training loss, and y-axis is prediction task loss. The lower left corner is the ideal position.