

MedLDA: Maximum Margin Supervised Topic Models

Jun Zhu

Amr Ahmed

Eric P. Xing

Presented by Zhengming Xing

Outline

- Introduction
- Background
- MedLDA
- Experiment result

Introduction

Main idea:

- Supervised topic modeling (Utilize the side information, such as rating and label information)
- Integrate the mechanism behind max-margin prediction algorithm with hierarchical Bayesian topic modeling.

Advantage:

- Discover sparse and highly discriminative topical representation
- Achieve state of art prediction performance
- More efficient than existing supervised model

Background (Topic modeling)

LDA

1. For a document d , draw a topic mixing proportion vector θ_d according to a K -dimensional Dirichlet prior: $\theta_d | \alpha \sim \text{Dir}(\alpha)$;
2. For the n -th word in document d , where $1 \leq n \leq N$,
 - (a) draw a topic assignment z_{dn} according to θ_d : $z_{dn} | \theta_d \sim \text{Mult}(\theta_d)$;
 - (b) draw the word instance w_{dn} according to z_{dn} : $w_{dn} | z_{dn}, \beta \sim \text{Mult}(\beta_{z_{dn}})$,

Joint distribution

$$p(\{\theta_d, \mathbf{z}_d\}, \mathbf{W} | \alpha, \beta) = \prod_{d=1}^D p(\theta_d | \alpha) \left(\prod_{n=1}^N p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right)$$

Supervised LDA

$$p(\{\theta_d, \mathbf{z}_d\}, y, \mathbf{W} | \alpha, \beta, \eta, \delta^2) = \prod_{d=1}^D p(\theta_d | \alpha) \left(\prod_{n=1}^N p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right) p(y_d | \eta^\top \bar{\mathbf{z}}_d, \delta^2)$$

Side information

$$y | \mathbf{z}_d, \eta, \delta^2 \sim \mathcal{N}(\eta^\top \bar{\mathbf{z}}_d, \delta^2) \quad \bar{\mathbf{z}}_d \triangleq \frac{1}{N} \sum_n z_{dn}$$

Background (Topic modeling)

marginal data likelihood

$$P(\mathcal{D}) = \int P(\mathcal{D}, \Omega) d\Omega$$

Variational bound

$$\ln P(\mathcal{D}) = \mathcal{L}(Q) + KL(Q||P) \quad \mathcal{L} = \int Q(\Omega) \ln \frac{P(\mathcal{D}, \Omega)}{Q(\Omega)} d\Omega$$

LDA

$$\begin{aligned} \mathcal{L}^u(q; \alpha, \beta) &\triangleq -\mathbb{E}_q[\log p(\{\theta_d, \mathbf{z}_d\}, \mathbf{W} | \alpha, \beta)] - \mathcal{H}(q(\{\theta_d, \mathbf{z}_d\})) \\ &\geq -\log p(\mathbf{W} | \alpha, \beta), \end{aligned}$$

approximate the marginal data likelihood

$$p(\mathbf{W} | \alpha, \beta)$$

sLDA

$$\begin{aligned} \mathcal{L}^s(q; \alpha, \beta, \eta, \delta^2) &\triangleq -\mathbb{E}_q[\log p(\{\theta_d, \mathbf{z}_d\}, \mathbf{y}, \mathbf{W} | \alpha, \beta, \eta, \delta^2)] - \mathcal{H}(q(\{\theta_d, \mathbf{z}_d\})) \\ &\geq -\log p(\mathbf{y}, \mathbf{W} | \alpha, \beta, \eta, \delta^2). \end{aligned}$$

approximate the marginal data likelihood

$$p(\mathbf{y}, \mathbf{W} | \alpha, \beta, \eta, \delta^2)$$

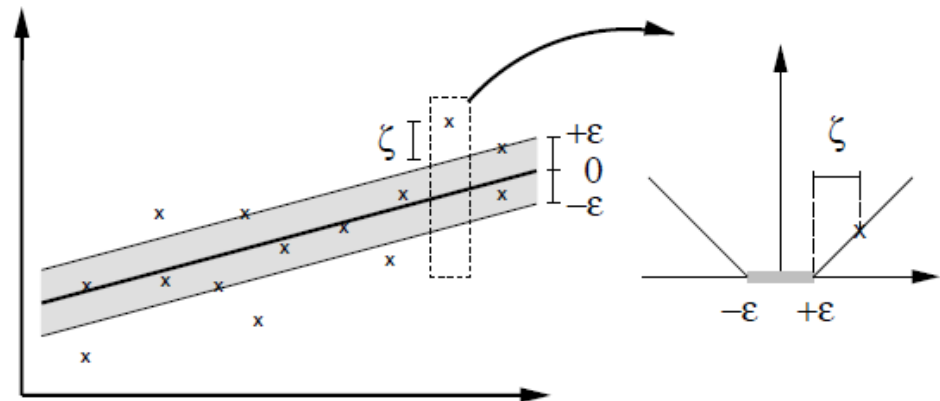
Background (SVR)

Given the training data

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_D, y_D)\}$$

Objective function

$$\text{P0(SVR)} : \min_{\boldsymbol{\eta}, \xi, \xi^*} \frac{1}{2} \|\boldsymbol{\eta}\|_2^2 + C \sum_{d=1}^D (\xi_d + \xi_d^*)$$
$$\forall d, \text{ s.t. : } \begin{cases} y_d - \boldsymbol{\eta}^\top \mathbf{f}(\mathbf{x}_d) \leq \epsilon + \xi_d \\ -y_d + \boldsymbol{\eta}^\top \mathbf{f}(\mathbf{x}_d) \leq \epsilon + \xi_d^* \\ \xi_d, \xi_d^* \geq 0 \end{cases}$$



Background (Med)

Maximum entropy discrimination

Find the distribution of all possible regression coefficient belong to a particular distribution class family

Objective function

$$\text{P1(MED}^{\text{r}}) : \min_{q(\boldsymbol{\eta}), \xi, \xi^*} KL(q(\boldsymbol{\eta}) || p_0(\boldsymbol{\eta})) + C \sum_{d=1}^D (\xi_d + \xi_d^*)$$
$$\forall d, \text{ s.t. : } \begin{cases} y_d - \mathbb{E}[\boldsymbol{\eta}]^{\text{T}} \mathbf{f}(\mathbf{x}_d) \leq \epsilon + \xi_d \\ -y_d + \mathbb{E}[\boldsymbol{\eta}]^{\text{T}} \mathbf{f}(\mathbf{x}_d) \leq \epsilon + \xi_d^* \\ \xi_d, \xi_d^* \geq 0 \end{cases},$$

$p_0(\boldsymbol{\eta})$

Prior distribution

$$KL(p||q) \triangleq \mathbb{E}_p[\log(p/q)]$$

Regression-MedLDA

Side information

Numerical rate of a document in the corpus

Regression case

Y is the response response

$$\hat{y} \triangleq \mathbb{E}[Y|\mathbf{w}, \alpha, \beta, \delta^2] = \mathbb{E}[\eta^\top \bar{Z}|\mathbf{w}, \alpha, \beta, \delta^2]$$

Take expectation with respect to $q(\eta, \mathbf{z})$

let

$$q(\eta, \mathbf{z}) = \int_{\theta} q(\eta)q(\mathbf{z}, \theta|\eta)$$

and

$q(\eta, \{\theta_d, \mathbf{z}_d\})$ be the approximation to $p(\eta, \{\theta_d, \mathbf{z}_d\}|\alpha, \beta, \delta^2, \mathbf{y}, \mathbf{W})$

Upper bound of negative log-likelihood

$$\begin{aligned} \mathcal{L}^{bs}(q; \alpha, \beta, \delta^2) &\triangleq -\mathbb{E}_q[\log p(\eta, \{\theta_d, \mathbf{z}_d\}, \mathbf{y}, \mathbf{W}|\alpha, \beta, \delta^2)] - \mathcal{H}(q(\eta, \{\theta_d, \mathbf{z}_d\})) \\ &= KL(q(\eta)||p_0(\eta)) + \mathbb{E}_{q(\eta)}[\mathcal{L}^s]. \end{aligned}$$

Regression-MedLDA

Objective function

$$\text{P2(MedLDA}^r) : \min_{q, \alpha, \beta, \delta^2, \xi, \xi^*} \mathbb{E}_{q(\boldsymbol{\eta})} [\mathcal{L}^s(q; \boldsymbol{\alpha}, \boldsymbol{\beta}, \delta^2)] + KL(q(\boldsymbol{\eta}) \| p_0(\boldsymbol{\eta})) + C \sum_{d=1}^D (\xi_d + \xi_d^*)$$

$$\forall d, \text{ s.t. : } \begin{cases} y_d - \mathbb{E}[\boldsymbol{\eta}^\top \bar{\mathbf{Z}}_d] \leq \epsilon + \xi_d, & (\mu_d) \\ -y_d + \mathbb{E}[\boldsymbol{\eta}^\top \bar{\mathbf{Z}}_d] \leq \epsilon + \xi_d^*, & (\mu_d^*) \\ \xi_d \geq 0, & (v_d) \\ \xi_d^* \geq 0, & (v_d^*) \end{cases}$$

$$\begin{aligned} \mathcal{L}^s(q; \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}, \delta^2) &\triangleq -\mathbb{E}_q[\log p(\{\boldsymbol{\theta}_d, \mathbf{z}_d\}, \mathbf{y}, \mathbf{W} | \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}, \delta^2)] - \mathcal{H}(q(\{\boldsymbol{\theta}_d, \mathbf{z}_d\})) \\ &\geq -\log p(\mathbf{y}, \mathbf{W} | \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\eta}, \delta^2). \end{aligned}$$

Classificational-MedLDA

Side information

Discrete label of a document in the corpus

Multi-class case

Y is discrete response and take value from $\mathcal{C} \triangleq \{1, \tilde{2}, \dots, J\}$.

Linear discriminant function

$$F(y, \mathbf{z}, \boldsymbol{\eta}; \mathbf{w}) = \boldsymbol{\eta}_y^\top \bar{\mathbf{z}}.$$

Effective discriminant function

$$F(y; \mathbf{w}) = \mathbb{E}[F(y, \mathbf{Z}, \boldsymbol{\eta}; \mathbf{w})] = \mathbb{E}[\boldsymbol{\eta}^\top \mathbf{f}(y, \bar{\mathbf{Z}}) | \boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{w}].$$

Take expectation with respect to

$$q(\boldsymbol{\eta}, \mathbf{z})$$

Prediction rule

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{C}} F(y; \mathbf{w}) = \operatorname{argmax}_{y \in \mathcal{C}} \mathbb{E}[\boldsymbol{\eta}^\top \mathbf{f}(y, \bar{\mathbf{Z}}) | \boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{w}]$$

Classificational-MedLDA

Objective function

$$\text{P3(MedLDA}^c\text{)} : \min_{q, q(\boldsymbol{\eta}), \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\xi}} \mathcal{L}^u(q; \boldsymbol{\alpha}, \boldsymbol{\beta}) + KL(q(\boldsymbol{\eta}) || p_0(\boldsymbol{\eta})) + \frac{C}{D} \sum_{d=1}^D \xi_d$$
$$\forall d, y \neq y_d, \text{ s.t. : } \begin{cases} \mathbb{E}[\boldsymbol{\eta}^\top \Delta \mathbf{f}_d(y)] \geq \Delta \ell_d(y) - \xi_d \\ \xi_d \geq 0, \end{cases}$$

$\Delta \ell_d(y)$ Cost function measure the different between prediction and the true label
eg.(0/1)

$$\mathbb{E}[\boldsymbol{\eta}^\top \Delta \mathbf{f}_d(y)] = F(y_d; \mathbf{w}_d) - F(y; \mathbf{w}_d)$$

$$\begin{aligned} \mathcal{L}^u(q; \boldsymbol{\alpha}, \boldsymbol{\beta}) &\triangleq -\mathbb{E}_q[\log p(\{\boldsymbol{\theta}_d, \mathbf{z}_d\}, \mathbf{W} | \boldsymbol{\alpha}, \boldsymbol{\beta})] - \mathcal{H}(q(\{\boldsymbol{\theta}_d, \mathbf{z}_d\})) \\ &\geq -\log p(\mathbf{W} | \boldsymbol{\alpha}, \boldsymbol{\beta}), \end{aligned}$$

A general framework

Objective function

$$\text{P5(MedTM)} : \min_{q(\Upsilon, H), \Psi, \xi} \mathbb{E}_{q(\Upsilon)} \left[\mathcal{L}^t(q(H|\Upsilon); \Psi, \Upsilon) \right] + KL(q(\Upsilon) \| p_0(\Upsilon)) + U(\xi)$$

s.t. : $q(\Upsilon, H)$ satisfies the expected margin constraints.

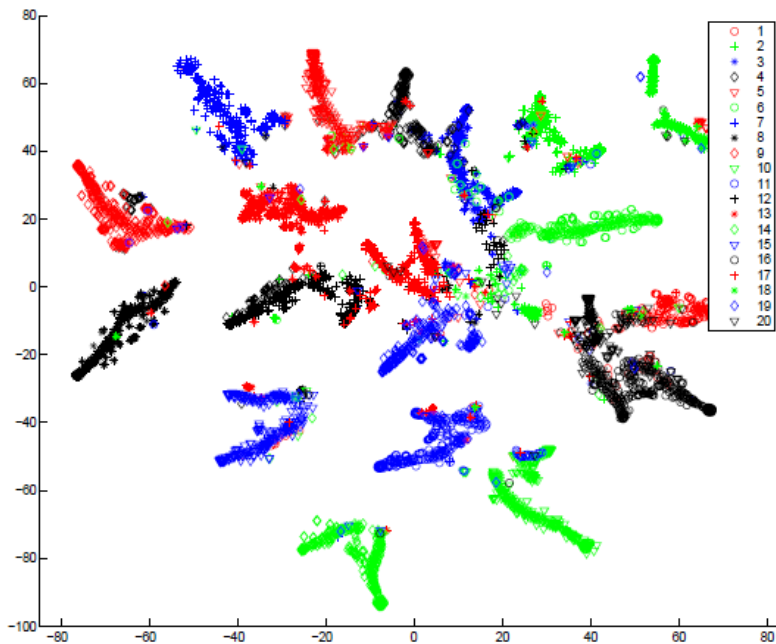
Υ	Denote the parameters of the model pertinent to the prediction task	η in sLDA)
H	Topic assignment and mixing variable	\mathbf{z} and θ
Ψ	Parameter of the underlying topic model	α and topics β

Experiment

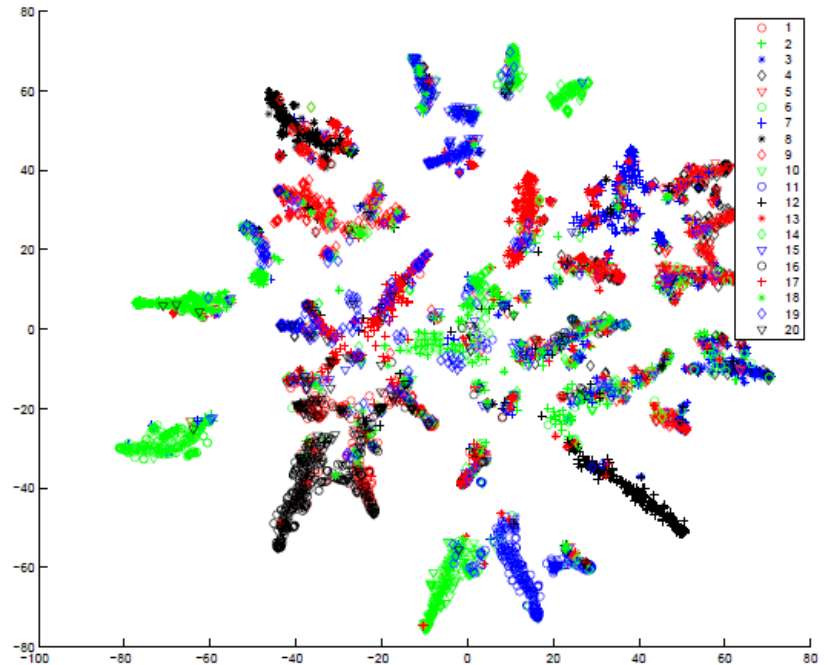
Dataset:

20 news group dataset, 20000 documents in 20 categories. Set the topic numbers to 110.

Embedding:



Med-LDA



LDA

Experiment

Top topic

Class	MedLDA			LDA			Average θ per class	
comp.graphics	T 69	T 11	T 80	T 59	T 104	T 31		
	image jpeg gif file color files bit images format program	graphics image data ftp software pub mail package fax images	db key chip encryption clipper system government keys law escrow	image jpeg color file gif images format bit files display	ftp pub graphics mail version tar file information send server	card monitor dos video apple windows drivers vga cards graphics		
sci.electronics	T 32	T 95	T 46	T 30	T 84	T 44		
	ground wire power wiring don current circuit neutral writes work	audio output input signal chip high data mhz time good	source rs time john cycle low dixie dog weeks face	power ground wire circuit supply voltage current wiring signal cable	water energy air nuclear loop hot cold cooling heat temperature	sale price offer shipping sell interested mail condition email cd		

Experiment

Top topic

politics.mideast	T 30	T 40	T 51	T 42	T 78	T 47	
	israel israeli jews arab writes people article jewish state rights	turkish armenian armenians armenia people turks greek turkey government soviet	israel lebanese israeli lebanon people attacks soldiers villages peace writes	israel israeli peace writes article arab war lebanese lebanon people	israel jewish israel israeli arab people arabs center jew nazi	armenian turkish armenians armenia turks genocide russian soviet people muslim	
misc.forsale	T 109	T 110	T 84	T 44	T 94	T 49	
	sale price shipping offer mail condition interested sell email dos	drive scsi mb drives controller disk ide hard bus system	mac apple monitor bit mhz card video speed memory system	sale price offer shipping sell interested mail condition email cd	don mail call package writes send number ve hotel credit	drive scsi disk hard mb drives ide controller floppy system	

Experiment(classification)

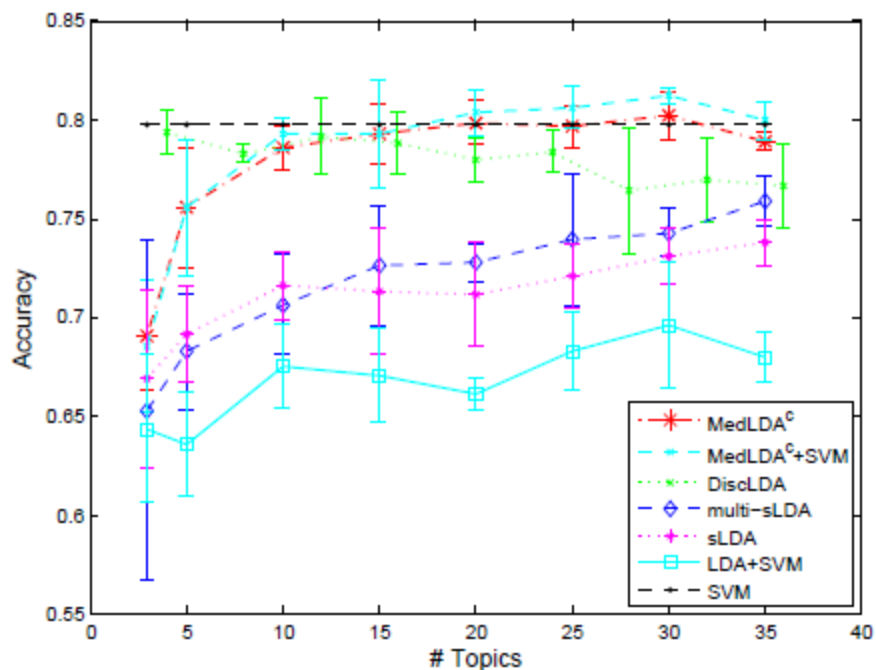
Dataset:

Binary classification:

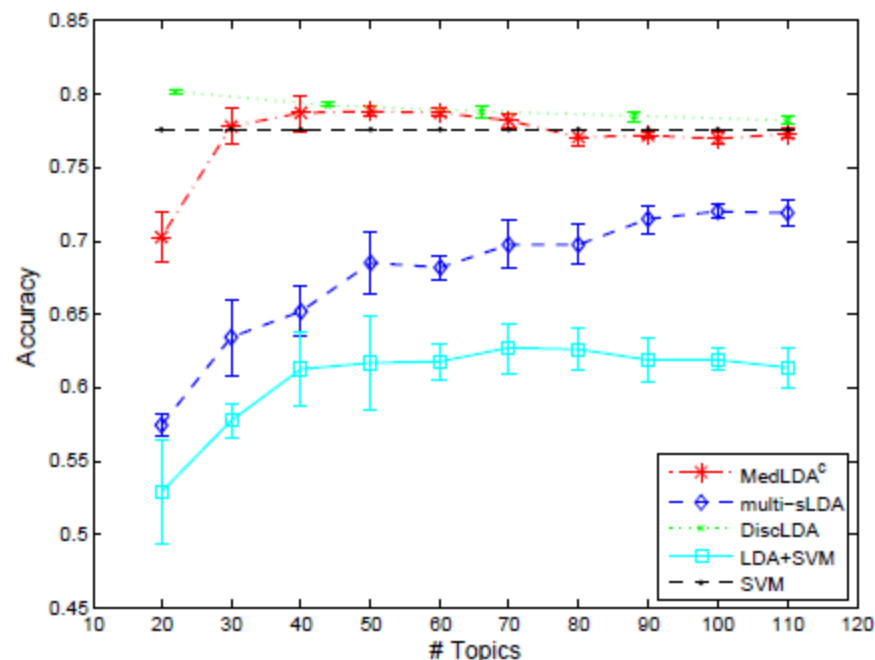
Training: 856 documents with 480 and 376 over two categories

Multi-class classification (20 categories):

Training: 11269 documents
Testing: 7505 documents



binary



Multi-class

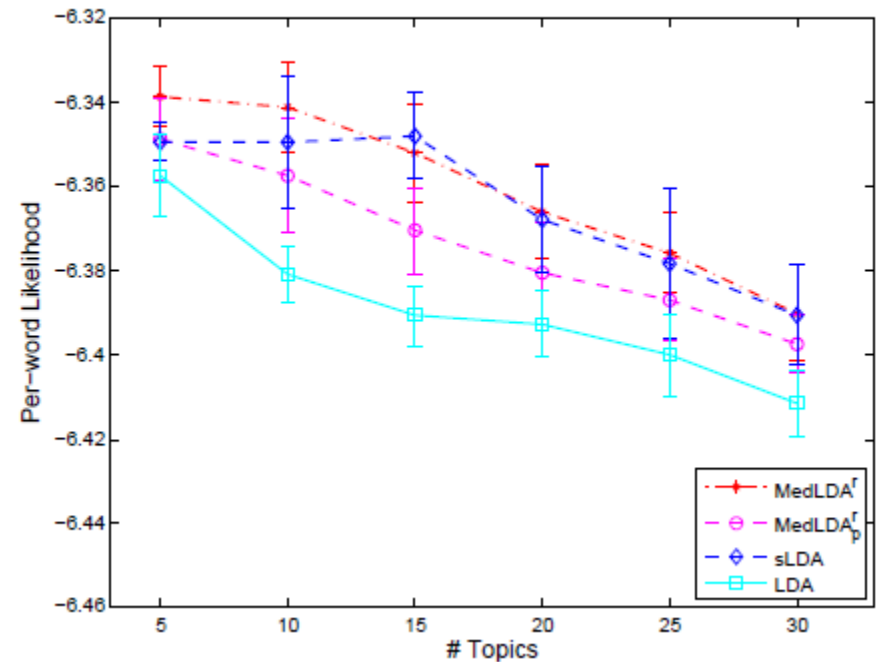
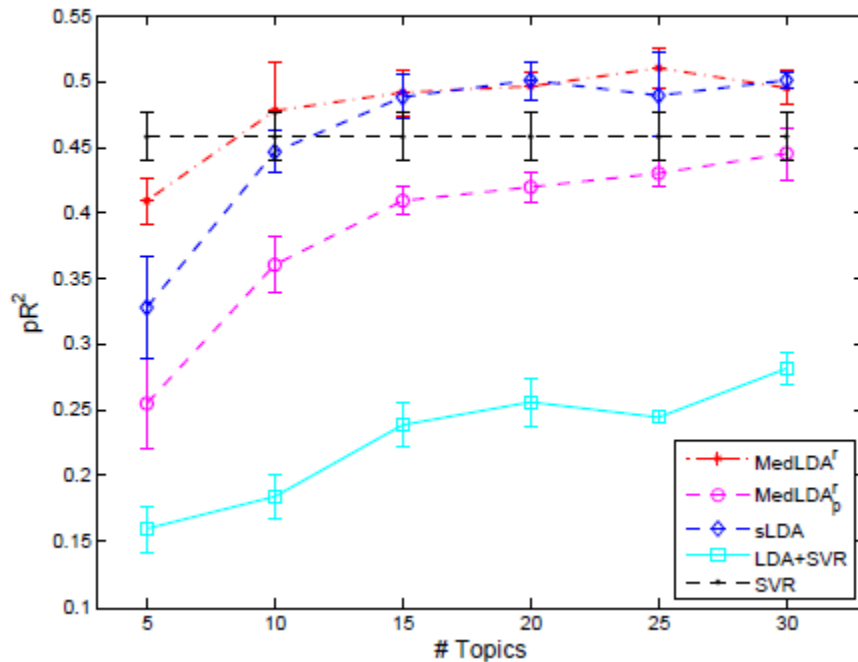
Experiment (regression)

Dataset:

Movie review data: 5006 documents

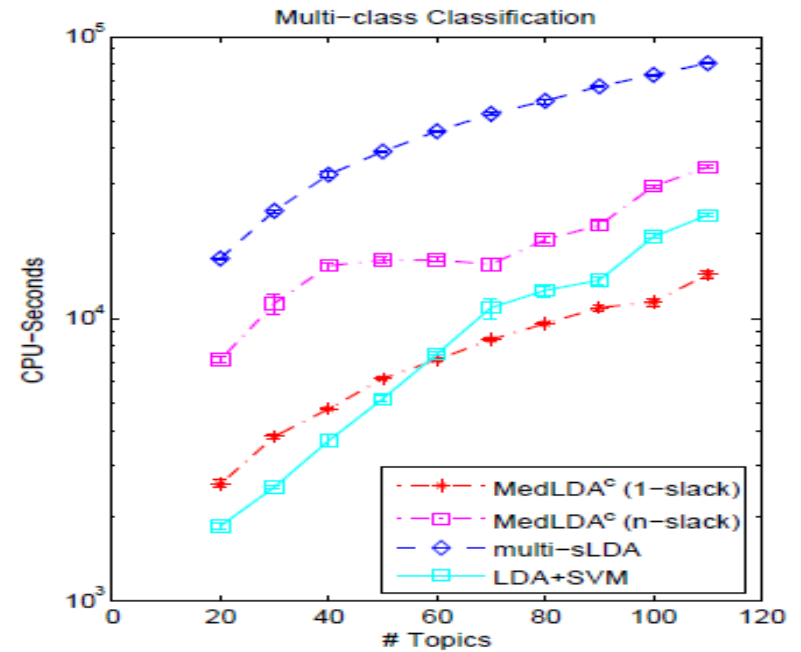
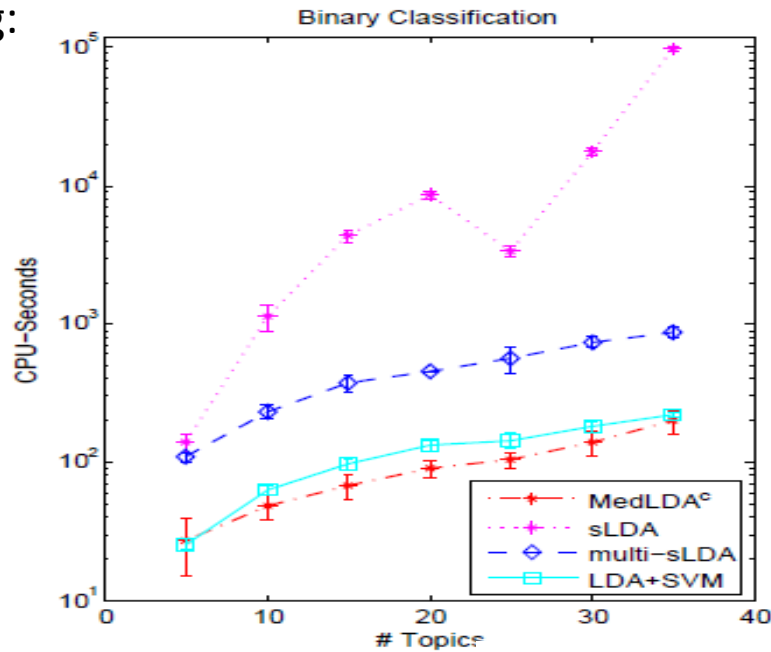
Measure criteria

$$pR^2 = 1 - \frac{\sum_{d=1}^D (y_d - \hat{y}_d)^2}{\sum_{d=1}^D (y_d - \bar{y})^2},$$



Experiment (time efficiency)

Training:



Testing:

