

Sparse Latent Semantic Analysis

Xi Chen¹, Yanjun Qi², Bing Bai², Qihang Lin¹, Jaime G. Carbonell¹

1. Machine Learning Department, Carnegie Mellon University
2. Machine Learning Department, NEC Lab America

The work is done during the
internship at NEC Lab America

Background

❖ Vector Space Model:

Document: $\mathbf{x} = [w_1, \dots, w_M] \in \mathbb{R}^M$ M : vocabulary size

w_i : normalized weight (tf-idf) of the i -th word

N Documents: $\mathbf{X} = [\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N] \in \mathbb{R}^{N \times M}$: Document-Word matrix

❖ Latent Semantic Analysis:

D latent topics (dimensionality of the latent space)

LSA applies SVD to construct a *rank- D* approximation:

$$\mathbf{X} \approx \mathbf{U}_{N \times D} \mathbf{S}_{D \times D} (\mathbf{V}_{M \times D})^T, \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}, \mathbf{V}^T \mathbf{V} = \mathbf{I}$$

Projection Matrix: $\mathbf{A} = \mathbf{S}^{-1} \mathbf{V}^T \in \mathbb{R}^{D \times M}$

Dimension reduction for a new document q : $q \in \mathbb{R}^M \Rightarrow \hat{q} = \mathbf{A}q \in \mathbb{R}^D$

Optimization Formulation for LSA

❖ Latent Semantic Analysis

$$\mathbf{X} \approx \mathbf{U}_{N \times D} \mathbf{S}_{D \times D} (\mathbf{V}_{M \times D})^T, \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}, \mathbf{V}^T \mathbf{V} = \mathbf{I}$$

❖ Relaxed Optimization Formulation:

[K. Yu et al. 05]

$$\begin{aligned} \min_{\mathbf{U}, \mathbf{A}} \quad & \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 \\ \text{subject to:} \quad & \mathbf{U}^T \mathbf{U} = \mathbf{I} \end{aligned}$$

❖ Sparse Latent Semantic Analysis:

Add *sparsity* constraint on the project matrix \mathbf{A} :

$$\begin{aligned} \min_{\mathbf{U}, \mathbf{A}} \quad & \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 + \lambda \|\mathbf{A}\|_1 \quad \Rightarrow \quad \|\mathbf{A}\|_1 = \sum_{d=1}^D \sum_{j=1}^M |a_{dj}| \\ \text{subject to:} \quad & \mathbf{U}^T \mathbf{U} = \mathbf{I} \end{aligned}$$

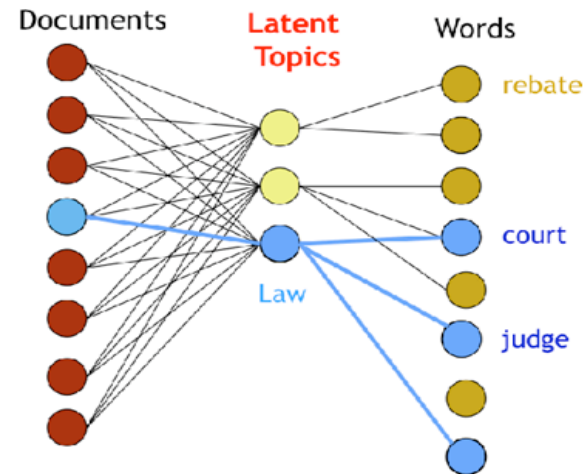
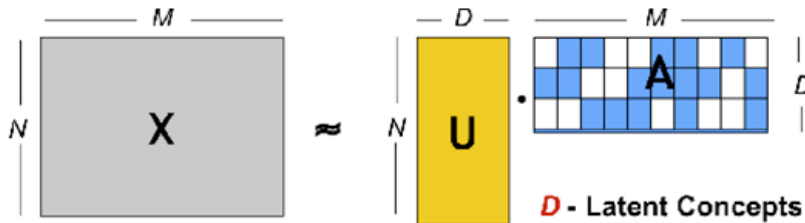
ℓ_1 -regularization

Sparse LSA

❖ Sparse LSA

$$\min_{\mathbf{U}, \mathbf{A}} \quad \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 + \lambda \|\mathbf{A}\|_1$$

subject to: $\mathbf{U}^T \mathbf{U} = \mathbf{I}$



New Document q : $\hat{q} = \mathbf{A}q \in \mathbb{R}^D$

Simple Projection, Computational Efficient



❖ Comparison to Sparse Coding

$$\min_{\mathbf{U}, \mathbf{A}} \quad \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 + \lambda \|\mathbf{U}\|_1$$

subject to: $\|A_j\|_2^2 \leq c, \quad j = 1, \dots, M$

New Document q : $\hat{q} = \arg \min_{\hat{q}} \frac{1}{2} \|q - \mathbf{A}^T \hat{q}\| + \lambda \|\hat{q}\|_1.$

Lasso Problem
More Computation Time 😞

Advantage of Sparse LSA

❖ Better Interpretability:

Sparse LSA selects most relevant words for each topic ($a_{dj} \neq 0$)
Compact representation of topic-word relationship

❖ Efficient Projection:

Sparse $\mathbf{A} \implies$ Efficient Projection for new documents: $\hat{q} = \mathbf{A}q \in \mathbb{R}^D$

❖ Cheap Storage:

Cheap storage for sparse \mathbf{A}

❖ Sparse Projected Documents: sparse $\hat{q} = \mathbf{A}q$

❖ Document-Topic Relationship: $\hat{q}_d = 0 \Leftrightarrow \hat{q}$ not belong to d -th topic

❖ Advantage as compared to PCA:

Do not need to centralize $\mathbf{X} \implies$ destroy the sparsity of \mathbf{X}


Do not need the covariance matrix $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{M \times M} \implies$ may not fit in the memory for large vocabulary size

Optimization Method

❖ Alternating Approach

$$\begin{aligned} \min_{\mathbf{U}, \mathbf{A}} \quad & \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 + \lambda \|\mathbf{A}\|_1 \\ \text{subject to:} \quad & \mathbf{U}^T \mathbf{U} = \mathbf{I} \end{aligned}$$

Fix \mathbf{U} and optimize with respect to \mathbf{A} :

$$\min_{\mathbf{A}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 + \lambda \|\mathbf{A}\|_1$$

$$\min_{A_j} \frac{1}{2} \|X_j - \mathbf{U}A_j\|_2^2 + \lambda \|A_j\|_1; \quad j = 1, \dots, M. \quad A_j: j\text{-th column of } \mathbf{A}$$

[J. Friedman et al. 10]

M independent *lasso* problem : *Solved via Coordinate Descent*

Fix \mathbf{A} and optimize with respect to \mathbf{U} :

$$\begin{aligned} \min_{\mathbf{U}} \quad & \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 \Leftrightarrow \text{tr}(\mathbf{U}^T \mathbf{X} \mathbf{A}^T) \\ \text{subject to:} \quad & \mathbf{U}^T \mathbf{U} = \mathbf{I} \end{aligned}$$

Closed-form Solution:

Let $\mathbf{V} = \mathbf{X}\mathbf{A}^T$ (projected documents onto the latent space)

Perform SVD on \mathbf{V} : $\mathbf{V} = \mathbf{P}\Delta\mathbf{Q} \Rightarrow \mathbf{U}^* = \mathbf{P}\mathbf{Q}$

Note: SVD on $\mathbf{V} \in \mathbb{R}^{M \times D}$ is much *cheaper* than that on $\mathbf{X} \in \mathbb{R}^{N \times M}$



Optimization Summary

Algorithm 1 Optimization Algorithm for Sparse LSA

Input: \mathbf{X} , dimensionality of the latent space D , regularization parameter λ

Initialization: $\mathbf{U}^0 = \begin{pmatrix} \mathbf{I}_D \\ \mathbf{0} \end{pmatrix}$,

Iterate until convergence of \mathbf{U} and \mathbf{A} :

1. Compute \mathbf{A} by solving M independent lasso problems via coordinate descent
2. Project \mathbf{X} onto the latent space: $\mathbf{V} = \mathbf{X}\mathbf{A}^T$.
3. Compute the SVD of \mathbf{V} : $\mathbf{V} = \mathbf{P}\Delta\mathbf{Q}$ and let $\mathbf{U} = \mathbf{P}\mathbf{Q}$.

Output: Sparse projection matrix \mathbf{A} .

Extension I : Nonnegative Sparse LSA

Constraint: $\mathbf{A} \geq 0$:

$$\begin{aligned} \min_{\mathbf{U}, \mathbf{A}} \quad & \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 + \lambda \|\mathbf{A}\|_1 \\ \text{subject to:} \quad & \mathbf{U}^T \mathbf{U} = \mathbf{I}, \quad \mathbf{A} \geq 0. \end{aligned}$$

Simulate the *probability* of the word w_j given the topic t_d :

$$\text{Normalize each row: } \tilde{a}_{dj} = \frac{a_{dj}}{\sum_{j=1}^M a_{dj}} \sim \mathbb{P}(w_j | t_d)$$

Optimization with respect to \mathbf{A} :

$$\min_{\mathbf{A}_j \geq \mathbf{0}} f(\mathbf{A}_j) = \frac{1}{2} \|\mathbf{X}_j - \mathbf{U}\mathbf{A}_j\|_F^2 + \lambda \sum_{d=1}^D a_{dj}. \quad j = 1, \dots, M$$

Optimize via the *coordinate descent* approach:

Iterating over d : fix $a_{\hat{d}j}$ for $\hat{d} \neq d$ and optimize over a_{dj}

$$a_{dj}^* = \begin{cases} \frac{b_d - \lambda}{c_d} & b_d > \lambda \\ 0 & b_d \leq \lambda \end{cases},$$

$$c_d = \sum_{i=1}^N u_{id}^2, b_d = \sum_{i=1}^N u_{id} (x_{ij} - \sum_{k \neq d} u_{ik} a_{kj}).$$

Extension II : Group Structured Sparse LSA

❖ Application: latent gene-function identification: determine relevant pathways (groups of genes) to a latent gene function (topic)

❖ Group Structured Sparse LSA:

The set of groups of input features $\mathcal{G} = \{g_1, \dots, g_{|\mathcal{G}|}\}$ (available as *a priori*)

$$\min_{\mathbf{U}, \mathbf{A}} \quad \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{A}\|_F^2 + \lambda \sum_{d=1}^D \sum_{g \in \mathcal{G}} w_g \|\mathbf{A}_{dg}\|_2$$

$$\text{subject to: } \mathbf{U}^T \mathbf{U} = \mathbf{I}.$$

Optimization with respect to \mathbf{A} :

Optimize via the *coordinate descent* approach:

$$\mathbf{A}_{dg}^* = \begin{cases} \frac{\mathbf{B}_{dg} (\|\mathbf{B}_{dg}\|_2 - \lambda w_g)}{C_d \|\mathbf{B}_{dg}\|_2} & \|\mathbf{B}_{dg}\|_2 > \lambda w_g \\ \mathbf{0} & \|\mathbf{B}_{dg}\|_2 \leq \lambda w_g \end{cases}.$$

$$C_d = \sum_{i=1}^N u_{id}^2, (\mathbf{B}_{dg})_{j \in g} = \sum_{i=1}^N u_{id} (x_{ij} - \sum_{k \neq d} u_{ik} a_{kj}).$$

Experimental Setup

Methods Compared

Traditional LSA

Sparse Coding (Code from Lee et. al. 07)

Latent Dirichlet allocation (LDA) (Code from Blei et. al. 03)

Sparse LSA

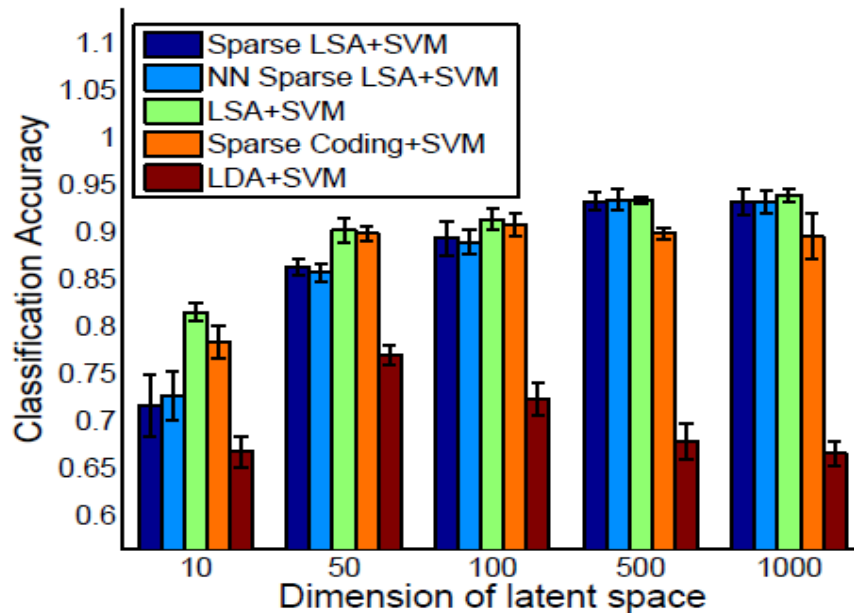
Nonnegative Sparse LSA (NN Sparse LSA)

Text Classification Data	N (No. of Documents)	M (Vocabulary Size)
20 news group (20NG) (alt.atheism vs talk.religion.misc)	1,425	17,390
RCV1 (20 classes)	15,564	7,413

Topic-Word Relationship Data	N (No. of Documents)	M (Vocabulary Size)
NIPS Proceedings from 98 to 99	1,714	13,649

Text Classification Performance

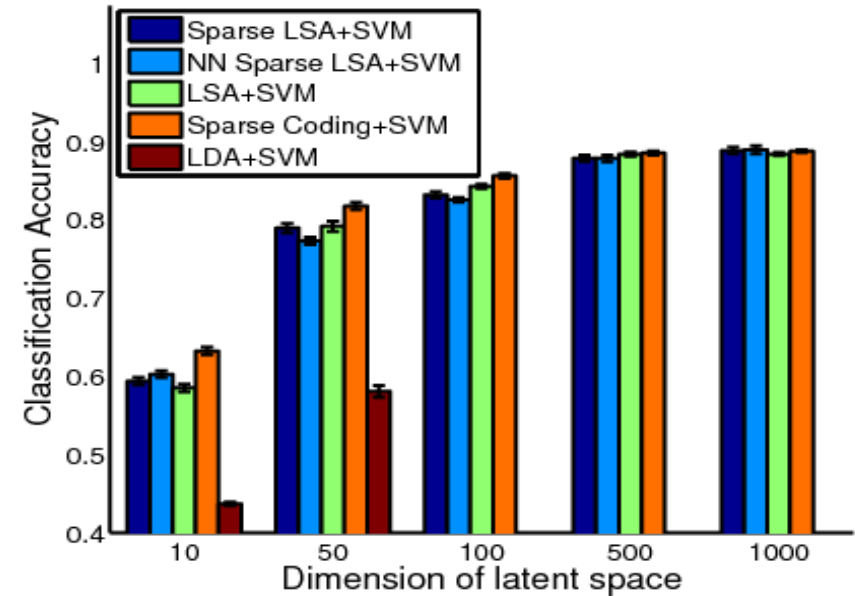
20NG



Dimension	10	50	100	500	1000
Sparse LSA	1.48	0.80	0.74	0.32	0.18
NN Sparse LSA	1.44	0.72	0.55	0.31	0.17
Other Methods	100	100	100	100	100

20NG: Density of \mathbf{A} (%) ($\lambda=0.05$)

RCV1



Dimension	10	50	100	500	1000
Sparse LSA	13.52	7.46	7.40	2.71	1.13
NN Sparse LSA	11.65	4.97	0.40	1.91	0.79
Other Methods	100	100	100	100	100

RCV1: Density of \mathbf{A} (%) ($\lambda=0.05$)

Conclusion: For large D , the classification performance of *Sparse LSA* is almost the same as *LSA* but with a much more *sparse* projection matrix \mathbf{A} .

Efficiency and Storage

20NG

	Proj. Time (ms)	Storage (MB)	Density of Proj. Doc. (%)
Sparse LSA	0.25 (4.05E-2)	0.6314	35.81 (15.39)
NN Sparse LSA	0.22 (2.78E-2)	0.6041	35.44 (15.17)
LSA	31.6 (1.10)	132.68	100 (0)
Sparse Coding	1711.1 (323.9)	132.68	86.94 (3.63)

RCV1

	Proj. Time (ms)	Storage (MB)	Density of Proj. Doc. (%)
Sparse LSA	0.59 (7.36E-2)	1.3374	55.38 (11.77)
NN Sparse LSA	0.46 (6.66E-2)	0.9537	46.47 (11.90)
LSA	13.2 (0.78)	113.17	100 (0)
Sparse Coding	370.5 (23.3)	113.17	83.88 (2.11)

Conclusion: Sparse LSA or NN Sparse LSA

- Efficient projection with less time
- Less storage for the projection matrix **A**
- Sparse projected documents: more efficient for subsequent retrieval tasks, e.g. ranking, text categorization, etc

$D = 1,000, \lambda = 0.05$ Table entry : mean (std)

Topic-Word Relationship

NIPS from 1988 to 1999

Nonnegative Sparse LSA

Topic 1	Topic 2	Topic 3	Topic 4
network neural networks system neurons neuron input output time systems	learning reinforcement algorithm function rule control learn weight action policy	network learning data neural training set function model input networks	model data models parameters mixture likelihood distribution gaussian em variables
Topic 5	Topic 6	Topic 7	
function functions approximation linear basis threshold theorem loss time systems	input output inputs chip analog circuit signal current action policy	image images recognition visual object system feature figure input networks	

LDA

Topic 1	Topic 2	Topic 3	Topic 4
learning data model training information number algorithm performance linear input	figure model output neurons vector networks state layer system order	algorithm method networks process learning input based function error parameter	single general sets time maximum paper rates features estimated neural
Topic 5	Topic 6	Topic 7	
rate unit data time estimation node set input neural properties	algorithms set problem weight temporal prior obtain parameter neural simulated	function neural hidden networks recognition output visual noise parameters references	

Conclusion: The topics learned by *NN Sparse LSA* are discriminative while the topics learned by *LDA* are all closely related to *neural network*.

Thank You !

