

Sparse Multinomial Logistic Regression via Approximate Message Passing



Evan Byrne and Philip Schniter

(Supported by NSF CCF-1018368 and CCF-1218754)

Problem Statement

- **Goal:** Infer the D -ary label y_0 from "test" feature vector $\mathbf{a}_0 \in \mathbb{R}^N$ given training $\{y_m, \mathbf{a}_m\}_{m=1}^M$.
- **Linear classification:** Estimate weight matrix $\widehat{\mathbf{X}} \in \mathbb{R}^{N \times D}$, then predict $\hat{y}_0 = \arg \max_d [\widehat{\mathbf{X}}^T \mathbf{a}_0]_d$.
- **Feature selection:** Determine which subset of N features is needed to accurately predict the label y_0 .
- We're especially interested in the case $M \ll N$ (MVPA, text-mining, micro-array gene expression).
 - Possible if "true" \mathbf{X} is K -row-sparse with $K \ll M$.

Multinomial Logistic Regression

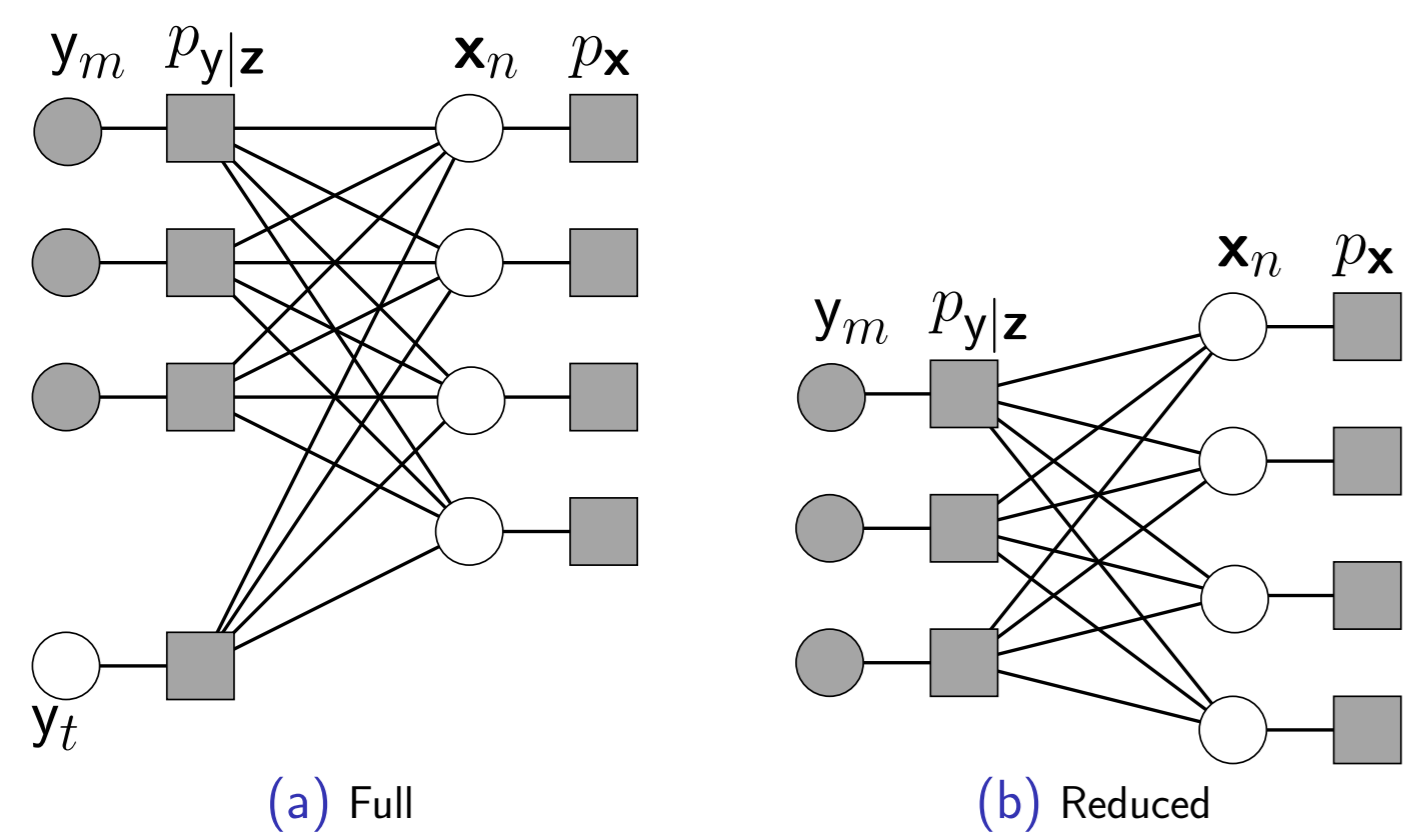
- One approach to designing \mathbf{X} is Multinomial Logistic Regression (MLR).
- In MLR, we use the multinomial logistic likelihood:

$$p_{y|z}(y_m | \mathbf{z}_m) = \frac{\exp(\mathbf{z}_m^T \mathbf{y}_m)}{\sum_{d=1}^D \exp(\mathbf{z}_m^T \mathbf{e}_d)}, \quad y_m \in \{1, \dots, D\} \quad \text{where} \quad \mathbf{z}_m \triangleq \mathbf{X}^T \mathbf{a}_m. \quad (1)$$

- Also, \mathbf{X} is regularized through some prior $p_{\mathbf{X}}(\mathbf{X})$.
- Existing approaches to sparse MLR include **SMLR** [Krishnapuram Carin Figueiredo Hartemink 05], **SBMLR** [Cawley Talbot Girolami 07], and **GLMNET** [Friedman Hastie Tibshirani 10], which all employ a Laplacian prior for $p_{\mathbf{X}}$ and MAP estimation to find $\widehat{\mathbf{X}}$.

HyGAMP for MLR

- Assuming a separable likelihood $p_{\mathbf{y}|\mathbf{z}}(\mathbf{y} | \mathbf{Z}) = \prod_m p_{y|z}(y_m | \mathbf{z}_m)$ and prior $p_{\mathbf{X}}(\mathbf{X}) = \prod_n p_{\mathbf{x}}(\mathbf{x}_n)$, $p_{\mathbf{y},\mathbf{X}}(\mathbf{y}, \mathbf{X})$ can be represented by the following factor graph:



- Through message passing, we break one large inference problem into many smaller inference problems.
- Under large i.i.d. \mathbf{A} and scalar z_m & \mathbf{x}_n , we can apply generalized approximate message passing (GAMP) [Rangan 11]. Has been used for binary logistic regression [Ziniel Schniter Sederberg 15].
- However, our z_m & \mathbf{x}_n are vector valued, so we instead apply hybrid GAMP (HyGAMP) [Rangan Fletcher Goyal Schniter 12].
 - MSA variant: computes MAP estimate of \mathbf{X} .
 - SPA variant: computes approximate marginal posteriors of y_0 and \mathbf{X} \Rightarrow approximately minimizes test-error rate!
 - Passes $O(M+N)$ messages in the form of D -dimensional Gaussian pdfs.

Algorithm Summary

- HyGAMP iteratively passes messages back and forth between the $p_{y|z}$ and $p_{\mathbf{x}}$ nodes until convergence.
- The algorithm can be divided into "linear" and "non-linear" steps.

Linear steps:

- Involve $N+M$ matrix inversions of size $D \times D$.
- Identical for SPA and MSA variants of HyGAMP.

Non-linear steps:

- At each node n and m , HyGAMP approximates the posterior distributions as:

$$p_{\mathbf{x}|r}(\mathbf{x}_n | \hat{\mathbf{r}}_n; \mathbf{Q}_n^r) \propto p_{\mathbf{x}}(\mathbf{x}_n) \mathcal{N}(\mathbf{x}_n; \hat{\mathbf{r}}_n, \mathbf{Q}_n^r) \quad (2)$$

$$p_{z|y,p}(\mathbf{z}_m | y_m, \hat{\mathbf{p}}_m; \mathbf{Q}_m^p) \propto p_{y|z}(y_m | \mathbf{z}_m) \mathcal{N}(\mathbf{z}_m; \hat{\mathbf{p}}_m, \mathbf{Q}_m^p), \quad (3)$$

for $\hat{\mathbf{p}}_m, \mathbf{Q}_m^p, \hat{\mathbf{r}}_n, \mathbf{Q}_n^r$ calculated in the linear steps.

- SPA variant: computes the means ($\hat{\mathbf{x}}_n$ and $\hat{\mathbf{z}}_m$) and covariances (\mathbf{Q}_n^x and \mathbf{Q}_m^z) of above posteriors.
- MSA variant: computes the modes ($\hat{\mathbf{x}}_n$ and $\hat{\mathbf{z}}_m$) and inverse Hessians of above log posteriors.
- For MLR likelihood and most sparsity-inducing priors, there are no closed-form solutions. Need approximations like numerical integration, importance sampling, Newtons method, minorize-maximization.
- Typically, to enforce sparsity, we use a Bernoulli-Gaussian prior in SPA-HyGAMP and a Laplacian prior in MSA-HyGAMP.

Simplified HyGAMP (SHyGAMP) for MLR

- Unfortunately, HyGAMP is not computationally competitive due to
 - expensive linear steps (e.g., matrix inversion)
 - expensive non-linear steps (e.g., iterative algorithms)
 - numerical instabilities
- **Our Solution:** Assume all matrices \mathbf{Q} are diagonal.
 - trivializes the linear steps (i.e., no matrix inversion)
 - drastically simplifies the non-linear steps
 - enables use of existing GAMPmatlab software framework [Rangan, Schniter, Parker, Ziniel, et al.]

SPA SHyGAMP

Non-linear z_m steps:

- Simplified posterior mean/variance computations:

$$\hat{z}_{md} = C_m^{-1} \int_{\mathbb{R}^D} z_d p_{y|z}(y_m | \mathbf{z}) \prod_{k=1}^D \mathcal{N}(z_k; \hat{p}_{mk}, q_{mk}^p) dz \quad (4)$$

$$q_{md}^z = C_m^{-1} \int_{\mathbb{R}^D} z_d^2 p_{y|z}(y_m | \mathbf{z}) \prod_{k=1}^D \mathcal{N}(z_k; \hat{p}_{mk}, q_{mk}^p) dz - \hat{z}_{md}^2 \quad (5)$$

- Investigated approaches based on numerical integration, importance sampling, Taylor series approximation
- Proposed novel Gaussian-mixture approximation with improved accuracy-runtime tradeoff

Non-linear x_n steps:

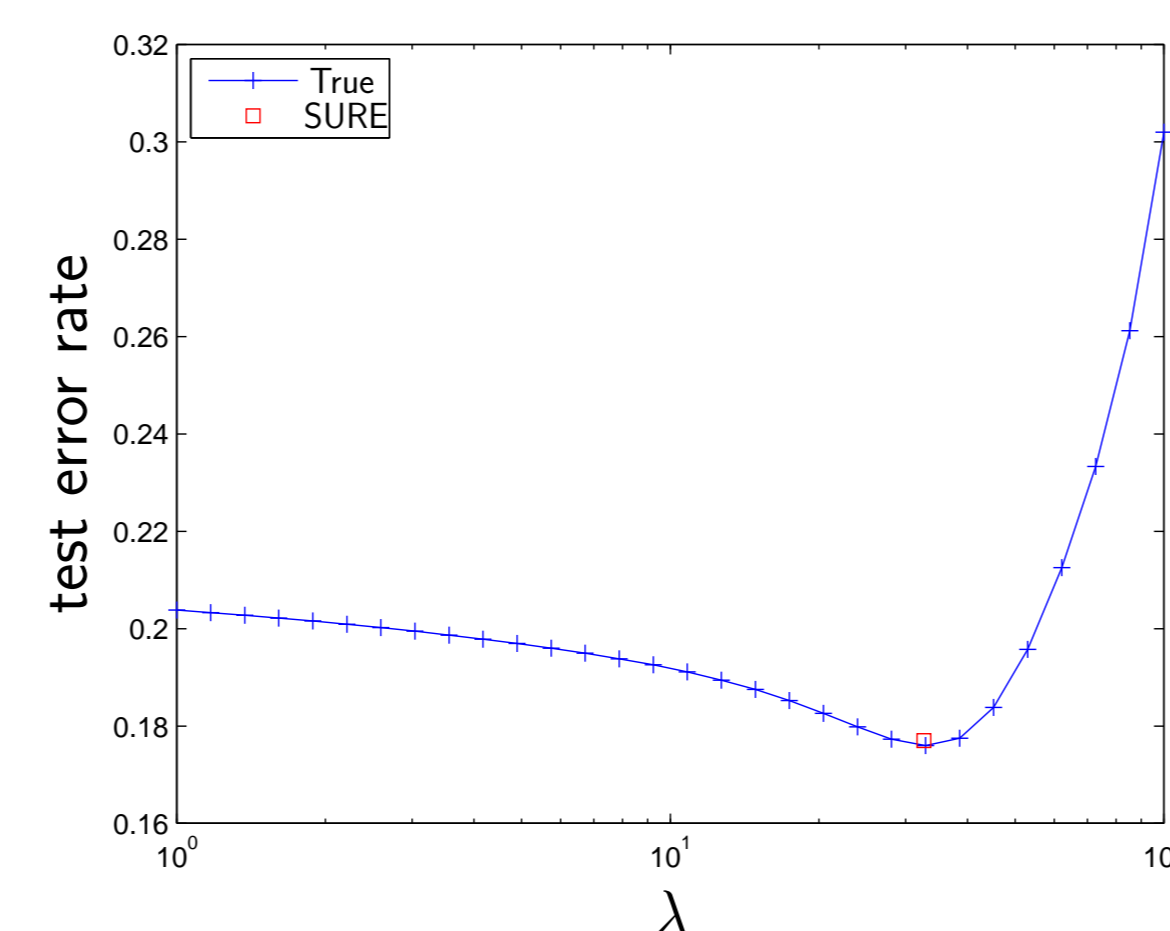
- Choosing separable prior allows further decoupling into D scalar inference problems
- Example: i.i.d. BG: $p_{\mathbf{x}}(x_{nd}) \triangleq \beta \mathcal{N}(x_{nd}; 0, \sigma_x^2) + (1 - \beta) \delta(x_{nd}) \quad \forall n, d$
- Parameters σ_x^2 and β can be tuned online via EM [Vila, Schniter 13].

MSA SHyGAMP

- Non-linear z_m steps: solved via component-wise Newton's method.
- Non-linear x_n steps: choose ℓ_1 regularization, solve via soft-thresholding.
- λ tuned online via variation on SURE procedure [Mousavi Maleki Baraniuk 13].

SURE tuning procedure

- Idea: at each GAMP iteration, choose λ to minimize the SURE of the thresholder.
- Challenge: objective function is highly non-smooth.
- Our Solution: approximate empirical data by GM distribution \Rightarrow smooth objective function. Minimize using conventional techniques (e.g., gradient descent, bisection search).



- $N = 30000$
- $M = 300$
- $K = 25$
- $D = 4$
- Bayes Error Rate = .1
- average of 12 trials

Test error rate vs λ for fixed- λ MSA SHyGAMP, with final error and λ of SURE-tuned MSA SHyGAMP superimposed

Selected References

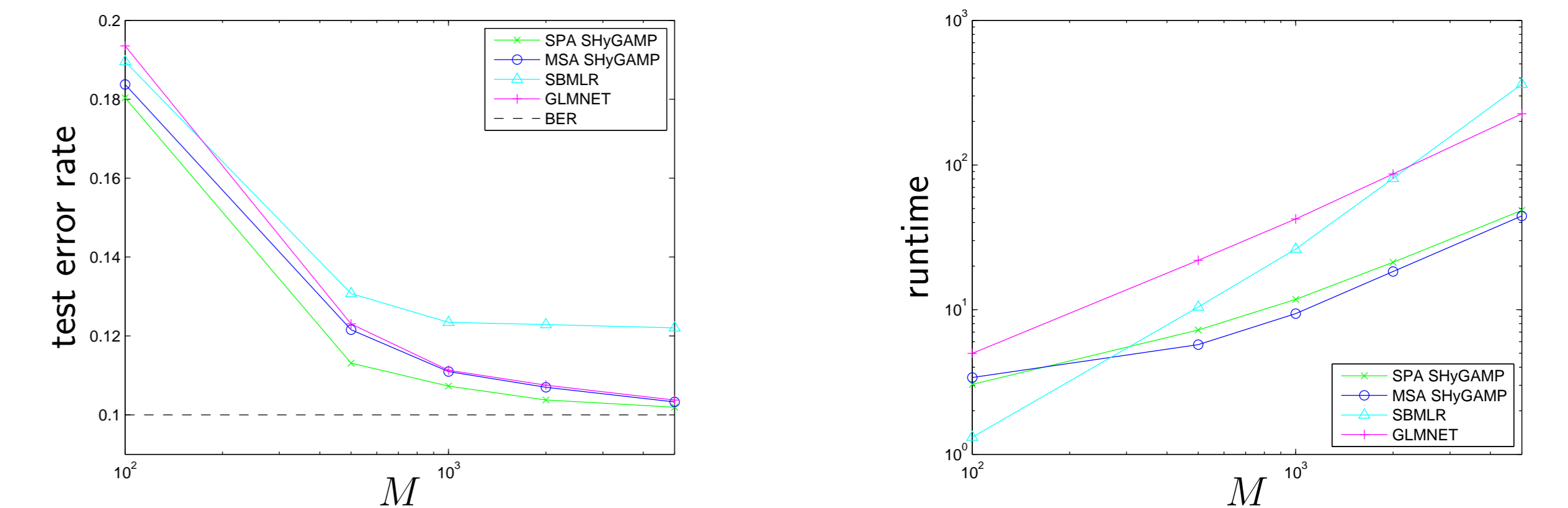
- E. M. Byrne, "Sparse multinomial logistic regression via approximate message passing," Master's thesis, The Ohio State University, July 2015.
- B. Krishnapuram, L. Carin, M. Figueiredo, and A. Hartemink, "Sparse multinomial logistic regression: Fast algorithms and generalization bounds," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, pp. 957-968, June 2005.
- G. C. Cawley, N. L. C. Talbot, and M. Girolami, "Sparse multinomial logistic regression via Bayesian L1 regularisation," in *Proc. Neural Inform. Process. Syst. Conf.*, pp. 209-216, 2007.
- J. Friedman, T. Hastie, and R. Tibshirani, "Regularization paths for generalized linear models via coordinate descent," *J. Statist. Softw.*, vol. 33, pp. 1-22, Jan. 2010.
- S. Rangan, "Generalized approximate message passing for estimation with random linear mixing," in *Proc. IEEE Int. Symp. Inform. Thy.* (Saint Petersburg, Russia), pp. 2168-2172, Aug. 2011. (full version at arXiv:1010.5141)
- J. Ziniel, P. Schniter, and P. Sederberg, "Binary classification and feature selection via generalized approximate message passing," *IEEE Trans. Signal Process.*, vol. 63, no. 8, pp. 2020-2032, 2015.
- S. Rangan, A. K. Fletcher, V. K. Goyal, and P. Schniter, "Hybrid generalized approximate message passing with applications to structured sparsity," in *Proc. IEEE Int. Symp. Inform. Thy.*, pp. 1236-1240, July 2012. (full version at arXiv:1111.2581).
- S. Rangan, P. Schniter, J. T. Parker, J. Ziniel, J. Vila, M. Borgerdig, and et al. GAMPmatlab. <https://sourceforge.net/projects/gampmatlab/>.
- J. P. Vila and P. Schniter, "Expectation-maximization Gaussian-mixture approximate message passing," *IEEE Trans. Signal Process.*, vol. 61, pp. 4658-4672, Oct. 2013.
- A. Mousavi, A. Maleki, and R. G. Baraniuk, "Parameterless optimal approximate message passing," Oct. 2013. (full version at arXiv:1311.0035)

Classification Performance on Synthetic Data

Data generation model:

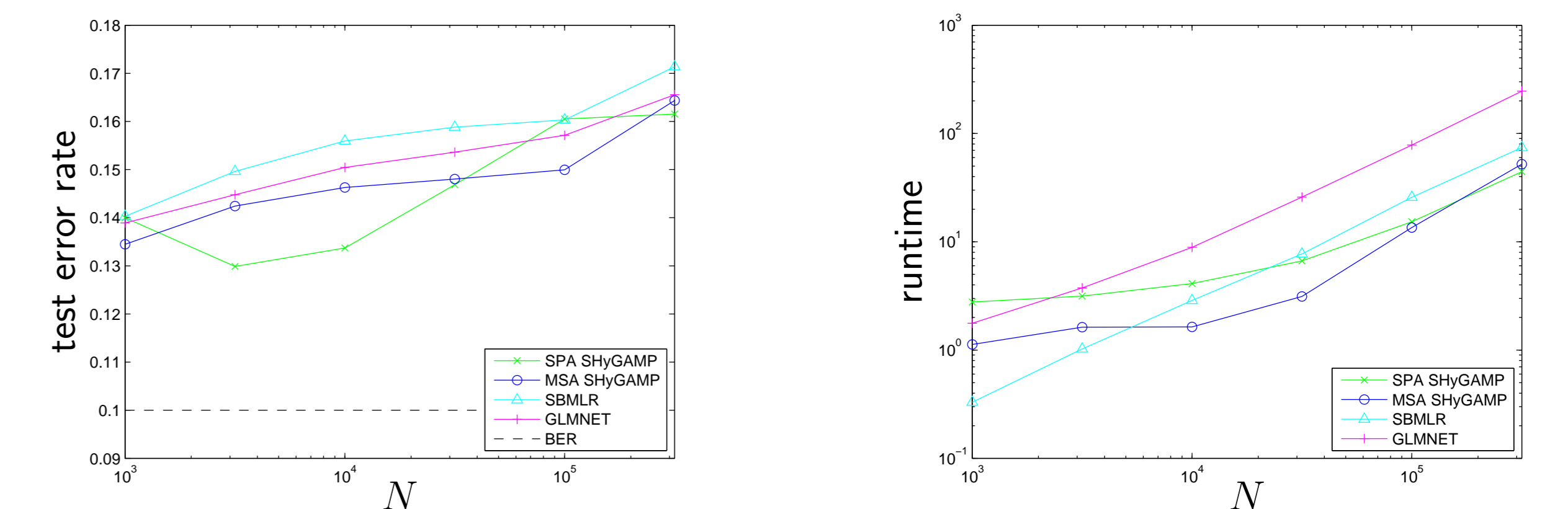
- features $\mathbf{a}_m | (y_m = d) \sim \mathcal{N}(\boldsymbol{\mu}_d, \sigma_a^2 \mathbf{I}_N)$
- feature means $\{\boldsymbol{\mu}_d\}_{d=1}^D$ orthonormal with K non-zero entries
- balanced training labels

Average classification error and runtime vs M for fixed $D = 4, N = 10000, K = 10, 12$ trials:



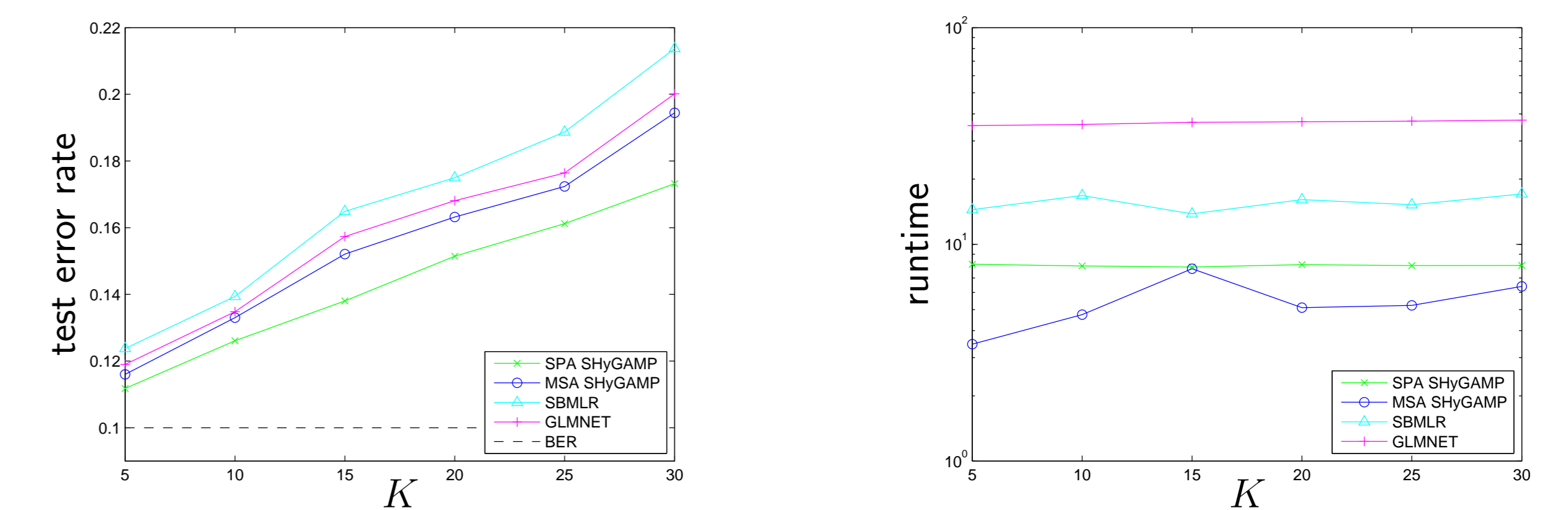
SPA-SHyGAMP wins in error. MSA-SHyGAMP beats SBMLR and GLMNET in both error and runtime (for large M).

Average classification error and runtime vs N for fixed $D = 4, M = 200, K = 10, 12$ trials:



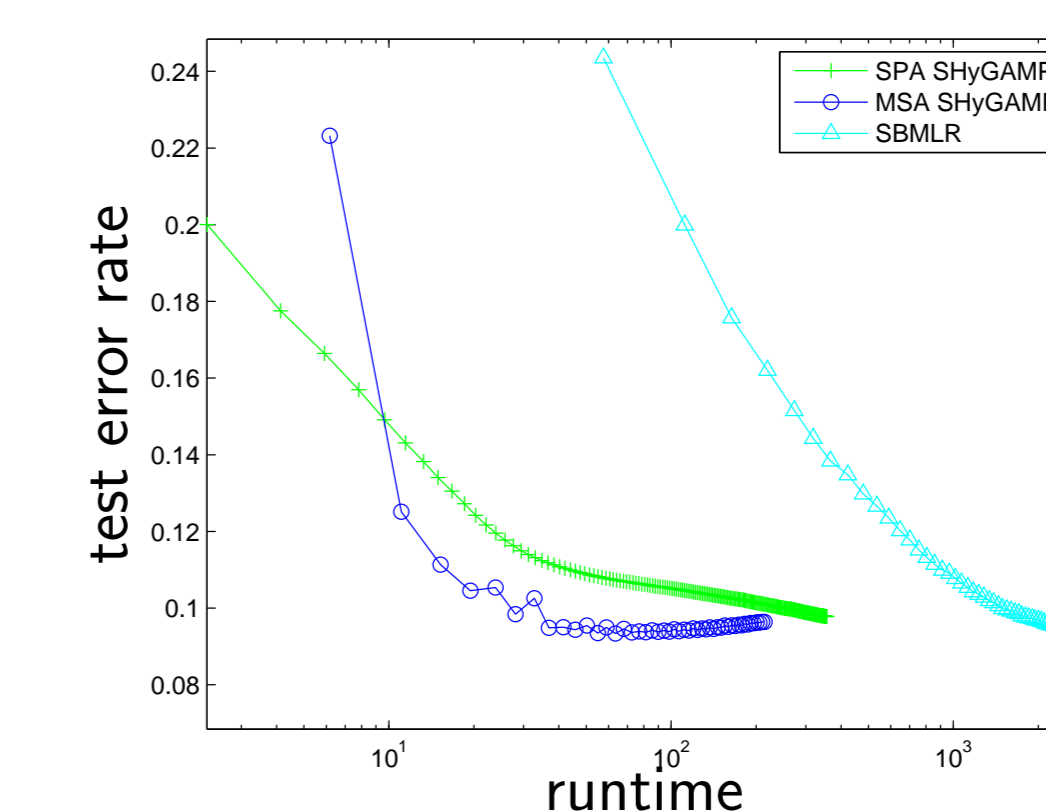
MSA-SHyGAMP beats SBMLR and GLMNET in both error and runtime (for large N).

Average classification error and runtime vs K for fixed $D = 4, M = 300, N = 30000, 12$ trials:



SPA-SHyGAMP wins overall in error. MSA-SHyGAMP beats SBMLR and GLMNET in both error and runtime.

Classification Performance on RCV1 Dataset



Both MSA and SPA SHyGAMP converge to the final error rate faster than SBMLR.

- Features are word frequency; labels are document subject (e.g., business).
- $D = 25, M_{\text{train}} = 14147, N = 47236,$ and $M_{\text{test}} = 469571$.
- Shown is test error rate vs training time for auto-tuned algorithms.