

# Sparse Bayesian Logistic Regression with Hierarchical Prior and Variational Inference

## Bayesian Sparse Classifiers

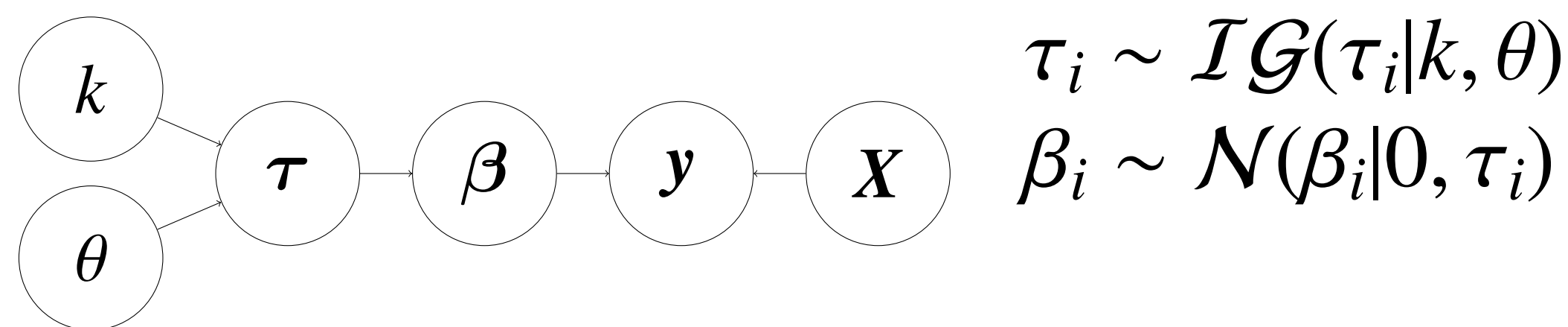
Logistic regression model:

$$p(y_j|\beta) = \sigma(\mathbf{x}_j^T \beta)^{y_j} (1 - \sigma(\mathbf{x}_j^T \beta))^{1-y_j},$$

$$\sigma(x) = 1 / (1 + e^{-x}).$$

Existing Bayesian sparse classifiers:

- Relevance Vector Machine [Tipping, 2001]



- Learning: ML or MAP estimate for  $\gamma$  and Laplace approximation
- Sparse Representation Prior [Serra et al., 2016]



- Learning: Majorize minimization and variational inference

## Proposed Model

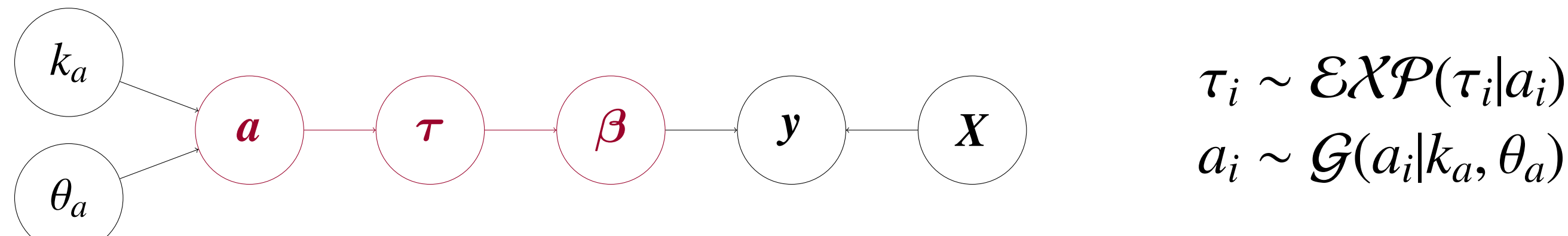
- Hierarchical prior

$$\beta_i \sim \mathcal{N}(\beta_i | 0, \tau_i), \quad \tau_i \sim \mathcal{GIG}(\tau_i | a_i, b_i, \rho)$$

$$(a_i, b_i) \sim \text{Conjugate prior for } p(\tau_i | a_i, b_i, \rho)$$

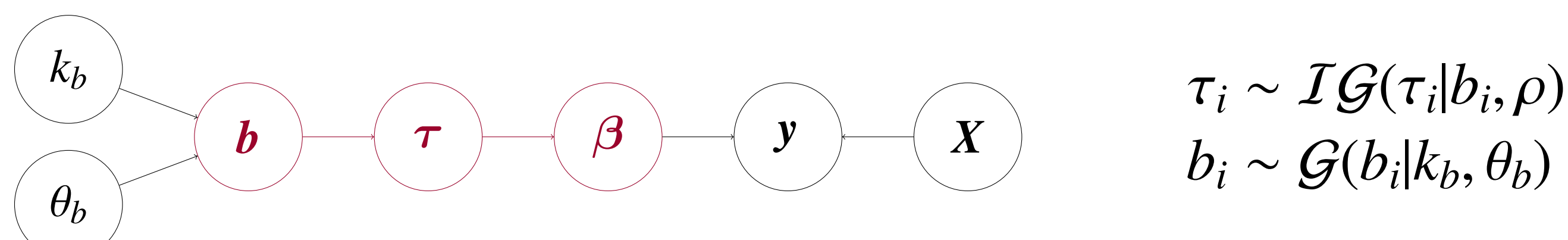
$\mathcal{GIG}$  is a **generalized inverse Gaussian**, which is **exponential** when  $b_i = \rho = 0$  and **inverse gamma** when  $a_i = 0, \rho < 0$ .

- Exponential mixing:



The marginal distribution  $p(\beta_i)$  is a **Laplace** distribution.

- Inverse gamma mixing:



The marginal distribution  $p(\beta_i)$  is a **Student's t** distribution.

## Learning Algorithm

- **Variational inference** with **Mean-field approximation** + **Majorize Minimization**

- Extension of [Jaakkola and Jordan, 1997]

$h(\beta, \xi)$ : a lower bound of  $p(y|\beta)$  and a quadratic function of  $\beta$

### Algorithm

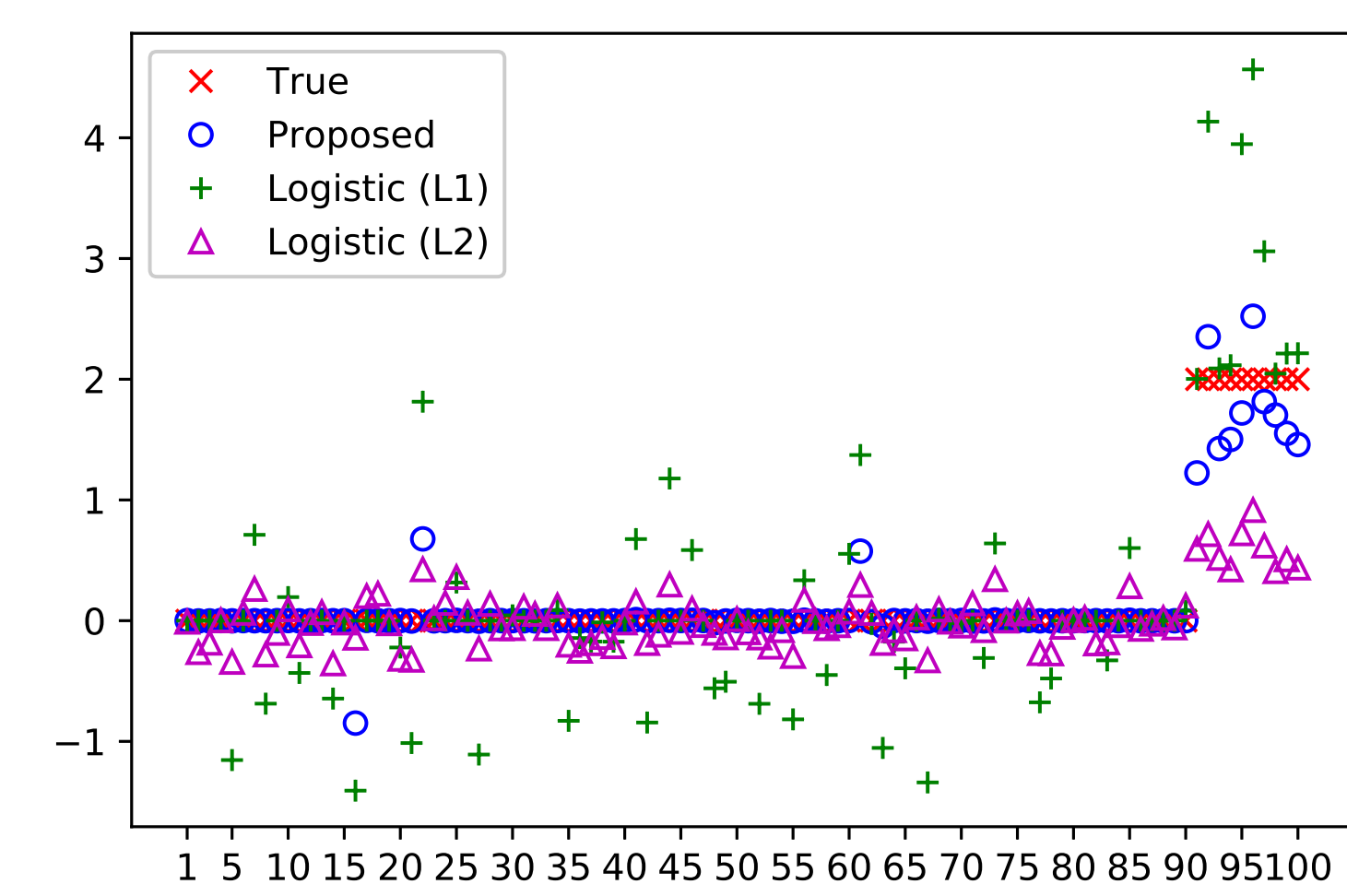
1.  $q^*(\beta) \propto \exp(\mathbb{E}_{q(\tau)} [\ln h(\beta, \xi) p(\beta|\tau)])$
2.  $q^*(\tau) \propto \exp(\mathbb{E}_{q(\beta)q(a,b)} [\ln p(\beta|\tau) p(\tau|a,b)])$
3.  $q^*(a,b) \propto \exp(\mathbb{E}_{q(\tau)} [\ln p(\tau|a,b) p(a,b)])$
4.  $\xi^* = \arg\max_{\xi} \mathbb{E}_{q(\beta)} [\ln h(\beta, \xi)]$

## Experiments on synthetic data

- True parameter  $\beta^* = (\mathbf{0}_{90}, \mathbf{2}_{10})$
- Parameter setting for the proposed algorithm:
  - Exponential mixing ( $b_i = \rho = 0$ )
  - $k_a = \theta_a = 10^{-6}$  (very flat prior)
- Compare with  $L_1$  regularized logistic regression and  $L_2$  regularized logistic regression

**Table:** MSE and prediction accuracy for synthetic data.

	MSE	Accuracy
Proposed	<b>0.1589 ± 0.1133</b>	<b>0.8195 ± 0.0477</b>
$L_1$	0.3974 ± 0.2939	0.7750 ± 0.0456
$L_2$	0.4391 ± 0.2597	0.7112 ± 0.0242



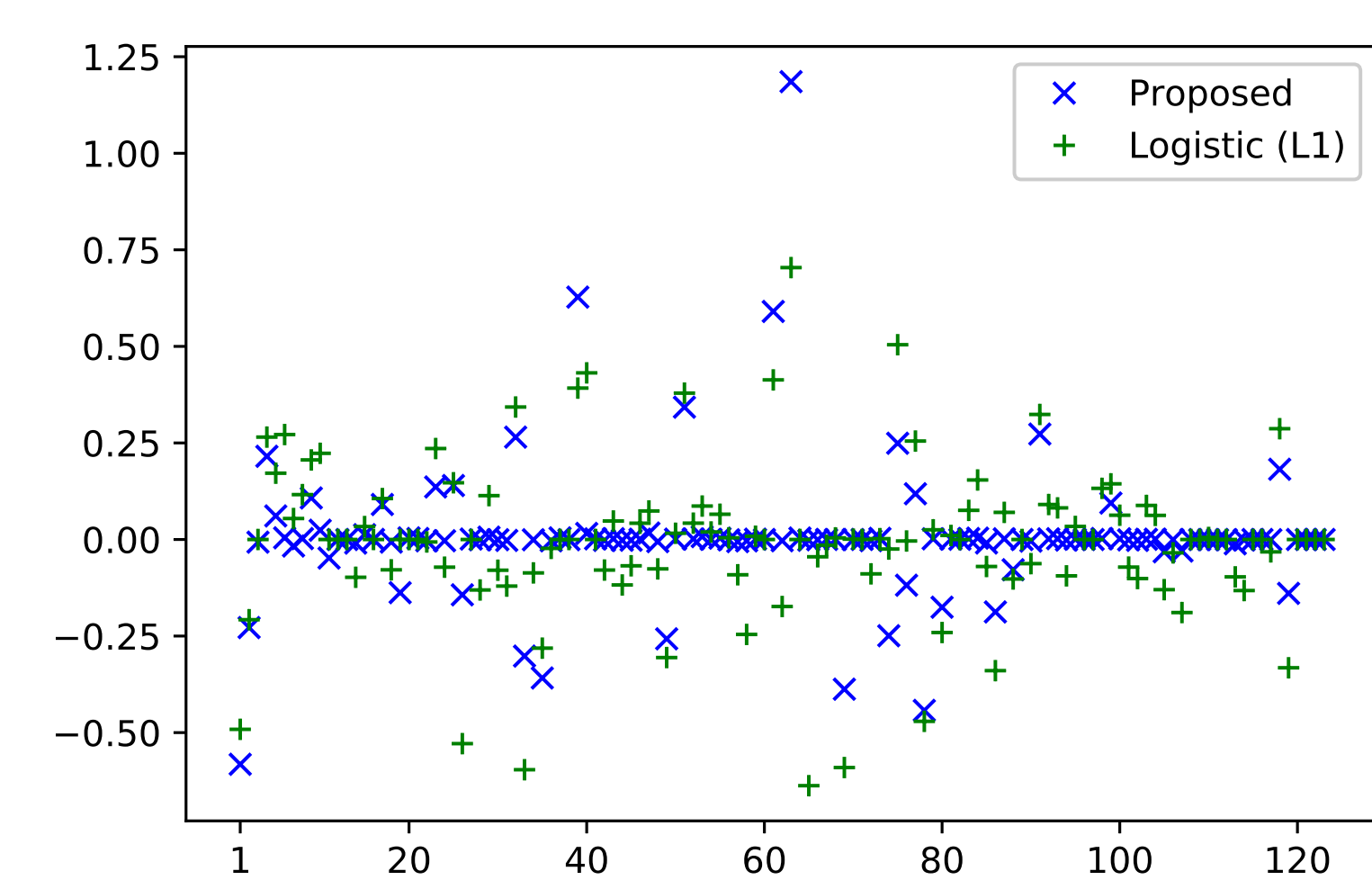
**Figure:** Estimation results on synthetic data.

## Experiments on real world data

- Parameter setting is the same with the case of synthetic data

**Table:** Prediction accuracy for real world data.

	a1a	w1a	covtype
Proposed	<b>0.8400</b>	0.9757	<b>0.7436</b>
$L_1$	<b>0.8400</b>	0.9702	0.7412
$L_2$	0.8386	<b>0.9779</b>	0.7398



**Figure:** Estimation results on 'a1a' data.

## References

- M. E. Tipping. **Sparse bayesian learning and the relevance vector machine**. *Journal of machine learning research*, pp. 211-244, 2001.
- J. G. Serra, P. Ruiz, R. Molina, and A. K. Katsaggelos. **Bayesian logistic regression with sparse general representation prior for multispectral image classification**. In *Image Processing (ICIP), 2016 IEEE International Conference on*, pp. 1893-1897, 2016.
- T. Jaakkola and M. Jordan. **A variational approach to bayesian logistic regression models and their extensions**. In *Sixth International Workshop on Artificial Intelligence and Statistics*, Vol. 82, page 4, 1997.

