

Statistical Inference for High-dimensional Logistic Regression

Zijian Guo

Department of Statistics, Rutgers University

CMStatistics 2021, London, UK

Joint with Prabrisha Rakshit, Daniel Herman and Jinbo Chen.

Motivation

Electronic health record (EHR)

EHR: Document for medical billing

1. **Limited** labeled data: 50 to 1000
2. **Binary** outcome (e.g. disease status.)
3. **Many** predictors: billing codes, demographics, disease histories, co-morbid conditions, laboratory test results, prescription codes, and concepts extracted from doctors' notes.

EHR phenotyping

High-dimensional inference for binary outcome labeling.

High-dimensional Logistic Regression

For $1 \leq i \leq n$, consider the model for $y_i \in \{0, 1\}, X_i. \in \mathbb{R}^p$

$$\mathbb{P}(y_i = 1 | X_i.) = h(X_i.^T \beta), \quad h(z) = \exp(z) / [1 + \exp(z)]$$

- ▶ $p \gg n$, β is sparse
- ▶ Case probability

$$\mathbb{P}(y_i = 1 | X_i. = x_{\text{new}}) \equiv h(x_{\text{new}}^T \beta)$$

EHR phenotyping

$$H_0 : h(x_{\text{new}}^T \beta) < 1/2.$$

The penalized log-likelihood estimator $\hat{\beta}$

$$\hat{\beta} = \arg \min_{\beta} \ell(\beta) + \lambda \|\beta\|_1,$$

where $\ell(\beta) = \sum_{i=1}^n [\log(1 + \exp(\mathbf{X}_i^T \beta)) - y_i \cdot (\mathbf{X}_i^T \beta)]$.

Debiasing Inference

- ▶ β_j in linear models (Zhang & Zhang '14, Javanmard & Montanari '14)
- ▶ β_j in GLM (van de Geer, Bühlmann, Ritov & Dezeure '14)

Bias Correction Intuition

$$\hat{\beta}_j + \hat{u}^\top \frac{1}{n} \sum_{i=1}^n X_i \cdot (y_i - h(X_i^\top \hat{\beta})) \quad (1)$$

$$\frac{1}{n} \sum_{i=1}^n X_i \cdot (y_i - h(X_i^\top \hat{\beta})) \approx \hat{H}(\hat{\beta})(\beta - \hat{\beta}) + \frac{1}{n} \sum_{i=1}^n \epsilon_i X_i,$$

with $\hat{H}(\beta) = \frac{1}{n} \sum_{i=1}^n h(X_i^\top \beta)(1 - h(X_i^\top \beta))X_i X_i^\top$.

Bias Correction Intuition

$$\hat{\beta}_j + \hat{u}^\top \frac{1}{n} \sum_{i=1}^n X_i \cdot (y_i - h(X_i^\top \hat{\beta})) \quad (1)$$

$$\frac{1}{n} \sum_{i=1}^n X_i \cdot (y_i - h(X_i^\top \hat{\beta})) \approx \hat{H}(\hat{\beta})(\beta - \hat{\beta}) + \frac{1}{n} \sum_{i=1}^n \epsilon_i X_i,$$

with $\hat{H}(\beta) = \frac{1}{n} \sum_{i=1}^n h(X_i^\top \beta)(1 - h(X_i^\top \beta))X_i X_i^\top$.

$$\begin{aligned} \hat{u}^\top \frac{1}{n} \sum_{i=1}^n X_i \cdot (y_i - h(X_i^\top \hat{\beta})) &\approx \hat{u}^\top \hat{H}(\hat{\beta})(\beta - \hat{\beta}) \\ &\approx \mathbf{e}_j^\top (\beta - \hat{\beta}). \end{aligned}$$

Challenge

For $\beta_j = \mathbf{e}_j^\top \beta$,

$$\hat{H}(\hat{\beta})\hat{u} \approx \mathbf{e}_j$$

- ▶ Sparse $[\mathbb{E}\hat{H}(\beta)]^{-1}\mathbf{e}_j$ (van de Geer et.al., 14)

For $\mathbf{x}_{\text{new}}^\top \beta$, we construct \hat{u} such that

$$\hat{H}(\hat{\beta})\hat{u} \approx \mathbf{x}_{\text{new}}.$$

Challenge

$[\hat{H}(\hat{\beta})]^{-1}\mathbf{x}_{\text{new}}$ can be **DENSE!**

Our Proposed Method

Existing

$$\hat{\beta}_j + \hat{u}^T \frac{1}{n} \sum_{i=1}^n X_{i.} (y_i - h(X_{i.}^T \hat{\beta}))$$

Linearization and **V**ariance **E**nhancement

$$\widehat{x_{\text{new}}^T \beta} = x_{\text{new}}^T \hat{\beta} + \hat{u}^T \frac{1}{n} \sum_{i=1}^n \underbrace{[h(X_{i.}^T \hat{\beta})(1 - h(X_{i.}^T \hat{\beta}))]^{-1}}_{\text{weight for } i\text{-th observation}} X_{i.} (y_i - h(X_{i.}^T \hat{\beta})).$$

Linearization: Logistic to Linear

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^T \hat{\beta})(1 - h(\mathbf{X}_i^T \hat{\beta}))]^{-1} \mathbf{X}_i (y_i - h(\mathbf{X}_i^T \hat{\beta})) \\ & \approx \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T (\beta - \hat{\beta}) + \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^T \hat{\beta})(1 - h(\mathbf{X}_i^T \hat{\beta}))]^{-1} \epsilon_i \mathbf{X}_i. \end{aligned}$$

Linearization: Logistic to Linear

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^T \hat{\beta})(1 - h(\mathbf{X}_i^T \hat{\beta}))]^{-1} \mathbf{X}_i (y_i - h(\mathbf{X}_i^T \hat{\beta})) \\ & \approx \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T (\beta - \hat{\beta}) + \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^T \hat{\beta})(1 - h(\mathbf{X}_i^T \hat{\beta}))]^{-1} \epsilon_i \mathbf{X}_i. \end{aligned}$$

$\widehat{\mathbf{x}}_{\text{new}}^T \beta - \mathbf{x}_{\text{new}}^T \beta$ is decomposed as

$$(\widehat{\Sigma} \hat{\mathbf{u}} - \mathbf{x}_{\text{new}})^T (\beta - \hat{\beta}) + \hat{\mathbf{u}}^T \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^T \hat{\beta})(1 - h(\mathbf{X}_i^T \hat{\beta}))]^{-1} \epsilon_i \mathbf{X}_i.$$

Variance Enhancement: Uniform for \mathbf{x}_{new}

$$(\widehat{\Sigma}\widehat{\mathbf{u}} - \mathbf{x}_{\text{new}})^\top(\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}) + \widehat{\mathbf{u}}^\top \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^\top \widehat{\boldsymbol{\beta}})(1 - h(\mathbf{X}_i^\top \widehat{\boldsymbol{\beta}}))]^{-1} \epsilon_i \mathbf{X}_i.$$

$$\widehat{\Sigma}\widehat{\mathbf{u}} \approx \mathbf{x}_{\text{new}}$$

Variance Enhancement: Uniform for \mathbf{x}_{new}

$$(\widehat{\Sigma}\widehat{\mathbf{u}} - \mathbf{x}_{\text{new}})^\top(\beta - \widehat{\beta}) + \widehat{\mathbf{u}}^\top \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^\top \widehat{\beta})(1 - h(\mathbf{X}_i^\top \widehat{\beta}))]^{-1} \epsilon_i \mathbf{X}_i.$$

$$\widehat{\Sigma}\widehat{\mathbf{u}} \approx \mathbf{x}_{\text{new}}$$

Variance enhancement projection direction.

$$\widehat{\mathbf{u}} = \arg \min_{\mathbf{u} \in \mathbb{R}^p} \mathbf{u}^\top \widehat{\Sigma} \mathbf{u}$$

$$\text{subject to } \|\widehat{\Sigma}\mathbf{u} - \mathbf{x}_{\text{new}}\|_\infty \leq \|\mathbf{x}_{\text{new}}\|_2 \lambda_n$$

$$|\mathbf{x}_{\text{new}}^\top \widehat{\Sigma}\mathbf{u} - \|\mathbf{x}_{\text{new}}\|_2^2| \leq \|\mathbf{x}_{\text{new}}\|_2^2 \lambda_n$$

where $\lambda_n \asymp (\log p/n)^{1/2}$.

What if no additional constraint?

$$\hat{u} = \arg \min_{u \in \mathbb{R}^p} u^T \hat{\Sigma} u$$

subject to $\|\hat{\Sigma} u - x_{\text{new}}\|_{\infty} \leq \|x_{\text{new}}\|_2 \lambda_n$

What if no additional constraint?

$$\hat{u} = \arg \min_{u \in \mathbb{R}^p} u^T \hat{\Sigma} u$$

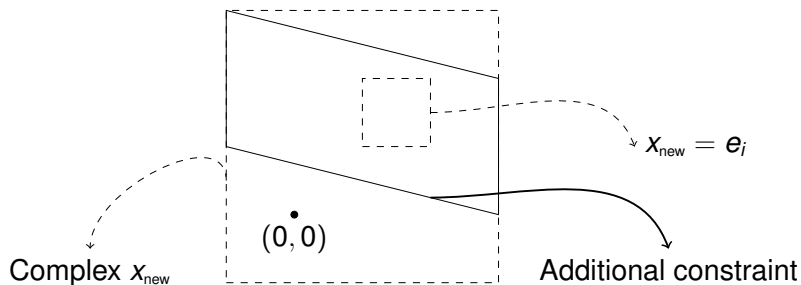
$$\text{subject to } \|\hat{\Sigma} u - x_{\text{new}}\|_{\infty} \leq \|x_{\text{new}}\|_2 \lambda_n$$

For a **dense** x_{new} : there is no bias correction,

$$\|x_{\text{new}}\|_{\infty} \leq \|x_{\text{new}}\|_2 \lambda_n \Rightarrow \hat{u} = 0!$$

Curse of dimensionality: too much flexibility of searching the direction.

Additional Constraint and Feasible Set



- ▶ Large dashed: dense x_{new} **without** additional constraint.
- ▶ Solid parallelogram: dense x_{new} **with** additional constraint.

$$\left| x_{\text{new}}^T \hat{\Sigma} u - \|x_{\text{new}}\|_2^2 \right| \leq \|x_{\text{new}}\|_2^2 \lambda$$

$$\widehat{\mathbf{x}}_{\text{new}}^{\top} \widehat{\beta} = \mathbf{x}_{\text{new}}^{\top} \widehat{\beta} + \widehat{\mathbf{u}}^{\top} \frac{1}{n} \sum_{i=1}^n \underbrace{[h(\mathbf{X}_i^{\top} \widehat{\beta})(1 - h(\mathbf{X}_i^{\top} \widehat{\beta}))]^{-1}}_{\text{weight for } i\text{-th observation}} \mathbf{X}_i (y_i - h(\mathbf{X}_i^{\top} \widehat{\beta})).$$

with the projection direction $\widehat{\mathbf{u}}$ defined as

$$\widehat{\mathbf{u}} = \arg \min_{\mathbf{u} \in \mathbb{R}^p} \mathbf{u}^{\top} \widehat{\Sigma} \mathbf{u}$$

$$\text{subject to } \|\widehat{\Sigma} \mathbf{u} - \mathbf{x}_{\text{new}}\|_{\infty} \leq \|\mathbf{x}_{\text{new}}\|_2 \lambda_n$$

$$|\mathbf{x}_{\text{new}}^{\top} \widehat{\Sigma} \mathbf{u} - \|\mathbf{x}_{\text{new}}\|_2^2| \leq \|\mathbf{x}_{\text{new}}\|_2^2 \lambda_n$$

where $\lambda_n \asymp (\log p/n)^{1/2}$.

Statistical Inference

We construct the CI for $\mathbb{P}(y_i = 1 | \mathbf{X}_i = \mathbf{x}_{\text{new}})$ as,

$$\text{CI}_\alpha(\mathbf{x}_{\text{new}}) = \left[h\left(\widehat{\mathbf{x}}_{\text{new}}^\top \widehat{\beta} - z_{\alpha/2} \widehat{\mathbf{V}}^{1/2}\right), h\left(\widehat{\mathbf{x}}_{\text{new}}^\top \widehat{\beta} + z_{\alpha/2} \widehat{\mathbf{V}}^{1/2}\right) \right],$$

with

$$\widehat{\mathbf{V}} = \widehat{\mathbf{u}}^\top \left[\frac{1}{n^2} \sum_{i=1}^n [h(\mathbf{X}_i^\top \widehat{\beta})(1 - h(\mathbf{X}_i^\top \widehat{\beta}))]^{-1} \mathbf{X}_i \mathbf{X}_i^\top \right] \widehat{\mathbf{u}}.$$

EHR phenotyping

$$\phi_\alpha(\mathbf{x}_{\text{new}}) = \mathbf{1} \left(\widehat{\mathbf{x}}_{\text{new}}^\top \widehat{\beta} - z_\alpha \widehat{\mathbf{V}}^{1/2} \geq 0 \right).$$

Theory and Optimality

Theoretical Justification

Theorem 1.

Under regularity conditions, if

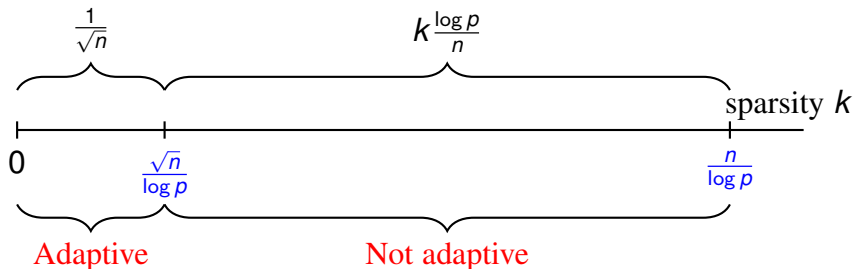
$$k \ll \sqrt{n}/[\log p(\log n)^{1/2}],$$

then

$$\mathbb{P} \left[\mathbf{V}^{-1/2} \left(\widehat{\mathbf{x}}_{\text{new}}^{\top} \beta - \mathbf{x}_{\text{new}}^{\top} \beta \right) \geq \mathbf{z}_{\alpha} \right] \rightarrow \alpha.$$

- ▶ No sparsity on Σ^{-1} and \mathbf{x}_{new} .
- ▶ Approximate $\widehat{\mathbf{u}}^{\top} \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^{\top} \widehat{\beta})(1 - h(\mathbf{X}_i^{\top} \widehat{\beta}))]^{-1} \mathbf{X}_i \epsilon_i$ by $\widehat{\mathbf{u}}^{\top} \frac{1}{n} \sum_{i=1}^n [h(\mathbf{X}_i^{\top} \beta)(1 - h(\mathbf{X}_i^{\top} \beta))]^{-1} \mathbf{X}_i \epsilon_i$
- ▶ Contraction Principle.

Discussion: Optimality of CI for β_j



Cai, T. Tony, Zijian Guo, and Rong Ma. "Statistical inference for high-dimensional generalized linear models with binary outcomes." JASA, to appear.

Numerical Results

Simulation Studies

1. $p = 501$
2. $n \in \{200, 400, 600\}$
3. $\beta_1 = 0$, $\beta_j = (j - 1)/20$ for $2 \leq j \leq 11$ and $\beta_j = 0$ for $12 \leq j \leq p$

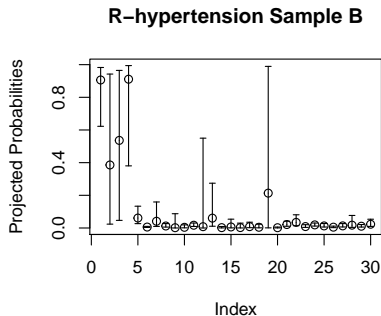
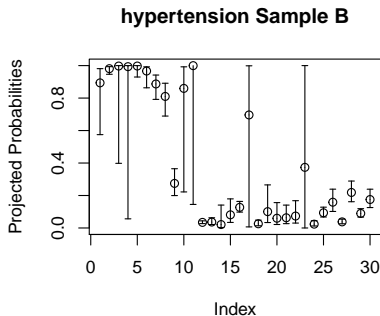
Simulation Studies

		LiVE				Post Selection				hdi				WLDP			
$\ x_{\text{new}}\ _2$	n	Cov	ERR	Len	t	Cov	ERR	Len	t	Cov	ERR	Len	t	Cov	ERR	Len	t
16.1	200	0.98	0.05	0.88	5	0.68	0.54	0.42	1	0.97	0.06	0.93	370	1.00	0.00	1.00	34
	400	0.97	0.10	0.81	14	0.71	0.57	0.38	2	0.96	0.10	0.87	751	1.00	0.00	1.00	56
	600	0.95	0.13	0.74	23	0.70	0.68	0.32	6	0.94	0.10	0.83	3212	1.00	0.00	1.00	118
1.90	200	0.96	0.62	0.34	5	0.80	0.77	0.31	1	0.92	0.86	0.31	371	1.00	0.36	0.58	34
	400	0.94	0.92	0.23	14	0.83	0.93	0.24	2	0.92	0.96	0.23	751	1.00	0.45	0.53	54
	600	0.95	0.95	0.19	22	0.82	0.95	0.20	5	0.95	0.97	0.19	3211	1.00	0.47	0.50	118

Real Data Applications

1. Data: extracted from the Penn Medicine clinical data repository, including demographics, laboratory results, medication prescriptions, vital signs, and encounter meta information.
2. 348 patients, 198 predictors in the final analyses
3. Goal: predicting hypertension, hypertension resistant to standard treatment ("R-hypertension").
4. Outcome prevalence: 39.4% and 8.1%

Real Data Results



We randomly sampled 30 patients as the test sample,

- ▶ Left, indexes 1 to 11 correspond to hypertension.
- ▶ Right, indices 1 to 4 correspond to R-hypertension.

Conclusion and Discussion

1. Non-linear outcome model: **Reweighting**
2. Uniform inference for x_{new} : **Additional constraint**
3. Optimality of CI construction

Future research

1. Outcome surrogates
2. Model misspecification

Reference and Acknowledgement

Guo, Z., Rakshit, P., Herman, D. S., and Chen, J. (2021). Inference for the case probability in high-dimensional logistic regression. *JMLR*, 22(254), 1-54.

Cai, T., Guo, Z., and Ma, R. (2021+). "Statistical inference for high-dimensional generalized linear models with binary outcomes." *JASA*, to appear.

Acknowledgement to NIH and NSF for fundings.

Thank you!