

# Scale-invariant Optimal Sampling for Rare-events data with Sparse Models

Jing Wang

University of Connecticut

---

<sup>1</sup>Collaborates: Haiying Wang<sup>1</sup>, Hao Helen Zhang<sup>2</sup>

<sup>2</sup>Institutions: University of Connecticut<sup>1</sup>, University of Arizona<sup>2</sup>

# Motivation

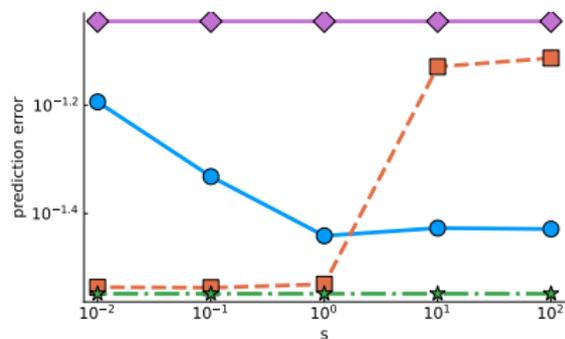
## Rare-events data with sparse models

- Rare-events data are **highly imbalanced** binary response data.
  - ▶ Rare diseases, click-data on recommendation system.
  - ▶ Massive, highly imbalanced.
  - ▶ Involve sparse models, e.g., limited number of key genes related to rare-diseases.
  - ▶ Variable selection is not studied.
  
- **Subsampling** is a popular approach for rare-events data analysis.
  - ▶ Data balancing, reducing computational burdens.
  - ▶ Usually done with strategy:
    - ① Keeping all ones.
    - ② Subsampling zeros according to an important function  $\varphi(\mathbf{x})$ .
  - ▶ Non-uniform subsampling reduces information loss.

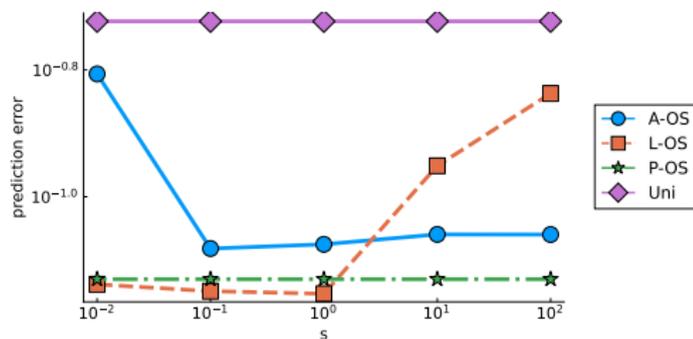
# Motivation

## Scale-dependent issue

- Existing optimal subsampling functions are scale-dependent.
- May lead to inefficient results.
- A wide concern in literature for various data types and models.



(a) Non-sparse parameter



(b) Sparse parameter

## Problem setup

- Rare-events Model:

$$p(\mathbf{x}; \boldsymbol{\theta}_t) := \mathbb{P}(y = 1|\mathbf{x}) = \frac{e^{\alpha_t + f(\mathbf{x}; \boldsymbol{\beta}_t)}}{1 + e^{\alpha_t + f(\mathbf{x}; \boldsymbol{\beta}_t)}} = \frac{e^{g(\mathbf{x}; \boldsymbol{\theta}_t)}}{1 + e^{g(\mathbf{x}; \boldsymbol{\theta}_t)}}.$$

Then,  $\alpha_t \rightarrow -\infty$  as  $N \rightarrow \infty$  implies that  $\frac{N_1}{N_0} \rightarrow 0$ .

- IPW Adaptive Lasso for Variable selection

$$\hat{\boldsymbol{\theta}}_w^{\text{adp}} := \arg \max_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^{N_{\text{sub}}^*} \frac{\ell_i^{\text{sub}}}{\pi(\mathbf{x}_i^{\text{sub}}, y_i^{\text{sub}})} - \lambda_N \sum_{j=1}^p \frac{|\beta_{(j)}|}{|\hat{\beta}_{\text{pl}(j)}|^\gamma} \right\}, \quad (1)$$

where  $\ell_i^{\text{sub}} = y_i^{\text{sub}} g(\mathbf{x}_i^{\text{sub}}; \boldsymbol{\theta}) - \log\{1 + e^{g(\mathbf{x}_i^{\text{sub}}; \boldsymbol{\theta})}\}$ .

- Both optimal probabilities and adaptive lasso requires a pilot estimator. It is natural to combine them into one unified framework.

# Theoretical analysis

## Asymptotic properties of $\hat{\theta}_w^{\text{adp}}$

- 1 Consistency in variable selection:  $\lim_{N \rightarrow \infty} \mathbb{P}(\hat{\mathcal{A}}_w = \mathcal{A}) = 1$
- 2 Asymptotic normality:  $\sqrt{N_1} \mathbf{V}_{w(\mathcal{A})}^{-1/2} (\hat{\theta}_{w(\mathcal{A})}^{\text{adp}} - \theta_{t(\mathcal{A})}) \rightsquigarrow \mathbb{N}(0, \mathbf{I})$ ,

$$\mathbf{V}_{w(\mathcal{A})} = \underbrace{\mathbf{V}_{\text{mle}(\mathcal{A})}^{-1}}_{\text{Full data}} + \underbrace{c \mathbf{V}_{\text{sub}(\mathcal{A})}}_{\text{Information loss}}, \text{ where}$$

$$c \mathbf{V}_{\text{sub}(\mathcal{A})} \propto c \mathbf{M}_{(\mathcal{A})}^{-1} \mathbb{E} \left\{ \frac{e^{2f(\mathbf{x}; \beta_t)}}{\varphi(\mathbf{x})} \dot{g}_{(\mathcal{A})}^{\otimes 2}(\mathbf{x}; \theta_t) \right\} \mathbf{M}_{(\mathcal{A})}^{-1}$$

- $c = \lim_{N \rightarrow \infty} \frac{e^{\alpha t}}{\rho}$  is the imbalance rate in the subsample.

## Message from theoretical analysis

- The asymptotic variances  $\mathbf{V}_{\text{mle}(\mathcal{A})}$  and  $\mathbf{V}_{w(\mathcal{A})}$  are of order  $\frac{1}{N_1}$ .
- If **remain enough 0's**, e.g.,  $c = 0$ , there will be **no information loss**.
- In case there is information loss  $c > 0$ , we can choose  $\varphi(\mathbf{x})$  to minimize the information loss.

# Optimal subsampling function

## Traditional optimal subsampling function and limitations

### 1 A-optimality:

$$\min \text{tr}(\mathbf{V}_{w(\mathcal{A})}) \Rightarrow \varphi_{\text{A-OS}}^{\text{adp}}(\mathbf{x}) \propto p(\mathbf{x}; \boldsymbol{\theta}_t) \|\mathbf{M}_{(\mathcal{A})}^{-1} \dot{g}_{(\mathcal{A})}(\mathbf{x}; \boldsymbol{\theta}_t)\|.$$

### 2 L-optimality:

$$\min \text{tr}(\mathbf{M}_{w(\mathcal{A})}) \Rightarrow \varphi_{\text{L-OS}}^{\text{adp}}(\mathbf{x}) \propto p(\mathbf{x}; \boldsymbol{\theta}_t) \|\dot{g}_{(\mathcal{A})}(\mathbf{x}; \boldsymbol{\theta}_t)\|.$$

- If  $g(\mathbf{x}, \boldsymbol{\theta}) = \alpha + \mathbf{x}^T \boldsymbol{\beta}$ , then  $\varphi_{\text{L-OS}}^{\text{adp}}(\mathbf{x}) \propto p(\mathbf{x}; \boldsymbol{\theta}_t)(1 + \|\mathbf{x}_{(\mathcal{A})}\|)$ .
  - ▶ Due to inaccurate pilot, scale of  $\mathbf{x}_{(\mathcal{A}^c)}$  will affect  $\hat{\varphi}_{\text{L-OS}}^{\text{adp}}(\mathbf{x})$ .
- Construct optimal function by focusing on **prediction error**:

$$\text{MSPE}(\hat{\boldsymbol{\theta}}) = \mathbb{E}_{\mathbf{x}} \left[ \left\{ p(\mathbf{x}; \hat{\boldsymbol{\theta}}) - p(\mathbf{x}; \boldsymbol{\theta}_t) \right\}^2 \right].$$

# Optimal subsampling function

## Scale-invariant optimal subsampling function

- 1 We prove that

$$N_1 e^{-2\alpha_t} \text{MSPE}(\hat{\boldsymbol{\theta}}_{w(\mathcal{A})}^{\text{adp}}) \rightsquigarrow \mathbb{E}^{-1} \left\{ e^{f(\mathbf{x}; \boldsymbol{\beta}_t)} \right\} \mathbf{Z}_{(\mathcal{A})}^{\text{T}} \mathbf{L}_{(\mathcal{A})}^{\text{T}} \boldsymbol{\Omega}_{(\mathcal{A})} \mathbf{L}_{(\mathcal{A})} \mathbf{Z}_{(\mathcal{A})}.$$

where  $\mathbf{Z}_{(\mathcal{A})} \sim \mathbb{N}(0, \mathbf{I})$ ,  $\boldsymbol{\Omega}_{(\mathcal{A})} = \mathbb{E} \left[ e^{2f(\mathbf{x}; \boldsymbol{\beta}_t)} \dot{\mathbf{g}}_{(\mathcal{A})}^{\otimes 2}(\mathbf{x}, \boldsymbol{\theta}_t) \right]$ , and

$$\mathbf{L}_{(\mathcal{A})} = \mathbf{M}_{(\mathcal{A})}^{-1} \mathbf{M}_{w(\mathcal{A})}^{1/2}.$$

- 2 The optimal function that minimizes the asymptotic mean is

$$\varphi_{\text{P-OS}}^{\text{adp}}(\mathbf{x}) \propto p(\mathbf{x}; \boldsymbol{\theta}_t) \left\| \boldsymbol{\Omega}_{(\mathcal{A})}^{\frac{1}{2}} \mathbf{M}_{(\mathcal{A})}^{-1} \dot{\mathbf{g}}_{(\mathcal{A})}(\mathbf{x}; \boldsymbol{\theta}_t) \right\|,$$

which is scale-invariant for a class of  $g$  including neural networks.

## Penalized MSCL estimator and practical algorithm

- The IPW assigns smaller weights for more informative data points
- To improve the estimation efficiency, let  $I_i^{\text{sub}} = -\log \{ \rho \varphi(\mathbf{x}_i^{\text{sub}}) \}$ ,

$$\hat{\boldsymbol{\theta}}_{\text{mscl}}^{\text{adp}} := \arg \max_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^{N_{\text{sub}}^*} \ell_{\text{mscl},i}^{\text{sub}} - \lambda_N \sum_{j=1}^p \frac{|\beta_{(j)}|}{|\hat{\beta}_{\text{pl}(j)}|^\gamma} \right\}, \quad (2)$$

where  $\ell_{\text{mscl},i}^{\text{sub}} = y_i^{\text{sub}} g(\mathbf{x}_i^{\text{sub}}; \boldsymbol{\theta}) - \log \{ 1 + e^{g(\mathbf{x}_i^{\text{sub}}; \boldsymbol{\theta}) + I_i^{\text{sub}}} \}$ .

- **Efficiency:**  $\mathbf{V}_{\text{mscl}(\mathcal{A})} \leq \mathbf{V}_{\text{w}(\mathcal{A})}$ , and  $\mathbf{V}_{\text{mscl}(\mathcal{A})} = \mathbf{V}_{\text{mle}(\mathcal{A})}$  if  $c = 0$ .

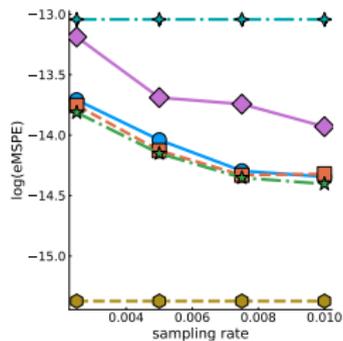
### Practical Algorithm

#### 1 First-stage screening:

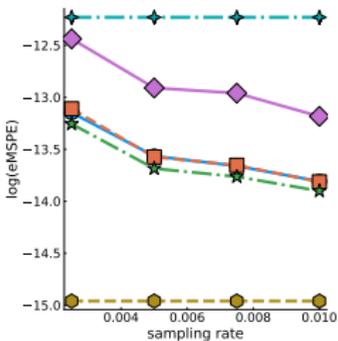
- 1 Take a pilot sample, and obtain a lasso estimator.
- 2 Estimate  $\hat{\varphi}(\mathbf{x}_i)$  for  $i = 1, \dots, N$  and  $\hat{\mathcal{A}}$ .

- #### 2 Second-stage screening:
- Subsampling from 0's with  $\hat{\varphi}(\mathbf{x}_i)$ , and compute adaptive lasso.

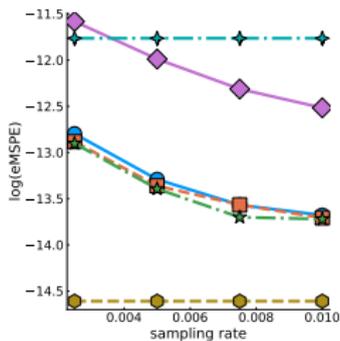
- 1 Case A: Small active effects.
- 2 Case B: Large and small active effects, different signs
- 3 Case C: Large and small active effects, same signs.



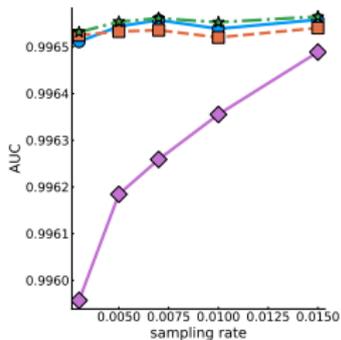
(a) Case A



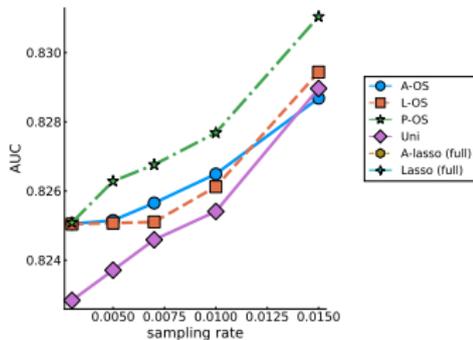
(b) Case B



(c) Case C



(d) covtype



(e) font

# Conclusion

## Conclusion

- 1 For rare-events data with sparse models, subsampling estimators can be as efficient as full data estimators under the true model
- 2 Traditional optimal functions are scale-dependent. The scale-invariant function based on prediction error is a better choice.

## Limitation and Future work

- 1 Optimal functions are based on asymptotic normality and asymptotic mean square error.
- 2 Optimal functions based on the quality of variable selection.
- 3 Non-asymptotic behaviors.

Thank you!