

# Adaptive Debiased SGD in High-dimensional GLMs with Streaming Data

Yuanhang Luo<sup>1</sup>

Joint work with Ruijian Han<sup>2</sup>, Lan Luo<sup>3</sup>, Yuanyuan Lin<sup>4</sup>, Jian Huang<sup>2</sup>

<sup>1</sup>Department of Applied Mathematics, The Hong Kong Polytechnic University

<sup>2</sup>Department of Data Science and Artificial Intelligence, The Hong Kong Polytechnic University

<sup>3</sup>Department of Biostatistics and Epidemiology, Rutgers University

<sup>4</sup>Department of Statistics, The Chinese University of Hong Kong

29th Aug 2024

# Table of Contents

1 Background

2 Online Method

3 Theoretical Results

4 Empirical Results

# Table of Contents

1 Background

2 Online Method

3 Theoretical Results

4 Empirical Results

# Big Data

Big data is about

*"datasets whose size is beyond the ability of typical database software tools to capture, store, manage, and analyze"* [Manyika et al., 2011]

McKinsey Global Institute

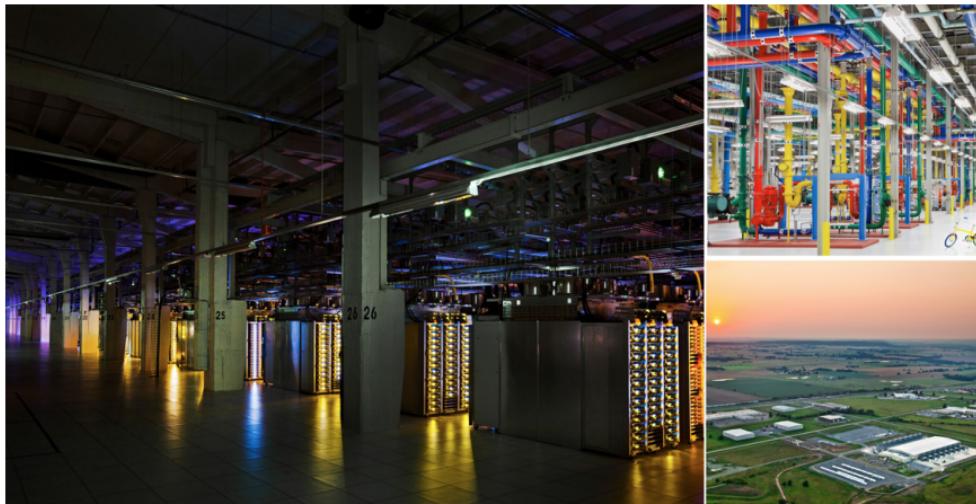


Figure: Google Data Center

# Streaming Data

Many applications must process **large streams** of live data and provide results in **near-real-time**

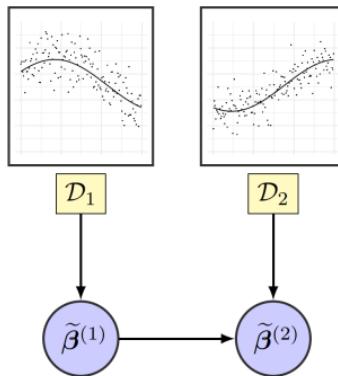
- Algorithmic trading
- Real-time operations management
- Online fraud detection
- Proximity/location tracking
- Intrusion detection systems
- Traffic management
- Real-time recommendations
- Social media/data analytics
- Precision agriculture
- Internet of things
- Gaming data feed
- ...

Some of the major challenges of handling streaming data.

- Historical data is no longer available (needs **one pass** process).
- Each data has **large number of features** (high-dimensional).
- Critical **memory** requirements.
- Limiting **computing power**.

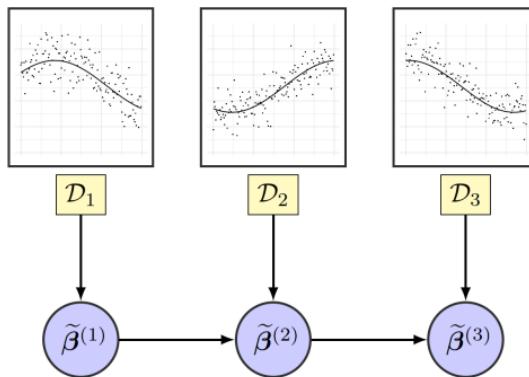
# Streaming Data

Proposed solution: **Online learning**



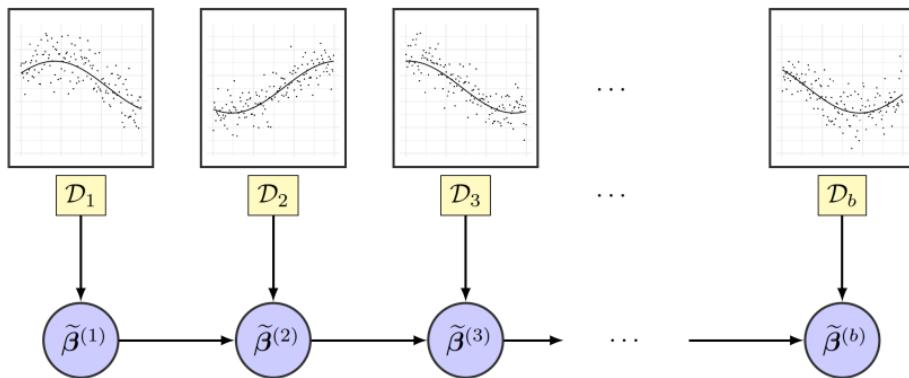
# Streaming Data

Proposed solution: **Online learning**



# Streaming Data

Proposed solution: **Online learning**



Objective: Develop an “online” estimator that

- can be **updated sequentially** using raw data in current batch only
- requires little **time and space complexities**
- still achieves desired **statistical property**

# High-dimensional GLM

A sequence of  $n$  data points  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  are independently and identically sampled from a generalized linear model (GLM):

$$\mathbb{P}(y | \mathbf{x}; \boldsymbol{\beta}^*) \propto \exp \left\{ \frac{y (\mathbf{x}^\top \boldsymbol{\beta}^*) - \Phi(\mathbf{x}^\top \boldsymbol{\beta}^*)}{c(\sigma)} \right\}, \quad (1)$$

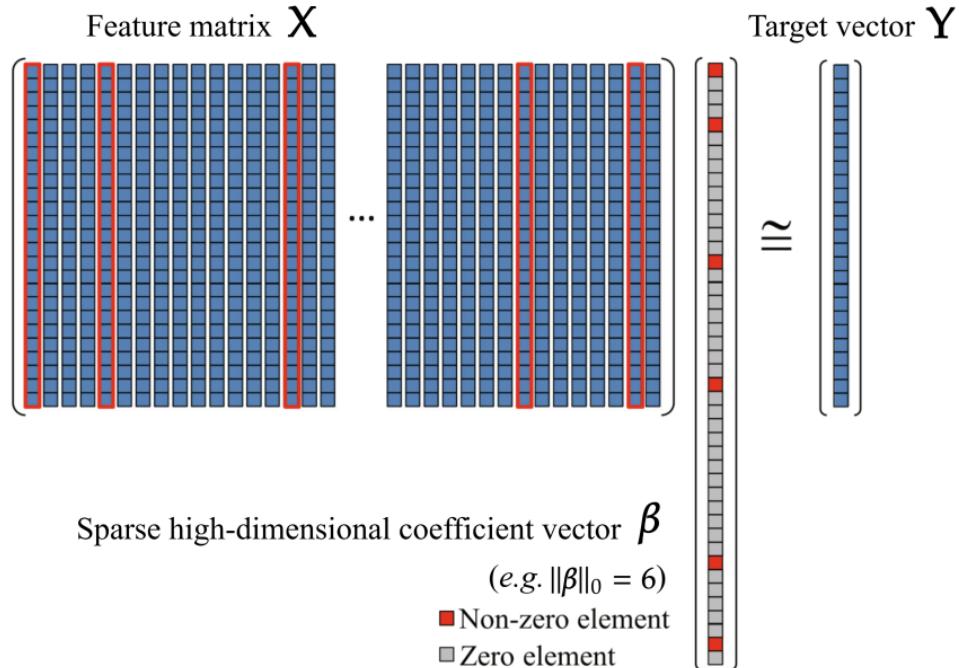
- $\mathbf{x} \in \mathbb{R}^p$ . High-dimensional case:  $n \ll p$ .
- Sparsity:  $\|\boldsymbol{\beta}^*\|_0 = s_0$ .

The ground truth  $\boldsymbol{\beta}^*$  is the minimum of the population loss function

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} \mathbb{E}[F_n(\boldsymbol{\beta})], \quad (2)$$

where  $F_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \{-y_i(\mathbf{x}_i^\top \boldsymbol{\beta}) + \Phi(\mathbf{x}_i^\top \boldsymbol{\beta})\}$ , the negative log-likelihood of model  $\mathbb{P}$ .

# Example: High-dimensional Linear Regression



# Statistical Inference in High Dimension

Traditionally, we estimate  $\beta^*$  using data  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  via two steps:

$$\tilde{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \{F_n(\beta) + \lambda_n \|\beta\|_1\}, \quad (3)$$

and

$$\hat{\beta}_{j,\text{de}} = \tilde{\beta}_j - \Theta_j^\top \nabla F_n(\tilde{\beta}), \quad (4)$$

where  $\Theta$  is the inverse of the covariance matrix (information matrix).  $\Theta_j$  can be estimated by [Zhang and Zhang, 2014, van de Geer et al., 2014]...

# Statistical Inference in High Dimension

Traditionally, we estimate  $\beta^*$  using data  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  via two steps:

$$\tilde{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \{F_n(\beta) + \lambda_n \|\beta\|_1\}, \quad (3)$$

and

$$\hat{\beta}_{j,\text{de}} = \tilde{\beta}_j - \Theta_j^\top \nabla F_n(\tilde{\beta}), \quad (4)$$

where  $\Theta$  is the inverse of the covariance matrix (information matrix).  $\Theta_j$  can be estimated by [Zhang and Zhang, 2014, van de Geer et al., 2014]...

By doing so, we have asymptotic normality

$$(\hat{\beta}_{j,\text{de}} - \beta_j^*) / \hat{\sigma}_j \rightarrow \mathcal{N}(0, 1) \text{ in distribution as } n \rightarrow \infty. \quad (5)$$

And we are able to make all kinds of statistical inference based on

$$\mathbb{P}\{\beta_j^* \in (\hat{\beta}_{j,\text{de}} - z_{\alpha/2} \hat{\sigma}_j, \hat{\beta}_{j,\text{de}} + z_{\alpha/2} \hat{\sigma}_j)\} \rightarrow 1 - \alpha \text{ as } n \rightarrow \infty. \quad (6)$$

To do

$$\tilde{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{n} \sum_{i=1}^n \{-y_i(\mathbf{x}_i^\top \beta) + \Phi(\mathbf{x}_i^\top \beta)\} + \lambda_n \|\beta\|_1 \right\},$$

and

$$\hat{\beta}_{j,\text{de}} = \tilde{\beta}_j - \Theta_j^\top \nabla F_n(\tilde{\beta}),$$

## Remark

we need to store data of size  $\mathcal{O}(n \times p)$ , and do calculations involving the covariance matrix of size  $\mathcal{O}(p \times p)$ . In large datasets (e.g. text data, web corpus, gene expression data... )  $p$  can be extensive. Let's say  $p = 3,000,000$ , then,

$$p \times p = 9,000,000,000,000. \quad (7)$$

To operate this  $p \times p$  matrix of double-precision floating-point (float64) numbers, you would need at least approximately **65.54 terabytes (TB)** of storage.

# Our Solution

## Question

Can we bring it "online"?

# Our Solution

## Question

Can we bring it "online"?

**Table:** A comparison in time and space complexities for statistical inference of one single coefficient. LM: linear model; GLM: generalized linear model.

Works	Model	One-pass?	Space complexity
[Chen et al., 2020]	LM	Yes	$\mathcal{O}(np)$
[Deshpande et al., 2023]	LM	No	$\mathcal{O}(p^2)$
[Han et al., 2023]	LM	Yes	$\mathcal{O}(p)$
[Shi et al., 2021]	GLM	No	$\mathcal{O}(np)$
[Luo et al., 2023]	GLM	No	$\mathcal{O}(p^2)$
Proposed Method	GLM	Yes	$\mathcal{O}(p)$

# Our Solution

## Question

Can we bring it "online"?

**Table:** A comparison in time and space complexities for statistical inference of one single coefficient. LM: linear model; GLM: generalized linear model.

Works	Model	One-pass?	Space complexity
[Chen et al., 2020]	LM	Yes	$\mathcal{O}(np)$
[Deshpande et al., 2023]	LM	No	$\mathcal{O}(p^2)$
[Han et al., 2023]	LM	Yes	$\mathcal{O}(p)$
[Shi et al., 2021]	GLM	No	$\mathcal{O}(np)$
[Luo et al., 2023]	GLM	No	$\mathcal{O}(p^2)$
Proposed Method	GLM	Yes	$\mathcal{O}(p)$

## Main Idea

Optimization + Statistics

# Table of Contents

1 Background

2 Online Method

3 Theoretical Results

4 Empirical Results

# Overview of the Algorithm

Upon the arrival of the  $k$ -th data, update  $\hat{\beta}^{(i)}$  and  $\hat{\gamma}_{\cdot j}^{(i)}$  with a variant of SGD, then construct the following summary statistics:

$$\mathbf{A}_1^{(k)} = \sum_{i=1}^k \mathbf{x}_i \{-y_i + \dot{\Phi}(\mathbf{x}_i^\top \hat{\beta}^{(i)})\}, \quad \mathbf{A}_2^{(k)} = \sum_{i=1}^k \mathbf{x}_i^\top \hat{\gamma}_{\cdot j}^{(i)} \ddot{\Phi}(\mathbf{x}_i^\top \hat{\beta}^{(i)}) \mathbf{x}_i,$$

$$a_3^{(k)} = \sum_{i=1}^k \mathbf{x}_i^\top \hat{\gamma}_{\cdot j}^{(i)} \ddot{\Phi}(\mathbf{x}_i^\top \hat{\beta}^{(i)}) \mathbf{x}_i^\top \hat{\beta}^{(i)}, \quad a_4^{(k)} = - \sum_{i=1}^k \mathbf{x}_i^\top \hat{\gamma}_{\cdot j}^{(i)} (\mathbf{x}_i)_j \{\ddot{\Phi}(\mathbf{x}_i^\top \hat{\beta}^{(i)})\}.$$

All are  **$p$ -dimensional vectors** and **scalars**. The approximated debiased lasso (ADL) estimator

$$\tilde{\beta}_{j,\text{de}}^{(k)} := \hat{\beta}_j^{(k)} - \frac{\{\mathbf{A}_1^{(k)}\}^\top \hat{\gamma}_{\cdot j}^{(k)} + \{\mathbf{A}_2^{(k)}\}^\top \hat{\beta}^{(k)} - a_3^{(k)}}{a_4^{(k)}}.$$

We update summary statistics **online**:

$$\mathbf{A}_1^{(k)} \leftarrow \mathbf{A}_1^{(k-1)} + \mathbf{x}_k \{-y_k + \dot{\Phi}(\mathbf{x}_k^\top \hat{\beta}^{(k)})\}.$$

# Overview of the Algorithm

Upon the arrival of the  $k$ -th data, update  $\hat{\beta}^{(i)}$  and  $\hat{\gamma}_{\cdot j}^{(i)}$  with a variant of SGD, then construct the following summary statistics:

$$\mathbf{A}_1^{(k)} = \sum_{i=1}^k \mathbf{x}_i \{-y_i + \dot{\Phi}(\mathbf{x}_i^\top \hat{\beta}^{(i)})\}, \quad \mathbf{A}_2^{(k)} = \sum_{i=1}^k \mathbf{x}_i^\top \hat{\gamma}_{\cdot j}^{(i)} \ddot{\Phi}(\mathbf{x}_i^\top \hat{\beta}^{(i)}) \mathbf{x}_i,$$

$$a_3^{(k)} = \sum_{i=1}^k \mathbf{x}_i^\top \hat{\gamma}_{\cdot j}^{(i)} \ddot{\Phi}(\mathbf{x}_i^\top \hat{\beta}^{(i)}) \mathbf{x}_i^\top \hat{\beta}^{(i)}, \quad a_4^{(k)} = - \sum_{i=1}^k \mathbf{x}_i^\top \hat{\gamma}_{\cdot j}^{(i)} (\mathbf{x}_i)_j \{\ddot{\Phi}(\mathbf{x}_i^\top \hat{\beta}^{(i)})\}.$$

All are  **$p$ -dimensional vectors** and **scalars**. The approximated debiased lasso (ADL) estimator

$$\tilde{\beta}_{j,\text{de}}^{(k)} := \hat{\beta}_j^{(k)} - \frac{\{\mathbf{A}_1^{(k)}\}^\top \hat{\gamma}_{\cdot j}^{(k)} + \{\mathbf{A}_2^{(k)}\}^\top \hat{\beta}^{(k)} - a_3^{(k)}}{a_4^{(k)}}.$$

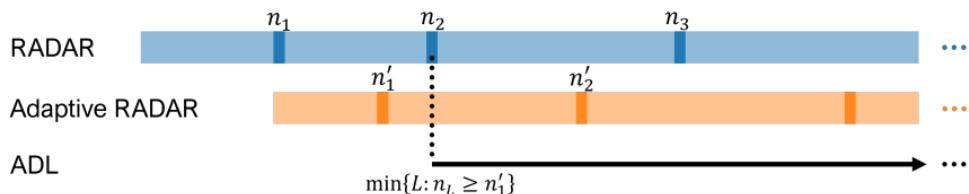
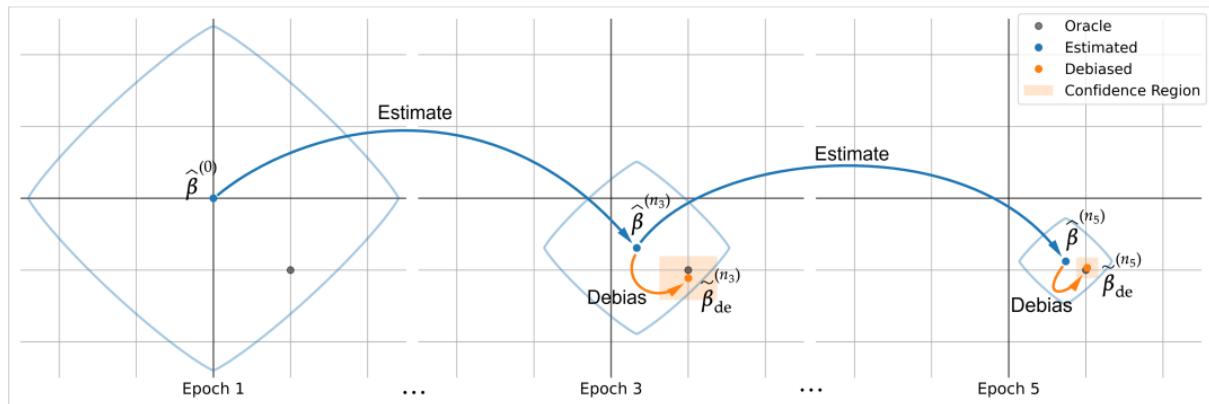
We update summary statistics **online**:

$$\mathbf{A}_1^{(k)} \leftarrow \mathbf{A}_1^{(k-1)} + \mathbf{x}_k \{-y_k + \dot{\Phi}(\mathbf{x}_k^\top \hat{\beta}^{(k)})\}.$$

## Remark

Now we can recurse  $\{\hat{\beta}^{(k)}, \hat{\gamma}_{\cdot j}^{(k)}, \mathbf{A}_1^{(k)}, \mathbf{A}_2^{(k)}, a_3^{(k)}, a_4^{(k)}\}$  and **clear all other variables in memory**.

# Overview of the Algorithm



# Complexities

Table: A comparison in time and space complexities for statistical inference of one single coefficient. LM: linear model; GLM: generalized linear model.

Works	Model	One-pass?	Space complexity
[Chen et al., 2020]	LM	Yes	$\mathcal{O}(np)$
[Deshpande et al., 2023]	LM	No	$\mathcal{O}(p^2)$
[Han et al., 2023]	LM	Yes	$\mathcal{O}(p)$
[Shi et al., 2021]	GLM	No	$\mathcal{O}(np)$
[Luo et al., 2023]	GLM	No	$\mathcal{O}(p^2)$
Proposed Method	GLM	Yes	$\mathcal{O}(p)$

# Table of Contents

1 Background

2 Online Method

3 Theoretical Results

4 Empirical Results

## Theorem (Convergence Rate)

With some regularity conditions and proper choice of hyperparameters, and initial error  $d_0 = \|\hat{\beta}^{(0)} - \beta^*\|_1$  the following events

$$\|\hat{\beta}^{(i)} - \beta^*\|_1 \leq C_1 s_0 d_0 \sqrt{\frac{(\log p)^3}{i}}, \quad \|\hat{\beta}^{(i)} - \beta^*\|_2 \leq C_1 d_0 \sqrt{\frac{s_0 (\log p)^3}{i}} \quad (8)$$

for  $i \geq n_1$  hold uniformly for a universal constant  $C_1$  with high probability.

# Theoretical Results

## Theorem (Convergence Rate)

With some regularity conditions and proper choice of hyperparameters, and initial error  $d_0 = \|\hat{\beta}^{(0)} - \beta^*\|_1$  the following events

$$\|\hat{\beta}^{(i)} - \beta^*\|_1 \leq C_1 s_0 d_0 \sqrt{\frac{(\log p)^3}{i}}, \quad \|\hat{\beta}^{(i)} - \beta^*\|_2 \leq C_1 d_0 \sqrt{\frac{s_0 (\log p)^3}{i}} \quad (8)$$

for  $i \geq n_1$  hold uniformly for a universal constant  $C_1$  with high probability.

## Theorem (Asymptotic Normality)

Further with

$$(s_0 s_1 d_0 d_1 + s_0^2 d_0^2 + s_0 s_1 d_0^2) (\log p)^4 \log_2(n) = o_{\mathbb{P}}(\sqrt{n}),$$

then, for any fixed  $j \in [p]$ ,

$$\left( \tilde{\beta}_{j,de}^{(n)} - \beta_j^* \right) / \tilde{\nu}_j^{(n)} \rightarrow \mathcal{N}(0, 1) \text{ as } n \rightarrow \infty \text{ in distribution.}$$

## Remark

Most existing methods rely on the bounded individual probability condition that

$$P(y_i = 1 | X_i) \in (c, 1 - c)$$

for all  $1 \leq i \leq n$  and some  $c \in (0, 1/2)$ . To the best of our knowledge, only [Guo et al., 2021, Cai et al., 2023] and ours remove this.

# Table of Contents

1 Background

2 Online Method

3 Theoretical Results

4 Empirical Results

# Simulation with Synthetic Data (1)

True coefficients:  $\beta_{D_1}^* = 1, \beta_{D_2}^* = -1$  and others are set to 0.

Table: Case 2:  $n = 200, p = 500, s_0 = 6, \Sigma = \{0.5^{|i-j|}\}_{i,j=1,\dots,p}$ . Simulation results are averaged over 200 replications. Other methods: deLasso [[van de Geer et al., 2014](#)], LSW [[Cai et al., 2023](#)], ODL [[Luo et al., 2023](#)].

$\beta_k^*$	deLasso	LSW	ODL			ADL		
Sample size $n$	200	200	40	120	200	40	120	200
Coverage probability	0	0.983	0.996	1.000	0.999	0.998	0.951	0.951
	1	0.598	0.936	1.000	0.880	0.705	0.891	0.955
	-1	0.610	0.946	1.000	0.861	0.670	0.881	0.961
Absolute bias	0	0.114	0.113	0.022	0.019	0.017	0.901	0.425
	1	0.325	0.260	0.657	0.590	0.539	1.008	0.423
	-1	0.326	0.259	0.660	0.598	0.548	1.034	0.426
Coverage length	0	0.753	1.207	3.423	1.691	1.287	4.419	2.053
	1	0.769	1.261	3.434	1.689	1.286	4.555	2.086
	-1	0.767	1.271	3.416	1.648	1.282	4.431	2.078
Time (s)	6.550	20.04	0.628			0.365		

# Simulation with Synthetic Data (1)

True coefficients:  $\beta_{D_1}^* = 1, \beta_{D_2}^* = -1$  and others are set to 0.

Table: Case 2:  $n = 200, p = 500, s_0 = 6, \Sigma = \{0.5^{|i-j|}\}_{i,j=1,\dots,p}$ . Simulation results are averaged over 200 replications. Other methods: deLasso [van de Geer et al., 2014], LSW [Cai et al., 2023], ODL [Luo et al., 2023].

$\beta_k^*$	deLasso	LSW	ODL			ADL		
Sample size $n$	200	200	40	120	200	40	120	200
Coverage probability	0	0.983	0.996	1.000	0.999	0.998	0.951	0.951
	1	0.598	0.936	1.000	0.880	0.705	0.891	0.955
	-1	0.610	0.946	1.000	0.861	0.670	0.881	0.961
Absolute bias	0	0.114	0.113	0.022	0.019	0.017	0.901	0.425
	1	0.325	0.260	0.657	0.590	0.539	1.008	0.423
	-1	0.326	0.259	0.660	0.598	0.548	1.034	0.426
Coverage length	0	0.753	1.207	3.423	1.691	1.287	4.419	2.053
	1	0.769	1.261	3.434	1.689	1.286	4.555	2.086
	-1	0.767	1.271	3.416	1.648	1.282	4.431	2.078
Time (s)	6.550	20.04	0.628			0.365		

## Remark

deLasso and ODL fail because they rely on the bounded individual probability condition.

## Simulation with Synthetic Data (2)

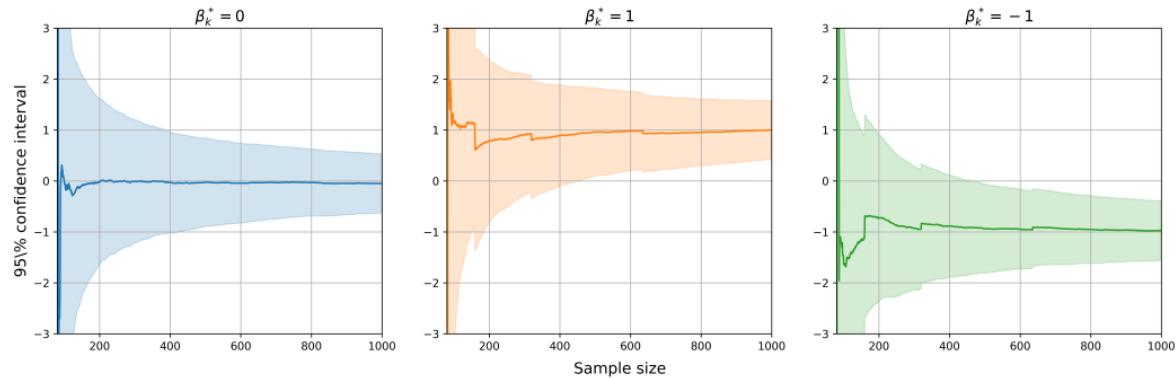


Figure:  $n = 1000$ ,  $p = 20000$ ,  $s_0 = 20$ ,  $\Sigma = \{0.5^{|i-j|}\}_{i,j=1,\dots,p}$ . Simulation results are averaged over 200 replications (around 110.4s for each replication).

# Large-scale URLs Data

Real data application: Detection of malicious websites [Ma et al., 2010].

- $x$ : 3,231,961 features containing lexical, host-based information...
- $y$ : Binary, is this URL a phishing site or not?
- Logistics regression.

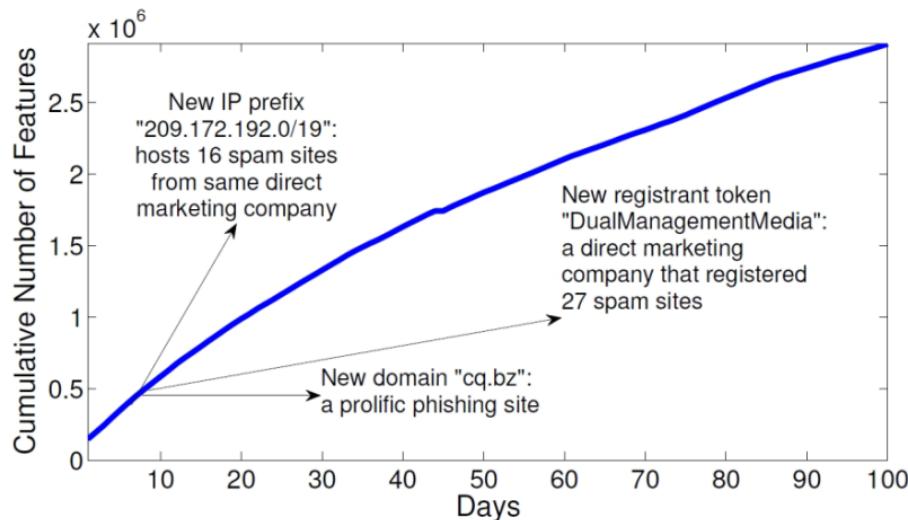
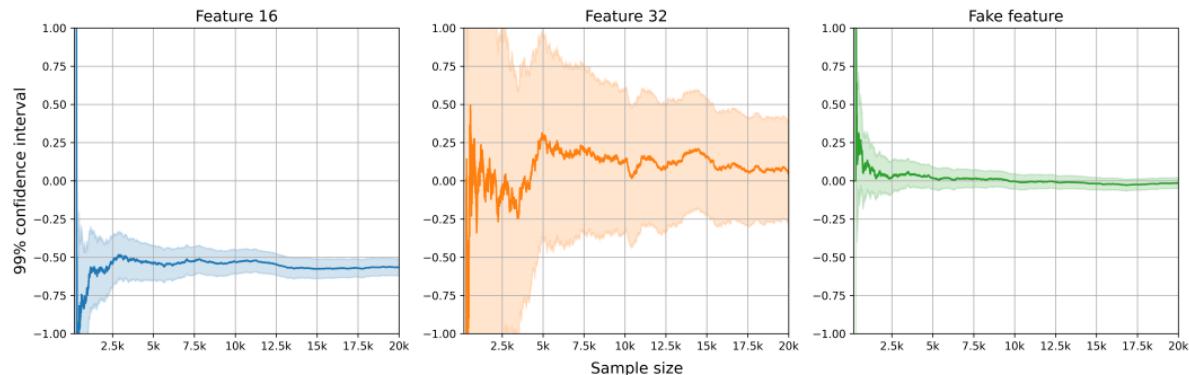


Figure: The number of features also grows with incoming data.

# Large-scale URLs Data

Real data application on logistic regression:

- $x$ : 3,231,961 features containing lexical, host-based information...
- $y$ : Binary, is this URL a phishing site or not?



- An incoming data needs around **2s** to arrive on average.
- ADL consumes around **0.3s** on each data.

# 10-K Financial Report Data

Data source: [Kogan et al., 2009]. We used more than 10,000 10-K forms from 2002 to 2004 to predict future market volatility.

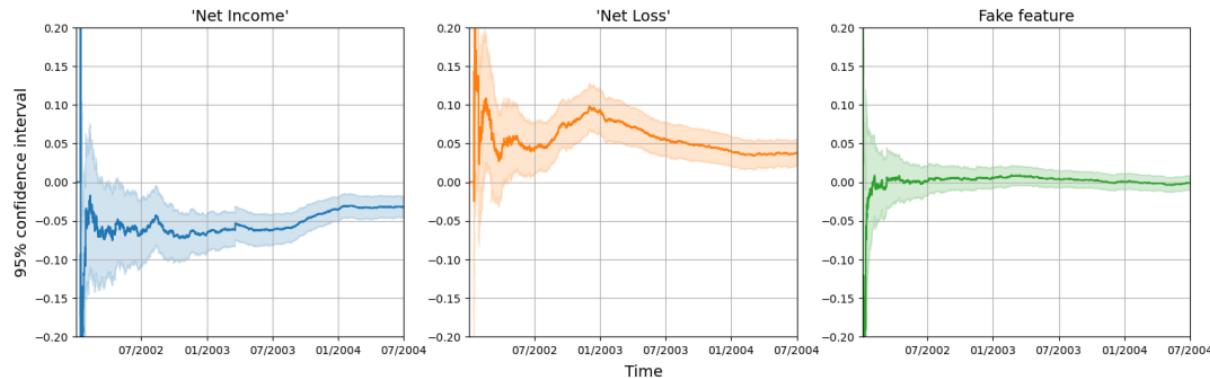
- $x$ : 4,500,000 features containing unigram and bigram frequency.
- $y$ : Increment of market volatility (log scale).
- Linear regression



**Figure:** Form 10-K is an annual report required by the U.S. Securities and Exchange Commission (SEC), that gives a comprehensive summary of a company's financial performance including, organizational structure, executive compensation, equity, subsidiaries, and audited financial statements....

# 10K Financial Report Data

- $x$ : 4,500,000 features containing unigram and bigram frequency.
- $y$ : Increment of market volatility (log scale).



- ADL consumes around **0.45s** on each data.

## Remark (Leverage/Announcement Effect)

“Bad news” has a larger impact on market volatility than “good news”.

## Conclusion & Remarks

In this research, we develop **online** statistical inference approach in high-dimensional GLM that

- sequentially updates **without retrieving** historical data.
- has a space complexity  $\mathcal{O}(p)$ , instead of  $\mathcal{O}(p^2)$  nor  $\mathcal{O}(np)$ .
- provides **theoretical guarantees** for convergence and asymptotic normality.
- achieves similar performance as offline counterparts.

Future works:

- non i.i.d. data, covariate shift
- non-convex objectives
- ...

Suggested reference:

- Han, R., Luo, L., Luo, Y., Lin, Y. and Huang, J. (2024). "Adaptive debiased SGD in high-dimensional GLMs with streaming data". Manuscript submitted for publication.

Thanks

Thanks for watching!

Q&A



Figure: Contact info: [chattelion.luo@connect.polyu.hk](mailto:chattelion.luo@connect.polyu.hk)

# Reference I

-  Cai, T. T., Guo, Z., and Ma, R. (2023). Statistical inference for high-dimensional generalized linear models with binary outcomes. *Journal of the American Statistical Association*, 118(542):1319–1332.
-  Chen, X., Lee, J. D., Tong, X. T., and Zhang, Y. (2020). Statistical inference for model parameters in stochastic gradient descent. *The Annals of Statistics*, 48(1):251 – 273.
-  Deshpande, Y., Javanmard, A., and Mehrabi, M. (2023). Online debiasing for adaptively collected high-dimensional data with applications to time series analysis. *Journal of the American Statistical Association*, 118(542):1126–1139.
-  Guo, Z., Rakshit, P., Herman, D. S., and Chen, J. (2021). Inference for the case probability in high-dimensional logistic regression. *Journal of Machine Learning Research*, 22(254):1–54.
-  Han, R., Luo, L., Lin, Y., and Huang, J. (2023). Online inference with debiased stochastic gradient descent. *Biometrika*, 111(1):93–108.
-  Kogan, S., Levin, D., Routledge, B. R., Sagi, J. S., and Smith, N. A. (2009). Predicting risk from financial reports with regression. *USA. Association for Computational Linguistics*.
-  Luo, L., Han, R., Lin, Y., and Huang, J. (2023). Online inference in high-dimensional generalized linear models with streaming data. *Electronic Journal of Statistics*, 17(2):3443 – 3471.
-  Ma, J., Kulesza, A., Dredze, M., Crammer, K., Saul, L., and Pereira, F. (2010). Exploiting feature covariance in high-dimensional online learning. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 493–500.
-  Manyika, J., Chui, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C., and Hung Byers, A. (2011). Big data: The next frontier for innovation, competition, and productivity.
-  Shi, C., Song, R., Lu, W., and Li, R. (2021). Statistical inference for high-dimensional models via recursive online-score estimation. *Journal of the American Statistical Association*, 116(535):1307–1318.

# Reference II



van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014).

On asymptotically optimal confidence regions and tests for high-dimensional models.  
*The Annals of Statistics*, 42(3):1166 – 1202.



Zhang, C.-H. and Zhang, S. S. (2014).

Confidence intervals for low dimensional parameters in high dimensional linear models.  
*Journal of the Royal Statistical Society Series B: Statistical Methodology*, 76(1):217–242.