

# Real-time Regression Analysis of Streaming Health Datasets

Lan Luo

This work is supervised by Professor Peter X.K. Song.



March 25, 2019

# Motivation

## Challenges:

- storage for cumulatively growing datasets;
- recomputation when new data batch arrives.

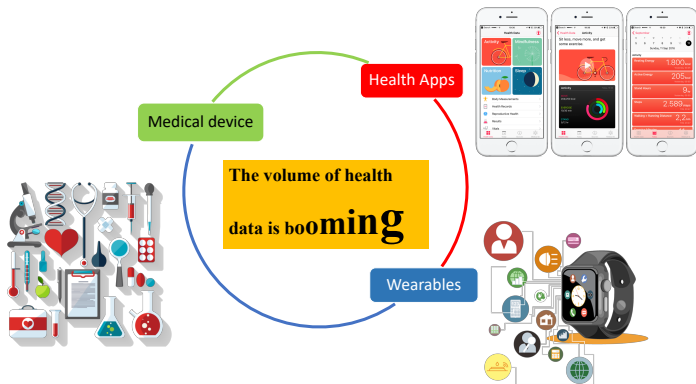
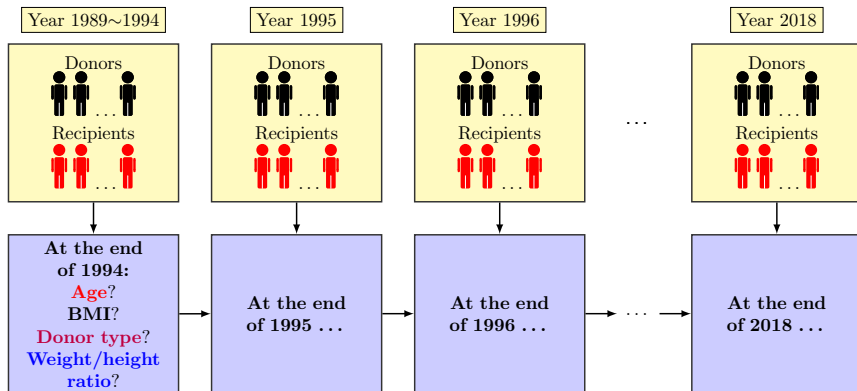


Image sources: <https://unixtitan.net/explore/see-clipart-smart-watch/>,  
<https://www.idownloadblog.com/2018/01/11/health-app-data-murder-investigation/>,  
<http://medtechasia.in/new-medical-devices-being-regulated-under-the-drugs-and-cosmetics-act/>.

# Motivation



Objective: updating the association between 5-year graft failure and demographic features of donors/recipients based on yearly updated kidney transplant data.



# Introduction

HMMcopy.pdf

- ① Rule: Updating the objective statistics without historical raw data but only summary statistics.

- ② A simple example:

$$M(D_1, D_2) = \frac{1}{n_1 + n_2} \left( \sum_{i=1}^{n_1} x_{1i} + \sum_{i=1}^{n_2} x_{2i} \right) = \frac{1}{n_1 + n_2} \left( n_1 M(D_1) + \sum_{i=1}^{n_2} x_{2i} \right).$$

- ③ Question: Can the maximum likelihood estimation (MLE) in generalized linear models (GLM) be updated sequentially like the sample mean?



# Model and Notations

- GLM: continuous, binary and count data;
- Suppose  $(y_i, \mathbf{x}_i) \sim f(y, \mathbf{x}; \beta_0, \phi_0)$  independently for  $i = 1, \dots, N_B$ ;
- Goal: fit a regression model with mean  $\mathbb{E}(y_i | \mathbf{x}_i) = g(\mathbf{x}_i^T \beta)$ .

notation.pdf



# Existing Methods

Methods	Pros	Cons
Oracle MLE	(i) asymptotically unbiased (ii) statistical efficient	(i) store large data (ii) refit model
AI-SGD	online point estimation and bootstrap inference	large bias and standard error
OLSE	online point estimation and inference	linear model only
CEE/CUEE	online point estimation and inference	restriction: $B \ll n_b$

**Table:** AI-SGD denotes the averaged implicit stochastic gradient descent (Toulis & Airolidi, 2015); OLSE, CEE and CUEE are three online meta-type estimation methods (Schifano *et al.* , 2016).



# Renewable Estimation

- A simple scenario with two data batches  $D_1$  and  $D_2$ , we aim to solve  $\hat{\beta}_2^*$  satisfying

$$\mathbf{U}_1(D_1; \hat{\beta}_2^*) + \mathbf{U}_2(D_2; \hat{\beta}_2^*) = \mathbf{0},$$

- Taking the first-order Taylor expansion of the first term around  $\hat{\beta}_1$  leads to

$$\mathbf{U}_1(D_1; \hat{\beta}_1) + \mathbf{J}_1(D_1; \hat{\beta}_1)(\hat{\beta}_1 - \hat{\beta}_2^*) + \mathbf{U}_2(D_2; \hat{\beta}_2^*) + O_p(\|\hat{\beta}_2^* - \hat{\beta}_1\|^2) = \mathbf{0}.$$

- The error term may be asymptotically ignored. We propose a new estimator  $\tilde{\beta}_2$  satisfies the following equation:

$$\mathbf{J}_1(D_1; \hat{\beta}_1)(\hat{\beta}_1 - \tilde{\beta}_2) + \mathbf{U}_2(D_2; \tilde{\beta}_2) = \mathbf{0}.$$

- $\tilde{\beta}_2$  is a **renewable estimator** of  $\beta_0$ , and the above equation is termed as an **incremental estimating equation**.



# Renewable Estimation

- Generalizing to an arbitrary time point  $B$ ,  $\tilde{\beta}_B$  is the solution to the following incremental estimating equation:

$$\sum_{b=1}^{B-1} \mathbf{J}_b(D_b; \tilde{\beta}_b)(\tilde{\beta}_{B-1} - \tilde{\beta}_B) + \mathbf{U}_B(D_B; \tilde{\beta}_B) = \mathbf{0}$$

- Solving via incremental updating algorithm:

$$\begin{aligned} \tilde{\beta}_B^{(r+1)} &= \tilde{\beta}_B^{(r)} + \left\{ \sum_{b=1}^{B-1} \mathbf{J}_b(D_b; \tilde{\beta}_b) + \mathbf{J}_B(D_B; \tilde{\beta}_{B-1}) \right\}^{-1} \\ &\quad \times \left\{ \sum_{b=1}^{B-1} \mathbf{J}_b(D_b; \tilde{\beta}_b)(\tilde{\beta}_{B-1} - \tilde{\beta}_B^{(r)}) + \mathbf{U}_B(D_B; \tilde{\beta}_B^{(r)}) \right\} \\ &= \tilde{\beta}_B^{(r)} + \left\{ \tilde{\mathbf{J}}_{B-1} + \mathbf{J}_B(D_B; \tilde{\beta}_{B-1}) \right\}^{-1} \tilde{\mathbf{U}}_B^{(r)} \end{aligned}$$

- Key components from historical data:  $\{\tilde{\beta}_{B-1}, \tilde{\mathbf{J}}_{B-1}\}$ .





## Main Theorem

Under some regularity conditions, as  $N_B = \sum_{b=1}^B n_b \rightarrow \infty$ ,

- ① consistency:  $\tilde{\beta}_B \xrightarrow{P} \beta_0$ ;
- ② asymptotic efficiency:  $\sqrt{N_B}(\tilde{\beta}_B - \beta_0) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathcal{I}^{-1}(\beta_0))$ ;
- ③ ignorable error to the oracle MLE:  $\|\tilde{\beta}_B - \hat{\beta}_B^*\|_2 = \mathcal{O}_p(1/N_B)$ .



# Implementation

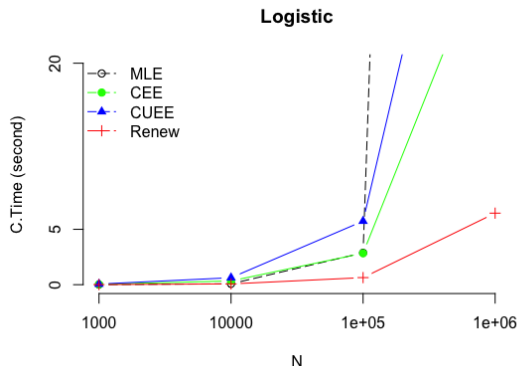
- No historical data storage;
- Real-time estimation and inference.

algorithmcol.pdf

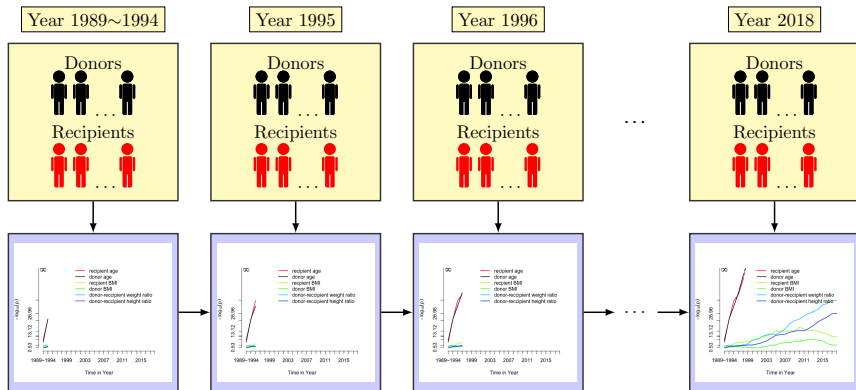


# Data Example

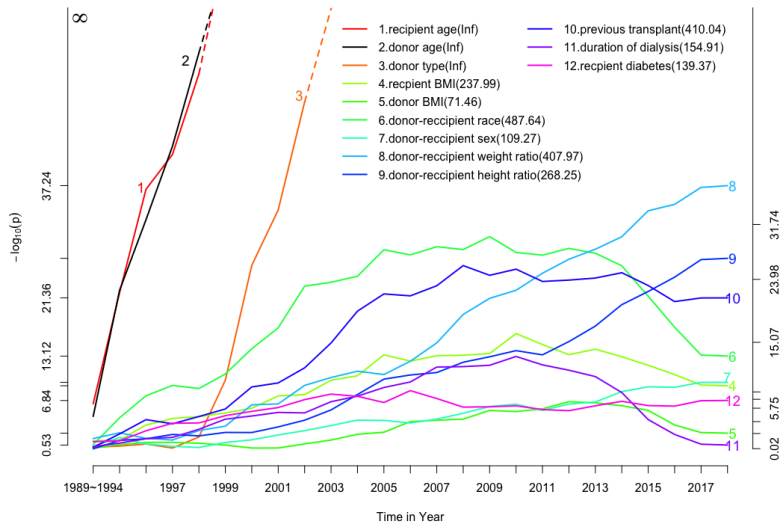
- Total sample size: 244,614 samples ( $N > 10^5$ );
- Model: logistic regression;
- Outcome: 1 for graft failure at the 5-th year after transplantation and 0 otherwise;
- Objective: update  $p$ -values based on yearly-arrived data batches.



# Data Example



# Data Example



# Conclusions and Future Directions

## ① Conclusions:

- A **novel renewable estimation method** in the generalized linear model;
- **Ignorable error** to the oracle MLE;
- Our method is **computationally efficient**;
- **Real-time regression analysis** without strong regularity conditions on relative scale of data batch size  $n_b$  and total number of batches  $B$ .

## ② Future Directions:

- Model homogeneity is assumed, and addressing partial homogeneous model is an undergoing direction.

## ③ Questions?



- Schifano, Elizabeth D, Wu, Jing, Wang, Chun, Yan, Jun, & Chen, Ming-Hui. 2016. Online updating of statistical inference in the big data setting. *Technometrics*, **58**(3), 393–403.
- Toulis, Panos, & Airoidi, Edoardo M. 2015. Scalable estimation strategies based on stochastic approximations: classical results and new insights. *Statistics and computing*, **25**(4), 781–795.

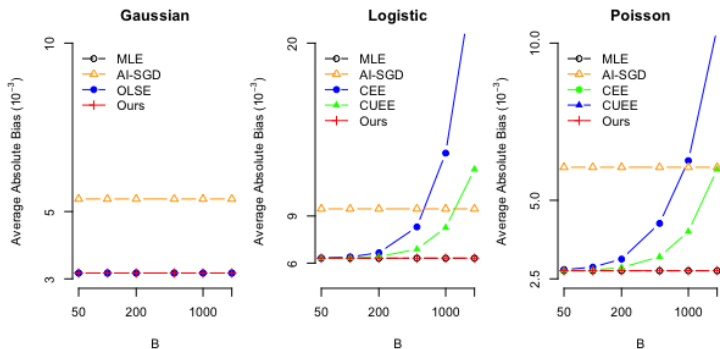


- $B$  data streams with  $N_B$  independent observations  $(y_i, \mathbf{x}_i)$  from GLMs with mean model  $\mathbb{E}(y_i | \mathbf{x}_i) = g(\mathbf{x}_i^T \boldsymbol{\beta}_0)$ ;
- $\boldsymbol{\beta}_0 = (0.2, -0.2, 0.2, -0.2, 0.2)'$ ,  $\mathbf{x}_{i[2:5]} \sim \mathcal{N}_4(\mathbf{0}, \mathbf{V}_4)$  where  $\mathbf{V}_4$  is a compound symmetry covariance matrix with correlation 0.5.
- **Scenarios:**
  - (1) divide a fixed  $N_B = 10^5$  evenly into  $B = 50, 100, 200, 500, 1000, 2000$  data streams;
  - (2) fix sub-sample size  $n_b = 100$ ,  $b = 1, \dots, B$  and let  $B$  increase from  $10^5$  to  $10^6$ .





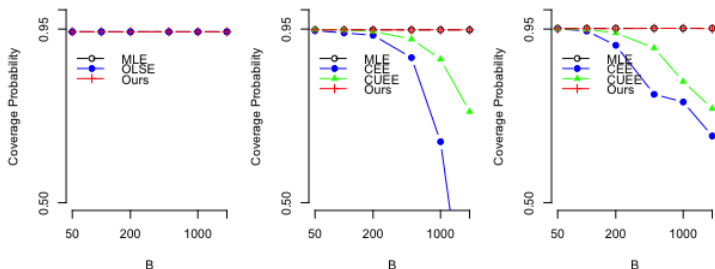
# Simulation 1 - Bias



- Average bias in AI-SGD is the largest;
- In linear model, OLSE = Renew = MLE;
- In logistic or Poisson model, as  $B$  increases,  $CEE > CUEE \gg \text{Renew} \approx \text{MLE}$ .



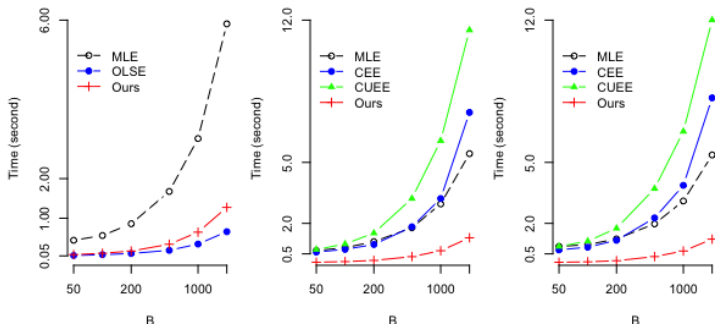
# Simulation 1 - Coverage probability



- In linear model, OLSE = Renew = MLE;
- In logistic or Poisson model,  $95\% \approx \text{MLE} \approx \text{Renew} \gg \text{CUEE} > \text{CEE} \approx 0$  as  $b$  increases;
- AI-SGD does not apply here.



# Simulation 1 - Computation time



- Computation time = data loading time + processing time;
- In linear model,  $MLE > Renew > OLSE$ ;
- In logistic or Poisson model,  $CUEE > CEE > MLE > Renew$ .



## Simulation 2

Fix  $n_b = 100$ , increase  $B$  from  $10^3$  to  $10^4$ :

Logistic Model ( $p = 5, n_b = 100$ )					
$B = 10^3, N_B = 10^5$					
	MLE	Renew	AI-SGD	CEE	CUEE
A.bias( $\times 10^{-3}$ )	5.023	5.020	<b>8.581</b>	<b>12.696</b>	6.218
Std. err.( $\times 10^{-3}$ )	6.538	6.539	-	6.576	6.581
Cov. prob.	0.956	0.958	-	<b>0.556</b>	<b>0.898</b>
C.Time (seconds)	2.951	0.663	-	3.022	5.763
R.Time (seconds)	0.471	0.487	0.165	2.846	5.587
$B = 10^4, N_B = 10^6$					
A.bias( $\times 10^{-3}$ )	1.626	1.626	<b>8.581</b>	<b>12.978</b>	<b>4.157</b>
Std. err.( $\times 10^{-3}$ )	2.067	2.067	-	2.136	2.081
Cov. prob.	0.958	0.956	-	<b>0</b>	<b>0.644</b>
C.Time (seconds)	323.566	8.149	-	34.483	68.642
R.Time (seconds)	98.362	5.257	1.119	31.594	65.750

- Increasing  $N_B$  helped to reduce bias in MLE and Renew, but not in AI-SGD;
- CEE and CUEE failed because  $B$  is much larger than  $n_b$ .

