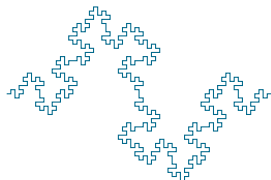


# Statistical analysis of stochastic gradient methods for generalized linear models

Panos Toulis  
ptoulis@fas.harvard.edu  
Harvard University,  
Department of Statistics



Joint work with Jason Rennie (Google) & Edoardo M Airoldi (Harvard)

# INTRODUCTION

## FOCUS OF THIS WORK

- ▶ Stochastic Gradient Descent (SGD): computationally attractive, slow convergence, great empirical performance.
- ▶ As a *statistical* estimation method, still not well-understood.
- ▶ Say  $\theta_n^{\text{sgd}}$  is the output of SGD given observed data  $\mathbf{y}$  generated by a model with true parameters  $\theta^*$ ; we ask:
  - ▶ What is the bias  $\mathbb{E}(\theta_n^{\text{sgd}} - \theta^*)$ ?
  - ▶ What is the variance  $\text{Var}(\theta_n^{\text{sgd}})$ ?
  - ▶ How to optimally set the learning rate?
- ▶ We also want to consider SGD with *explicit* and *implicit* updates and provide a meaningful comparison.

# PROBLEM AND NOTATION

## GLM FAMILY

At every time step indexed by  $n$ , assume the following data generating process:

- ▶  $\mathbf{x}_n \sim G$  sampled iid,  $\in \mathbb{R}^p$  (features)
- ▶  $y_n \sim f(y_n; \mathbf{x}_n^\top \boldsymbol{\theta}^*, \psi) \in \mathbb{R}$  (outcome)

such that  $f(\cdot)$  is a density in the *exponential family* and

$$\mathbb{E}(y_n | \mathbf{x}_n) = h(\mathbf{x}_n^\top \boldsymbol{\theta}^*) \quad (1)$$

where  $h(\cdot)$  is the (monotone) *link* function,  $\boldsymbol{\theta}^* \in \mathbb{R}^p$  are the unknown model parameters, and  $\psi > 0$  is the dispersion parameter.

Our **goal** is to *estimate*  $\boldsymbol{\theta}^*$  given observations  $(y_i, \mathbf{x}_i)$ , indexed by  $i = 1, \dots, N$ .

# KNOWN PROPERTIES OF EXPONENTIAL FAMILY/GLMs.

Let  $\ell(\boldsymbol{\theta}; y_n, \mathbf{x}_n)$  be the log-likelihood of  $\boldsymbol{\theta}$  for observation  $(y_n, \mathbf{x}_n)$ . The following hold for a GLM:

$$\nabla \ell(\boldsymbol{\theta}; y_n, \mathbf{x}_n) = \frac{1}{\psi} (y_n - h(\mathbf{x}_n^\top \boldsymbol{\theta})) \mathbf{x}_n \quad (2)$$

$$\mathcal{I}(\boldsymbol{\theta}) = -\mathbb{E}(\nabla \nabla \ell(\boldsymbol{\theta}; y_n, \mathbf{x}_n)) = \frac{1}{\psi} \mathbb{E}(h'(\mathbf{x}_n^\top \boldsymbol{\theta}) \mathbf{x}_n \mathbf{x}_n^\top) \quad (3)$$

- ▶ Equation (2) gives the gradient of the log-likelihood in terms of “**observed - expected**” of the sufficient statistics.
- ▶ Equation (3) gives the Fisher information matrix. The value  $\mathcal{I}(\boldsymbol{\theta}^*)^{-1}$  is the theoretically **best** possible variance we can achieve if we try to (unbiasedly) estimate  $\boldsymbol{\theta}^*$ .

# ITERATIVE ESTIMATION OF GLMs

## SGD PROCEDURES

The *explicit* SGD updates are given by

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + a_n (y_n - h(\mathbf{x}_n^\top \boldsymbol{\theta}_{n-1})) \mathbf{x}_n \quad (4)$$

The *implicit* SGD updates are given by

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + a_n (y_n - h(\mathbf{x}_n^\top \boldsymbol{\theta}_n)) \mathbf{x}_n \quad (5)$$

**Remark #1.** After  $n$  iterations, both procedures provide an *estimate* of  $\boldsymbol{\theta}^*$ :

- ▶  $\boldsymbol{\theta}_n^{\text{sgd}}$  of the explicit updates is the *explicit SGD estimator* of  $\boldsymbol{\theta}^*$ .
- ▶ Similarly,  $\boldsymbol{\theta}_n^{\text{im}}$  is the *implicit SGD estimator* of  $\boldsymbol{\theta}^*$ .

# ITERATIVE ESTIMATION OF GLMS

## SGD PROCEDURES

The *explicit* SGD updates are given by

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + a_n (y_n - h(\mathbf{x}_n^T \boldsymbol{\theta}_{n-1})) \mathbf{x}_n \quad (6)$$

The *implicit* SGD updates are given by

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + a_n (y_n - h(\mathbf{x}_n^T \boldsymbol{\theta}_n)) \mathbf{x}_n \quad (7)$$

**Remark #2.** Implicit methods are less well-studied.

- ▶ Similar ideas have been used in numerical analysis e.g., Crank-Nicolson method (1947), to solve PDE.
- ▶ The NLMS algorithm in signal processing (Nagumo & Noda, 1967) uses an implicit method for linear regression.
- ▶ Recent interest due to stability of the method (Kivinen, 1996), (Kivinen et al., 2006), (Kulis & Bartlett, 2010).

# EXAMPLE 1 : NORMAL MODEL

Assume

$$\mathbf{x}_n \sim G \text{ and } \mathbf{X}_n = \mathbf{x}_n \mathbf{x}_n^\top \text{ (possibly random)}$$

$$y_n \sim \mathcal{N}(\mathbf{x}_n^\top \boldsymbol{\theta}^*, \sigma^2)$$

The explicit SGD update is,

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + a_n \underbrace{(y_n - \mathbf{x}_n^\top \boldsymbol{\theta}_{n-1})}_{\text{observed - expected}} \mathbf{x}_n = (\mathbf{I} - a_n \mathbf{X}_n) \boldsymbol{\theta}_{n-1} + a_n y_n \mathbf{x}_n$$

The implicit SGD can be derived analytically as,

$$\boldsymbol{\theta}_n = (\mathbf{I} + a_n \mathbf{X}_n)^{-1} (\boldsymbol{\theta}_{n-1} + a_n y_n \mathbf{x}_n)$$

The latter update is known as the “Normalized Least Mean Squares” filter.

## EXAMPLE 2: POISSON REGRESSION

Assume

$$y_n \sim \text{Pois}(e^{\mathbf{x}_n^\top \boldsymbol{\theta}^*})$$

The explicit SGD is,

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + a_n \underbrace{(y_n - e^{\mathbf{x}_n^\top \boldsymbol{\theta}_{n-1}})}_{\text{observed - expected}} \mathbf{x}_n$$

The implicit SGD is,

$$\boldsymbol{\theta}_n = \boldsymbol{\theta}_{n-1} + a_n (y_n - e^{\mathbf{x}_n^\top \boldsymbol{\theta}_n}) \mathbf{x}_n$$

- ▶ Explicit SGD is problematic in this model because of the non-linearity of the score function.
- ▶ The implicit update cannot be derived analytically. However, we will show how it can be *computed efficiently*.



## EXAMPLE 2: POISSON REGRESSION (CONTINUED)

### INSTABILITY OF **explicit** SGD

$$\theta_n = \theta_{n-1} + a_n(y_n - e^{\mathbf{x}_n^\top \theta_{n-1}})\mathbf{x}_n$$

- ▶ Assume one-dimensional case, for which  $\theta_0 = 0, x_0 = x_1 = a_0 = 1, y_1 = 1001$ , then the update is:

$$\theta_1 = 0 + 1(1001 - 1)1 = 1000$$

- ▶ In the next iteration assume that  $y_2 = 500, a_1 = 0.5$ , then:

$$\theta_2 = 1000 + 0.5(500 - e^{1000})1 = -\infty$$

The problem is that the starting point ( $\theta_0 = 0$ ) was far away from the observation  $y_1 = 1001$  and the learning rate was not small enough to prevent a large update (misspecification).

## EXAMPLE 2: POISSON REGRESSION (CONTINUED)

### STABILITY OF **implicit** SGD

Implicit update for Poisson regression model

$$\theta_n = \theta_{n-1} + a_n(y_n - e^{\mathbf{x}_n^\top \theta_n})\mathbf{x}_n$$

- The implicit update would solve:

$$\theta_1 = 0 + 1(1001 - \underbrace{e^{\theta_1}}_{\text{implicit}})1$$

and so  $\theta_1 \approx \log(1001)$ .

In case of misspecification, the implicit update will try to “overfit” on the current data point but does not diverge like explicit SGD.

# BIAS

Let  $\mathbf{J} = \psi \mathcal{I}(\boldsymbol{\theta}^*)$ . It holds:

$$\|\mathbb{E}(\boldsymbol{\theta}_n^{\text{sgd}}) - \boldsymbol{\theta}^*\| \propto \prod_i^n \|(I - a_i \mathbf{J})\|$$

$$\|\mathbb{E}(\boldsymbol{\theta}_n^{\text{im}}) - \boldsymbol{\theta}^*\| \propto \prod_i^n \|(I + a_i \mathbf{J})^{-1}\|$$

- For large enough  $n$ ,

$$\|(I - a_n \mathbf{J})\| \leq \|(I + a_n \mathbf{J})^{-1}\|$$

and so the explicit SGD is converging *faster*. However, the asymptotic rates are equal.

- The spectra of  $(I - \epsilon \mathbf{J})$  and  $(I + \epsilon \mathbf{J})^{-1}$  are crucial for their stability properties.

# STABILITY

Recall that  $\mathbf{J} = \psi \mathcal{I}(\boldsymbol{\theta}^*) \geq 0$ . For  $\epsilon > 0$ ,

- ▶ The spectrum of  $(\mathbf{I} - \epsilon \mathbf{J})$  is equal to  $(1 - \epsilon \lambda_i(\mathbf{J}))$ . For stability, we thus need  $|1 - \epsilon \lambda_i(\mathbf{J})| < 1$ . The explicit updates are *conditionally stable*.
- ▶ The spectrum of  $(\mathbf{I} + \epsilon \mathbf{J})^{-1}$  is  $(1 + \epsilon \lambda_i(\mathbf{J}))^{-1} < 1$  and thus, the implicit updates are *unconditionally stable*.

# ASYMPTOTIC VARIANCE

Let  $\text{Var}(\boldsymbol{\theta}_n^{\text{sgd}}) = \boldsymbol{\Sigma}_n^{\text{sgd}}$  and  $\text{Var}(\boldsymbol{\theta}_n^{\text{im}}) = \boldsymbol{\Sigma}_n^{\text{im}}$ .

**Theorem 3.** *The asymptotic variance of the explicit SGD estimator is,*

$$n \cdot \boldsymbol{\Sigma}_n^{\text{sgd}} \rightarrow \alpha^2 \psi^2 (2\alpha\psi \mathbf{I}(\boldsymbol{\theta}^*) - \mathbf{I})^{-1} \mathbf{I}(\boldsymbol{\theta}^*) \quad (8)$$

*The asymptotic variance of the implicit SGD estimator is,*

$$n \cdot \boldsymbol{\Sigma}_n^{\text{im}} \rightarrow \alpha^2 \psi^2 (2\alpha\psi \mathbf{I}(\boldsymbol{\theta}^*) - \mathbf{I})^{-1} \mathbf{I}(\boldsymbol{\theta}^*) \quad (9)$$

*Assuming convergence, both SGD methods have the same asymptotic efficiency.*

# ASYMPTOTIC VARIANCE (CONTINUED)

- ▶ It holds that

$$a^2\psi^2(2a\psi\mathbf{I}(\boldsymbol{\theta}^*) - \mathbf{I})^{-1}\mathbf{I}(\boldsymbol{\theta}^*) \geq \underbrace{\mathbf{I}(\boldsymbol{\theta}^*)^{-1}}_{\text{MLE(theoretically optimal)}}, \forall \alpha, \psi > 0$$

Not surprisingly, both methods incur information loss.

- ▶ However, it is possible to **optimize** for the learning rate e.g., minimize the trace of the asymptotic variance:

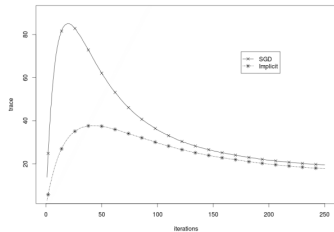
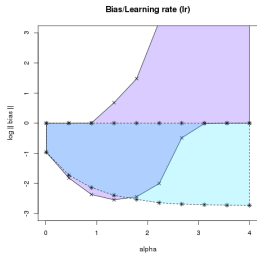
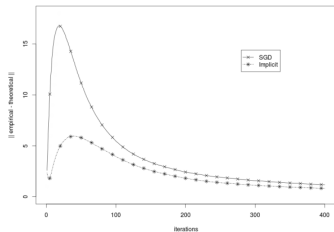
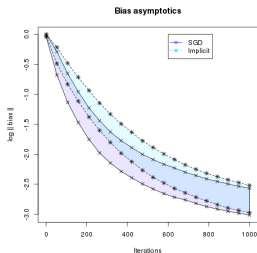
$$\hat{\alpha} = \arg \min_a \sum_i \frac{a^2 \lambda_i}{2a\lambda_i - 1} \quad (10)$$

for  $\lambda_i = \text{spectrum}(\psi\mathbf{I}(\boldsymbol{\theta}^*))$ .

$\theta^* = \mathbf{1} \in \mathbb{R}^{20}$ ,  $\mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{V})$ ,  $\lambda_i(\mathbf{V}) \in [0.2, 1]$ ,  $y_n \sim \mathcal{N}(\mathbf{x}_n^\top \theta^*, 1)$ .

(left) = log-bias (2.5-97.5% percentiles),

(right) trace of empirical covariance matrix over 2,000 samples



# EXPERIMENTS ON CLASSIFICATION TASKS (SVM)

## TESTING OUTSIDE GLM FAMILY

Table : Test errors of explicit and implicit SGD methods on the RCV1 dataset benchmark. Training times are roughly comparable. Best scores, for a particular loss and regularization, are bolded.

LOSS		REGULARIZATION ( $\lambda$ )		
		1E-5	1E-7	1E-12
HINGE	SGD	<b>4.65%</b>	<b>3.57%</b>	4.85%
	IMPLICIT	4.68%	3.6%	<b>3.46%</b>
LOG	SGD	5.23%	3.87%	5.42%
	IMPLICIT	<b>4.28%</b>	<b>3.69%</b>	<b>4.01%</b>



# SUMMARY

- ▶ Exact statistical analysis of SGD is possible in GLMs, both for explicit and implicit updates.
- ▶ Helps in optimizing for the learning rate.
- ▶ Implicit updates compare favorably to explicit ones:
  - ▶ They are easy to implement (Theorem 1).
  - ▶ They have the same asymptotic performance (bias and variance, Theorems 2,3).
  - ▶ They are unconditionally stable, and thus more robust to misspecification.
  - ▶ Vanilla implementation performs on par with standard SGD on large-scale optimization tasks.

Thank you!

(poster #T-76).