Influence Function and NLP Application

Jillian Fisher¹, Lang Liu¹, Krishna Pillutla², Yejin Choi^{3,4}, and Zaid Harchaoui¹

¹Department of Statistics, University of Washington, ²Google Research, ³Paul G. Allen School of Computer Science & Engineering, University of Washington, ⁴Allen Institute for Artificial Intelligence

UNIVERSITY of WASHINGTON

Motivation

We rely on models for important tasks...



But how do we know we can trust these models?



Outline

Background: Influential Points

- Statistical Finite Bound
- Computational Bound
 - Experiment: Is there always meaning?
- Most Influential Subset
 - Experiment: Are all statistics a lie?!
- NLP Connection
 - Will there be influence in your future?

Background: Notation

Setting: Consider $\theta \in \Theta$, constructed from i.i.d sample $z = \{(x_i, y_i)\}_{i=1}^n$





Background: Influence Function

Consider a prediction problem,



 $\mathbf W$ university of washington

Background: Notation



Outline

- Background: Influential Points
- Statistical Finite Bound
- Computational Bound
 - Experiment: Is there always meaning?
- Most Influential Subset
 - Experiment: Are all statistics a lie?!
- NLP Connection
 - Experiment: Will there be influence in your future?

Assumptions: Pseudo Self-Concordance

1. Simple definition if we assume *linear prediction models (i.e.* $\ell(\theta) = \ell(Y, X^T \theta)$). We consider $\ell(\theta)$ is pseudo self-concordant if

 $\left| \nabla^{3} \ell(z,\theta) \right| \leq \nabla^{2} \ell(z,\theta)$

Prevents $\nabla^2 \ell(z, \theta)$ from changing too quickly with θ

Consequence: Spectral Approximation of the Hessian $\frac{1}{2}H(\theta') \le H(\theta) \le 2H(\theta') \text{ for } \theta \text{ close to } \theta'$



Black curve: population function f(x); colored dot: reference point x_i ; colored dashed curve: quadratic approximation at the corresponding reference point $Q(x; x_i)$.

Assumptions

2. Normalized gradient $H(\theta_{\star})^{-1/2} \nabla \ell(Z, \theta_{\star})$ at θ_{\star} is sub-Gaussian with parameter K_1

Since $\mathbb{E}[\nabla \ell(Z, \theta_{\star})] = 0$, then Assumption 2 gives a high prob. bound on $\|\nabla \ell(Z, \theta_{\star})\|_{H_{\star}}^{-1}$

3. There exist $K_2 > 0$ such that the **standardized Hessian at** θ_{\star} **satisfies a Bernstein condition** with parameter K_2

Moreover,

$$\sigma_H^2 := \|\operatorname{Var}(H(\theta_\star)^{-1/2} \nabla^2 \mathscr{C}(Z, \theta_\star) H(\theta_\star)^{-1/2})\|_2 \text{ is finite.}$$

Assumption 3 gives spectral concentration

 $(1/2)H(\theta) \prec H_n(\theta) \prec 2H(\theta)$

Generalized Linear Models satisfy these assumptions

Results: Statistical Bound

Theorem 1. Suppose the assumptions¹ hold and

$$n \ge C\left(\frac{p}{\mu_{\star}}\log\frac{1}{\delta} + \log\frac{p}{\delta}\right)$$

where $\mu_{\star} = \lambda_{\min}(H(\theta_{\star}))$.

Then, with probability at least $1 - \delta$, we have $\frac{1}{4}H(\theta_{\star}) \leq H_n(\theta_n) \leq 3H(\theta_{\star})$ and

$$\|I_n(z) - I(z)\|_{H_{\star}}^2 \le C \frac{p_{\star}^2}{\mu_{\star} n} \text{poly} \log\left(\frac{p}{\delta}\right)$$

- Only logarithmic dependence on p (dim. of param.)
- p_{\star} is the degrees of freedom (model misspecification)
- Rate of 1/n

1. Assumptions met by Generalized Linear Models

Experiment: Simulation

Simulation

- $x \sim N(0,1)$
- Linear (Ridge) Regression Logistic Regression
- *X-axis: Training Sample Size Y-axis: Difference in empirical vs. population IF*

Results

- •See 1/n of our bound observed
- •Straight line in log-log scale
- Hard to approximate classification population



Experiment: Real Dataset

- **Real Dataset**
- Cash Transfer
 - X: Socio-economic covariates
 - Y: Total consumption (regression)
- Oregon Medicaid
 - X: Health-related covariates
 - 1. Y: Estimate overall health (classification)
 - 2. Y: Number of good days (regression)
- X-axis: Training Sample Size Y-axis: Difference in empirical vs. population IF

Results

- •See 1/n of our bound observed
- •Straight line in log-log scale
- Hard to approximate classification population



Outline

- Background: Influential Points
- Statistical Finite Bound

Computational Bound

- Experiment: Is there always meaning?
- Most Influential Subset
 - Experiment: Are all statistics a lie?!
- NLP Connection
 - Will there be influence in your future?

Computational Challenge

Second derivative (p x p) p = dim of parameter

Cook and Weisberg Formula

$$I_n(z) = -H_n(\theta_n)^{-1} \nabla \ell(z, \theta_n)$$

Can't be computed for large values of p

Instead use iterative algorithms to approximately minimize

$$g_n(\mu) \coloneqq \frac{1}{2} \langle \mu, H_n(\theta_n) \mu \rangle + \langle \nabla \ell(z, \theta_n), \mu \rangle$$

Algorithms

- > Conjugate Gradient (CG)
- > Stochastic Gradient Descent (SGD)
- > Stochastic Variance Reduced Gradient (SVRG)
- > Arnoldi Low Rank



Result: Computational Bound

Proposition 1. Consider the setting of Theorem 1, and let \mathscr{G} denote the event under which its conclusions hold. Let $\hat{I}_n(\theta)$ be an estimate of $I_n(\theta)$ that satisfies $\mathbb{E}_{Z_{1:n}} \left[\left\| \hat{I}_n(z) - I_n(z) \right\|_{H_n(\theta_n)}^2 \right] \le \epsilon.$ Then $\mathbb{E}_{\mathscr{G}} \left[\left\| \hat{I}_n(z) - I(z) \right\|_{H(\theta_\star)}^2 \right] \le 8\epsilon + C \frac{p_\star^2}{\mu_\star n} \text{poly } \log \frac{p}{\delta}$ • Using an ϵ -approximate minimizer of the empirical influence approximation

- Translating approx. error in $H_n(\theta_n)$ -norm to the H_{\star} -norm under \mathcal{G} (Theorem 1)
- **Total Error** under $O(\epsilon)$ is $O(n(\epsilon)T(\epsilon))$





Result: Computational Bound

Proposition 1. Consider the setting of Theorem 1, and let \mathscr{G} denote the event under which its conclusions hold. Let $\hat{I}_n(\theta)$ be an estimate of $I_n(\theta)$ that satisfies

$$\mathbb{E}_{Z_{1:n}}\left[\left\|\hat{I}_n(z) - I_n(z)\right\|_{H_n(\theta_n)}^2\right] \leq \epsilon.$$

Then

$$\mathbb{E}_{\mathscr{G}}\left[\left\|\hat{I}_{n}(z) - I_{n}(z)\right\|_{H(\theta_{\star})}^{2}\right] \leq 8\epsilon + C \frac{p_{\star}^{2}}{\mu_{\star}n} \text{poly} \log \frac{p}{\delta}$$

Example: Stochastic Variance Reduction Gradient (SVRG)

- Requires $T_n(\epsilon) = C(n + \kappa_n) \log\left(\frac{\kappa_n \|u_0 u_\star\|_{H_n(\theta_n)}}{\epsilon}\right)$ iterations to return an ϵ -approximate minimizer.
- Each iteration requires *n* Hessian-vector products
- To make statistical error to be smaller than ϵ , $n \ge n(\epsilon) = \tilde{O}\left(\frac{p_{\star}^2}{\mu_{\star}\epsilon}\right)$ from **Theorem 1**
- Total error under $O(\epsilon)$ is $O(n(\epsilon)T(\epsilon))$ by Proposition 1

 κ_{\star} is the condition number $\Delta_{\star} = \|I_n(z)\|_{H(\theta_{\star})}^2$

Total Error =
$$\kappa_{\star} \left(1 + \frac{p_{\star}}{\epsilon} \right) \log \left(\frac{\kappa_{\star} \Delta_{\star}}{\epsilon} \right)$$



Result: Global Bounds

Method	Computational Error	Total Error
Conjugate Gradient	$n\sqrt{\kappa_n}$	$\frac{\kappa_{\star}^{3/2} p_{\star}^2}{\epsilon}$
Stochastic Gradient Descent	$\frac{\sigma_n^2}{\epsilon} + \kappa_n$	$\frac{\sigma_{\star}^2}{\epsilon} + \kappa_{\star}$
Stochastic Variance Reduction Gradient	$(n + \kappa_n)$	$\kappa_{\star} \left(1 + \frac{p_{\star}^2}{\epsilon} \right)$
Accelerated Stochastic Variance Reduction Gradient	$(n + \sqrt{n\kappa_n})$	$\kappa_{\star} \left(\sqrt{\frac{p_{\star}^2}{\epsilon}} + \frac{p_{\star}^2}{\epsilon} \right)$



Experiment: Is there always meaning?

Question Answering

- Input: question
- Response: factual correct answer
 - X= What country did The Laughing Cow originate?
 - Y= France
- zsRE dataset (Levy et. al., 2017)/BART-base model
- Average over 5 data points

Question Answering



Experiment: Is there always meaning?

Text Continuation

- Input: Start of paragraph
- Response: 10 tokens continuation
 - X = "The interchange is considered by Popular Mechanics to be one of the...",
 - y = "World's 18 Strangest Roadways because of its height"
- WikiText (Merity et. al., 2017)/GPT2
- Averaged over 5 data points



Outline

- Background: Influential Points
- Statistical Finite Bound
- Computational Bound
 - Experiment: Is there always meaning?
- Most Influential Subset
 - Experiment: Are all statistics a lie?!
- NLP Connection
 - Experiment: Will there be influence in your future?



Most Influential Subset

•Given an $\alpha \in (0,1)$, and a test function $h : \mathbb{R}^p \to \mathbb{R}$

<u>Most influential set</u> is the subset of data (size at most αn), which when removed leads to largest increase in the test function.



MIS: Definition

Most Influential Subset

•Given n $\alpha \in (0,1)$, and a test function $h : \mathbb{R}^p \to \mathbb{R}$

<u>Most influential subset</u> is the subset of data (size at most αn), which when removed leads to largest increase in the test function.

Mathematically,



MIS: Definition

First-order Taylor expansion: $f(x) = f(a) + \frac{f'(a)}{(x - a)}$



Instead Broderick et al. (2020) use first-order Taylor expansion in $h(\theta_{n,w})$ around w = 1

$$h(\theta_{n,w}) \approx h(\theta_n, \frac{1}{n}) + \left\langle \nabla_w h(\theta_n, w) \big|_{w=\frac{1}{n}}, w - \frac{1}{n} \right\rangle$$

1 is a vector of all 1's

MIS: Definition

Instead Broderick et al. (2020) use linear approximation

$$h(\theta_{n,w}) \approx h(\theta_n) + \left\langle w - \frac{\mathbf{1}_n}{n}, \nabla_w h(\theta_n, w) \big|_{w = \mathbf{1}_n/n} \right\rangle$$

Which leads to the influence of the most influential subset,

$$I_{\alpha,n}(h) := \max_{w \in W_{\alpha}} \left\langle w, \nabla_{w} h(\theta_{n}, w) \big|_{w = \mathbf{1}_{n}/n} \right\rangle$$

Which can be simplified using the implicit function theorem and the chain rule to a closed form

$$I_{\alpha,n}(h) := \max_{w \in W_{\alpha}} \sum_{i=1}^{n} w_i v_i$$

Where $v_i = -\langle \nabla h(\theta_n), H_n(\theta_n)^{-1} \nabla \ell(Z_i, \theta_n) \rangle$

$$I_n(Z_i, \theta_n)$$

Main Results: Most Influential Subset

Theorem 2. Suppose the added assumptions hold and the sample size n satisfies the condition in Theorem 1. Then with probability at least $1-\delta$

$$\frac{\left(I_{\alpha,n}(h)-I_{\alpha}(h)\right)^{2}}{\left(1-\alpha\right)^{2}} \leq \frac{C_{M_{1},M_{2},M_{1}',M_{2}'}}{\left(1-\alpha\right)^{2}} \frac{R^{2}p_{\star}}{\mu_{\star}n} \log \frac{n \vee p}{\delta}$$

- Only logarithmic dependence on *p*
- p_{\star} is affine-invariant

$$\frac{1}{-rat}$$

$$-$$
 rate n

Experiment: Real Dataset

Oregon Medicaid study (Finkestein et al., 2012)

- Lottery from 90,000 people to sign up for Medicaid = randomization into treatment (Medicaid) and control (no Medicaid) groups
- Measured outcomes one year after treatment group received Medicaid ($\underline{n} \approx 22,000$)

 $y = \beta_0 + \beta_1 LOTTERY + \beta_2 X_{covariates}$

• Test function, h(x) : is β_1 significant?

/	le thic case	Original estimate	Target change	Refit estimate	Observations dropped	
	surprising?	2m 0.133 (0.026)*	Sign change Significance change Significant sign change	-0.006 (0.025) 0.044 (0.026) -0.043 (0.024)	$275 = 1.18\% \\ 162 = 0.69\% \\ 381 = 1.63\%$	On average
	Health not for 12	2m 0.099 (0.018)*	Sign change Significance change Significant sign change	$\begin{array}{c} -0.003 \ (0.015) \\ 0.027 \ (0.016) \\ -0.030 \ (0.015)^* \end{array}$	155 = 0.66% 100 = 0.43% 219 = 0.94%	the removal of <.05% of
(Do you think large da (like the ones we use	ntasets 113 (0.023)* in LLM	Sign change Significance change Significant sign change	-0.006 (0.022) 0.039 (0.022) -0.049 (0.022)*	197 = 0.84% 106 = 0.45% 291 = 1.24%	the data
	pretraining) are this ser change?	nsitive to 317 (0.563)*	Sign change Significance change Significant sign change	-0.023 (0.535) 1.078 (0.558) -1.009 (0.521)	73 = 0.33% 10 = 0.05% 144 = 0.66%	Of significance II
	Not bad days physica	al 12m 585 (0.606)*	Sign change Significance change Significant sign change	$\begin{array}{c} -0.040 (0.577) \\ 1.131 (0.597) \\ -1.141 (0.566)^* \end{array}$	$87 = 0.41\% \\ 20 = 0.09\% \\ 164 = 0.77\%$	
	Not bad days menta	l 12m 2.082 (0.640)*	Sign change Significance change Significant sign change	$\begin{array}{c} -0.062 \ (0.607) \\ 1.171 \ (0.625) \\ -1.201 \ (0.594)^* \end{array}$	$\begin{array}{l} 123 = 0.57\% \\ 42 = 0.19\% \\ 212 = 0.98\% \end{array}$	A7

Experiment: Most Influential Subset

MIS (Question Answering)

- 4 different test points (questions/answer)
- $\alpha = 0.05, .0.1$ (size of subset)

Alpha = 0.05

Arnoldi method was used to approximate influence

Why is there such variety in slopes?

Alpha = 0.01



Downward trend —> similar to influence of 1 datapoint

Experiment: Most Influential Subset



Outline

- Background: Influential Points
- Statistical Finite Bound
- Computational Bound
 - Discussion: Is there always meaning?
- Most Influential Subset
 - Discussion: Are all statistics a lie?!
- NLP Connection
 - Will there be influence in your future?

Related Work in NLP

Influential points

- Leave one out training (data point importance)
- Saliency maps (token importance)
- Self-influence (Bejan et al., 2023)
- Influence function for NLP.... Still in development

Machine Unlearning

- Quark reinforcement learning (Lu et al., 2022)
- SISA Training (Kumar et al., 2022)

Explaining Black Box Predictions and Unveiling Data Artifacts through Influence Functions

Xiaochuang Han, Byron C. Wallace, Yulia Tsvetkov

INFLUENCE FUNCTIONS IN DEEP LEARNING ARE FRAGILE

Samyadeep Basu, Phillip Pope & Soheil Feizi Department of Computer Science University of Maryland, College Park {sbasul2, pepope, sfeizi}@cs.umd.edu Influence Functions Do Not Seem to Predict Usefulness in NLP Transfer Learning

AUG 27 2020

Author(s): Vid Kocijan and Samuel R. Bowman Publication date: August 27 2020 Reviewer: Alex Wang Editor: Kyunghyun Cho

Conclusion and Future Extensions

Conclusion

- Presented statistical and computational guarantees for influence functions for generalized linear models
- •Established the statistical consistency of most influential subsets method (Broderick et at., 2020) together with non-asymptotic bounds
- •Illustrated our results on simulated and real datasets

Future Extension

- Non-convex/Non-smooth penalized M-estimation
- Application for toxicity/bias removal in NLP

Thank You!

Full Paper





References

R. Cook and S. Weisberg. Residuals and influence in regression.New York: Chapman and Hall, New York: Chapman Hall, 1982.

T. Broderick, R. Giordano, and R. Meager. An Automatic Finite-Sample Robustness Metric: When Can Dropping a Little Data Make a Big Difference? arXiv Preprint, 2020

D. M. Ostrovskii and F. Bach. Finite-sample analysis of M-estimators using self-concordance. Electronic Journal of Statistics, 15(1), 2021



Appendix Slides

UNIVERSITY of WASHINGTON



Algorithms: Conjugate Gradient

Algorithm 1 Conjugate Gradient Method to Compute the Influence Function

Input: vector v, batch Hessian vector product oracle $HVP_n(u) = H_n(\theta_n)u$, number of iterations T 1: $u_0 = 0, r_0 = -v - HVP_n(u_0), d_0 = r_0$ 2: for t = 0, ..., T - 1 do 3: $\alpha_t = \frac{d_t^\top r_t}{d_t^\top HVP_n(d_t)}$ 4: $u_{t+1} = u_t + \alpha_t d_t$ 5: $r_{t+1} = -v - HVP_n(u_{t+1})$ 6: $\beta_t = \frac{r_{t+1}^\top r_{t+1}}{r_t^\top r_t}$ 7: $d_{t+1} = r_{t+1} + \beta_t d_t$ 8: return u_T

Algorithms: Stochastic Gradient Descent

Algorithm 2 Stochastic Gradient Descent Method to Compute the Influence Function

Input: vector v, Hessian vector product oracle $HVP(i, u) = \nabla^2 \ell(z_i, \theta_n) u$, number of iterations T, learning rate γ 1: $u_0 = 0$ 2: for t = 0, ..., T - 1 do 3: Sample $i_t \sim Unif([n])$ 4: $u_{t+1} = u_t - \gamma(HVP(i_t, u_t) + v)$ 5: return u_T

UNIVERSITY of WASHINGTON

Algorithms: Stochastic Variance Reduction Gradient

Algorithm 4 Stochastic Variance Reduced Gradient Method to Compute the Influence Function

Input: vector v, Hessian vector product oracle HVP $(i, u) = \nabla^2 \ell(z_i, \theta_n) u$, number of epochs S, number of iterations per epoch T, learning rate γ 1: $u_T^{(0)} = 0$ 2: for s = 1, 2, ..., S do 3: $u_0^{(s)} = u_T^{(s-1)}$ 4: $\tilde{u}_0^{(s)} = \frac{1}{n} \sum_{i=1}^n \text{HVP}(u_0^{(s)}) - v$ 5: for t = 0, ..., T - 1 do 6: Sample $i_t \sim \text{Unif}([n])$ 7: $u_{t+1}^{(s)} = u_t^{(s)} - \gamma(\text{HVP}(i_t, u_t^{(s)}) - \text{HVP}(i_t, u_0^{(s)}) + \tilde{u}_0^{(s)})$ 8: return $u_T^{(S)}$

Algorithms: Arnoldi

Algorithm 5 Arnoldi Method to Compute the Influence Function (Schioppa et al., 2022)

Input: vector v, test function h, initial guess u_0 , batch Hessian vector product oracle HVP_n $(u) = H_n(\theta_n)u$, number of top eigenvalues k, number of iterations T**Output:** An estimate of $\langle \nabla h(\theta), H_n(\theta_n)^{-1}v \rangle$ 1: Obtain $\Lambda, G = \text{ARNOLDI}(u_0, T, k)$ \triangleright Cache the results for future calls 2: return $\langle G \nabla h(\theta), \Lambda^{-1} G v \rangle$ 3: procedure ARNOLDI (u_0, T, k) $w_0 = 1 = u_0 / \|u_0\|_2$ 4: $A = \mathbf{0}_{T+1 \times T}$ 5: for t = 1, ..., T do 6: Set $u_t = \text{HVP}_n(w_t) - \sum_{j=1}^t \langle u_t, w_j \rangle w_j$ 7: Set $A_{i,t} = \langle u_t, w_i \rangle$ for j = 1, ..., t and $A_{t+1,t} = ||u_t||_2$ 8: Update $w_{t+1} = u_t / ||u_t||$ 9: Set $\tilde{A} = A[1:T, :] \in \mathbb{R}^{T \times T}$ (discard the last row) 10: Compute an eigenvalue decomposition $\tilde{A} = \sum_{j=1}^{T} \lambda_j e_j e_j^{\top}$ with λ_j 's in descending order 11: Define $G: \mathbb{R}^p \to \mathbb{R}^k$ as the operator $Gu = (\langle u, W^\top e_1 \rangle, \cdots, \langle u, W^\top e_k \rangle)$, where $W = (w_1^\top; \cdots; w_T^\top) \in \mathbb{R}^{T \times p}$ 12: **return** diagonal matrix $\Lambda = \text{Diag}(\lambda_1, \dots, \lambda_k)$ and the operator G 13:

UNIVERSITY of WASHINGTON

Computational Results: CG

Proposition 1. Consider the setting of Theorem 1, and let **3** denote the event under which its

conclusions hold. Let $\hat{I}_n(\theta)$ be an estimate of $I_n(\theta)$ that satisfies $\mathbb{E}\left[\left\|\hat{I}_n(z) - I_n(z)\right\|_{H_n(\theta_n)}^2\right] Z_{1:n}\right] \leq \epsilon$.

Then

$$\mathbb{E}\left[\left\|\hat{I}_{n}(z)-I_{n}(z)\right\|_{H_{\star}}^{2}\right] \leq 8\epsilon + C\frac{R^{2}p_{\star}^{2}}{\mu_{\star}n}\log^{3}\left(\frac{p}{\delta}\right)$$

Example: Conjugate Gradient

- Requires $T_n(\epsilon) \coloneqq \sqrt{k_n} \log(\left\|I_n(z)\right\|_{H_n(\theta_n)}^2 / \epsilon)$ iterations to return an ϵ -approximate minimizer.
- Each iteration requires *n* Hessian-vector products

To make statistical error to be smaller than ϵ , $n \ge n(\epsilon) = \tilde{O}\left(\frac{R^2 p_{\star}^2}{\mu_{\star} \epsilon}\right)$ Total error under $O(\epsilon)$ is $O(n(\epsilon)T(\epsilon))$ – by Proposition 1

Experiment: Most Influential Subset

MIS Test Questions

- 1. What position did Víctor Vázquez Solsona play? midfielder
- 2. The nationality of Jean-Louis Laya was what? French
- 3. Where is Venera 9 found? Venus
- 4. Who set the standards for ISO 3166-1 alpha-2? International Organization for Standardization
- *5.* In which language Nintendo La Rivista Ufficiale monthly football magazine reporting? *Italian*



Experiment: Most Influential Subset

- 1. What position did Víctor Vázquez Solsona play? Midfielder
- 2. Was Goldmoon male or female? Female
- 3. Where is Venera 9 found? Venus
- 4. In which language Nintendo La Rivista Ufficiale monthly football magazine reporting? Italian

