

JSM 2022

Model Editing in Language Models Using Influence Functions

Jillian Fisher
08/10/2022



Liwei Jiang



Krishna Pillutla



Swabha Swayamdipta



Yejin Choi



Zaid Harchaoui

Language models trained on un-curated data cause issues...

On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?

En e
Unive
S
Angeli
a
Unive
S

UW NEWS

NEWS RELEASES

March 10, 2021

Large computer language models carry environmental, social risks

[Jackson Holtz](#)

ru*
ai.org
USA
hitchell
@gmail.com

The New York Times Magazine

A.I. Is Mastering Language. Should We Trust What It Says?

OpenAI's GPT-3 and other neural nets can now write original

Challenges in Detoxifying Language




Johannes Welbl* Amelia Glaese* Jonathan Uesato*
John Mellor* Lisa Anne Hendricks Kiran
Pushmeet Kohli Ben Coppin Po-Sen Huang*
DeepMind
{welbl, glamia, juesato, sdathath, johnme, posenhuang}@deepmind.com

PUBLISHED ON AUGUST 19, 2021 IN OPINIONS

With A Rush To Create Larger Language Models, Are We Beating Their Purpose

Language Models trained on large, uncurated, static datasets from the Web encode hegemonic views that are harmful to marginalised populations.

By [Kumar Gandharv](#)

 Adrian Yijie Xu
Jun 7, 2020 · 11 min read · Member-only ·  

Language Models and Fake News: the Democratization of Propaganda

Deepfaking information with the OpenAI GPT- 3 model




03-16-22

A new Stanford study still has a bias problem

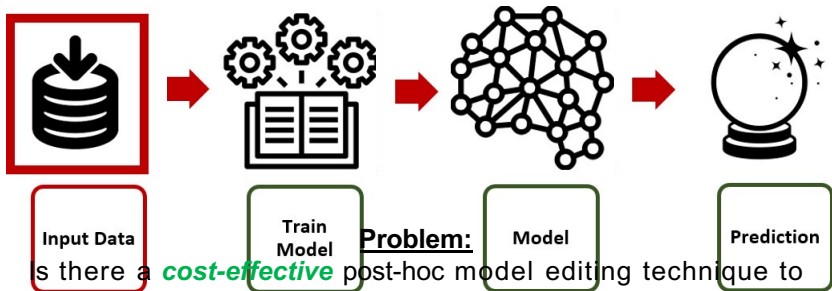
Stanford's annual report gives a snapshot of both the good and bad of its AI forms.

The Machine
Making sense of AI

Language models that can search the web hold promise — but also raise concerns

Kyle Wiggers
[@Kyle_L_Wiggers](#)
March 19, 2022 6:21 AM
  

Motivation



Is there a *cost-effective* post-hoc model editing technique to remove (or edit) behaviors in trained language models, *without re-training* the model?

Solution:

Influence Functions

Outline

1. Theory

- Background
- Influence Functions

2. Useful for Language Models?

- Implementation Challenges

3. Influence Functions for Model Editing

- Methods
- Experiment
- Results

Background: Empirical Risk Minimizer

Set-Up:

- $\theta \in \Theta$, constructed from i.i.d sample $z = \{(x_i, y_i)\}_{i=1}^n$
- loss function $L(z_i, \theta)$

Empirical Risk Minimizer (ERM):

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n L(z_i, \theta)$$

Maximum Likelihood Estimation (MLE):

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n -\log p_{\theta}(y|x)$$

Theory: Influence Functions

Functional = function that maps a **distribution** to a **real number**.

Example: the sample mean, $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, as a functional, $T(F_n)$,

$$\bar{x} = \int x dF_n = T(F_n)$$

where F_n is the empirical distribution.

Why Important?

Easily examine estimator under different distributions

Example: “Contaminated Distribution”

$$F_\epsilon = (1 - \epsilon)F(x) + \epsilon G(x) \text{ for } \epsilon \in [0,1]$$



Dog

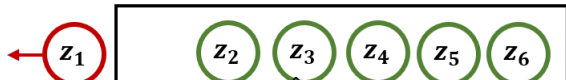


Cat

Theory: Influence Functions

Consider a prediction problem,

Training Set



Using a linear approximation¹ of $\hat{\theta}_{\epsilon, z}$ around $\epsilon \equiv 0$, this leads to the following approximation $z_i = (x_i, y_i) \in X \times Y$

Derivative of contamination distribution

Parameter of Interest

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \frac{1}{6} \sum_{i=1}^6 L(z_i, \theta)$$

Parameters with all training data (known)

$$\hat{\theta}_{\epsilon, z} \in \arg \min_{\theta \in \Theta} (1 - \epsilon) \frac{1}{6} \sum_{i=1}^6 L(z_i, \theta) + \epsilon L(z_1, \theta)$$

Change in parameters due to removing 1 training point (unknown)

$$\hat{\theta}_{\epsilon, z} - \hat{\theta} \approx \frac{d \hat{\theta}_{\epsilon, z}}{d \epsilon}$$

As a Functional

$$\hat{\theta}_{\epsilon, z} = T \left((1 - \epsilon) F_{\theta} + \epsilon \delta_{z_1} \right) \rightarrow \text{“Contaminated Distribution”}$$

¹ $f(x) \approx f(a) + f'(a)(x - a)$

Theory: Influence Functions

Definition 2: Influence Function

The influence functions, $IF(x; T, F_\theta)$ of T , is

Derivative of
"Contaminated Distribution"

$$IF(x; T, F_\theta) = \lim_{\epsilon \rightarrow 0} \frac{T((1 - \epsilon)F_\theta + \epsilon\delta_x) - T(F_\theta)}{\epsilon}$$

where δ_x is the probability measure that places a point mass 1 at x .

Theory: Influence Functions

Cook and Weisberg (1982) classical result

$$IF(x; T, F_{\hat{\theta}}) = -H_{\hat{\theta}}^{-1} \nabla_{\theta} L(z, \hat{\theta})$$

where¹ $H_{\theta} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \nabla_{\theta}^2 L(z_i, \theta)$ and assumed positive definite.

Parameter of Interest

To remove a training examples $z \rightarrow \epsilon = \frac{1}{n}$

$$\hat{\theta}_{-z} \approx \hat{\theta} - \frac{1}{n} H_{\hat{\theta}}^{-1} \nabla_{\theta} L(z, \hat{\theta})$$

¹ H_{θ} is also called the observed Fisher's Information Matrix

² For simplicity allow $\theta_{-z} = \theta_{\vec{1}, z}$

Outline

1. Theory
 - Background: ERM and GLM
 - Influence Functions
- 2. Useful for Language Models?**
 - Implementation Challenges
3. Influence Functions for Model Editing
 - Methods
 - Experiment
 - Results

Theory: Implementation Challenges

$$\hat{\theta}_{-z} \approx \hat{\theta} - \frac{1}{n} H_{\hat{\theta}}^{-1} \nabla_{\theta} L(z, \hat{\theta})$$

Challenge #1:

Computing the **inverse Hessian** of the empirical risk alone with large parameters = COSTLY

Solution:

Hessian Vector Products (HVP) to efficiently approximate,

$$s \approx H_{\hat{\theta}}^{-1} \nabla_{\theta} L(z, \hat{\theta})$$

Theory: Implementation Challenges

HVP:

The approximation of $\hat{\theta}_z$ still requires the calculation of

$$s \approx H_{\hat{\theta}}^{-1} \nabla_{\theta} L(z, \hat{\theta})$$

Solution:

Frame as solving a linear system

$$Hx = v \rightarrow H^{-1}v = x$$

Use first or second order stochastic methods

- Conjugate Gradient Descent
- CURVEBALL (Henriques et al. [2019](#))

Method: Conjugate Gradient Descent

Algorithm 1: Conjugate Gradient Descent

for $t = 1, \dots, T$ **do**

$$x_0 \leftarrow 0, r_0 = b - Ax_0, p_0 = r_0$$

for $k = 1, \dots, K$ **do**

$$\alpha_k = \frac{p_k^T r_k}{p_k^T A p_k}$$

$$x_{k+1} = x_k + \alpha_k p_k$$

$$r_{k+1} = b - Ax_{k+1}$$

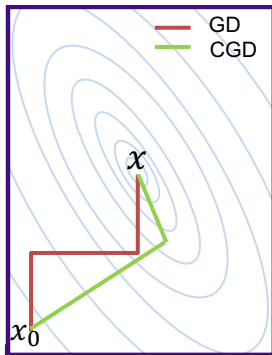
$$p_{k+1} = r_{k+1} - \sum_{i < k} \frac{p_i^T A r_k}{A p_i} p_i$$

$$\theta_{t+1} = \theta_t + x_K$$

Two "For"
loops
needed

New direction
conjugate to all
previous directions

Conjugate Gradient Descent



- First-order method

➤ $y_t = x_t - \eta \nabla f(x_t)$

- Accelerates convergence rate of gradient descent by not repeating a direction

Method: CURVEBALL [Henriques et al. 2019]

Algorithm 2: CURVEBALL [Henriques et al., 2019]

$x_0 \leftarrow 0$

for $k = 1, \dots, K$ **do**

$$\nabla_x = \hat{A}x_k + J$$

$$x_{k+1} = \rho x_k + \alpha \nabla_x$$

$$\theta_{k+1} = \theta_k + x_{k+1}$$

Only one
"For" loop

Second-order newton method

Uses momentum parameter

- Second-order method
 - $y_t = x_t - \eta[\nabla^2 f(x_t)]^{-1}\nabla f(x_t)$
- Adds curvature term (second order)
- Faster and less memory use
 - Reuses previous direction
 - Interweaves search direction and parameter update (only 1 "For" loop)
- Specifically tailored for deep-learning-scale stochastic optimization problems

Implementation Challenges

Challenge #2:

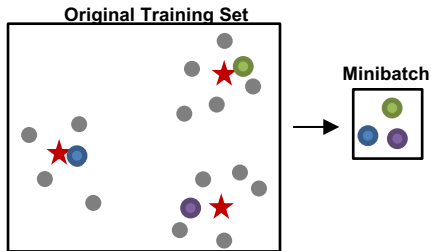
The HVP, requires the loss over the original data set which can be computationally expensive.

Solution:

Use clustering techniques to find a mini-batch that is representative of the original data set.

Steps:

1. Start with all features from the Original data set
2. Perform PCA to reduce dimensionality
3. Perform clustering on the features
4. Close data point to each cluster center to include in the minibatch



Outline

1. Theory
 - Background: ERM and GLM
 - Influence Functions
2. Useful for Language Models?
 - Implementation Challenges
- 3. Influence Functions for Model Editing**
 - Methods
 - Experiment
 - Results

Model Editing: Goal

Task: Model Editing

Goal: Develop a cost-effective, post-hoc model editing technique to edit knowledge in a trained language models.

Data: Zero Shot Relation Extraction (zsRE)

Dataset (D)

Edited (D_{ed})

Q: What is the name of Another Side of Bob Dylan's record label?
A: Capitol Records

Non-edited (D_{non-ed})

Q: What country did The Laughing Cow originate?
A: France

Forget (D_F)

Q: What is the name of Another Side of Bob Dylan's record label?
A: Capitol Records

Remember (D_R)

Q: What is the name of Another Side of Bob Dylan's record label?
A: Colombia Records

Model Editing: Baseline

Notation

- θ_0 = parameters of the original model
- $\hat{\theta}$ = parameters of the edited model
- T = number of updates
- L_R = loss over DR
- η = learn. rate
- P = Distribution of reg. subset under θ_0
- Q = Distribution of reg. subset under $\hat{\theta}$

Algorithm 1: Baseline

Initialize $\hat{\theta}_0 = \theta_0$.

for $t = 0, \dots, T - 1$ **do**

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \eta \nabla L_r(\hat{\theta}_t) + D_{KL}(P || Q)$$

end for

Newton
step

Regularization
term

Model Editing: Influence Function Method

Notation

- T = number of updates for step 1
- S = number of updates for step 2
- L_F = loss over D_F

Algorithm 2 Influence Function Method

Step 1: Forgetting

Initialize $\hat{\theta}_0 = \theta_0$.

For $t = 0, \dots, T - 1$ do

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \underbrace{H_{\theta_0}^{-1} \nabla_{\theta} L_F(\hat{\theta}_t)}_{\text{Approximation using IF}}$$

end for

Step 2: Remembering

Initialize $\hat{\theta}_0 = \hat{\theta}_T$.

for $s = 0, \dots, S - 1$ do

$$\hat{\theta}_{s+1} = \hat{\theta}_s - \underbrace{\eta \nabla L_r(\hat{\theta}_s) + D_{KL}(P||Q)}_{\text{Baseline}}$$

end for

Baseline

Experimentation: Details

Details:

- # of Non-Edits: 10,000
- # of Edits: 40
- Repetition: 10
- Epochs Baseline: 60
- Epochs Step 1 (forg.): 4
- Epochs Step 2 (rem.): 56

Metrics:

1. *Reliability*: made edits successfully

Accuracy over D_R (increase)

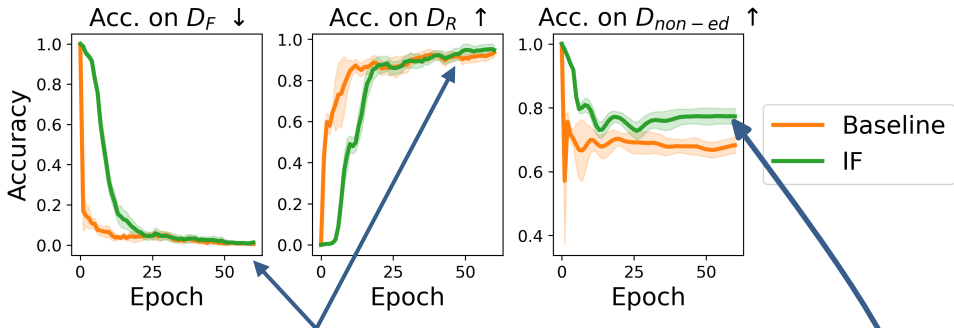
Accuracy over D_F (decrease)

2. *Generality*: did not change non-edited input/out

Accuracy over D_{non-ed} (increase)

Results: Accuracy

Results: zsRE 40 Edits



- *Reliability*: Accuracy over D_F and D_R are similar
- *Generality*: Retention of Non-edited is better with IF method
 - “Forgetting” is more targeted

Conclusion

- Theoretically influence functions good solution to problem
- However, approximation techniques → fragile or inaccurate
- Promising results in application to model editing indicates more experimentation



Thank you!

Questions?

References

- 
- Agarwal, Naman, Brian Bullins, and Elad Hazan (2016). “Second-order stochastic optimization in linear time”. In: *stat* 1050, p. 15.
- 
- Basu, Samyadeep, Philip Pope, and Soheil Feizi (2020). “Influence functions in deep learning are fragile”. In: *arXiv preprint arXiv:2006.14651*.
- 
- Cook, R Dennis and Sanford Weisberg (1982). *Residuals and influence in regression*. New York: Chapman and Hall.
- 
- De Cao, Nicola, Wilker Aziz, and Ivan Titov (2021). “Editing factual knowledge in language models”. In: *arXiv preprint arXiv:2104.08164*.
- 
- Han, Xiaochuang, Byron C. Wallace, and Yulia Tsvetkov (2020). “Explaining Black Box Predictions and Unveiling Data Artifacts through Influence Functions”. In: *arXiv: [2005.06676](https://arxiv.org/abs/2005.06676) [cs.CL]*.
- 
- Henriques, João F et al. (2019). “Small steps and giant leaps: Minimal newton solvers for deep learning”. In: *Proceedings of*