

Scaling Up Influence Functions

Andrea Schioppa, Artem Sokolov, Polina Zablotskaia, David Vilar {arischio, artemsok, polinaz, vilar}@google.com

Introduction

Influence Functions

Questions that Influence Methods address:

- Which training examples are most responsible for the prediction at a given test point?
- Which training examples advocate / oppose the prediction at a test point?
- Which part of the training data is corrupt / of lower quality?
- Which training examples are easier / harder to memorize?

The formal answer [1]:

$$\mathcal{I}_{\text{true}}(x, z) = L(z|\Theta : x \notin \mathcal{D}) - L(z|\Theta : x \in \mathcal{D})$$

x is a training example; z is a test example; Θ is the parameters of a model; L is the loss function used in training; \mathcal{D} is the corpus of all training examples.

However removing each training example and retraining the model is not feasible!

- Instead [1] compute the effect of the perturbation on the current optimum:

$$\mathcal{I}_H(x, z) = \langle \nabla_{\Theta} L(z), H^{-1} \nabla_{\Theta} L(x) \rangle$$

This is faster, but for DL models still not feasible due to the inverse of the Hessian $H = \nabla_{\Theta}^2 L$

- [2] (LISSA) suggest iterative inversion of Hessian with:

$$H_r^{-1} v = v + (I - H) H_{r-1}^{-1} v$$

Hessian vector products can be estimated in $O(b \cdot p)$ -time, where $b = |D|, p = |\Theta|$. The total $O(b \cdot p \cdot r)$ -time complexity will be incurred at every x in whose influence on z we are interested.

Method is slow: either small models, reduction of parameters to the last layer, small datasets or a subset of points.

Algorithm 1 Arnoldi

```

1: procedure ARNOLDI( $v, n$ )
2:  $\triangleright$  Build orthonormal basis for the Krylov subspaces  $K_n$ .
3:  $w_0 \leftarrow \frac{v}{\|v\|_2}$ 
4:  $A_{l,m} \leftarrow 0$  for  $0 \leq l \leq n$  and  $0 \leq m < n$ 
5: for  $l \leftarrow 1, n$  do
6:  $w_l \leftarrow H \cdot w_{l-1}$   $\triangleright$  HVP in fwd-over-rev mode
7: Set  $A_{l,j} = \langle w_l, w_j \rangle$  for  $j < l$ 
8:  $w_l \leftarrow w_l - \sum_{j < l} A_{l,j} w_j$   $\triangleright$  Orthogonalize
9:  $A_{l+1,l} \leftarrow \|w_l\|_2, w_l \leftarrow \frac{w_l}{\|w_l\|_2}$ 
10: return:  $A, \{w_i\}$ 
11: procedure DISTILL( $A, \{w_i\}, \bar{p}$ )
12:  $\triangleright$  Distill  $A, \{w_i\}$  to its top- $\bar{p}$  eigenvalues.
13: Discard the last row of  $A$  and the last  $w_n$ 
14: Obtain  $A$ 's eigenvalues  $\{\lambda_i\}$  and eigenvectors  $\{e_i\}$ 
15: Set  $\{\lambda'_i\}$  to the  $\bar{p}$ -largest (in absolute value) of  $\{\lambda_i\}$ 
16: Set  $\{e'_i\}$  to the corresponding eigenvectors
17: Set  $G$  to the projection onto the spans  $\{e'_i\}$  in  $\{w_i\}$ -basis
18: return:  $\{\lambda'_i\}, G$ 
19: procedure INFLUENCE( $x, z, n, \bar{p}$ )
20:  $\triangleright$  Influence of  $x$  on  $z$  with  $n$  iterations and top- $\bar{p}$  eigenvalues.
21:  $v \leftarrow \mathcal{N}(0, 1)$ 
22:  $A, \{w_i\} = \text{ARNOLDI}(v, n)$ 
23:  $\{\lambda'_i\}, G = \text{DISTILL}(A, \{w_i\}, \bar{p})$   $\triangleright$  Executed once and cached
24:  $g_x \leftarrow G \cdot \nabla_{\Theta} L(x)$   $\triangleright$  fwd JVP for  $x$  over  $G$ -rows
25:  $g_z \leftarrow G \cdot \nabla_{\Theta} L(z)$   $\triangleright$  fwd JVP for  $z$  over  $G$ -rows
26:  $g_x \leftarrow g_x / \|\lambda'_i\|$   $\triangleright$  Multiply with diagonalized  $H^{-1}$ 
27: return:  $\langle g_x, g_z \rangle$ 
    
```

Method

Scalable Influence Functions

$O(p)$ complexity is the major bottleneck for efficient implementation of IFs.

We can use Random Projection $G \in \mathbb{R}^{p \times \bar{p}}$ to create a subspace of size \bar{p} such that $\mathbb{E}[G^T G] = I$

$$\mathcal{I}_{\tilde{H}}(x, z) = \langle G \nabla_{\Theta} L(z), \tilde{H}^{-1} G \nabla_{\Theta} L(x) \rangle$$

$$\tilde{H} = G \cdot H \cdot G^T$$

However the image of G is not H -invariant. To solve this we propose using the standard the Arnoldi iteration [3] technique of building an approximately H -invariant subspace and constructing the n -th order Krylov subspace:

$$K_n(H; v) = \text{Span}\{v, H v, H^2 v, \dots, H^n v\}$$

The procedure additionally builds an orthonormal basis for Krylov subspace so that the diagonalization of the restriction \tilde{H} to H yields an approximation of the largest (in absolute value) eigenvalues of H and of the corresponding eigenvectors.

The matrix H gets replaced with now diagonal \tilde{H} simplifying the matrix inversion appearing in the definition of \mathcal{I}_H and dispensing with the expensive LISSA [2] procedure.

Results

MNIST Benchmark

Corrupted MNIST -- Small Model (Figure 1 & 2)

- [2] considers only a 10% subset of data and a CNN with 3k parameters
- We evaluate ability to recall synthetically mislabeled examples [4].
- IFs outperform gradient methods and RandProj does not capture correct eigenvalues.

Corrupted MNIST -- Large Model (Table 1)

- CNN with 800k parameters..
- LISSA is prohibitively slow; Arnoldi & RandProj scale well; no advantage in using curvature information.

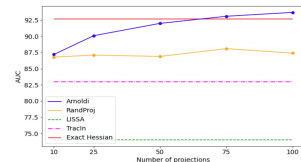


Figure 1

Method	\bar{p}	T, secs	AUC	AP
LISSA, $r = 10$	100	4900	98.9	95.0
LISSA, $r = 100$	100	32300	98.8	94.8
TraceIn [1]	-	5	98.7	94.0
TraceIn [100]	-	42	99.7	98.7
RandProj	10	0.2	97.2	87.7
RandProj	100	1.9	98.6	93.9
Arnoldi	10	0.2	95.0	84.0
Arnoldi	100	1.9	98.2	92.9

Table 1

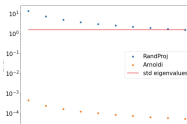


Figure 2

Conclusion

- We proposed a new way of calculating influence scores of [2] for large DNNs by approximate diagonalization of their Hessians and avoiding re-estimating them on every training example.
- We demonstrated finding influential or noisy examples in datasets of up to 100M training examples and models with up to 300M parameters.
- We showed IFs can be either superior or match random projection and gradient-based approaches, when measured on the benchmark of retrieving synthetically mislabeled data with self-influence scores.
- We will make our code publicly available.

References

[1] Cook, R. D.; and Weisberg, S. 1980. Characterizations of an Empirical Influence Function for Detecting Influential Cases in Regression. *Technometrics* 22(4): 495–508.
 [2] Koh, P. W.; and Liang, P. 2017. Understanding Black-box Predictions via Influence Functions. In *ICML*.
 [3] Arnoldi, W. E. 1951. The principle of minimized iterations in the solution of the matrix eigenvalue problem. *Quarterly of Applied Mathematics* 9(1): 17–29.
 [4] Pruthi, G.; Liu, F.; Sundararajan, M.; and Kale, S. 2020. Estimating Training Data Influence by Tracing Gradient Descent. In *NeurIPS*.

Machine Translation

Synthetic Noise (Table 2)

- A model on WMT17 German -> English.
- Noisy data where references were permuted (with 6% probability).
- Subsetting the layers leads to worse performance.

Natural Noise [Quality] (Table 3)

- Bad quality Paracrawl (100M pairs, ~1B words) from WMT'18 data cleaning task..
- Applying Arnoldi on top of Microsoft winning submission Scores to get 50% of the selected clean data gives 37.2.
- Method has high throughput: 2.2M examples per hour using TPUv2 (32 cores)

Layers	Method	\bar{p}	AUC	AP
all layers	RandProj	10	85.6	31.3
	RandProj	20	85.2	28.0
	Arnoldi	10	92.0	47.8
	Arnoldi	20	93.8	54.4
last 3 decoder layers	RandProj	10	80.6	23.5
	RandProj	20	83.0	25.8
	Arnoldi	10	82.0	28.1
	Arnoldi	20	83.7	28.5

Table 2

Method	% selected	BLEU@10k	BLEU@200k
None	100	9.9	17.8
Pre-cleaning	14	11.7	22.6
RandProj	1	7.7	8.7
Arnoldi	1	17.8	19.6
RandProj	5	18.7	27.9
Arnoldi	5	24.0	30.3
RandProj	10	21.3	28.6
Arnoldi	10	25.0	30.8

Table 3