

Weekly meetings

If Influence Functions are the Answer, Then What is the Question?

Setting

Consider a prediction task (regression problem) with:

- Input space \mathcal{X} ;
- Output space \mathcal{Y} ;
- Training set $\mathcal{D}^n = \{z_i\}_{i=1}^n$ where $z_i = (x_i, y_i)$ for all $i = 1, \dots, n$;
 $X = (x_1, \dots, x_n)$, $Y = (y_1, \dots, y_n)$;
- Parameter $\theta \in \Theta := \mathbb{R}^d$;
- $f(\theta; x)$ estimator of $\mathcal{Y} \mid \mathcal{X}$;
- Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ (e.g., $\ell(y', y) \mapsto \|y' - y\|^2$).

We aim to minimize the training error (empirical risk):

$$L(\theta; \mathcal{D}^n) = \frac{1}{n} \sum_{i=1}^n \ell(f(\theta; x_i), y_i).$$

Call

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} L(\theta; \mathcal{D}^n).$$

How different is it from

$$\hat{\theta}_{\epsilon, -z} = \arg \min_{\theta \in \Theta} (L(\theta; \mathcal{D}^n) - \epsilon \ell(f(\theta; x), y))$$

when $\epsilon = 1/n$ and $z = z_i$ for some $i \in [n]$?

To answer this question, we can re-train the whole model on $\mathcal{D}^n \setminus \{z\}$ (Leave-One-Out method), or use influence functions.

In the context where the influence functions are well defined, they are a powerful tool. However, when applied to multi-layer perceptrons, for example, their capability of approximating the effect of the LOO decreases drastically. This paper presents a new point of view: IFs do not approximate the LOO retraining, but instead the effect of another method they present, called PBRF.

Influence functions

Definition 1. Given $(\bar{x}, \bar{y}) = \bar{z} \in \mathcal{D}^n$, the *influence loss difference* relative to \bar{z} is

$$\mathcal{Q}(\bar{z}; \hat{\theta}) = \frac{d}{d\varepsilon} \left[\ell(f(\hat{\theta}_{\varepsilon, -z}), \bar{y}) \right] \Big|_{\varepsilon=1/n}.$$

Interpretation: It measures how much the training error changes when the data point \bar{z} is removed.

Definition 2. Given $(\bar{x}, \bar{y}) = \bar{z} \in \mathcal{D}^n$, the *influence function* relative to \bar{z} is

$$\mathcal{I}(\bar{z}; \hat{\theta}) = \lim_{\varepsilon \rightarrow 0} \frac{\hat{\theta}_{\varepsilon+1/n, -\bar{z}} - \hat{\theta}_{1/n, -\bar{z}}}{\varepsilon}.$$

Interpretation: It represents the direction in which the optimal parameter moves when the training objective is perturbed by removing \bar{z} .

Assuming that L is strongly convex in θ , we can rewrite the previous quantities in the closed forms:

$$\mathcal{I}(z; \hat{\theta}) = H_{\hat{\theta}}^{-1} \nabla \ell(f(\hat{\theta}; x), y), \quad \mathcal{Q}(z; \hat{\theta}) = \nabla \ell(f(\hat{\theta}; x), y)^\top H_{\hat{\theta}}^{-1} \nabla \ell(f(\hat{\theta}; x), y),$$

where $H_{\hat{\theta}}$ is the Hessian of L at $\hat{\theta}$.

Remark. This part is not really clear: referring to the original reference [2], the derivation should be as follow, with our notation.

Let $L_i(\theta) = (L(\hat{\theta}; \mathcal{D}^n) - 1/n \ell(f(\hat{\theta}; x_i), y_i))$. Assuming that L is twice differentiable, we can use Taylor near $\hat{\theta}$:

$$L_i(\theta) = L_i(\hat{\theta}) + (\theta - \hat{\theta})^\top L'_i(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^\top L''_i(\hat{\theta})(\theta - \hat{\theta}) + O(\|\theta - \hat{\theta}\|^3)$$

Then, by minimizing the difference $L_i(\theta) - L_i(\hat{\theta})$ for $\theta \neq \hat{\theta}$, we get:

$$\theta = \hat{\theta} - (HL_i(\hat{\theta}))^{-1} JL_i(\hat{\theta}).$$

From this expression is also clear that the Influence Function estimator represents the parameter after one step of the Newton algorithm starting from $\hat{\theta}$ trying to get to $\hat{\theta}_{1/n, -z_i}$.

Unfortunately, there is some problems with these formulae:

1. the strong convexity is essential. If at a minimum point H has any 0-eigenvalue, we cannot invert it.
2. Even when L is strongly convex, the problem is not trivial because computing the inverse of the Hessian and the matrix-vector product are heavy computations.

Solution 1. For iHVP, there exist efficient approximations requiring $O(nd)$ flops instead of $O(n^3)$.

Solution 2. Change point of view: Influence functions are not approximators of LOO retraining, but instead of the proximal Bregman response function (PBRF).

Response functions

We define the response function as:

$$\hat{r}_z(\varepsilon) = \arg \min_{\theta \in \Theta} (L(\theta; \mathcal{D}^n) - \varepsilon \ell(f(\theta; x), y)).$$

Observe that $\hat{r}_z(\varepsilon) = \hat{\theta}_{\varepsilon, -z}$ and $\hat{r}_z(0) = \hat{\theta}$.

Since \hat{r} is differentiable at 0, we can expand with Taylor at first order and we get:

$$\hat{r}_{z, \text{lin}}(1/n) \approx \hat{\theta} + \frac{1}{n} H_{\hat{\theta}}^{-1} \nabla \ell(f(\hat{\theta}; x), y).$$

We require $H_{\hat{\theta}}$ positive definite to invert it; thus $\hat{\theta}$ must be a minimizer.

To compute influence functions for MLPs, we can approximate $H_{\hat{\theta}}$ using the Gauss–Newton Hessian (GNH) and add damping:

$$\mathcal{I}^\dagger(z; \hat{\theta}) = \left(J_{y, \hat{\theta}}^\top H_{\ell, \hat{\theta}} J_{y, \hat{\theta}} + \lambda I \right)^{-1} \nabla \ell(f(\hat{\theta}; x), y),$$

where $J_{y, \hat{\theta}}$ is the Jacobian of $F(\theta) = (f(\theta; x_1), \dots, f(\theta; x_n))$ and $H_{\ell, \hat{\theta}}$ is the hessian of $\ell(f(\theta; x), y)$ in $\theta = \hat{\theta}$.

Note that the damped GNH is always positive definite, as long as $H_{\hat{\theta}}$ is SPD.

We can get the previous formula by linearizing the response function of the regularized loss:

$$\hat{r}_{z, \text{damp}}(\varepsilon) = \arg \min_{\theta \in \Theta} \left(L(\theta; \mathcal{D}^n) - \varepsilon \ell(f(\theta; x), y) + \frac{\lambda}{2} \|\theta - \hat{\theta}\|^2 \right),$$

$$\hat{r}_{z, \text{damp}, \text{lin}}(1/n) \approx \hat{\theta} + \frac{1}{n} \mathcal{I}^\dagger(z; \hat{\theta}).$$

Another issue is that, in practice, the parameter we utilise is not a minimizer of L . Thus, we want to consider a risk for which the early-stopped point θ^s is optimal:

$$\mathcal{L}(\theta; \theta^s, \mathcal{D}^n) = \frac{1}{n} \sum_{i=1}^n D_{\ell(i)}(f(\theta; x_i), f(\theta^s; x_i)),$$

with $D_{\ell(i)}$ being the Bregman difference:

$$D_{\ell(i)}(y, y') = \ell(y, y_i) - \ell(y', y_i) - \nabla_1 \ell(y', y_i)^\top (y - y').$$

Intuitively, this quantity measures the difference between the evaluation of ℓ on y and the first order Taylor expansion of ℓ around y' computed on y .

Observation. This quantity is always non-negative as long as ℓ is convex. To exemplify, choose $\ell(y, y') = \|y - y'\|^2/2$. This yields:

$$D_{\ell(i)}(y, y') = \|y - y_i\|^2/2 - \|y' - y_i\|^2/2 - \langle \nabla \|y' - y_i\|^2/2, y - y_i \rangle = \|y - y'\|^2/2.$$

Consequently, we can define the *Proximal Bregman Response Function* (PBRF) as:

$$r_{z,\text{damp}}^b(\varepsilon) = \arg \min_{\theta \in \Theta} \left(\mathcal{L}(\theta; \theta^s, \mathcal{D}^n) - \varepsilon \ell(f(\theta; x), y) + \frac{\lambda}{2} \|\theta - \theta^s\|^2 \right).$$

The interesting property of this object is that the linearized PBRF satisfies

$$r_{z,\text{damp},\text{lin}}^b(1/n) = \theta^s + \frac{1}{n} \mathcal{I}^\dagger(z; \theta^s).$$

As a consequence, influence functions are *not* approximating LOO retraining under the original loss L . Instead, they approximate the effect of training from θ^s under the modified objective

$$\mathcal{L}(\theta; \theta^s, \mathcal{D}^n) - \frac{1}{n} \ell(f(\theta; x), y) + \frac{\lambda}{2} \|\theta - \theta^s\|^2.$$

This fact makes PBRF a more suitable benchmark when testing the performances of IFs.

Concrete examples show that actually PBRF achieves good results in tasks such as mislabeled example detection, making it a viable alternative to LOO retraining.

Error decomposition

The approximation error of influence functions decomposes into:

- **Warm-start gap:** LOO starts from a random parameter (cold start), while IF are related to θ^s ; we can then converge to another "optimal" point;
- **Proximity gap:** the factor $\|\theta - \theta^s\|$ induces the warm start not to move far away from θ^s ;
- **Non-convergence gap:** in practice we almost never start from a fully trained network;
- **Linearization error:** produced by approximating the Taylor expansion at first order;
- **Solver error:** inexact iHVP computation.

Remark. The PBRF formulation eliminates the first three gaps.

Influence Functions vs Leverage

For simplicity, let's consider the Ordinary Least Squares problem.

Let $X = (x_1 | \dots | x_n)^\top \in \mathbb{R}^{(n \times d)}$, $y = (y_1, \dots, y_n) \in \mathbb{R}^n$, $e = (e_1, \dots, e_n) \in \mathbb{R}^n$, and $\theta \in \mathbb{R}^d$. In this setting, assume the dataset is generated by $y = f(\theta^*; x) + e = X \theta^* + e$ where e is the observational error. We can estimate θ^* with $\hat{\theta} = (X^\top X)^{-1} X^\top y$. Consequently, the predicted data with our model will be $\hat{y} = X(X^\top X)^{-1} X^\top y$ and we call $P = X(X^\top X)^{-1} X^\top$, as it is an orthogonal projection on $\text{ran}(X)$.

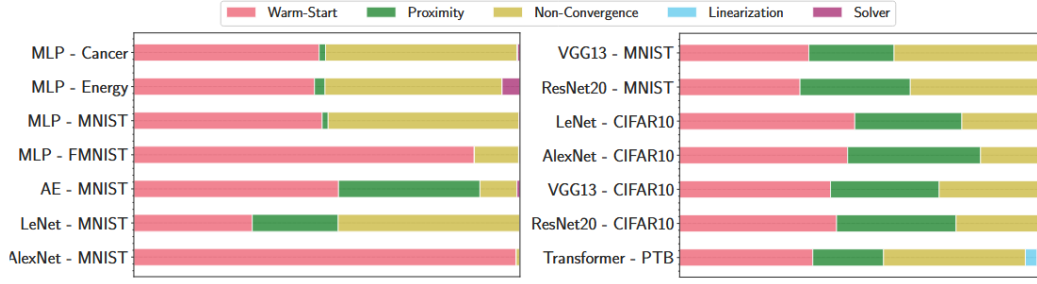


Figure 1: Visual representation of the error decomposition on different datasets and models. The main focus is that the largest components are the first three.

Definition 3. The *leverage score* of the i -th sample data is $P_{ii} = x_i(X^\top X)^{-1}x_i^\top$.

This quantity describes how much the i -th sample data affects the i -th prediction of our model. The bigger it is, the more probable it is that the i -th sample point is an outlier.

Indeed, if we consider $X = JF(\theta)$ as the X in our case, then the influence loss difference is approximately ¹ the same as the leverage. However, the interpretation of the problem would become different.

In any case, an expression of interest that involves both leverage and influence functions is the following (cf. [4]):

$$\hat{\theta} - \hat{\theta}_{1/n, -z} = \frac{(X^\top X)^{-1}x_i\hat{e}_i}{1 - P_{ii}},$$

where $\hat{e}_i = y_i - x_i^\top \hat{\theta}$.

Observe that the $1 - P_{ii}$ at denominator means that outliers also have high influence in the training.

Example: IFs for OLS

Let's compute all the quantities we defined so far in the case of OLS with the euclidean loss.

We have:

- Loss function: $\ell(f(\theta; x_i), y_i) = \frac{1}{2}(y_i - x_i^\top \theta)^2$;
- Gradient: $\nabla \ell(f(\theta; x_i), y_i) = -x_i(y_i - x_i^\top \theta)$;
- Hessian: $H_\theta = \sum_{i=1}^n x_i x_i^\top = X^\top X$;
- Influence function: $\mathcal{I}(z_i; \hat{\theta}) = (X^\top X)^{-1}x_i(y_i - x_i^\top \hat{\theta})$;

¹we can write $HF(\theta) = JF(\theta)^\top JF(\theta) + \sum \dots$. If we omit the second term, we have the sought approximation. This is acceptable when the parameter is close to optimal, since the sum annihilates for optimal parameters.

- PBRF: $\frac{1}{2} \arg \min_{\theta \in \Theta} \left(\frac{1}{n} \sum_{i=1}^n (x_i^\top (\theta - \hat{\theta}))^2 - \frac{1}{n} \|y_i - x_i^\top \theta\|^2 + \lambda \|\theta - \hat{\theta}\|^2 \right)$; note that in this case $\theta^s = \hat{\theta}$ since we can compute it explicitly.

Reformulation of influence functions

Remark. The object similar to what we are going to discuss, which is present in the paper, is $\hat{r}_z(\varepsilon)$.

Let \mathcal{X}, \mathcal{Y} be two measurable spaces and fix $D_n \subset (\mathcal{X} \times \mathcal{Y})^n$ such that $D_n = (z_i)_{i=1}^n$ with $z_i = (x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$.

Given two random variables X and Y respectively on \mathcal{X} and \mathcal{Y} , we are interested in studying the distribution of $Y|X$. To do so, we choose a Banach space Θ and a parametric function $f : \Theta \times \mathcal{X} \mapsto \mathcal{Y}$ as an estimator of such distribution.

In order to evaluate the estimators, we give the following definitions:

Definition 4. Given $\ell \in C^2(\mathcal{Y} \times \mathcal{Y}; \mathbb{R}_+)$ (called *loss function*), we define the *empirical risk* for the estimator $f(\theta; x)$ on the dataset D_n as:

$$R(\theta) := \sum_{i=1}^n \ell(f(\theta; x_i); y_i). \quad (0.1)$$

Given $\varepsilon \in [0, 1]$ and $j \in [n]$, we are interested in solving the minimization problem:

$$\min_{\theta \in \Theta} \sum_{i=1}^n \ell(f(\theta; x_i); y_i) + (\varepsilon - 1) \ell(f(\theta; x_j); y_j).$$

For simplicity, assume that for any ε the objective function has a unique minimum. We call $\theta(\varepsilon)$ such minimum. Let us introduce the notation:

$$\begin{aligned} \ell_j(\theta) &= \ell(f(\theta; x_j); y_j), \\ R_j(\theta) &= \sum_{\substack{i \in [n] \\ i \neq j}}^n \ell_i(\theta). \end{aligned}$$

Thus, we can rewrite:

$$\theta(\varepsilon) = \arg \min_{\theta \in \Theta} R_j(\theta) + \varepsilon \ell_j(\theta).$$

Notice that $\theta(\varepsilon)$ symbolizes the parameter for which our model best fits the data while we change the weight of one data point in the training. It is particularly of interest to understand how

$$\theta(1) = \theta^* = \arg \min R(\theta)$$

and

$$\theta(0) = \theta_{-j}^* = \arg \min R_j(\theta)$$

are different, since they represent respectively the optimal parameter obtained while training on the whole data and on the dataset minus one particular point.

One way to analyze this, is by the means of influence functions.

Definition 5. The *influence function* of z_j is:

$$I(j) := \left. \frac{d}{d\varepsilon} \theta(\varepsilon) \right|_{\varepsilon=0} = \dot{\theta}(0). \quad (0.2)$$

Remark. The previous quantity is well defined because $\theta(\varepsilon)$ is C^1 thanks to the Implicit function theorem (applied to $\nabla \mathcal{L}(\theta, \varepsilon) = L(\theta) + \varepsilon \ell_j(\theta)$ we get that there exists $\theta(\varepsilon)$ differentiable such that $\nabla \mathcal{L}(\theta(\varepsilon), \varepsilon) = 0$ in a neighbourhood of 1 at least, since by definition $\nabla \mathcal{L}(\theta^*, 1) = \nabla R(\theta^*) = 0$).

It is not clear yet if we are more interested in $\dot{\theta}(0)$ or $\dot{\theta}(1)$.

Proposition 0.0.1. We can write $I(j)$ explicitly as:

$$I(j) = -H_{R_j}^{-1} \nabla_{\theta} \ell_j(\theta_{-j}^*) \quad (0.3)$$

where H_{R_j} is the Hessian of R_j in θ_{-j}^* .

Proof. By definition, $\theta(\varepsilon)$ satisfies:

$$\nabla_{\theta}(R_j(\theta(\varepsilon)) + \varepsilon \ell_j(\theta(\varepsilon))) = 0.$$

Taking another derivative in ε yields:

$$\nabla_{\theta, \varepsilon}^2(R_j(\theta(\varepsilon)) + \varepsilon \ell_j(\theta(\varepsilon))) = 0 \iff \quad (0.4)$$

$$\nabla_{\theta\theta}^2 R_j(\theta(\varepsilon)) \dot{\theta}(\varepsilon) + \nabla_{\theta} \ell_j(\theta(\varepsilon)) + \varepsilon \nabla_{\theta\theta}^2 \ell_j(\theta(\varepsilon)) \dot{\theta}(\varepsilon) = 0 \iff \quad (0.5)$$

$$\dot{\theta}(\varepsilon) = -(\nabla_{\theta\theta}^2(R_j(\theta(\varepsilon)) + \varepsilon \ell_j(\theta(\varepsilon))))^{-1} \nabla_{\theta} \ell_j(\theta(\varepsilon)). \quad (0.6)$$

Evaluating in $\varepsilon = 0$ concludes the proof:

$$\dot{\theta}(0) = -(\nabla_{\theta\theta}^2 R_j(\theta(0)))^{-1} \nabla_{\theta} \ell_j(\theta(0)).$$

■

Remark. If we wanted to compute $\dot{\theta}(1)$, it would also have a nice form:

$$\dot{\theta}(1) = -H_R^{-1} \nabla_{\theta} \ell_j(\theta^*).$$

where H_R is the Hessian of R in θ^* . Indeed, if our goal is to *unlearn* a data point, this formulation does not require us to retrain the model, as we have already access to θ^* (but not to θ_{-j}^*).

It could also be of interest to consider:

$$Q(z) := \left. \frac{d}{d\varepsilon} (R_j(\theta(\varepsilon)) + \varepsilon \ell_j(\theta(\varepsilon))) \right|_{\varepsilon=0}. \quad (0.7)$$

Corollary 0.0.1. The following formulae hold:

$$\begin{aligned} \left. \frac{d}{d\varepsilon} (R_j(\theta(\varepsilon)) + \varepsilon \ell_j(\theta(\varepsilon))) \right|_{\varepsilon=0} &= \ell_j(\theta_{-j}^*), \\ \left. \frac{d}{d\varepsilon} (\ell_j(\theta(\varepsilon))) \right|_{\varepsilon=0} &= -\nabla_{\theta} \ell_j(\theta_{-j}^*)^{\top} H_{R_j}^{-1} \nabla_{\theta} \ell_j(\theta_{-j}^*), \end{aligned} \quad (0.8)$$

Proof. For the first equality, taking the derivative yields:

$$\nabla_{\theta} R_j(\theta(\varepsilon)) \dot{\theta}(\varepsilon) + \ell_j(\theta(\varepsilon)) + \varepsilon \nabla_{\theta} \ell_j(\theta(\varepsilon)) \dot{\theta}(\varepsilon)$$

and recalling that on $\theta(\varepsilon)$ it holds $\nabla(R_j(\theta(\varepsilon)) + \varepsilon \ell_j(\theta(\varepsilon))) = 0$ we can conclude. For the second one it suffices to substitute the definition of $\dot{\theta}$ from the previous equation after taking the derivative. ■

Remark. As discussed in [5], it is also possible to consider $\epsilon = (\varepsilon_1, \dots, \varepsilon_n) \in [0, 1]^n$ and $R_{\epsilon}(\theta) = \sum_{i=1}^n \varepsilon_i \ell_i(\theta)$. Note also that in this case, their notation figures “ $\theta(0)$ ”, but indeed in our notation it would be $\theta(1)$ (which makes sense since we are considering a group of points).

In such work, they provide formal statements about when the influence function approximations are accurate, taking into account also the difference between the influence function estimation and the *Newton estimation*.

Definition 6. We call Newton estimation the quantity:

$$I_{Nt}(j) = -(HR_j(\theta(1)))^{-1} \nabla R_j(\theta(1))$$

The name for this stems from the fact that this formula estimates the parameter after one step of Newton method while minimizing ∇R_j (equivalently, we are considering the second order Taylor expansion of R_j in $\theta(1)$):

$$\theta_{Nt} = \theta(1) - (HR_j(\theta(1)))^{-1} \nabla R_j(\theta(1)).$$

Unfortunately, $I(j)$ and $I_{Nt}(j)$ are close when the smallest eigenvalue of the hessian is big, which is not always the case. One way to ensure this difference is small, is to add the regularization term $\lambda \|\theta\|$ in the risk. This way, we can decrease the aforementioned error by choosing a large λ .

When working with more than one index, it makes sense to define the following quantity:

$$\theta^j(\varepsilon) = \theta(\underbrace{1, \dots, \varepsilon, \dots, 1}_{j\text{-th place}}).$$

Question: Is there a link between $\dot{\theta}^j$ and $J\theta$?

0.1 Gradient Descent case

Let us specialize the analysis to the instance where we are using the gradient descent algorithm to pursue the minimization task. Furthermore, consider the gradient flow dynamics:

$$\begin{cases} \frac{d}{dt} \theta(t, \varepsilon) = -\nabla R_j(\theta(t, \varepsilon)) - \varepsilon \nabla \ell_j(\theta(t, \varepsilon)) \\ \theta(0, \varepsilon) = \theta_0 \quad \forall \varepsilon \in [0, 1] \end{cases}, \quad (\text{GF})$$

where we consider fixed $j \in [n]$ and $\theta_0 \in \Theta$.

Let us call $\xi(t_0, t_1, \varepsilon; \theta_0)$ the flux of (GF), i.e., $\theta(t, \varepsilon) = \xi(0, t, \varepsilon; \theta_0)$. By the differential equation theory, since we are assuming $\ell \in C^2$, we know that at least $\xi \in C_t^2 \times C_\varepsilon^0$ (even Lipschitz in ε). Therefore, $\theta : [0, +\infty] \times [0, 1] \rightarrow \Theta$ can be seen as a homotopy between the learning trajectories in the parameters' space (taking as definition $\theta(+\infty, \varepsilon) := \theta(\varepsilon)$).

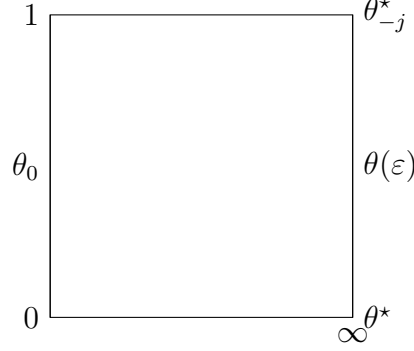


Figure 2: Visual representation of the homotopy $\theta(t, \varepsilon)$ where on the horizontal axis we have evolution in time and on vertical axis the change of parameter ε .

Question: How different really is $\dot{\theta}(1)$ from $\dot{\theta}(0)$?

Assuming $\theta(\varepsilon)$ is C^2 , we can give the bound:

$$\dot{\theta}(1) - \dot{\theta}(0) = \int_0^1 \ddot{\theta}(\varepsilon) d\varepsilon \leq \|\ddot{\theta}\|_{\infty, [0, 1]},$$

but its computation requires a third derivative of R_j ...

I would like to do some experiments in a simple setting for visualizing $\theta(\varepsilon)$. This may give insights.

Some experiments are summarized in Figure 3. What I found interesting about them is that the distance from the target function does not affect the linearity of the trajectory. In fact, in the figure you can see that the first considered point is close to the real value, but the trajectory of its parameter is not a rect, unlike the second point (which is more distant).

Question: Consider now the map $\theta(t, s) : [0, +\infty] \times [0, T]$, with $T \in (0, +\infty]$ fixed a priori, such that:

$$\begin{cases} \theta(t, s) = \xi(0, t, 1; \theta_0) & t \in [0, s] \\ \theta(t, s) = \xi(s, t, 0; \xi(0, s, 1; \theta_0)) & t \in [s, T] \end{cases}.$$

What can we say about $\theta(T, s)$?

If $T = \infty$, then for any $s < \infty$ it holds that $\theta(\infty, s) = \theta_{-j}^*$, but for finite T the answer is not clear.

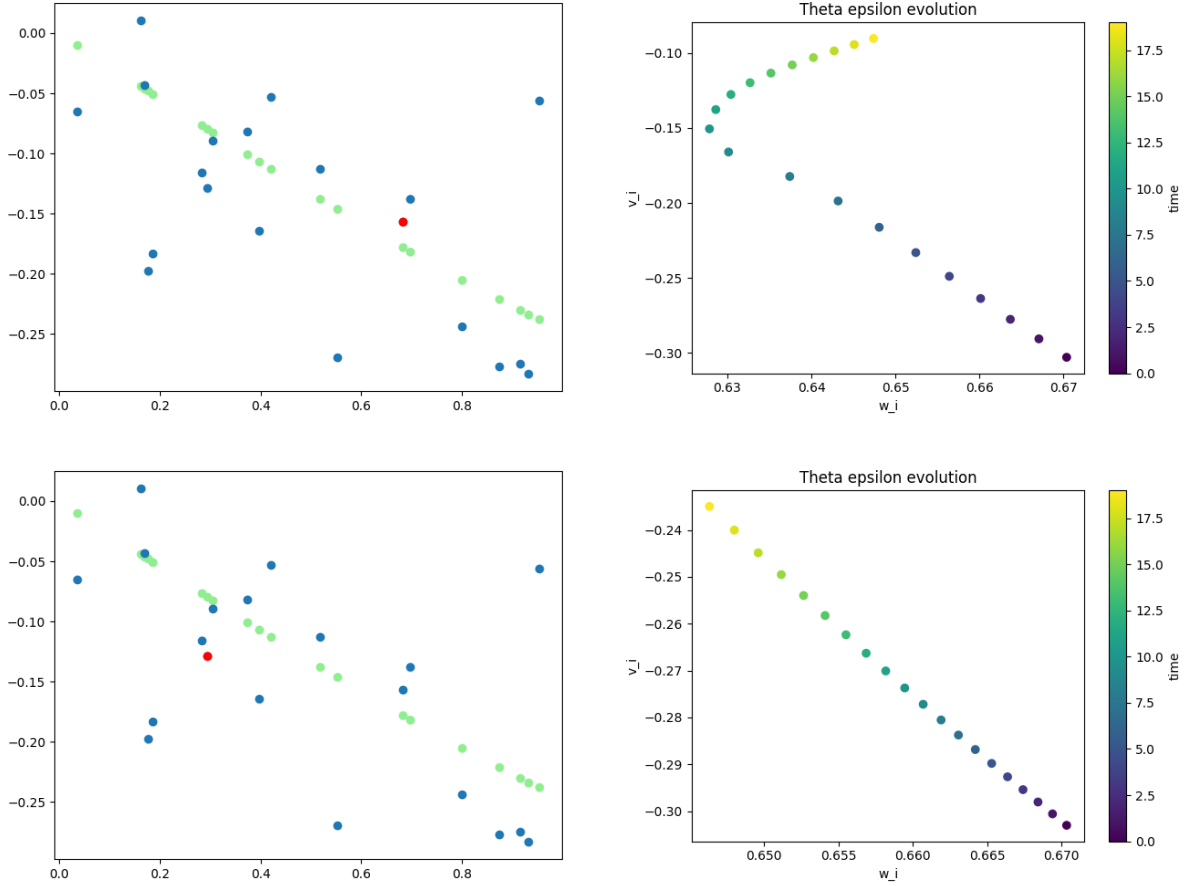


Figure 3: I used a 3-neuron 1-hidden layer MLP trained on 20 points generated by $y_i = f_*(x_i) + x_i$ where f_* is another MLP of the same type and $\xi_i \sim \mathcal{N}(0, 0.1)$. On the left column we can observe the considered data point we are deleting in red, in blue the other data points and in green the target function. On the right side, I represented the trajectory of the first neuron as a function of ε (I trained the model for 1000 iterations, it should be enough)

Algorithm implementations: D2D vs R2D

An application for which we may want to use influence functions is unlearning a data point (for example for privacy reasons). However, as we have seen, IFs may not work well in the MLP and Deep Learning setting. Therefore, in the literature some algorithms have been developed to unlearn training points without using the highly inefficient LOO retraining. For example, in [9], the authors take inspiration from the *Descend-to-Delete* algorithm (D2D) and create a more efficient version called *Rewind-to-Delete* (R2D).

The former algorithm is very intuitive and consists on simply continuing the training of the model but on the updated dataset where we deleted the target data point. By continuing for enough iteration, there have been proven some theoretical guarantees that the final parameter distribution is in some sense undistinguishable from the one of LOO

method.

The problem with this algorithm is that such guarantees only hold in the convex case and cannot be extended to the more general setting (e.g. ReLU and tanh functions are not convex, so this result does not hold).

On the other hand, the R2D algorithm saves a check-point parameter during the training (let's say after the K -th iteration) and then keeps going until iteration T . The idea is that when we receive the request of unlearning the target point, we will continue the training on the updated dataset starting from $\theta(K)$, instead of $\theta(T)$.

The authors of the paper show that also this method yields indistinguishability, but in this case the result holds also for non-convex functions and the iterations required are much less than the total training time T (required by LOO).

In [8], the same researchers extend the previous results also to the (projected) SGD case. In particular, in addition to all previous observations, they highlight that D2D provides tighter bounds for the undistinguishability due to its reliance on the convergence to a unique global minimum, while R2D has more loose estimates as it only counts on the underlying contractivity of gradient systems.

0.2 Other interesting things we might want to explore

- If our goal is to gain a better understanding of what data points are the most informative during the training, instead of trying to unlearn certain data, the results in [3] might be interesting. In their paper, the authors prove that their Shapley values-based method performs better on this task rather than influence function methods.
- It may be of interest to study Bayesian Influence Functions (BIFs) as well as frequentistic ones. In [7], the researchers present an unlearning method that uses BIFs instead of IFs. The reason for this is we don't need to compute the iHVP to evaluate BIFs, therefore this method works better with more singular loss landscapes. As another consequence, we don't need to evaluate these quantities on local minima (we don't need the hessian to be invertible), which is one of the less realistic hypotheses for IFs.

On the other hand, computations are not always faster than IFs (here, the leading cost is estimating the covariance between two elements) and achieving good results requires more hyperparameter tuning than classical methods.

Bibliography

- [1] Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger Grosse. If Influence Functions are the Answer, Then What is the Question?, September 2022. arXiv:2209.05364 [cs].
- [2] R. D. Cook and S. Weisberg. *Residuals and Influence in Regression*. Springer, 1983.
- [3] Amirata Ghorbani and James Zou. Data Shapley: Equitable Valuation of Data for Machine Learning, June 2019. arXiv:1904.02868 [stat].
- [4] Fumio Hayashi. *Econometrics*. Princeton University Press, 2000.
- [5] Pang Wei Koh, Kai-Siang Ang, Hubert H. K. Teo, and Percy Liang. On the Accuracy of Influence Functions for Measuring Group Effects, November 2019. arXiv:1905.13289 [cs].
- [6] Pang Wei Koh and Percy Liang. Understanding Black-box Predictions via Influence Functions, December 2020. arXiv:1703.04730 [stat].
- [7] Philipp Alexander Kreer, Wilson Wu, Maxwell Adam, Zach Furman, and Jesse Hoogland. Bayesian Influence Functions for Hessian-Free Data Attribution, September 2025. arXiv:2509.26544 [cs].
- [8] Siqiao Mu and Diego Klabjan. Descend or Rewind? Stochastic Gradient Descent Unlearning, November 2025. arXiv:2511.15983 [cs].
- [9] Siqiao Mu and Diego Klabjan. Rewind-to-Delete: Certified Machine Unlearning for Nonconvex Functions, October 2025. arXiv:2409.09778 [cs].