# Stochastic Gradient Descent in NPIV estimation

**Anonymous Author(s)**
Affiliation
Address
`email`

## 1 Problem setup

### 1.1 Basic definitions

Fix a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. Given $X \in L^2(\Omega, \mathcal{A}, \mathbb{P}; \mathcal{X} \subseteq \mathbf{R}^p)$, we define

$$L^2(X) \triangleq \left\{ h : \mathcal{X} \to \mathbf{R} \ : \ \mathbb{E}[h(X)^2] < \infty \right\},$$

that is, $L^2(X) = L^2(\mathcal{X}, \mathcal{B}(\mathcal{X}), \nu_X)$, where we denote by $\nu_X$ the distribution of the r.v. $X$ and by $\mathcal{B}(\mathcal{X})$ the Borel $\sigma$-algebra in $\mathcal{X}$. This is a Hilbert space equipped with the inner product $\langle h, g \rangle_{L^2(X)} = \mathbb{E}[h(X)g(X)]$. The regression problem we are interested in has the form

$$Y = h^\star(X) + \varepsilon, \tag{1}$$

where $h^\star \in L^2(X)$ and $\varepsilon$ is an square-integrable r.v. such that $\mathbb{E}[\varepsilon \mid X] \neq 0$. We assume there exists $Z \in L^2(\Omega, \mathcal{A}, \mathbb{P}; \mathcal{Z} \subseteq \mathbf{R}^q)$ such that

    i) $Z$ influences $X$, that is, $\nu_{X|Z}(\cdot \mid Z) \neq \nu_X(\cdot)$;

    ii) $Z$ influences $Y$ only through $Z$;

    iii) $Z$ and $\varepsilon$ are uncorrelated, that is, $\mathbb{E}[\varepsilon \mid Z] = 0$.

The space $L^2(Z)$ is defined accordingly. This variable is called the *instrumental variable*. The problem consists of estimating $h^\star$ based on independent joint samples from $X, Z$ and $Y$.

Conditioning (1) in $Z$, we find

$$\mathbb{E}[Y \mid Z] = \mathbb{E}[h^\star(X) \mid Z]. \tag{2}$$

This motivates us to introduce the operator $\mathcal{P} : L^2(X) \to L^2(Z)$ defined by

$$\mathcal{P}[h](z) \triangleq \mathbb{E}[h(X) \mid Z = z].$$

Clearly $\mathcal{P}$ is linear and, using Jensen's inequality, one may prove that it's bounded. It's also interesting to notice that its adjoint $\mathcal{P}^* : L^2(Z) \to L^2(X)$ satisfies

$$\mathcal{P}^*[g](x) = \mathbb{E}[g(Z) \mid X = x]. \tag{3}$$

Define $r_0 : \mathcal{Z} \to \mathbf{R}$ by $r_0(Z) = \mathbb{E}[Y \mid Z]$. Again by Jensen's inequality, we have $r_0 \in L^2(Z)$, and thus we can rewrite (2) as

$$\mathcal{P}[h^\star] = r_0. \tag{4}$$

Hence, (1) can be formulated as an inverse problem, where we wish to invert the operator $\mathcal{P}$. ⎤ Discuss the other implication, that if $h$ satisfies $\mathcal{P}[h] = r_0$, then $h = h^\star$. This is false, but the reason can be connected to the strength of the instrument $Z$.

### 1.2 Risk measure

Let $\ell : \mathbf{R} \times \mathbf{R} \to \mathbf{R}$ be a pointwise loss function, which, with respect to its second argument, is convex and differentiable. We use the symbol $\partial_2$ to denote a derivative with respect to the second

argument. The example to keep in mind is the quadratic loss function $\ell(y, y') = \frac{1}{2}(y - y')^2$. Given $h \in L^2(X)$, we define the *populational risk* associated with it to be

$$\mathcal{R}(h) \triangleq \mathbb{E}[\ell(r_0(Z), \mathcal{P}h(Z))].$$

We would like to solve

$$\inf_{h \in \mathcal{F}} \mathcal{R}(h),$$

where $\mathcal{F} \subseteq L^2(X)$ is a bounded, closed, convex set such that $h^\star \in \mathcal{F}$. A possible choice for the set $\mathcal{F}$ is

$$\mathcal{F} = \left\{ h \in L^2(X) : \|h\|_\infty \leq A \right\},$$

where $A > 0$ is a constant chosen *a priori*. This set is obviously closed, convex and bounded in the $L^2(X)$ norm. Furthermore, the projection operator $\pi_\mathcal{F}$ is very easy to compute, as $\pi_\mathcal{F}[h]$ is obtained by cropping $h$ inside $[-A, A]$. More formally,

$$\pi_\mathcal{F}[h] = h^+ \wedge A - h^- \wedge A.$$

We now state all the assumptions needed about the function $\ell$ for future reference:

**Assumption 1** (Regularity of $\ell$).

    *1. The function $\ell : \mathbf{R} \times \mathbf{R} \to \mathbf{R}$ is convex and $C^2$ with respect to its second argument;*

    *2. There exists $\theta_0 > 0$ such that for all $f, g \in L^2(X)$*

$$\sup_{|\theta| < \theta_0} \mathbb{E}\left[ \partial_2^2 \ell(r_0(Z), \mathcal{P}[g + \theta f](Z)) \cdot \mathcal{P}[f](Z)^2 \right] < \infty; \tag{5}$$

Assumption 1.2 is a mild integrability condition which can be easily shown to hold in the quadratic case.

# 2 Gradient computation

We'd like to compute $\nabla \mathcal{R}(h)$ for $h \in L^2(X)$. We start by computing the directional derivative of $\mathcal{R}$ at $h$ in the direction $f$, denoted by $D\mathcal{R}[h](f)$:

$$\begin{aligned}
D\mathcal{R}[h](f) &= \lim_{\delta \to 0} \frac{1}{\delta}\left[ \mathcal{R}(h + \delta f) - \mathcal{R}(f) \right] \\
&= \lim_{\delta \to 0} \frac{1}{\delta} \mathbb{E}\left[ \ell(r_0(Z), \mathcal{P}[h + \delta f](Z)) - \ell(r_0(Z), \mathcal{P}[h](Z)) \right] \\
&= \lim_{\delta \to 0} \frac{1}{\delta} \mathbb{E}\left[ \ell(r_0(Z), \mathcal{P}[h](Z) + \delta \mathcal{P}[f](Z)) - \ell(r_0(Z), \mathcal{P}[h](Z)) \right] \\
&= \lim_{\delta \to 0} \frac{1}{\delta} \mathbb{E}\left[ \delta \partial_2 \ell(r_0(Z), \mathcal{P}[h](Z)) \cdot \mathcal{P}[f](Z) \right. \\
&\qquad\qquad \left. + \frac{\delta^2}{2} \partial_2^2 \ell(r_0(Z), \mathcal{P}[h + \theta f](Z)) \cdot \mathcal{P}[f](Z)^2 \right] \\
&= \mathbb{E}\left[ \partial_2 \ell(r_0(Z), \mathcal{P}[h](Z)) \cdot \mathcal{P}[f](Z) \right] \\
&\qquad + \lim_{\delta \to 0} \mathbb{E}\left[ \frac{\delta}{2} \partial_2^2 \ell(r_0(Z), \mathcal{P}[h + \theta f](Z)) \cdot \mathcal{P}[f](Z)^2 \right] \\
&= \mathbb{E}\left[ \partial_2 \ell(r_0(Z), \mathcal{P}[h](Z)) \cdot \mathcal{P}[f](Z) \right],
\end{aligned}$$

where $\theta \in \mathbf{R}$ is due to Taylor's formula and can be assumed to be inside a fixed interval $(-\theta_0, \theta_0)$, with $\theta_0$ arbitrarily small. The last step is then due to Assumption 1.2.

We can in fact expand the calculation a bit more, as follows:

$$\begin{aligned}
D\mathcal{R}[h](f) &= \mathbb{E}\left[ \partial_2 \ell(r_0(Z), \mathcal{P}[h](Z)) \cdot \mathcal{P}[f](Z) \right] \\
&= \langle \partial_2 \ell(r_0(\cdot), \mathcal{P}[h](\cdot)), \mathcal{P}[f] \rangle_{L^2(Z)} \\
&= \langle \mathcal{P}^*[\partial_2 \ell(r_0(Z), \mathcal{P}[h](\cdot))], f \rangle_{L^2(X)},
\end{aligned}$$

45 where we are assuming that $\partial_2 \ell(r_0(\cdot), \mathcal{P}[h](\cdot)) \in L^2(Z)$. This shows that $\mathcal{R}$ is Gateux-differentiable,
46 with Gateux derivative at $h$ given by

$$D\mathcal{R}[h] = \mathcal{P}^*[\partial_2 \ell(r_0(\cdot), \mathcal{P}[h](\cdot))].$$

47 If we assume[1] that $h \mapsto D\mathcal{R}[h]$ is a continuous mapping from $L^2(Z)$ to $L^2(Z)$, then $\mathcal{R}$ is also
48 Fréchet-differentiable, and both derivatives coincide. Therefore, under this assumption, which we
49 henceforth make, $\nabla \mathcal{R}(h) = \mathcal{P}^*[\partial_2 \ell(r_0(\cdot), \mathcal{P}[h](\cdot))]$.

## 50  3   Estimating the gradient

51 We have found that

$$\nabla \mathcal{R}(h)(x) = \mathcal{P}^*[\partial_2 \ell(r_0(\cdot), \mathcal{P}[h](\cdot))](x) = \mathbb{E}[\partial_2 \ell(r_0(Z), \mathcal{P}[h](Z)) \mid X = x].$$

52 This turns out to be hard to estimate in practice, as we have two nested conditional expectation
53 operators. Our objective in this section is to write $\nabla \mathcal{R}(h)(x) = \mathbb{E}[\Phi(x, Z)\partial_2 \ell(r_0(Z), \mathcal{P}[h](Z))]$,
54 for some suitable kernel $\Phi$. Then, for a given sample of $Z$, the function $\Phi(\cdot, Z)\partial_2 \ell(r_0(Z), \mathcal{P}[h](Z))$
55 acts as an stochastic estimate for $\nabla \mathcal{R}(h)$. To ease the notation, define $\Psi_h(z) \triangleq \partial_2 \ell(r_0(z), \mathcal{P}[h](z))$.
56 Assuming that $X$ and $Z$ have a joint distribution which is absolutely continuous with respect to
57 Lebesgue measure in $\mathbf{R}^{p+q}$, we can write

$$\nabla \mathcal{R}(h)(x) = \mathbb{E}[\Psi_h(Z) \mid X = x]$$

$$= \int_{\mathbb{Z}} p(z \mid x)\Psi_h(z) \, \mathrm{d}z$$

$$= \int_{\mathbb{Z}} p(z)\frac{p(z \mid x)}{p(z)}\Psi_h(z) \, \mathrm{d}z$$

$$= \mathbb{E}\left[\frac{p(Z \mid x)}{p(Z)}\Psi_h(Z)\right].$$

58 Thus, we must take

$$\Phi(x, z) = \frac{p(z \mid x)}{p(z)} = \frac{p(x \mid z)}{p(x)} = \frac{p(x, z)}{p(x)p(z)}.$$

59 With this choice, setting $u_h(x) = \Phi(x, Z)\Psi_h(Z)$ we clearly have $\mathbb{E}[u_h(x)] = \nabla \mathcal{R}(h)(x)$.

60 An obvious obstacle for this approach is that we don't know how to analytically compute $\Phi, r_0$ nor $\mathcal{P}$,
61 se we will proceed with estimators $\widehat{\Phi}, \widehat{r_0}$ and $\widehat{\mathcal{P}}$. In what follows, we will remain agnostic to the exact
62 form taken by these estimators and will present the algorithm assuming we know how to compute
63 them. Later, we will show how the individual convergence rates of these three pieces come together
64 to determine the convergence rate of our method.

## 65  4   Algorithm

66 Having an estimator of the gradient, we can construct Functional GD algorithm for estimating $h^\star$.

---

**Algorithm 1:** SGD-NPIV

**input**  : Datasets $\mathcal{D}_{r_0}, \mathcal{D}_\Phi$ and $\mathcal{D}_\mathcal{P}$ for estimating $r_0, \Phi$ and $\mathcal{P}$, respectively. Samples
$\{(\boldsymbol{z}_m)\}_{m=1}^M$ for the gradient descent loop. Discretization $\{\boldsymbol{x}_k\}_{k=1}^K$ of $\mathcal{X}$ which contains
the observed values of $X$. Sequence of learning rates $(\alpha_m)_{m=1}^M$.

**output** : $\widehat{h}$

67  Compute $\widehat{r_0}, \widehat{\Phi}, \widehat{\mathcal{P}}$ using $\mathcal{D}_{r_0}, \mathcal{D}_\Phi, \mathcal{D}_\mathcal{P}$, respectively ;

**for** $1 \le m \le M$ **do**

$\quad$ Set $u_m = \widehat{\Phi}(\cdot, \boldsymbol{z}_m)\partial_2 \ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)$ ;

$\quad$ Set $\widehat{h}_m(\boldsymbol{x}_k) = \pi_\mathcal{F}\left[\widehat{h}_{m-1} - \alpha_m u_m\right](\boldsymbol{x}_k) \quad$ for $1 \le k \le K$ ;

**end**

Set $\widehat{h} = \frac{1}{M}\sum_{m=1}^M \widehat{h}_m$ ;

---

---

[1]It is if $\ell$ is quadratic.

# 5 Proof of convergence

69 The first problem is proving our sequence of estimates is, in fact, contained in $L^2(X)$. This amounts
70 to proving $u_m \in L^2(X)$ for every $m$. It's not even immediate why $u_h(x) = \Phi(x, Z)\xi_h(Z)$ (the
71 unbiased gradient when we know $r_0, \Phi$ and $\mathcal{P}$) belongs to $L^2(X)$

> We'll need to bound the norm of $u_m$ by a constant later in the proof.

72 After doing this, we check that $\mathcal{R}$ is convex in $\mathcal{F}$: if $h, g \in \mathcal{F}$ and $\lambda \in [0, 1]$, then

$$\mathcal{R}(\lambda h + (1 - \lambda)g) = \mathbb{E}[\ell(r_0(Z), \mathcal{P}[\lambda h + (1 - \lambda)g](Z))]$$
$$= \mathbb{E}[\ell(r_0(Z), \lambda \mathcal{P}[h](Z) + (1 - \lambda)\mathcal{P}[g](Z))]$$
$$\leq \lambda \mathbb{E}[\ell(r_0(Z), \mathcal{P}[h](Z))] + (1 - \lambda)\mathbb{E}[\ell(r_0(Z), \mathcal{P}[g](Z))]$$
$$= \lambda \mathcal{R}(h) + (1 - \lambda)\mathcal{R}(g).$$

73 To lighten the notation, the symbols $\|\cdot\|$ and $\langle \cdot, \cdot \rangle$, when written without a subscript to specify which
74 space they refer to, will act as the norm and inner product, respectively, of $L^2(X)$. By the Algorithm
75 1 procedure, we have

$$\frac{1}{2}\left\|\widehat{h}_m - h^\star\right\|^2 = \frac{1}{2}\left\|\pi_\mathcal{F}\left[\widehat{h}_{m-1} - \alpha_m u_m\right] - h^\star\right\|^2$$
$$\leq \frac{1}{2}\left\|\widehat{h}_{m-1} - \alpha_m u_m - h^\star\right\|^2$$
$$= \frac{1}{2}\left\|\widehat{h}_{m-1} - h^\star\right\|^2 - \alpha_m\langle u_m, \widehat{h}_{m-1} - h^\star\rangle + \frac{\alpha_m^2}{2}\|u_m\|^2.$$

76 After adding and subtracting $\alpha_m\langle \nabla\mathcal{R}(\widehat{h}_{m-1}), \widehat{h}_{m-1} - h^\star\rangle$, we are left with

$$\frac{1}{2}\left\|\widehat{h}_{m-1} - h^\star\right\|^2 - \alpha_m\langle u_m - \nabla\mathcal{R}(\widehat{h}_{m-1}), \widehat{h}_{m-1} - h^\star\rangle + \frac{\alpha_m^2}{2}\|u_m\|^2 - \alpha_m\langle \nabla\mathcal{R}(\widehat{h}_{m-1}), \widehat{h}_{m-1} - h^\star\rangle.$$

77 Applying the basic convexity inequality on the last term give us, in total,

$$\frac{1}{2}\left\|\widehat{h}_m - h^\star\right\|^2 \leq \frac{1}{2}\left\|\widehat{h}_{m-1} - h^\star\right\|^2 - \alpha_m\langle u_m - \nabla\mathcal{R}(\widehat{h}_{m-1}), \widehat{h}_{m-1} - h^\star\rangle$$
$$+ \frac{\alpha_m^2}{2}\|u_m\|^2 - \alpha_m(\mathcal{R}(\widehat{h}_{m-1}) - \mathcal{R}(h^\star)).$$

78 Rearranging terms, we get

$$\mathcal{R}(\widehat{h}_{m-1}) - \mathcal{R}(h^\star) \leq \frac{1}{2\alpha_m}\left(\left\|\widehat{h}_{m-1} - h^\star\right\|^2 - \left\|\widehat{h}_m - h^\star\right\|^2\right)$$
$$+ \frac{\alpha_m}{2}\|u_m\|^2 - \langle u_m - \nabla\mathcal{R}(\widehat{h}_{m-1}), \widehat{h}_{m-1} - h^\star\rangle.$$

79 Finally, summing over $1 \leq m \leq M$ leads to

$$\sum_{n=1}^{M}\left[\mathcal{R}(\widehat{h}_{m-1}) - \mathcal{R}(h^\star)\right] \leq \sum_{m=1}^{M}\frac{1}{2\alpha_m}\left(\left\|\widehat{h}_{m-1} - h^\star\right\|^2 - \left\|\widehat{h}_m - h^\star\right\|^2\right)$$
$$+ \sum_{m=1}^{M}\frac{\alpha_m}{2}\|u_m\|^2$$
$$- \sum_{m=1}^{M}\langle u_m - \nabla\mathcal{R}(\widehat{h}_{m-1}), \widehat{h}_{m-1} - h^\star\rangle.$$

80 We then treat each of the three terms in the RHS of the inequality above separately:

81 **First term** By assumption, we have $\operatorname{diam}\mathcal{F} = D < \infty$. Hence

$$\sum_{m=1}^{M}\frac{1}{2\alpha_m}\left(\left\|\widehat{h}_{m-1} - h^\star\right\|^2 - \left\|\widehat{h}_m - h^\star\right\|^2\right) = \sum_{m=2}^{M}\left(\frac{1}{2\alpha_m} - \frac{1}{2\alpha_{m-1}}\right)\left\|\widehat{h}_{m-1} - h^\star\right\|^2$$
$$+ \frac{1}{2\alpha_1}\left\|\widehat{h}_0 - h^\star\right\|^2 - \frac{1}{2\alpha_M}\left\|\widehat{h}_M - h^\star\right\|^2$$
$$\leq \sum_{m=2}^{M}\left(\frac{1}{2\alpha_m} - \frac{1}{2\alpha_{m-1}}\right)D^2 + \frac{1}{2\alpha_1}D^2 = \frac{D^2}{2\alpha_M}.$$

**Second term** We are fixing the offline data $\mathcal{D}_{\Phi,\mathcal{P},r_0}$ and averaging with respect to the other samples of the instrumental variable. Therefore, what we wish to compute is

$$\mathbb{E}_{\boldsymbol{z}_{1:M}}\left[\|u_m\|^2 \mid \mathcal{D}_{\Phi,\mathcal{P},r_0}\right] = \mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\mathbb{E}_X\left[\widehat{\Phi}(X,\boldsymbol{z}_m)^2 \partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2\right] \;\Big|\; \mathcal{D}_{\Phi,\mathcal{P},r_0}\right]$$

$$= \mathbb{E}_X\left[\mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\widehat{\Phi}(X,\boldsymbol{z}_m)^2 \partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2 \;\Big|\; \mathcal{D}_{\Phi,\mathcal{P},r_0}\right]\right].$$

Since $\boldsymbol{z}_{1:m}$ is independent from $\mathcal{D}_{\Phi,\mathcal{P},r_0}$, this is equal to

$$\mathbb{E}_X\left[\mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\widehat{\Phi}(X,\boldsymbol{z}_m)^2 \partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2\right]\right].$$

Reversing back the expectations, we get

$$\mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\mathbb{E}_X\left[\widehat{\Phi}(X,\boldsymbol{z}_m)^2 \partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2\right]\right]$$

$$= \mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\mathbb{E}_X\left[\widehat{\Phi}(X,\boldsymbol{z}_m)^2\right] \partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2\right].$$

Now we use Assumption 14.5.1 in [1], which states that

$$\sup_{\boldsymbol{w}\in\mathbb{W}} k(\boldsymbol{w},\boldsymbol{w}) \leq 1,$$

where $\mathbb{W} = \mathbb{X}\times\mathbb{Z}$, $\boldsymbol{w} = (\boldsymbol{x},\boldsymbol{z})$ and $k : \mathbb{W}\times\mathbb{W}\to\mathbf{R}$ is the kernel corresponding to the RKHS used to estimate $\Phi$, which we denote by $\mathcal{R}_{\mathbb{W}}$. This assumption implies

$$\widehat{\Phi}(\boldsymbol{w}) = \langle\widehat{\Phi}, k(\boldsymbol{w},\cdot)\rangle_{\mathcal{R}_{\mathbb{W}}} \leq \left\|\widehat{\Phi}\right\|_{\mathcal{R}_{\mathbb{W}}}\|k(\boldsymbol{w},\cdot)\|_{\mathcal{R}_{\mathbb{W}}} = \left\|\widehat{\Phi}\right\|_{\mathcal{R}_{\mathbb{W}}}\sqrt{\langle k(\boldsymbol{w},\cdot),k(\boldsymbol{w},\cdot)\rangle} =$$

$$= \left\|\widehat{\Phi}\right\|_{\mathcal{R}_{\mathbb{W}}}\sqrt{k(\boldsymbol{w},\boldsymbol{w})} \leq \left\|\widehat{\Phi}\right\|_{\mathcal{R}_{\mathbb{W}}}$$

for all $\boldsymbol{w}\in\mathbb{W}$. Therefore,

$$\mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\mathbb{E}_X\left[\widehat{\Phi}(X,\boldsymbol{z}_m)^2\right] \partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2\right]$$

$$\leq \mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\mathbb{E}_X\left[\left\|\widehat{\Phi}\right\|^2_{\mathcal{R}_{\mathbb{W}}}\right] \partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2\right]$$

$$= \left\|\widehat{\Phi}\right\|^2_{\mathcal{R}_{\mathbb{W}}} \mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2\right].$$

To bound the expectation, we assume the loss is quadratic and then

<span style="border:1px solid green;border-radius:8px;padding:2px 6px;">Assumption</span>

$$\mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\left(\widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m) - \widehat{r_0}(\boldsymbol{z}_m)\right)^2\right]$$

$$= \mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\left(\left(\widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m) - \mathcal{P}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right) + (r_0(\boldsymbol{z}_m) - \widehat{r_0}(\boldsymbol{z}_m))\right.\right.$$

$$\left.\left. + \left(\mathcal{P}[\widehat{h}_{m-1}](\boldsymbol{z}_m) - r_0(\boldsymbol{z}_m)\right)\right)^2\right]$$

$$\leq 3\mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\left(\widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m) - \mathcal{P}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2 + (r_0(\boldsymbol{z}_m) - \widehat{r_0}(\boldsymbol{z}_m))^2\right.$$

$$\left. + \left(\mathcal{P}[\widehat{h}_{m-1}](\boldsymbol{z}_m) - r_0(\boldsymbol{z}_m)\right)^2\right]$$

$$= 3\left\{\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|(\widehat{\mathcal{P}} - \mathcal{P})[\widehat{h}_{m-1}]\right\|^2_{L^2(\mathbb{Z})}\right] + \mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\|r_0 - \widehat{r_0}\|^2_{L^2(\mathbb{Z})}\right]\right.$$

$$\left. + \mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|\mathcal{P}[\widehat{h}_{m-1}] - r_0\right\|^2_{L^2(\mathbb{Z})}\right]\right\}.$$

5

We treat each part of this expression separately. Firstly,

$$\left\|(\widehat{\mathcal{P}} - \mathcal{P})[\widehat{h}_{m-1}]\right\|_{L^2(\mathbb{Z})}^2 \leq \left\|\widehat{\mathcal{P}} - \mathcal{P}\right\|_{\mathrm{op}}^2 \left\|\widehat{h}_{m-1}\right\|_{L^2(\mathbb{X})}^2 \leq M^2 \left\|\widehat{\mathcal{P}} - \mathcal{P}\right\|_{\mathrm{op}}^2.$$

We leave the second part as $\|r_0 - \widehat{r_0}\|_{L^2(\mathbb{Z})}^2$. Finally, for the third part, we have

$$
\begin{aligned}
\left\|\mathcal{P}[\widehat{h}_{m-1}] - r_0\right\|_{L^2(\mathbb{Z})}^2 &= \mathbb{E}_Z\left[\left(\mathcal{P}[\widehat{h}_{m-1}](Z) - r_0(Z)\right)^2\right] \\
&= \mathbb{E}_Z\left[\left(\mathbb{E}\left[\widehat{h}_{m-1}(X) - Y \mid Z\right]\right)^2\right] \\
&\leq \mathbb{E}_{(X,Y)}\left[\left(\widehat{h}_{m-1}(X) - Y\right)^2\right] \\
&\leq 2\left(\mathbb{E}_X\left[\widehat{h}_{m-1}(X)^2\right] + \mathbb{E}\left[Y^2\right]\right) \\
&= 2\left(\left\|\widehat{h}_{m-1}\right\|_{L^2(\mathbb{X})}^2 + \mathbb{E}\left[Y^2\right]\right) \\
&\leq 2\left(M^2 + \mathbb{E}\left[Y^2\right]\right).
\end{aligned}
$$

Putting everything together, what we conclude is

$$\mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\|u_m\|_{L^2(\mathbb{X})}^2 \mid \mathcal{D}_{\Phi,\mathcal{P},r_0}\right] \leq 3\left\|\widehat{\Phi}\right\|_{\mathcal{R}_\mathbb{W}}^2\left(M^2\left\|\widehat{\mathcal{P}} - \mathcal{P}\right\|_{\mathrm{op}}^2 + \|r_0 - \widehat{r_0}\|_{L^2(\mathbb{Z})}^2 + 2\left(M^2 + \mathbb{E}[Y^2]\right)\right).$$

<span style="color:red">We still have to use convergence results for $\widehat{\mathcal{P}}$ and $\widehat{r_0}$ to finish this bound. It doesn't need to be good, we only need to bound this by something which remains bounded as $|\mathcal{D}_{\Phi,\mathcal{P},r_0}|$ and the number of iterations grow. Another idea is to simply say that this whole thing is $\mathcal{O}_p(1)$, that is, almost surely finite, and rely on the (fast enough) decay of the learning rate to achieve convergence.</span>

**Third term**

Our goal is to open up the inner product and make explicit the estimation errors of our model's different components, like we did before. Here, we define $\Psi_m(Z) \triangleq \partial_2\ell(r_0(Z), \mathcal{P}[\widehat{h}_{m-1}](Z))$. The hat version $\widehat{\Psi}_m$ is defined accordingly, replacing $r_0$ and $\mathcal{P}$ by their estimators.

$$
\begin{aligned}
\mathbb{E}_{\boldsymbol{z}_{1:m}}&\left[\langle\nabla\mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^\star\rangle \mid \mathcal{D}_{\Phi,\mathcal{P},r_0}\right] \\
&= \mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\langle\nabla\mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^\star\rangle\right] && (\boldsymbol{z}_{1:m} \perp\!\!\!\perp \mathcal{D}_{\Phi,\mathcal{P},r_0}) \\
&= \mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_{\boldsymbol{z}_m}\left[\langle\nabla\mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^\star\rangle\right]\right] && (\boldsymbol{z}_m \perp\!\!\!\perp \boldsymbol{z}_{1:m-1}) \\
&= \mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\langle\nabla\mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\boldsymbol{z}_m}[u_m], \widehat{h}_{m-1} - h^\star\rangle\right] \\
&\leq \mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|\nabla\mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\boldsymbol{z}_m}[u_m]\right\|\left\|\widehat{h}_{m-1} - h^\star\right\|\right] \\
&\leq D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|\nabla\mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\boldsymbol{z}_m}[u_m]\right\|\right] && (\mathrm{diam}\,\mathcal{F} = D) \\
&\leq D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_X\left[\left(\nabla\mathcal{R}(\widehat{h}_{m-1})(X) - \mathbb{E}_{\boldsymbol{z}_m}[u_m]\right)^2\right]\right]^{\frac{1}{2}} && (\text{Jensen}) \\
&= D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\Bigg[\mathbb{E}_X\Bigg[\Big(\mathbb{E}_Z\left[\Phi(X,Z)\Psi_m(Z)\right] \\
&\qquad\qquad\qquad - \mathbb{E}_{\boldsymbol{z}_m}\left[\widehat{\Phi}(X,\boldsymbol{z}_m)\widehat{\Psi}_m(\boldsymbol{z}_m)\right]\Big)^2\Bigg]\Bigg]^{\frac{1}{2}} \\
&= D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_X\left[\mathbb{E}_Z\left[\Phi(X,Z)\Psi_m(Z) - \widehat{\Phi}(X,Z)\widehat{\Psi}_m(Z)\right]^2\right]\right]^{\frac{1}{2}} && (Z \overset{\mathrm{iid}}{\sim} \boldsymbol{z}_m)
\end{aligned}
$$

$$= D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_X\left[\mathbb{E}_Z\left[\Psi_m(Z)\left(\Phi(X,Z)-\widehat{\Phi}(X,Z)\right)\right.\right.\right.$$

$$\left.\left.\left.+\,\widehat{\Phi}(X,Z)\left(\Psi_m(Z)-\widehat{\Psi}_m(Z)\right)\right]^2\right]\right]^{\frac{1}{2}}$$

$$\leq D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_X\left[\left(\|\Psi_m\|_{L^2(Z)}\left\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}\right.\right.\right.$$

$$\left.\left.\left.+\left\|\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}\left\|\Psi_m-\widehat{\Psi}_m\right\|_{L^2(Z)}\right)^2\right]\right]^{\frac{1}{2}}$$

$$\leq \sqrt{2}D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_X\left[\|\Psi_m\|_{L^2(Z)}^2\left\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\right.\right.$$

$$\left.\left.+\left\|\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\left\|\Psi_m-\widehat{\Psi}_m\right\|_{L^2(Z)}^2\right]\right]^{\frac{1}{2}}$$

$$= \sqrt{2}D\left(\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\|\Psi_m\|_{L^2(Z)}^2\mathbb{E}_X\left[\left\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\right]\right]\right.$$

$$\left.+\,\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|\Psi_m-\widehat{\Psi}_m\right\|_{L^2(Z)}^2\mathbb{E}_X\left[\left\|\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\right]\right]\right)^{\frac{1}{2}}$$

$$= \sqrt{2}D\left(\mathbb{E}_X\left[\left\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\right]\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\|\Psi_m\|_{L^2(Z)}^2\right]\right.$$

$$\left.+\,\mathbb{E}_X\left[\left\|\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\right]\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|\Psi_m-\widehat{\Psi}_m\right\|_{L^2(Z)}^2\right]\right)^{\frac{1}{2}}$$

$$=: \sqrt{2}D(A+B)^{\frac{1}{2}}.$$

We proceed to analyze each term separately:

- To bound $A$, first notice that

$$\mathbb{E}_X\left[\left\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\right\|^2\right]=\mathbb{E}_X\left[\mathbb{E}_Z\left[\left(\Phi(X,Z)-\widehat{\Phi}(X,Z)\right)^2\right]\right]=\left\|\Phi-\widehat{\Phi}\right\|_{L^2(X\otimes Z)}^2,$$

where $L^2(X\otimes Z)$ is the space of square integrable functions with respect to the measure induced by independent copies of $X$ and $Z$. If we estimate $\widehat{\Phi}$ using the uLSIF algorithm described in [1], under some regularity conditions, and decreasing the regularization parameter according to a specific rate, we have the following estimate:

> Create section describing how we are estimating each term.

$$\left\|\Phi-\widehat{\Phi}\right\|_{L^2(X\otimes Z)}^2=\mathcal{O}_p\left(\left(\frac{\log|\mathcal{D}_\Phi|}{|\mathcal{D}_\Phi|}\right)^{\frac{2}{2+\gamma}}\right).$$

Furthermore, we can bound $\|\Psi_m\|_{L^2(Z)}^2$ as follows:

$$\|\Phi_m\|_{L^2(Z)}^2=\left\|r_0-\mathcal{P}[\widehat{h}_{m-1}]\right\|_{L^2(Z)}^2$$

$$\leq 2\left(\|r_0\|_{L^2(Z)}^2+\left\|\mathcal{P}[\widehat{h}_{m-1}]\right\|_{L^2(Z)}^2\right)$$

$$\leq 2\left(\mathbb{E}[Y^2]+\|\mathcal{P}\|_{\text{op}}^2\left\|\widehat{h}_{m-1}\right\|_{L^2(Z)}^2\right)$$

$$\leq 2\left(\mathbb{E}[Y^2]+M^2\right) \qquad\qquad (\|\mathcal{P}\|_{\text{op}}\leq 1).$$

109    In total, what we have is

$$
\begin{aligned}
A &= \mathbb{E}_X\left[\left\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\right]\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\|\Psi_m\|_{L^2(Z)}^2\right] \\
&\le \left\|\Phi-\widehat{\Phi}\right\|_{L^2(Z)}^2\cdot 2(\mathbb{E}[Y^2]+M^2) \\
&= \mathcal{O}_p\left(\left(\frac{\log|\mathcal{D}_\Phi|}{|\mathcal{D}_\Phi|}\right)^{\frac{2}{2+\gamma}}\right).
\end{aligned}
$$

110    • To bound $B$, notice that, by Assumption 14.15 of [1], we have

$$
\mathbb{E}_X\left[\left\|\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\right]=\mathbb{E}_X\left[\mathbb{E}_Z\left[\widehat{\Phi}(X,Z)^2\right]\right]\le\left\|\widehat{\Phi}\right\|_{\mathcal{R}_{\mathbb{W}}}^2.
$$

112    Furthermore, we also have

$$
\begin{aligned}
\left\|\Psi_m-\widehat{\Psi}_m\right\|_{L^2(Z)}^2 &= \left\|\left(\mathcal{P}[\widehat{h}_{m-1}]-r_0\right)-\left(\widehat{\mathcal{P}}[\widehat{h}_{m-1}]-\widehat{r}_0\right)\right\|_{L^2(Z)}^2 \\
&= \left\|\left(\mathcal{P}[\widehat{h}_{m-1}]-\widehat{\mathcal{P}}[\widehat{h}_{m-1}]\right)-(r_0-\widehat{r}_0)\right\|_{L^2(Z)}^2 \\
&\le 2\left(\left\|\mathcal{P}[\widehat{h}_{m-1}]-\widehat{\mathcal{P}}[\widehat{h}_{m-1}]\right\|_{L^2(Z)}^2+\|r_0-\widehat{r}_0\|_{L^2(Z)}^2\right) \\
&\le 2\left(\left\|\mathcal{P}-\widehat{\mathcal{P}}\right\|_{\mathrm{op}}^2\left\|\widehat{h}_{m-1}\right\|_{L^2(Z)}^2+\|r_0-\widehat{r}_0\|_{L^2(Z)}^2\right) \\
&\le 2\left(M^2\left\|\mathcal{P}-\widehat{\mathcal{P}}\right\|_{\mathrm{op}}^2+\|r_0-\widehat{r}_0\|_{L^2(Z)}^2\right).
\end{aligned}
$$

113    Therefore,

$$
\begin{aligned}
B &= \mathbb{E}_X\left[\left\|\widehat{\Phi}(X,\cdot)\right\|_{L^2(Z)}^2\right]\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|\Psi_m-\widehat{\Psi}_m\right\|_{L^2(Z)}^2\right] \\
&\le \left\|\widehat{\Phi}\right\|_{\mathcal{R}_{\mathbb{W}}}^2\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[2\left(M^2\left\|\mathcal{P}-\widehat{\mathcal{P}}\right\|_{\mathrm{op}}^2+\|r_0-\widehat{r}_0\|_{L^2(Z)}^2\right)\right] \\
&= 2\left\|\widehat{\Phi}\right\|_{\mathcal{R}_{\mathbb{W}}}^2\left(M^2\left\|\mathcal{P}-\widehat{\mathcal{P}}\right\|_{\mathrm{op}}^2+\|r_0-\widehat{r}_0\|_{L^2(Z)}^2\right).
\end{aligned}
$$

- Find way to estimate $r_0$ which gives estimate on $\|r_0 - \widehat{r_0}\|_{L^2(Z)}$. Maybe use the same estimation technique we have for $\mathcal{P}$ as an operator from $L^2(Y) \to L^2(Z)$ applied to the identity and employ the same bound?

For the rest of the paper:

- Create section which describes, in detail, how we are estimating $\Phi$, $\mathcal{P}$ and $r_0$, lists all the references, states the main convergence theorems and lists all of the assumptions that are being made.
- Adapt the algorithm section to use the KIV first stage, which directly estimates $\mathcal{P}$.
- Find better letter for either the number of iterations or the upper bound for the set $\mathcal{F}$. Right now, both are being denoted by the letter $M$.

# References

[1] Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density Ratio Estimation in Machine Learning*. Cambridge University Press, 2012. DOI: 10.1017/CBO9781139035613.