# 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 $L^2(Z) = L^2(\mathcal{Z}, \mathcal{B}(\mathcal{Z}), \nu_Z)$ 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

Assumption

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

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 $\ell$:

**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. The function $\ell$ has Lipschitz first derivative with respect to the second coordinate, i.e., there exists $L \geq 0$ such that, for all $y, y', u, u' \in \mathbf{R}$ we have*

$$|\partial_2\ell(y, y') - \partial_2\ell(u, u')| \leq L(|y - u| + |y' - u'|).$$

Some useful facts which follow immediately from these assumptions are:

**Proposition 1.** *Under Assumption 1 we have:*

    *1. Setting $C_0 = |\partial_2\ell(0, 0)|$ we have*

$$|\partial_2\ell(y, y')| \leq C_0 + L(|y| + |y'|)$$

    *for all $y, y' \in \mathbf{R}$;*

    *2. The map $f \mapsto \partial_2\ell(r_0(\cdot), f(\cdot))$ from $L^2(Z)$ to $L^2(Z)$ is well-defined and $L$-Lipschitz.*

    *3. The second derivative with respect to the second coordinate is bounded: $\left|\partial_2^2\ell(y, y')\right| \leq L$ for all $y, y' \in \mathbf{R}$;*

*Proof.*

    1. Write $\partial_2\ell(y, y') = \partial_2\ell(y, y') - \partial_2\ell(0, 0) + \partial_2\ell(0, 0)$ and apply the triangle inequality as well as Assumption 1.2.

    2. From the previous item we know this map is well-defined. If $f$ and $g$ belong to $L^2(Z)$, we have

$$\begin{aligned}
\|\partial_2\ell(r_0(\cdot), f(\cdot)) - \partial_2\ell(r_0(\cdot), g(\cdot))\|^2_{L^2(Z)} &= \mathbb{E}\left[|\partial_2(r_0(Z), f(Z)) - \partial_2(r_0(Z), g(Z))|^2\right] \\
&\leq L^2\mathbb{E}\left[|f(Z) - g(Z)|^2\right] \\
&= L^2\|f - g\|^2_{L^2(Z)}.
\end{aligned}$$

    3. Follows from the definition of derivative and Assumption 1.2.

$\square$

## 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 $D\mathcal{R}[h](f)$:

$$
\begin{aligned}
D\mathcal{R}[h](f) &= \lim_{\delta\to0}\frac{1}{\delta}\left[\mathcal{R}(h+\delta f) - \mathcal{R}(f)\right] \\
&= \lim_{\delta\to0}\frac{1}{\delta}\mathbb{E}\left[\ell(r_0(Z), \mathcal{P}[h+\delta f](Z)) - \ell(r_0(Z), \mathcal{P}[h](Z))\right] \\
&= \lim_{\delta\to0}\frac{1}{\delta}\mathbb{E}\left[\ell(r_0(Z), \mathcal{P}[h](Z) + \delta\mathcal{P}[f](Z)) - \ell(r_0(Z), \mathcal{P}[h](Z))\right] \\
&= \lim_{\delta\to0}\frac{1}{\delta}\mathbb{E}\left[\delta\partial_2\ell(r_0(Z), \mathcal{P}[h](Z))\cdot\mathcal{P}[f](Z)\right. \\
&\qquad\qquad\left. + \frac{\delta^2}{2}\partial_2^2\ell(r_0(Z), \mathcal{P}[h+\theta f](Z))\cdot\mathcal{P}[f](Z)^2\right] \\
&= \mathbb{E}\left[\partial_2\ell(r_0(Z), \mathcal{P}[h](Z))\cdot\mathcal{P}[f](Z)\right] \\
&\qquad\qquad + \lim_{\delta\to0}\mathbb{E}\left[\frac{\delta}{2}\partial_2^2\ell(r_0(Z), \mathcal{P}[h+\theta f](Z))\cdot\mathcal{P}[f](Z)^2\right] \\
&= \mathbb{E}\left[\partial_2\ell(r_0(Z), \mathcal{P}[h](Z))\cdot\mathcal{P}[f](Z)\right],
\end{aligned}
$$

where $\theta \in \mathbf{R}$ is due to Taylor's formula. The last step is then due to Proposition 1.3.

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

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

This shows that $\mathcal{R}$ is Gateux-differentiable, with Gateux derivative at $h$ given by

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

By Proposition 1.2 we have that $h \mapsto D\mathcal{R}[h]$ is a continuous mapping from $L^2(X)$ to $L^2(X)$, which implies that $\mathcal{R}$ is also Fréchet-differentiable, and both derivatives coincide. Therefore,

<div style="text-align:right">Cite a reference for this.</div>

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

## 3  Estimating the gradient

We have found that

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

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

<div style="text-align:right">Assumption</div>

$$
\begin{aligned}
\nabla\mathcal{R}(h)(x) &= \mathbb{E}[\Psi_h(Z) \mid X = x] \\
&= \int_{\mathbb{Z}} p(z \mid x)\Psi_h(z)\,\mathrm{d}z \\
&= \int_{\mathbb{Z}} p(z)\frac{p(z \mid x)}{p(z)}\Psi_h(z)\,\mathrm{d}z \\
&= \mathbb{E}\left[\frac{p(Z \mid x)}{p(Z)}\Psi_h(Z)\right].
\end{aligned}
$$

69 Thus, we must take

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

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

71 An obvious obstacle for this approach is that we don't know how to analytically compute $\Phi, r_0$ nor $\mathcal{P}$,
72 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
73 form taken by these estimators and will present the algorithm assuming we know how to compute
74 them. Later, we will show how the individual convergence rates of these three pieces come together
75 to determine the convergence rate of our method.

76 We state here all the assumptions which we need from these estimators to bound the excess risk:

77 **Assumption 2.**

78     *1. $\widehat{r_0} \in L^2(Z)$;*

79     *2. $\widehat{\mathcal{P}} : L^2(X) \to L^2(Z)$ is a bounded linear operator;*

80     *3. Letting $\mathcal{W} = \mathcal{X} \times \mathcal{Z}$, we have*

$$\left\|\widehat{\Phi}\right\|_\infty \triangleq \sup_{\boldsymbol{w} \in \mathcal{W}} |\Phi(\boldsymbol{w})| < \infty.$$

# 4 Algorithm

82 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 :** $\widehat{h}$

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

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

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

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

**end**

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

---

# 5 Proof of convergence

85 We start by proving that our sequence of estimates is, in fact, contained in $L^2(X)$. This is clear, since,
86 by Assumption 2 we have:

$$
\begin{aligned}
\|u_m\|_{L^2(X)}^2 &= \left\|\widehat{\Phi}(\cdot, \boldsymbol{z}_m)\partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)\right\|_{L^2(X)}^2 \\
&= \mathbb{E}_X\left[\left|\widehat{\Phi}(X, \boldsymbol{z}_m)\partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)\right|^2\right] \\
&\leq \partial_2\ell\left(\widehat{r_0}(\boldsymbol{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\boldsymbol{z}_m)\right)^2\left\|\widehat{\Phi}\right\|_\infty^2 \qquad (5) \\
&< \infty.
\end{aligned}
$$

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

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

88 To lighten the notation, the symbols $\|\cdot\|$ and $\langle \cdot, \cdot \rangle$, when written without a subscript to specify which
89 space they refer to, will act as the norm and inner product, respectively, of $L^2(X)$. By the Algorithm
90 1 procedure, we have

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

91 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

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

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

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

93 Rearranging terms, we get

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

94 Summing over $1 \leq m \leq M$ leads to

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

95 Define $\varepsilon_m \triangleq u_m - \nabla \mathcal{R}(\widehat{h}_{m-1})$, so what we have is

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

96 The next step is to take the average of both sides with respect to $z_{1:M}$, taking advantage of the
97 independence between $z_{1:M}$ and $\mathcal{D}_{r_0, \Phi, \mathcal{P}}$. To better organize the argument, we treat each of the three
98 summations in the RHS of the inequality above separately:

5

**First summation** This one admits a deterministic bound. 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 summation** The computation in equation (5) shows that

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

Together with Proposition 1.1, this implies

$$\mathbb{E}_{\boldsymbol{z}_{1:M}}\left[\|u_m\|^2\right] \leq 3\left\|\widehat{\Phi}\right\|_\infty^2\left(C_0^2 + L^2\left(\|\widehat{r_0}\|_{L^2(Z)}^2 + \left\|\widehat{\mathcal{P}}[\widehat{h}_{m-1}]\right\|_{L^2(Z)}^2\right)\right)$$

$$\leq 3\left\|\widehat{\Phi}\right\|_\infty^2\left(C_0^2 + L^2\left(\|\widehat{r_0}\|_{L^2(Z)}^2 + \left\|\widehat{\mathcal{P}}\right\|_{\mathrm{op}}^2\left\|\widehat{h}_{m-1}\right\|^2\right)\right)$$

$$\leq 3\left\|\widehat{\Phi}\right\|_\infty^2\left(C_0^2 + L^2\left(\|\widehat{r_0}\|_{L^2(Z)}^2 + D^2\left\|\widehat{\mathcal{P}}\right\|_{\mathrm{op}}^2\right)\right).$$

Therefore,

$$\mathbb{E}_{\boldsymbol{z}_{1:M}}\left[\sum_{m=1}^{M}\frac{\alpha_m}{2}\|u_m\|^2\right] \leq \frac{3}{2}\left\|\widehat{\Phi}\right\|_\infty^2\left(C_0^2 + L^2\left(\|\widehat{r_0}\|_{L^2(Z)}^2 + D^2\left\|\widehat{\mathcal{P}}\right\|_{\mathrm{op}}^2\right)\right)\sum_{m=1}^{M}\alpha_m.$$

**Third summation**

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.

$$\mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\langle\nabla\mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^\star\rangle \mid \mathcal{D}_{\Phi,\mathcal{P},r_0}\right]$$

$$= \mathbb{E}_{\boldsymbol{z}_{1:m}}\left[\langle\nabla\mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^\star\rangle\right] \qquad (\boldsymbol{z}_{1:m} \perp\!\!\!\perp \mathcal{D}_{\Phi,\mathcal{P},r_0})$$

$$= \mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_{\boldsymbol{z}_m}\left[\langle\nabla\mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^\star\rangle\right]\right] \qquad (\boldsymbol{z}_m \perp\!\!\!\perp \boldsymbol{z}_{1:m-1})$$

$$= \mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\langle\nabla\mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\boldsymbol{z}_m}[u_m], \widehat{h}_{m-1} - h^\star\rangle\right]$$

$$\leq \mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|\nabla\mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\boldsymbol{z}_m}[u_m]\right\|\left\|\widehat{h}_{m-1} - h^\star\right\|\right]$$

$$\leq D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\left\|\nabla\mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\boldsymbol{z}_m}[u_m]\right\|\right] \qquad (\operatorname{diam}\mathcal{F} = D)$$

$$\leq D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_X\left[\left(\nabla\mathcal{R}(\widehat{h}_{m-1})(X) - \mathbb{E}_{\boldsymbol{z}_m}[u_m]\right)^2\right]\right]^{\frac{1}{2}} \qquad (\text{Jensen})$$

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

$$\left.\left.\left. - \mathbb{E}_{\boldsymbol{z}_m}\left[\widehat{\Phi}(X,\boldsymbol{z}_m)\widehat{\Psi}_m(\boldsymbol{z}_m)\right]\right)^2\right]\right]^{\frac{1}{2}}$$

$$= D\mathbb{E}_{\boldsymbol{z}_{1:m-1}}\left[\mathbb{E}_X\left[\mathbb{E}_Z\left[\Phi(X,Z)\Psi_m(Z) - \widehat{\Phi}(X,Z)\widehat{\Psi}_m(Z)\right]^2\right]\right]^{\frac{1}{2}} \qquad (Z \overset{\mathrm{iid}}{\sim} \boldsymbol{z}_m)$$

$$
= D\mathbb{E}_{\mathbf{z}_{1:m-1}}\Bigg[\mathbb{E}_X\Big[\mathbb{E}_Z\Big[\Psi_m(Z)\Big(\Phi(X,Z)-\widehat{\Phi}(X,Z)\Big)
$$

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

$$
\leq D\mathbb{E}_{\mathbf{z}_{1:m-1}}\Bigg[\mathbb{E}_X\Big[\Big(\|\Psi_m\|_{L^2(Z)}\big\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\big\|_{L^2(Z)}
$$

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

$$
\leq \sqrt{2}D\mathbb{E}_{\mathbf{z}_{1:m-1}}\Bigg[\mathbb{E}_X\Big[\|\Psi_m\|^2_{L^2(Z)}\big\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\big\|^2_{L^2(Z)}
$$

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

$$
= \sqrt{2}D\Bigg(\mathbb{E}_{\mathbf{z}_{1:m-1}}\Big[\|\Psi_m\|^2_{L^2(Z)}\mathbb{E}_X\Big[\big\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\big\|^2_{L^2(Z)}\Big]\Big]
$$

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

$$
= \sqrt{2}D\Bigg(\mathbb{E}_X\Big[\big\|\Phi(X,\cdot)-\widehat{\Phi}(X,\cdot)\big\|^2_{L^2(Z)}\Big]\mathbb{E}_{\mathbf{z}_{1:m-1}}\Big[\|\Psi_m\|^2_{L^2(Z)}\Big]
$$

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

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

We proceed to analyze each term separately:

- To bound $A$, first notice that

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

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.

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

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

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

In total, what we have is

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

$$\leq \left\| \Phi - \widehat{\Phi} \right\|_{L^2(Z)}^2 \cdot 2(\mathbb{E}[Y^2] + M^2)$$

$$= \mathcal{O}_p \left( \left( \frac{\log |\mathcal{D}_\Phi|}{|\mathcal{D}_\Phi|} \right)^{\frac{2}{2+\gamma}} \right).$$

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

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

<span style="color:red">We still need to bound this norm somehow.</span>

Furthermore, we also have

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

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

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

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

$$\leq 2 \left( M^2 \left\| \mathcal{P} - \widehat{\mathcal{P}} \right\|_{op}^2 + \| r_0 - \widehat{r_0} \|_{L^2(Z)}^2 \right).$$

Therefore,

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

$$\leq \left\| \widehat{\Phi} \right\|_{\mathcal{R}_\mathbb{W}}^2 \mathbb{E}_{\boldsymbol{z}_{1:m-1}} \left[ 2 \left( M^2 \left\| \mathcal{P} - \widehat{\mathcal{P}} \right\|_{op}^2 + \| r_0 - \widehat{r_0} \|_{L^2(Z)}^2 \right) \right]$$

$$= 2 \left\| \widehat{\Phi} \right\|_{\mathcal{R}_\mathbb{W}}^2 \left( M^2 \left\| \mathcal{P} - \widehat{\mathcal{P}} \right\|_{op}^2 + \| r_0 - \widehat{r_0} \|_{L^2(Z)}^2 \right).$$

<span style="color:red">What's left to do:</span>

- <span style="color:red">Bound $\left\| \widehat{\Phi} \right\|_{\mathcal{R}_\mathbb{W}}$. (May not be strictly necessary. This is finite, and since it multiplies something which is $\mathcal{O}_p$ of something which goes to zero, we may not need to further bound it.)</span>

- <span style="color:red">Use some estimate on $\left\| \mathcal{P} - \widehat{\mathcal{P}} \right\|_{op}$ (Adapt notation and setup in the KIV paper).</span>

<span style="color:red">Conclusion (20/08/2023): We might need the extra hypothesis that $\mathrm{Im}(\mathrm{id}_{L^2(X)} - \iota_X \iota_X^*) \subseteq \ker \mathcal{P}$, where $\iota_X : \mathcal{H}_X \to L^2(X)$ is the inclusion operator, whose adjoint is given by</span>

$$\iota_X^*(f) = (x \mapsto \mathbb{E}_X[f(X) k_X(X, x)]),$$

<span style="color:red">with $k_X : \mathbb{X} \times \mathbb{X} \to \mathbf{R}$ being the kernel associated with $\mathcal{H}_X$. Then $\mathcal{P} = \mathcal{P} \circ \iota_X \iota_X^*$ and we can directly apply the result on KIV's paper, since $\mathcal{P} \circ i_X$ can be seen as the restriction of $\mathcal{P}$ to $\mathcal{H}_X$. We then also need the further hypothesis that $\mathrm{Im}(\mathcal{P} \circ \iota_X) \subseteq \mathcal{H}_Z$, or something like this (because, rigorously speaking, $\mathcal{P}f$ is an equivalence class of functions, so in what way can we say that this equivalence class is "in $\mathcal{H}_Z$"?). This hypothesis is implicitly made in the KIV paper, when they say that $E : \mathcal{H}_\mathcal{X} \to \mathcal{H}_\mathcal{Z}$ without providing any assumptions on $\mathcal{H}_X$ and $\mathcal{H}_Z$, other than saying that they are RKHS. Who can guarantee that $(z \mapsto \mathbb{E}[f(X) \mid Z = z]) \in \mathcal{H}_Z$ for every $f \in \mathcal{H}_X$?</span>

8

- 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.