
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

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

that is, $L^2(X) = L^2(\mathcal{X}, \mathcal{B}(\mathcal{X}), \nu_X)$, where we denote by ν_X the distribution of the r.v. X and by $\mathcal{B}(\mathcal{X})$ the Borel σ -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

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

where $h^* \in L^2(X)$ and ε is an square-integrable r.v. such that $\mathbb{E}[\varepsilon | X] \neq 0$. 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, $\nu_{X|Z}(\cdot | Z) \neq \nu_X(\cdot)$;

ii) Z influences Y only through X ;

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

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^* based on independent joint samples from X, Z and Y .

Conditioning (1) in Z , we find

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

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

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

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 $\mathcal{P}^* : L^2(Z) \rightarrow L^2(X)$ satisfies

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

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

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

Hence, (1) can be formulated as an inverse problem, where we wish to invert the operator \mathcal{P} .

1.2 Risk measure

Let $\ell : \mathbf{R} \times \mathbf{R} \rightarrow \mathbf{R}$ be a pointwise loss function, which, with respect to its second argument, is convex and differentiable. We use the symbol ∂_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 $h^* \in \mathcal{F}$. We also assume that $D \triangleq \text{diam } \mathcal{F} < \infty$ and that $0 \in \mathcal{F}$, so that $\|h\| \leq D$ if $h \in \mathcal{F}$. A possible choice for the set \mathcal{F} is

Assumption

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

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 ℓ :

Assumption 1 (Regularity of ℓ).

1. The function $\ell : \mathbf{R} \times \mathbf{R} \rightarrow \mathbf{R}$ is convex and C^2 with respect to its second argument;
2. The function ℓ has Lipschitz first derivative with respect to the second argument, i.e., there exists $L \geq 0$ such that, for all $y, y', u, u' \in \mathbf{R}$ we have

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

38

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 argument is bounded: $|\partial_2^2 \ell(y, y')| \leq L$ 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

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

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

□

53 2 Gradient computation

54 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}
 55 at h in the direction f , denoted by $D\mathcal{R}[h](f)$:

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

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

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

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

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

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

59 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
 60 implies that \mathcal{R} is also Fréchet-differentiable, and both derivatives coincide. Therefore,

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

Cite a reference for this.

61 3 Estimating the gradient

62 We have found that

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

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

Assumption

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

69 Thus, we must take

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

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

Must discuss why $u_h \in L^2(X)$.

71 An obvious obstacle for this approach is that we don't know how to analytically compute Φ , r_0 nor \mathcal{P} ,
72 se we will proceed with estimators $\hat{\Phi}$, \hat{r}_0 and $\hat{\mathcal{P}}$. 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.

Must we? Since we end up not using u_h , but an approximation which we know is in $L^2(X)$.

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

77 **Assumption 2.**

- 78 1. $\hat{r}_0 \in L^2(Z)$;
- 79 2. $\hat{\mathcal{P}} : L^2(X) \rightarrow L^2(Z)$ is a bounded linear operator;
- 80 3. Letting $\mathcal{W} = \mathcal{X} \times \mathcal{Z}$, we have

$$\|\hat{\Phi}\|_{\infty} \triangleq \sup_{\mathbf{w} \in \mathcal{W}} |\Phi(\mathbf{w})| < \infty.$$

81 4 Algorithm

82 Having an estimator of the gradient, we can construct Functional GD algorithm for estimating h^* .

Discuss everything we don't know and must estimate.

Algorithm 1: SGD-NPIV

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

Comment on exactly what is needed to estimate each unknown (samples from which r.v.'s).

output : \hat{h}

83 Compute $\hat{r}_0, \hat{\Phi}, \hat{\mathcal{P}}$ using $\mathcal{D}_{r_0}, \mathcal{D}_{\Phi}, \mathcal{D}_{\mathcal{P}}$, respectively ;

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

Set $u_m = \hat{\Phi}(\cdot, z_m) \partial_2 \ell(\hat{r}_0(z_m), \hat{\mathcal{P}}[\hat{h}_{m-1}](z_m))$;
Set $\hat{h}_m(\mathbf{x}_k) = \pi_{\mathcal{F}}[\hat{h}_{m-1} - \alpha_m u_m](\mathbf{x}_k)$ for $1 \leq k \leq K$;

end

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

Discuss necessity of discretizing \mathcal{X} .

84 5 Proof of convergence

85 To lighten the notation, the symbols $\|\cdot\|$ and $\langle \cdot, \cdot \rangle$, when written without a subscript to specify which
86 space they refer to, will act as the norm and inner product, respectively, of $L^2(X)$.

87 **Lemma 1.** In the procedure of Algorithm 1 we have $u_m \in L^2(X)$ for all m and, furthermore,

$$\mathbb{E}_{\mathbf{z}_{1:M}}[\|u_m\|^2] \leq \rho(\hat{\Phi}, \hat{r}_0, \hat{\mathcal{P}}),$$

88 where

$$\rho(\hat{\Phi}, \hat{r}_0, \hat{\mathcal{P}}) = 3\|\hat{\Phi}\|_{\infty}^2 \left(C_0^2 + L^2 \|\hat{r}_0\|_{L^2(Z)}^2 + L^2 D^2 \|\hat{\mathcal{P}}\|_{\text{op}}^2 \right).$$

89 *Proof.* By Assumption 2 we have:

$$\begin{aligned}
\|u_m\|_{L^2(X)}^2 &= \left\| \widehat{\Phi}(\cdot, \mathbf{z}_m) \partial_2 \ell \left(\widehat{r}_0(\mathbf{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) \right) \right\|_{L^2(X)}^2 \\
&= \mathbb{E}_X \left[\left| \widehat{\Phi}(X, \mathbf{z}_m) \partial_2 \ell \left(\widehat{r}_0(\mathbf{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) \right) \right|^2 \right] \\
&\leq \partial_2 \ell \left(\widehat{r}_0(\mathbf{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) \right)^2 \left\| \widehat{\Phi} \right\|_{\infty}^2 \\
&< \infty.
\end{aligned} \tag{5}$$

90 Hence, $u_m \in L^2(X)$ for all m . This computation and Proposition 1.1 then imply

$$\begin{aligned}
\mathbb{E}_{\mathbf{z}_{1:M}} \left[\|u_m\|^2 \right] &\leq 3 \left\| \widehat{\Phi} \right\|_{\infty}^2 \left(C_0^2 + L^2 \left(\|\widehat{r}_0\|_{L^2(Z)}^2 + \left\| \widehat{\mathcal{P}}[\widehat{h}_{m-1}] \right\|_{L^2(Z)}^2 \right) \right) \\
&\leq 3 \left\| \widehat{\Phi} \right\|_{\infty}^2 \left(C_0^2 + L^2 \left(\|\widehat{r}_0\|_{L^2(Z)}^2 + \left\| \widehat{\mathcal{P}} \right\|_{\text{op}}^2 \left\| \widehat{h}_{m-1} \right\|^2 \right) \right) \\
&\leq 3 \left\| \widehat{\Phi} \right\|_{\infty}^2 \left(C_0^2 + L^2 \left(\|\widehat{r}_0\|_{L^2(Z)}^2 + D^2 \left\| \widehat{\mathcal{P}} \right\|_{\text{op}}^2 \right) \right) \\
&= 3 \left\| \widehat{\Phi} \right\|_{\infty}^2 \left(C_0^2 + L^2 \|\widehat{r}_0\|_{L^2(Z)}^2 + L^2 D^2 \left\| \widehat{\mathcal{P}} \right\|_{\text{op}}^2 \right) \triangleq \rho \left(\widehat{\Phi}, \widehat{r}_0, \widehat{\mathcal{P}} \right). \quad \square
\end{aligned}$$

91 **Lemma 2.** In the procedure of Algorithm 1 we have

$$\left\| \mathbb{E}_{\mathbf{z}_m} \left[\nabla \mathcal{R}(\widehat{h}_{m-1}) - u_m \right] \right\| \leq \kappa \left(\widehat{\Phi} \right) \left(\left\| \Phi - \widehat{\Phi} \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 + \|r_0 - \widehat{r}_0\|_{L^2(Z)}^2 + \left\| \mathcal{P} - \widehat{\mathcal{P}} \right\|_{\text{op}}^2 \right)^{\frac{1}{2}},$$

Comment on how this is the step that is different from the other article, since in the simpler scenario, this difference would vanish.

92 where

$$\kappa^2 \left(\widehat{\Phi} \right) \triangleq 2 \max \left\{ 3(C_0^2 + L^2 \mathbb{E}[Y^2] + L^2 D^2), 2L^2 \left\| \widehat{\Phi} \right\|_{\infty}^2, 2L^2 D^2 \left\| \widehat{\Phi} \right\|_{\infty}^2 \right\}.$$

93 *Proof.* To ease the notation, we define

$$\begin{aligned}
\Psi_m(Z) &\triangleq \partial_2 \ell(r_0(Z), \mathcal{P}[\widehat{h}_{m-1}](Z)), \\
\widehat{\Psi}_m(Z) &\triangleq \partial_2 \ell(\widehat{r}_0(Z), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](Z)).
\end{aligned}$$

94 Let's expand the definition of $\|\cdot\|$:

$$\begin{aligned}
\left\| \mathbb{E}_{\mathbf{z}_m} \left[\nabla \mathcal{R}(\widehat{h}_{m-1}) - u_m \right] \right\| &= \mathbb{E}_X \left[\mathbb{E}_{\mathbf{z}_m} \left[\nabla \mathcal{R}(\widehat{h}_{m-1})(X) - u_m(X) \right]^2 \right]^{\frac{1}{2}} \\
&= \mathbb{E}_X \left[\left(\nabla \mathcal{R}(\widehat{h}_{m-1})(X) - \mathbb{E}_{\mathbf{z}_m} [u_m(X)] \right)^2 \right]^{\frac{1}{2}} \\
&= \mathbb{E}_X \left[\left(\mathbb{E}_Z [\Phi(X, Z) \Psi_m(Z)] - \mathbb{E}_{\mathbf{z}_m} [\widehat{\Phi}(X, \mathbf{z}_m) \widehat{\Psi}_m(\mathbf{z}_m)] \right)^2 \right]^{\frac{1}{2}} \\
&= \mathbb{E}_X \left[\left(\mathbb{E}_Z [\Phi(X, Z) \Psi_m(Z) - \widehat{\Phi}(X, Z) \widehat{\Psi}_m(Z)] \right)^2 \right]^{\frac{1}{2}},
\end{aligned}$$

95 Now we add and subtract $\widehat{\Phi}(X, Z)\Psi_m(Z)$, so that

$$\begin{aligned}
& \mathbb{E}_X \left[\left(\mathbb{E}_Z \left[\Phi(X, Z)\Psi_m(Z) - \widehat{\Phi}(X, Z)\widehat{\Psi}_m(Z) \right] \right)^2 \right]^{\frac{1}{2}} \\
&= \mathbb{E}_X \left[\left(\mathbb{E}_Z \left[\Psi_m(Z) \left(\Phi(X, Z) - \widehat{\Phi}(X, Z) \right) + \widehat{\Phi}(X, Z) \left(\Psi_m(Z) - \widehat{\Psi}_m(Z) \right) \right] \right)^2 \right]^{\frac{1}{2}} \\
&\leq \mathbb{E}_X \left[\left(\left\| \Psi_m \right\|_{L^2(Z)} \left\| \Phi(X, \cdot) - \widehat{\Phi}(X, \cdot) \right\|_{L^2(Z)} + \left\| \widehat{\Phi}(X, \cdot) \right\|_{L^2(Z)} \left\| \Psi_m - \widehat{\Psi}_m \right\|_{L^2(Z)} \right)^2 \right]^{\frac{1}{2}} \\
&\leq \sqrt{2} \mathbb{E}_X \left[\left\| \Psi_m \right\|_{L^2(Z)}^2 \left\| \Phi(X, \cdot) - \widehat{\Phi}(X, \cdot) \right\|_{L^2(Z)}^2 + \left\| \widehat{\Phi}(X, \cdot) \right\|_{L^2(Z)}^2 \left\| \Psi_m - \widehat{\Psi}_m \right\|_{L^2(Z)}^2 \right]^{\frac{1}{2}} \\
&= \sqrt{2} \left(\left\| \Psi_m \right\|_{L^2(Z)}^2 \left\| \Phi - \widehat{\Phi} \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 + \left\| \widehat{\Phi} \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 \left\| \Psi_m - \widehat{\Psi}_m \right\|_{L^2(Z)}^2 \right)^{\frac{1}{2}},
\end{aligned}$$

96 where

$$\left\| \Phi \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 = \int_{\mathcal{X} \times \mathcal{Z}} \Phi(x, z)^2 p(x) p(z) \, dx dz$$

97 is the norm with respect to the independent coupling of the distributions of X and Z . By Proposition
98 1.1 we have

$$\begin{aligned}
\left\| \Psi_m \right\|_{L^2(Z)}^2 &= \mathbb{E}_Z \left[\partial_2 \ell(r_0(Z), \mathcal{P}[\widehat{h}_{m-1}](Z))^2 \right] \\
&\leq \mathbb{E}_Z \left[\left(C_0 + L \left(|r_0(Z)| + \left| \mathcal{P}[\widehat{h}_{m-1}](Z) \right| \right) \right)^2 \right] \\
&\leq 3 \left(C_0^2 + L^2 \|r_0\|_{L^2(Z)}^2 + L^2 \left\| \mathcal{P}[\widehat{h}_{m-1}] \right\|_{L^2(Z)}^2 \right) \\
&\leq 3 (C_0^2 + L^2 \mathbb{E}[Y^2] + L^2 D^2).
\end{aligned}$$

99 It is also clear that, by Assumption 2,

$$\left\| \widehat{\Phi} \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 \leq \left\| \widehat{\Phi} \right\|_{\infty}^2.$$

100 Finally, by Assumption 1.2 we also have

$$\begin{aligned}
\left\| \Psi_m - \widehat{\Psi}_m \right\|_{L^2(Z)}^2 &= \mathbb{E}_Z \left[\left(\partial_2 \ell(r_0(Z), \mathcal{P}[\widehat{h}_{m-1}](Z)) - \partial_2 \ell(\widehat{r}_0(Z), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](Z)) \right)^2 \right] \\
&\leq 2L^2 \left(\|r_0 - \widehat{r}_0\|_{L^2(Z)}^2 + \left\| (\mathcal{P} - \widehat{\mathcal{P}})[\widehat{h}_{m-1}] \right\|_{L^2(Z)}^2 \right) \\
&\leq 2L^2 \left(\|r_0 - \widehat{r}_0\|_{L^2(Z)}^2 + D^2 \left\| \mathcal{P} - \widehat{\mathcal{P}} \right\|_{\text{op}}^2 \right).
\end{aligned}$$

101 To combine all terms, we first define

$$\kappa^2(\widehat{\Phi}) \triangleq 2 \max \left\{ 3(C_0^2 + L^2 \mathbb{E}[Y^2] + L^2 D^2), 2L^2 \left\| \widehat{\Phi} \right\|_{\infty}^2, 2L^2 D^2 \left\| \widehat{\Phi} \right\|_{\infty}^2 \right\}.$$

102 Then, it's easy to see that

$$\left\| \mathbb{E}_{\mathbf{z}_m} \left[\nabla \mathcal{R}(\widehat{h}_{m-1}) - u_m \right] \right\| \leq \kappa(\widehat{\Phi}) \left(\left\| \Phi - \widehat{\Phi} \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 + \|r_0 - \widehat{r}_0\|_{L^2(Z)}^2 + \left\| \mathcal{P} - \widehat{\mathcal{P}} \right\|_{\text{op}}^2 \right)^{\frac{1}{2}},$$

103 as we wanted to show. \square

104 **Theorem 1.** Assume that $\widehat{h}_0, \dots, \widehat{h}_{M-1}$ are generated according to Algorithm 1. If we let $\widehat{h} =$
105 $\sum_{m=1}^M \widehat{h}_{m-1}$, the following excess risk bound holds:

$$\begin{aligned}
\mathbb{E}_{\mathbf{z}_{1:M}} \left[\mathcal{R}(\widehat{h}) - \mathcal{R}(h^*) \right] &\leq \frac{D^2}{2M\alpha_m} + \alpha \left(\widehat{\Phi}, \widehat{r}_0, \widehat{\mathcal{P}} \right) \frac{1}{M} \sum_{m=1}^M \alpha_m \\
&\quad + \tau(\widehat{\Phi}) \left(\left\| \Phi - \widehat{\Phi} \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 + \|r_0 - \widehat{r}_0\|_{L^2(Z)}^2 + \left\| \mathcal{P} - \widehat{\mathcal{P}} \right\|_{\text{op}}^2 \right)^{\frac{1}{2}},
\end{aligned}$$

106 where

$$\begin{aligned}\alpha(\hat{\Phi}, \hat{r}_0, \hat{\mathcal{P}}) &= \frac{3}{2} \|\hat{\Phi}\|_\infty^2 \left(C_0^2 + L^2 \|\hat{r}_0\|_{L^2(Z)}^2 + L^2 D^2 \|\hat{\mathcal{P}}\|_{\text{op}}^2 \right), \\ \tau(\hat{\Phi}) &= 2D \max \left\{ 3(C_0^2 + L^2 \mathbb{E}[Y^2] + L^2 D^2), 2L^2 \|\hat{\Phi}\|_\infty^2, 2L^2 D^2 \|\hat{\Phi}\|_\infty^2 \right\}.\end{aligned}$$

107 *Proof.* We start by checking that \mathcal{R} is convex in \mathcal{F} : if $h, g \in \mathcal{F}$ and $\lambda \in [0, 1]$, then

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

108 By the Algorithm 1 procedure, we have

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

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

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

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

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

111 Rearranging terms, we get

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

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

$$\begin{aligned}\sum_{n=1}^M [\mathcal{R}(\hat{h}_{m-1}) - \mathcal{R}(h^*)] &\leq \sum_{m=1}^M \frac{1}{2\alpha_m} \left(\|\hat{h}_{m-1} - h^*\|^2 - \|\hat{h}_m - h^*\|^2 \right) \\ &\quad + \sum_{m=1}^M \frac{\alpha_m}{2} \|u_m\|^2 \\ &\quad + \sum_{m=1}^M \langle \nabla \mathcal{R}(\hat{h}_{m-1}) - u_m, \hat{h}_{m-1} - h^* \rangle.\end{aligned}\tag{6}$$

113 The next step is to take the average of both sides with respect to $\mathbf{z}_{1:M}$, taking advantage of the
114 independence between $\mathbf{z}_{1:M}$ and $\mathcal{D}_{r_0, \Phi, \mathcal{P}}$. Each summation in the RHS is then bounded separately.

115 The first summation admits a deterministic bound: By assumption, we have $\text{diam } \mathcal{F} = D < \infty$.
 116 Hence

$$\begin{aligned} \sum_{m=1}^M \frac{1}{2\alpha_m} \left(\|\hat{h}_{m-1} - h^*\|^2 - \|\hat{h}_m - h^*\|^2 \right) &= \sum_{m=2}^M \left(\frac{1}{2\alpha_m} - \frac{1}{2\alpha_{m-1}} \right) \|\hat{h}_{m-1} - h^*\|^2 \\ &\quad + \frac{1}{2\alpha_1} \|\hat{h}_0 - h^*\|^2 - \frac{1}{2\alpha_M} \|\hat{h}_M - h^*\|^2 \\ &\leq \sum_{m=2}^M \left(\frac{1}{2\alpha_m} - \frac{1}{2\alpha_{m-1}} \right) D^2 + \frac{1}{2\alpha_1} D^2 \\ &= \frac{D^2}{2\alpha_M}. \end{aligned} \quad (7)$$

117 The second summation can be bounded with the aid of Lemma 1:

$$\mathbb{E}_{\mathbf{z}_{1:M}} \left[\sum_{m=1}^M \frac{\alpha_m}{2} \|u_m\|^2 \right] = \frac{\mathbb{E}_{\mathbf{z}_{1:M}} [\|u_m\|^2]}{2} \sum_{m=1}^M \alpha_m \leq \frac{\rho(\hat{\Phi}, \hat{r}_0, \hat{\mathcal{P}})}{2} \sum_{m=1}^M \alpha_m. \quad (8)$$

118 Finally, the third summation can be bounded using Lemma 2. Let $\mathbb{E}_{\mathbf{z}_{-m}}$ denote the expectation with
 119 respect to $\mathbf{z}_1, \dots, \mathbf{z}_{m-1}, \mathbf{z}_{m+1}, \dots, \mathbf{z}_M$ and notice that

$$\begin{aligned} \mathbb{E}_{\mathbf{z}_{1:M}} \left[\langle \nabla \mathcal{R}(\hat{h}_{m-1}) - u_m, \hat{h}_{m-1} - h^* \rangle \right] &= \mathbb{E}_{\mathbf{z}_{-m}} \left[\mathbb{E}_{\mathbf{z}_m} \left[\langle \nabla \mathcal{R}(\hat{h}_{m-1}) - u_m, \hat{h}_{m-1} - h^* \rangle \right] \right] \\ &= \mathbb{E}_{\mathbf{z}_{-m}} \left[\langle \mathbb{E}_{\mathbf{z}_m} [\nabla \mathcal{R}(\hat{h}_{m-1}) - u_m], \hat{h}_{m-1} - h^* \rangle \right] \\ &= \mathbb{E}_{\mathbf{z}_{-m}} \left[\left\| \mathbb{E}_{\mathbf{z}_m} [\nabla \mathcal{R}(\hat{h}_{m-1}) - u_m] \right\| \left\| \hat{h}_{m-1} - h^* \right\| \right] \\ &\leq D \mathbb{E}_{\mathbf{z}_{-m}} \left[\left\| \mathbb{E}_{\mathbf{z}_m} [\nabla \mathcal{R}(\hat{h}_{m-1}) - u_m] \right\| \right]. \end{aligned}$$

120 Then, applying Lemma 2 and setting $\tau \triangleq D\kappa$ we get

$$\begin{aligned} \mathbb{E}_{\mathbf{z}_{1:M}} \left[\langle \nabla \mathcal{R}(\hat{h}_{m-1}) - u_m, \hat{h}_{m-1} - h^* \rangle \right] \\ \leq \tau(\hat{\Phi}) \left(\left\| \Phi - \hat{\Phi} \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 + \|r_0 - \hat{r}_0\|_{L^2(Z)}^2 + \left\| \mathcal{P} - \hat{\mathcal{P}} \right\|_{\text{op}}^2 \right)^{\frac{1}{2}}. \end{aligned} \quad (9)$$

121 All that is left to do is to apply equations (6), (7), (8) and (9) along with a basic convexity inequality.

122 Let $\hat{h} \triangleq \frac{1}{M} \sum_{m=1}^M \hat{h}_{m-1}$ and $\alpha = \rho/2$. Then:

$$\begin{aligned} \mathbb{E}_{\mathbf{z}_{1:M}} \left[\mathcal{R}(\hat{h}) - \mathcal{R}(h^*) \right] \\ \leq \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{\mathbf{z}_{1:M}} \left[\mathcal{R}(\hat{h}_m) - \mathcal{R}(h^*) \right] \\ \leq \frac{D^2}{2M\alpha_m} + \alpha(\hat{\Phi}, \hat{r}_0, \hat{\mathcal{P}}) \frac{1}{M} \sum_{m=1}^M \alpha_m \\ + \tau(\hat{\Phi}) \left(\left\| \Phi - \hat{\Phi} \right\|_{L^2(\nu_X \otimes \nu_Z)}^2 + \|r_0 - \hat{r}_0\|_{L^2(Z)}^2 + \left\| \mathcal{P} - \hat{\mathcal{P}} \right\|_{\text{op}}^2 \right)^{\frac{1}{2}}. \quad \square \end{aligned}$$

123 What's left to do:

124 • Use some estimate on $\left\| \mathcal{P} - \hat{\mathcal{P}} \right\|_{\text{op}}$ (Adapt notation and setup in the KIV paper).

125 Conclusion (20/08/2023): We might need the extra hypothesis that $\text{Im}(\text{id}_{L^2(X)} - \iota_X \iota_X^*) \subseteq$
 126 $\ker \mathcal{P}$, where $\iota_X : \mathcal{H}_X \rightarrow L^2(X)$ is the inclusion operator, whose adjoint is given by

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

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

- 135 • Find way to estimate r_0 which gives estimate on $\|r_0 - \hat{r}_0\|_{L^2(Z)}$. Maybe use the same
 136 estimation technique we have for \mathcal{P} as an operator from $L^2(Y) \rightarrow L^2(Z)$ applied to the
 137 identity and employ the same bound?

138 For the rest of the paper:

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