
Stochastic Gradient Descent in NPIV estimation

Anonymous Author(s)

Affiliation

Address

email

1 Problem setup

2 1.1 Basic definitions

3 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

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

4 that is, $L^2(X) = L^2(\mathcal{X}, \mathcal{B}(\mathcal{X}), \nu_X)^1$, a Hilbert space equipped with the inner product $\langle h, g \rangle_{L^2(X)} =$
5 $\mathbb{E}[h(X)g(X)]$. The regression problem we are interested in has the form

$$Y = h^*(X) + \varepsilon, \quad (1)$$

6 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
7 $Z \in L^2(\Omega, \mathcal{A}, \mathbb{P}; \mathcal{Z} \subseteq \mathbf{R}^q)$ such that

- 8 i) Z influences X , that is, $\nu_{X|Z}(\cdot | Z) \neq \nu_X(\cdot)$;
- 9 ii) Z influences Y only through X ;
- 10 iii) Z and ε are uncorrelated, that is, $\mathbb{E}[\varepsilon | Z] = 0$.

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

13 Conditioning (1) in Z , we find

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

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

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

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

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

17 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
18 thus we can rewrite (2) as

$$\mathcal{P}[h^*] = r_0. \quad (4)$$

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

¹We denote by ν_X the distribution of the r.v. X and by $\mathcal{B}(\mathcal{X})$ the Borel σ -algebra in \mathcal{X} .

Discuss the other implication, that if h satisfies $\mathcal{P}[h] = r_0$, then $h = h^*$. This is false, but the reason can be connected to the strength of the instrument Z .

20 1.2 Risk measure

21 Let $\ell : \mathbf{R} \times \mathbf{R} \rightarrow \mathbf{R}$ be a pointwise loss function, which, with respect to its second argument, is
 22 convex and differentiable. We use the symbol ∂_2 to denote a derivative with respect to the second
 23 argument. The example to keep in mind is the quadratic loss function $\ell(y, y') = \frac{1}{2}(y - y')^2$. Given
 24 $h \in L^2(X)$, we define the *populational risk* associated with it to be

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

25 We would like to solve

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

26 where $\mathcal{F} \subseteq L^2(X)$ is a bounded, closed, convex set such that $h^* \in \mathcal{F}$.

Assumption

27 We now state all the assumptions needed about the function ℓ for future reference:

28 **Assumption 1** (Regularity of ℓ).

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

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

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

31

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

34 2 Gradient computation

35 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}
 36 at h in the direction f , denoted by $DR[h](f)$:

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

37 where $\theta \in \mathbf{R}$ is due to Taylor's formula and can be assumed to be inside a fixed interval $(-\theta_0, \theta_0)$,
 38 with θ_0 arbitrarily small. The last step is then due to Assumption 1.2.

Assumption

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

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

40 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,
 41 with Gateux derivative at h given by

Assumption

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

42 If we assume² that $h \mapsto D\mathcal{R}[h]$ is a continuous mapping from $L^2(Z)$ to $L^2(Z)$, then \mathcal{R} is also
 43 Fréchet-differentiable, and both derivatives coincide. Therefore, under this assumption, which we
 44 henceforth make, $\nabla \mathcal{R}(h) = \mathcal{P}^*[\partial_2 \ell(r_0(\cdot), \mathcal{P}[h](\cdot))]$.

Assumption

Talk about which conditions ℓ can satisfy so that this is continuous.

45 3 Estimating the gradient

46 We have found that

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

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

Assumption

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

53 Thus, we must take

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

54 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 $u_h \in L^2(X)$.

55 An obvious obstacle for this approach is that we don't know how to analytically compute Φ , r_0 nor
 56 \mathcal{P} , so we will proceed with estimators $\hat{\Phi}$, \hat{r}_0 and $\hat{\mathcal{P}}$. In what follows, we remain agnostic to the exact
 57 form of these estimators and present the algorithm assuming we know how to compute them. Later,
 58 we'll show how the individual convergence rates of these three pieces come together to determine the
 59 convergence rate of our method.

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

60 4 Algorithm

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

Discuss everything we don't know and must estimate.

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

Discuss necessity of discretizing \mathcal{X} .

²It is if ℓ is quadratic.

Algorithm 1: SGD-NPIV

input : Datasets $\mathcal{D}_{r_0} = \{(y_i, z_i)\} \stackrel{\text{iid}}{\sim} \nu_{YZ}$, $\mathcal{D}_\Phi = \{(\mathbf{x}_i, z_i)\} \stackrel{\text{iid}}{\sim} \nu_{XZ}$,
 $\mathcal{D}_\mathcal{P} = \{(\mathbf{x}_i, z_i)\} \stackrel{\text{iid}}{\sim} \nu_{XZ}$, discretization $\{\mathbf{x}_k\}_{k=1}^K$ of \mathcal{X} which contains the observed
values of X , sequence of learning rates $(\alpha_m)_{m=1}^M$.

output : $\{\hat{h}(\mathbf{x}_k)\}_{k=1}^K$

Compute $\{\hat{r}_0(z_m; \mathcal{D}_{r_0})\}_{m=1}^M$;

62 Compute $\hat{\Phi}(\mathbf{x}, z; \mathcal{D}_\Phi)$;

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

 Compute $\hat{\mathcal{P}}[\hat{h}_{m-1}](z_m; \mathcal{D}_\mathcal{P})$;

 Set $u_m(\mathbf{x}_k) = \hat{\Phi}(\mathbf{x}_k, z_m) \partial_2 \ell(\hat{r}_0(z_m, \mathcal{D}_{r_0}), \hat{\mathcal{P}}[\hat{h}_{m-1}](z_m; \mathcal{D}_\mathcal{P}))$ for $1 \leq k \leq K$;

 Set $\hat{h}_m(\mathbf{x}_k) = \hat{h}_{m-1}(\mathbf{x}_k) - \alpha_m u_m(\mathbf{x}_k)$ for $1 \leq k \leq K$;

end

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

63 An option we have is to project onto the closed, convex, bounded set \mathcal{F} after applying the stochastic
64 gradient, that is, constructing the new estimate as

Should we do this?

$$\hat{h}_m = P_{\mathcal{F}}[\hat{h}_{m-1} - \alpha_m u_m].$$

65 From what I can see, this would require minor changes to the proof and would justify the assumption
66 that $\hat{h}_m \in \mathcal{F}$ for all m .

67 A possible choice for the set \mathcal{F} is

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

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

$$P_{\mathcal{F}}[h] = h^+ \wedge M - h^- \wedge M.$$

71 5 Proof of convergence

72 The first problem is proving our sequence of estimates is, in fact, contained in $L^2(X)$. This amounts
73 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
74 unbiased gradient when we know r_0, Φ 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.

75 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

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

76 To lighten the notation, the symbols $\|\cdot\|$ and $\langle \cdot, \cdot \rangle$, when written without a subscript to specify which
77 space they refer to, will act as the norm and inner product, respectively, of $L^2(X)$. By the Algorithm
78 1 procedure, we have

$$\begin{aligned} \frac{1}{2} \|\hat{h}_m - h^*\|^2 &= \frac{1}{2} \|\hat{h}_{m-1} - \alpha_m u_m - h^*\|^2 \\ &= \frac{1}{2} \|\hat{h}_{m-1} - h^*\|^2 - \alpha_m \langle u_m, \hat{h}_{m-1} - h^* \rangle + \frac{\alpha_m^2}{2} \|u_m\|^2. \end{aligned}$$

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

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

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

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

81 Rearranging terms, we get

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

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

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

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

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

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

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

$$\begin{aligned} \mathbb{E}_{\mathbf{z}_{1:M}} \left[\|u_m\|^2 \mid \mathcal{D}_{\Phi, \mathcal{P}, r_0} \right] &= \mathbb{E}_{\mathbf{z}_{1:m}} \left[\mathbb{E}_X \left[\hat{\Phi}(X, \mathbf{z}_m)^2 \partial_2 \ell \left(\hat{r}_0(\mathbf{z}_m), \hat{\mathcal{P}}[\hat{h}_{m-1}](\mathbf{z}_m) \right)^2 \mid \mathcal{D}_{\Phi, \mathcal{P}, r_0} \right] \right] \\ &= \mathbb{E}_X \left[\mathbb{E}_{\mathbf{z}_{1:m}} \left[\hat{\Phi}(X, \mathbf{z}_m)^2 \partial_2 \ell \left(\hat{r}_0(\mathbf{z}_m), \hat{\mathcal{P}}[\hat{h}_{m-1}](\mathbf{z}_m) \right)^2 \mid \mathcal{D}_{\Phi, \mathcal{P}, r_0} \right] \right]. \end{aligned}$$

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

$$\mathbb{E}_X \left[\mathbb{E}_{\mathbf{z}_{1:m}} \left[\hat{\Phi}(X, \mathbf{z}_m)^2 \partial_2 \ell \left(\hat{r}_0(\mathbf{z}_m), \hat{\mathcal{P}}[\hat{h}_{m-1}](\mathbf{z}_m) \right)^2 \right] \right].$$

88 Reversing back the expectations, we get

$$\begin{aligned} &\mathbb{E}_{\mathbf{z}_{1:m}} \left[\mathbb{E}_X \left[\hat{\Phi}(X, \mathbf{z}_m)^2 \partial_2 \ell \left(\hat{r}_0(\mathbf{z}_m), \hat{\mathcal{P}}[\hat{h}_{m-1}](\mathbf{z}_m) \right)^2 \right] \right] \\ &= \mathbb{E}_{\mathbf{z}_{1:m}} \left[\mathbb{E}_X \left[\hat{\Phi}(X, \mathbf{z}_m)^2 \right] \partial_2 \ell \left(\hat{r}_0(\mathbf{z}_m), \hat{\mathcal{P}}[\hat{h}_{m-1}](\mathbf{z}_m) \right)^2 \right]. \end{aligned}$$

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

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

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

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

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

$$\begin{aligned}\mathbb{E}_{\mathbf{z}_{1:m}} \left[\mathbb{E}_X \left[\widehat{\Phi}(X, \mathbf{z}_m)^2 \right] \partial_2 \ell \left(\widehat{r}_0(\mathbf{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) \right)^2 \right] \\ \leq \mathbb{E}_{\mathbf{z}_{1:m}} \left[\mathbb{E}_X \left[\left\| \widehat{\Phi} \right\|_{\mathcal{R}_{\mathbb{W}}}^2 \right] \partial_2 \ell \left(\widehat{r}_0(\mathbf{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) \right)^2 \right] \\ = \left\| \widehat{\Phi} \right\|_{\mathcal{R}_{\mathbb{W}}}^2 \mathbb{E}_{\mathbf{z}_{1:m}} \left[\partial_2 \ell \left(\widehat{r}_0(\mathbf{z}_m), \widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) \right)^2 \right].\end{aligned}$$

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

Assumption

$$\begin{aligned}\mathbb{E}_{\mathbf{z}_{1:m}} \left[\left(\widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) - \widehat{r}_0(\mathbf{z}_m) \right)^2 \right] \\ = \mathbb{E}_{\mathbf{z}_{1:m}} \left[\left(\left(\widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) - \mathcal{P}[\widehat{h}_{m-1}](\mathbf{z}_m) \right) + (r_0(\mathbf{z}_m) - \widehat{r}_0(\mathbf{z}_m)) \right. \right. \\ \left. \left. + \left(\mathcal{P}[\widehat{h}_{m-1}](\mathbf{z}_m) - r_0(\mathbf{z}_m) \right) \right)^2 \right] \\ \leq 3 \mathbb{E}_{\mathbf{z}_{1:m}} \left[\left(\widehat{\mathcal{P}}[\widehat{h}_{m-1}](\mathbf{z}_m) - \mathcal{P}[\widehat{h}_{m-1}](\mathbf{z}_m) \right)^2 + (r_0(\mathbf{z}_m) - \widehat{r}_0(\mathbf{z}_m))^2 \right. \\ \left. + \left(\mathcal{P}[\widehat{h}_{m-1}](\mathbf{z}_m) - r_0(\mathbf{z}_m) \right)^2 \right] \\ = 3 \left\{ \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\left\| (\widehat{\mathcal{P}} - \mathcal{P})[\widehat{h}_{m-1}] \right\|_{L^2(\mathbb{Z})}^2 \right] + \mathbb{E}_{\mathbf{z}_{1:m}} \left[\|r_0 - \widehat{r}_0\|_{L^2(\mathbb{Z})}^2 \right] \right. \\ \left. + \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\left\| \mathcal{P}[\widehat{h}_{m-1}] - r_0 \right\|_{L^2(\mathbb{Z})}^2 \right] \right\}.\end{aligned}$$

94 We treat each part of this expression separately. Firstly,

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

95 We leave the second part as $\|r_0 - \widehat{r}_0\|_{L^2(\mathbb{Z})}^2$. Finally, for the third part, we have

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

96 Putting everything together, what we conclude is

$$\mathbb{E}_{\mathbf{z}_{1:m}} \left[\|u_m\|_{L^2(\mathbb{X})}^2 \mid \mathcal{D}_{\Phi, \mathcal{P}, r_0} \right] \leq 3 \left\| \widehat{\Phi} \right\|_{\mathcal{R}_{\mathbb{W}}}^2 \left(M^2 \left\| \widehat{\mathcal{P}} - \mathcal{P} \right\|_{\text{op}}^2 + \|r_0 - \widehat{r}_0\|_{L^2(\mathbb{Z})}^2 + 2 \left(M^2 + \mathbb{E}[Y^2] \right) \right).$$

97 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,
 98 we only need to bound this by something which remains bounded as $|\mathcal{D}_{\Phi, \mathcal{P}, r_0}|$ and the number of
 99 iterations grow. Another idea is to simply say that this whole thing is $\mathcal{O}_p(1)$, that is, almost surely
 100 finite, and rely on the (fast enough) decay of the learning rate to achieve convergence.

101 Third term

102 Our goal is to open up the inner product and make explicit the estimation errors of our model's
 103 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
 104 hat version $\widehat{\Psi}_m$ is defined accordingly, replacing r_0 and \mathcal{P} by their estimators.

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{z}_{1:m}} \left[\langle \nabla \mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^* \rangle \mid \mathcal{D}_{\Phi, \mathcal{P}, r_0} \right] \\
 &= \mathbb{E}_{\mathbf{z}_{1:m}} \left[\langle \nabla \mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^* \rangle \right] \quad (\mathbf{z}_{1:m} \perp\!\!\!\perp \mathcal{D}_{\Phi, \mathcal{P}, r_0}) \\
 &= \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\mathbb{E}_{\mathbf{z}_m} \left[\langle \nabla \mathcal{R}(\widehat{h}_{m-1}) - u_m, \widehat{h}_{m-1} - h^* \rangle \right] \right] \quad (\mathbf{z}_m \perp\!\!\!\perp \mathbf{z}_{1:m-1}) \\
 &= \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\langle \nabla \mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\mathbf{z}_m} [u_m], \widehat{h}_{m-1} - h^* \rangle \right] \\
 &\leq \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\left\| \nabla \mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\mathbf{z}_m} [u_m] \right\| \left\| \widehat{h}_{m-1} - h^* \right\| \right] \\
 &\leq D \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\left\| \nabla \mathcal{R}(\widehat{h}_{m-1}) - \mathbb{E}_{\mathbf{z}_m} [u_m] \right\| \right] \quad (\text{diam } \mathcal{F} = D) \\
 &\leq D \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\mathbb{E}_X \left[\left(\nabla \mathcal{R}(\widehat{h}_{m-1})(X) - \mathbb{E}_{\mathbf{z}_m} [u_m] \right)^2 \right] \right]^{\frac{1}{2}} \quad (\text{Jensen}) \\
 &= D \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\mathbb{E}_X \left[\left(\mathbb{E}_Z [\Phi(X, Z) \Psi_m(Z)] \right. \right. \right. \\
 &\quad \left. \left. \left. - \mathbb{E}_{\mathbf{z}_m} [\widehat{\Phi}(X, \mathbf{z}_m) \widehat{\Psi}_m(\mathbf{z}_m)] \right)^2 \right] \right]^{\frac{1}{2}} \\
 &= D \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\mathbb{E}_X \left[\mathbb{E}_Z [\Phi(X, Z) \Psi_m(Z) - \widehat{\Phi}(X, Z) \widehat{\Psi}_m(Z)]^2 \right] \right]^{\frac{1}{2}} \quad (Z \stackrel{\text{iid}}{\sim} \mathbf{z}_m) \\
 &= D \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\mathbb{E}_X \left[\mathbb{E}_Z [\Psi_m(Z) (\Phi(X, Z) - \widehat{\Phi}(X, Z)) \right. \right. \right. \\
 &\quad \left. \left. \left. + \widehat{\Phi}(X, Z) (\Psi_m(Z) - \widehat{\Psi}_m(Z)) \right]^2 \right] \right]^{\frac{1}{2}} \\
 &\leq D \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\mathbb{E}_X \left[\left(\left\| \Psi_m \right\|_{L^2(Z)} \left\| \Phi(X, \cdot) - \widehat{\Phi}(X, \cdot) \right\|_{L^2(Z)} \right. \right. \right. \\
 &\quad \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}_{\mathbf{z}_{1:m-1}} \left[\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 \right. \right. \right. \\
 &\quad \left. \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}_{\mathbf{z}_{1:m-1}} \left[\left\| \Psi_m \right\|_{L^2(Z)}^2 \mathbb{E}_X \left[\left\| \Phi(X, \cdot) - \widehat{\Phi}(X, \cdot) \right\|_{L^2(Z)}^2 \right] \right. \right. \right. \\
 &\quad \left. \left. \left. + \mathbb{E}_{\mathbf{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] \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}_{\mathbf{z}_{1:m-1}} \left[\left\| \Psi_m \right\|_{L^2(Z)}^2 \right] \right. \\
 &\quad \left. + \mathbb{E}_X \left[\left\| \widehat{\Phi}(X, \cdot) \right\|_{L^2(Z)}^2 \right] \mathbb{E}_{\mathbf{z}_{1:m-1}} \left[\left\| \Psi_m - \widehat{\Psi}_m \right\|_{L^2(Z)}^2 \right] \right)^{\frac{1}{2}}
 \end{aligned}$$

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

105 We proceed to analyze each term separately:

106 • To bound A , first notice that

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

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

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

Create section describ-
ing how we are esti-
mating each term.

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

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

112 In total, what we have is

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

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

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

114 **We still need to bound this norm somehow.**

115 Furthermore, we also have

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

116

Therefore,

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

117 What's left to do:

118 • Bound $\left\| \widehat{\Phi} \right\|_{\mathcal{R}_W}$. (May not be strictly necessary. This is finite, and since it multiplies
 119 something which is \mathcal{O}_p of something which goes to zero, we may not need to further bound
 120 it.)

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

122 Conclusion (20/08/2023): We might need the extra hypothesis that $\text{Im}(\text{id}_{L^2(X)} - \iota_X \iota_X^*) \subseteq$
 123 $\ker \mathcal{P}$, where $\iota_X : \mathcal{H}_X \rightarrow L^2(X)$ is the inclusion operator, whose adjoint is given by

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

124 with $k_X : \mathbb{X} \times \mathbb{X} \rightarrow \mathbf{R}$ being the kernel associated with \mathcal{H}_X . Then $\mathcal{P} = \mathcal{P} \circ \iota_X \iota_X^*$
 125 and we can directly apply the result on KIV's paper, since $\mathcal{P} \circ \iota_X$ can be seen as the
 126 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
 127 something like this (because, rigorously speaking, $\mathcal{P}f$ is an equivalence class of functions,
 128 so in what way can we say that this equivalence class is “in \mathcal{H}_Z ”?). This hypothesis is
 129 implicitly made in the KIV paper, when they say that $E : \mathcal{H}_X \rightarrow \mathcal{H}_Z$ without providing
 130 any assumptions on \mathcal{H}_X and \mathcal{H}_Z , other than saying that they are RKHS. Who can guarantee
 131 that $(z \mapsto \mathbb{E}[f(X) \mid Z = z]) \in \mathcal{H}_Z$ for every $f \in \mathcal{H}_X$?

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

135 For the rest of the paper:

136 • Create section which describes, in detail, how we are estimating Φ , \mathcal{P} and r_0 , lists all the
 137 references, states the main convergence theorems and lists all of the assumptions that are
 138 being made.

139 • Adapt the algorithm section to use the KIV first stage, which directly estimates \mathcal{P} .

140 • Find better letter for either the number of iterations or the upper bound for the set \mathcal{F} . Right
 141 now, both are being denoted by the letter M .

142 **References**

143 [1] Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density Ratio Estimation in Machine*
 144 *Learning*. Cambridge University Press, 2012. DOI: 10.1017/CB09781139035613.