

Worksheet 4: Hilbert Projection and Estimation with Polynomials

Due October 9, 2023

Name: Aaron Zoll

The Hilbert Projection Theorem provides a way of concretely interpreting conditional expectation. In this worksheet, you will be working out a computational method for using Hilbert Projection, and this exercise will serve as foundation for implementation in the next programming assignment.

Consider joint probability space $(\mathcal{X} \times \mathcal{Y} = \mathbb{R}^2, \mathbb{P}_{\mathcal{X} \times \mathcal{Y}})$ and suppose that we would like to estimate random variable $\pi_{\mathcal{Y}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Y}$ by some polynomial in x of degree at most n , namely $\hat{y}(x) = \sum_{j=0}^n a_j x^j$ with coefficients $a_j \in \mathbb{R}$.

1. (Prepping for Regression) Let $\mathcal{H} = \left\{ \sum_{j=0}^n a_j x^j : a_j \in \mathbb{R} \right\}$ be $n+1$ dimensional space of degree (at most) n polynomials.

Verify that $\mathcal{H} \subset \{ \hat{y} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Y} \}$ is a linear subspace. Can you also show that \mathcal{H} is also closed?

$$1) f \equiv 0 \in \mathcal{H} \quad \text{since } \deg(f) = -\infty < n \quad \forall n \in \mathbb{N}$$

$$2) f+g \in \mathcal{H} \Rightarrow f+g = \sum_{j=0}^n a_j^{(f)} x^j + \sum_{j=0}^n a_j^{(g)} x^j = \sum_{j=0}^n (a_j^{(f)} + a_j^{(g)}) x^j \in \mathcal{H}$$

$$3) \lambda \in \mathbb{R}, f \in \mathcal{H} \Rightarrow \lambda f = \lambda \sum_{j=0}^n a_j x^j = \sum_{j=0}^n (\lambda a_j) x^j \in \mathcal{H}$$

\mathcal{H} is closed by completeness of \mathbb{R} and finite dimensionality.
If $\{f^n\}_{n \in \mathbb{N}} \subset \mathcal{H}$, $f^n \rightarrow f$
then $a_j^{(n)} \rightarrow a_j \quad \forall j$
and $f = \sum_{j=0}^n a_j x^j \in \mathcal{H}$

2. (Linear Regression) Replicating computation in class for the linear case, enumerate the separate conditions which allow us to find optimal $y^*(x) = \sum_{j=0}^n a_j^* x^j$, the point in \mathcal{H} closest to $\pi_{\mathcal{Y}}$. If notation is cumbersome, try your hand first at the quadratic case (a pattern should emerge for generalizing to arbitrary degree). This problem is still called *linear* regression even when $n > 1$. Why?

$$\text{Want } \forall k=0,1,\dots,n \quad \mathbb{E}((y^* - y)x^k) = 0 \Rightarrow \sum_{j=0}^n a_j^* \mathbb{E}(x^{j+k}) = \mathbb{E}(y x^k)$$

$$\begin{aligned} &= \mathbb{E}\left(\left(\sum_{j=0}^n a_j^* x^j - y\right)x^k\right) = 0 \quad \text{i.e.} \quad a_0^* + a_1^* \mathbb{E}(x) + \dots + a_n^* \mathbb{E}(x^n) = \mathbb{E}(y) \\ & \quad a_0^* \mathbb{E}(x) + a_1^* \mathbb{E}(x^2) + \dots + a_n^* \mathbb{E}(x^{n+1}) = \mathbb{E}(xy) \\ & \quad a_0^* \mathbb{E}(x^2) + a_1^* \mathbb{E}(x^3) + \dots + a_n^* \mathbb{E}(x^{n+2}) = \mathbb{E}(x^2 y) \\ & \quad \vdots \\ & \quad a_0^* \mathbb{E}(x^n) + a_1^* \mathbb{E}(x^{n+1}) + \dots + a_n^* \mathbb{E}(x^{2n}) = \mathbb{E}(x^n y) \end{aligned}$$

note: linear b/c problem/solution is linear. I.e. matrices

3. In problem #2, you should have written $n+1$ equations in $n+1$ unknowns a_0^*, \dots, a_n^* . Using matrix notation, express the solution for a^* .

$$\underbrace{\begin{bmatrix} 1 & \mathbb{E}(x) & \mathbb{E}(x^2) & \dots & \mathbb{E}(x^n) \\ \mathbb{E}(x) & \mathbb{E}(x^2) & \mathbb{E}(x^3) & \dots & \mathbb{E}(x^{n+1}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbb{E}(x^n) & \mathbb{E}(x^{n+1}) & \mathbb{E}(x^{n+2}) & \dots & \mathbb{E}(x^{2n}) \end{bmatrix}}_A \underbrace{\begin{bmatrix} a_0^* \\ a_1^* \\ \vdots \\ a_n^* \end{bmatrix}}_a = \underbrace{\begin{bmatrix} \mathbb{E}(y) \\ \mathbb{E}(xy) \\ \vdots \\ \mathbb{E}(x^n y) \end{bmatrix}}_b$$

$$a^* = A^{-1} b$$

(1)

4. Now suppose you are provided data $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \subset \mathcal{X} \times \mathcal{Y}$. Express an approximation for the terms in eq. (1) in terms of S . You do not need to write them all out individually; write an approximation for the general term and the corresponding approximation equation for a_s of a^* .

$$\mathbb{E}(x^j y^l) \approx \frac{1}{m} \sum_{k=1}^m x_k^j y_k^l =: a_{j,l}$$

$$j = 0, 1, \dots, 2n \\ l = 0, 1$$

$$a_s = \begin{bmatrix} a_{0,0} & a_{1,0} & \dots & a_{n,0} \\ a_{0,1} & a_{1,1} & \dots & a_{n,1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{0,n} & a_{1,n} & \dots & a_{n,n} \end{bmatrix}^{-1} \begin{bmatrix} a_{0,1} \\ a_{1,1} \\ \vdots \\ a_{n,1} \end{bmatrix}$$

(2)

5. Problem #4 (eq. (2) in particular) allows you to find an a_S providing *approximate* solution of eq. (1). Under what condition(s) does—and which (probabilistic) principle justifies that— $a_S \rightarrow a^*$?

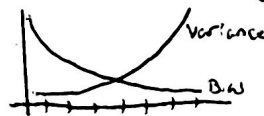
Finite moments + LLN
(up to $2n$)

6. (Bias-Variance) Let us write $y_H^* = E_H(y|x)$ even when $H \neq \{g: \mathcal{X} \rightarrow \mathcal{Y}\}$ (and ensure that we are clear on what H is!). For arbitrary $\tilde{y} \in H$, we can decompose the *mean squared error* $E((\tilde{y} - y)^2)^*$ as

$$E((\tilde{y} - y)^2) = \underbrace{E((\tilde{y} - y_H^*)^2)}_{\text{variance}} + \underbrace{E((y_H^* - y)^2)}_{\text{bias}}.$$

Rederive this decomposition and for $H_0 \subsetneq H_1 \subsetneq \dots \subsetneq H_n \subsetneq \dots$, plot a heuristic of each term against N (indexing H_j). Be careful: $\tilde{y} \in H_j$ for each j ; therefore explain why you expect the variance term to behave as you've drawn it.

$$\begin{aligned} E((\tilde{y} - y)^2) &= E((\tilde{y} - y_H^* + y_H^* - y)^2) = E((\tilde{y} - y_{H_1}^*)^2 - 2(\tilde{y} - y_{H_1}^*)(y - y_{H_1}^*) + (y - y_{H_1}^*)^2) \\ &= E((\tilde{y} - y_{H_1}^*)^2) - 2E(\tilde{y} - y_{H_1}^*)E(y - y_{H_1}^*) + E((y - y_{H_1}^*)^2) \end{aligned}$$



Variance does this b/c adding more moments/terms to sum, and the square make it more positive

7. (RKHS: Recall HP) Let (\mathcal{X}, P_X) be a probability space, and $\mathcal{V} \subset \{f: \mathcal{X} \rightarrow \mathbb{R} : \|f\|^2 < \infty\}$ a Hilbert space of squared integrable random variables. We say that a subspace $\mathcal{H} \subset \mathcal{V}$ is a *reproducing kernel Hilbert space* if there is function $k: \mathcal{X}^2 \rightarrow \mathbb{R}$ satisfying:

- (a) $k(\cdot, x) \in \mathcal{H}$ for each $x \in \mathcal{X}$ and
(b) $f(x) = \langle f, k(\cdot, x) \rangle$ for each $f \in \mathcal{H} \Rightarrow f = \langle f, k \rangle$

Suppose that $\mathcal{H} \subset \mathcal{V}$ is a *closed* subspace. Show that $\langle f, k \rangle = \pi_{\mathcal{H}}(f)$ for $f \in \mathcal{V}$, where $\pi_{\mathcal{H}}(f) := \arg \min_{h \in \mathcal{H}} \|f - h\|^2$.

$$\begin{aligned} \forall h \in \mathcal{H}, \quad \langle f - \langle f, k \rangle, h \rangle &= \int_{\mathcal{X}} f(x) \int_{\mathcal{X}} h(y) k(x, y) dP_X(y) dP_X(x) - \int_{\mathcal{X}} f(x) k(x, y) \int_{\mathcal{X}} h(y) dP_X(y) dP_X(x) \\ &= \int_{\mathcal{X}} f(x) \int_{\mathcal{X}} h(y) k(x, y) dP_X(y) dP_X(x) - \int_{\mathcal{X}} f(x) k(x, y) h(y) dP_X(y) dP_X(x) \\ &= 0 \Rightarrow f - \langle f, k \rangle \in \mathcal{H}^\perp \Rightarrow \langle f, k \rangle = \arg \min_{h \in \mathcal{H}} \|f - h\|^2 =: \pi_{\mathcal{H}}(f) \end{aligned}$$

8. (Symmetrization) For random variable (\mathcal{X}, P_X) , we have seen that $\mathcal{X} \times \mathcal{S} \rightarrow \mathbb{R}$ mapping $(x, s) \mapsto sx$ —for $\mathcal{S} = \{1, -1\}$ with $P_S(1) = 1/2$ —is symmetric. Show the same for $\mathcal{X}^2 \times \mathcal{S}^2 \rightarrow \mathbb{R}$ sending $(x_1, x_2, s_1, s_2) \mapsto s_1 x_1 + s_2 x_2$, namely that

$$P_R(s_1 x_1 + s_2 x_2 > t) = P_R(s_1 x_1 + s_2 x_2 < -t).$$

You may use the single symmetry version, and may have to cite more than one application of LTP/LTE.

$$\begin{aligned} P_{X^2 \times S^2}(s_1 x_1 + s_2 x_2 > t) &= \int_{\mathcal{S}^2} \int_{\mathcal{X}^2} \mathbb{1}_{s_1 x_1 + s_2 x_2 > t} dP_{X^2}(x_1, x_2) dP_{S^2}(s_1, s_2) \\ &= \frac{1}{4} P_{X^2 \times S^2}(x_1, x_2 > t | s_1 = 1, s_2 = 1) + \frac{1}{4} P_{X^2 \times S^2}(x_1, x_2 > t | s_1 = 1, s_2 = -1) \\ &\quad + \frac{1}{4} P_{X^2 \times S^2}(-x_1, x_2 > t | s_1 = -1, s_2 = 1) + \frac{1}{4} P_{X^2 \times S^2}(-x_1, x_2 > t | s_1 = -1, s_2 = -1) \\ &= \frac{1}{4} P_{X^2 \times S^2}(x_1, x_2 > t | s_1 = 1, s_2 = 1) + \frac{1}{4} P_{X^2 \times S^2}(x_1, x_2 < -t | s_1 = 1, s_2 = -1) \\ &\quad + \frac{1}{4} P_{X^2 \times S^2}(x_1, x_2 < -t | s_1 = -1, s_2 = 1) + \frac{1}{4} P_{X^2 \times S^2}(x_1, x_2 > t | s_1 = -1, s_2 = -1) \\ &= P_{X^2 \times S^2}(s_1 x_1 + s_2 x_2 < -t) \end{aligned}$$

*Of course, we are still operating with standard diagram: y may be read as evaluation of π_y at (x, y) and \tilde{y} as the composition $\tilde{y} \circ \pi_X$.

*Part of this exercise is to carefully define domains / maps: $\langle \cdot, \cdot \rangle: \mathcal{V} \rightarrow \mathbb{R}$ takes inputs in \mathcal{V} and values in \mathbb{R} . Apparently there is some looseness with notation, which part of this exercise is to tighten up.