

# Finite-Sample Analysis of Prediction-Powered Linear Regression

## 1 Preliminaries

We consider a linear regression model:

$$Y = X\theta^* + \varepsilon,$$

where  $\theta^* \in \mathbb{R}^d$  is the target parameter and  $\varepsilon$  is a noise term.

We observe independent labeled and unlabeled samples:

$$\{(X_i, Y_i)\}_{i=1}^n, \quad \{\tilde{X}_j\}_{j=1}^N, \quad \text{and} \quad f(X_i), f(\tilde{X}_j),$$

where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is a black-box predictor which is fixed or independent of both labeled and unlabeled samples (e.g. pre-trained).

Let  $X \in \mathbb{R}^{n \times d}$  and  $\tilde{X} \in \mathbb{R}^{N \times d}$  denote the design matrices with rows  $X_i^\top$  and  $\tilde{X}_j^\top$ , and let  $Y \in \mathbb{R}^n$  be the vector of labeled responses. We adopt the following empirical PPI loss with tuning parameter  $\lambda \in \mathbb{R}$ :

$$L_{PP,n,N}^\lambda(\theta) := \frac{1}{2n} \|Y - X\theta\|_2^2 + \lambda \left( \frac{1}{2N} \|f(\tilde{X}) - \tilde{X}\theta\|_2^2 - \frac{1}{2n} \|f(X) - X\theta\|_2^2 \right). \quad (1)$$

Here  $f(X) \in \mathbb{R}^n$  and  $f(\tilde{X}) \in \mathbb{R}^N$  are understood componentwise.

**PPI estimator.** The PPI estimator is any minimizer of the empirical loss:

$$\hat{\theta} := \hat{\theta}_{PP}^\lambda \in \arg \min_{\theta \in \mathbb{R}^d} L_{PP,n,N}^\lambda(\theta).$$

By definition of the argmin, for every  $\theta \in \mathbb{R}^d$  and in particular for  $\theta^*$ ,

$$L_{PP,n,N}^\lambda(\hat{\theta}) \leq L_{PP,n,N}^\lambda(\theta^*).$$

This is the starting point of our basic inequality.

## 2 Finite-Sample Error Bound for the PPI Estimator

Let  $\theta^* := \arg \min_\theta \mathbb{E}[\frac{1}{2}(X^\top \theta - Y)^2]$  denote the population risk minimizer under squared loss, and define the residual and prediction bias

$$\varepsilon := Y - X^\top \theta^*, \quad r(X) := f(X) - X^\top \theta^*.$$

Write  $\Delta := \hat{\theta}_\lambda - \theta^*$ .

**Assumptions.** Throughout the finite-sample analysis, we assume:

- (A1) The design  $X \in \mathbb{R}^d$  is mean-zero sub-Gaussian with  $\|X\|_{\psi_2} \leq K$  and covariance matrix  $\Sigma_X$  satisfying  $\lambda_{\min}(\Sigma_X) \geq \kappa_0 > 0$ .
- (A2) The noise  $\varepsilon$  satisfies  $\mathbb{E}[\varepsilon | X] = 0$  and  $\|\varepsilon\|_{\psi_2} \leq K_\varepsilon$ .
- (A3) The prediction residual  $r(X) = f(X) - X^\top \theta^*$  satisfies  $\|r(X)\|_{\psi_2} \leq K_r$ .

**Step 1: Algebraic expansion.** Recalling the optimality condition  $L_{PP,n,N}^\lambda(\hat{\theta}) \leq L_{PP,n,N}^\lambda(\theta^*)$  and the definition:

$$L_{PP,n,N}^\lambda(\theta) := \frac{1}{2n} \|Y - X\theta\|_2^2 + \lambda \left( \frac{1}{2N} \|f(\tilde{X}) - \tilde{X}\theta\|_2^2 - \frac{1}{2n} \|f(X) - X\theta\|_2^2 \right). \quad (2)$$

We can derive the basic inequality by algebraically expanding the differences in the squared norms:

$$\begin{aligned} \|Y - X\hat{\theta}\|_2^2 - \|Y - X\theta^*\|_2^2 &= \|X\theta^* + \varepsilon - X\hat{\theta}\|_2^2 - \|\varepsilon\|_2^2 \\ &= \|X\Delta\|_2^2 - 2\varepsilon^\top X\Delta, \end{aligned}$$

$$\begin{aligned} \|f(\tilde{X}) - \tilde{X}\hat{\theta}\|_2^2 - \|f(\tilde{X}) - \tilde{X}\theta^*\|_2^2 &= \|\tilde{X}\theta^* + r(\tilde{X}) - \tilde{X}\hat{\theta}\|_2^2 - \|r(\tilde{X})\|_2^2 \\ &= \|\tilde{X}\Delta\|_2^2 - 2r(\tilde{X})^\top \tilde{X}\Delta, \end{aligned}$$

$$\|f(X) - X\hat{\theta}\|_2^2 - \|f(X) - X\theta^*\|_2^2 = \|X\Delta\|_2^2 - 2r(X)^\top X\Delta.$$

Substituting into  $L_{PP,n,N}^\lambda(\hat{\theta}) - L_{PP,n,N}^\lambda(\theta^*) \leq 0$ , we obtain:

$$\frac{1-\lambda}{2n} \|X\Delta\|_2^2 + \frac{\lambda}{2N} \|\tilde{X}\Delta\|_2^2 \leq \frac{1}{n} \varepsilon^\top X\Delta + \frac{\lambda}{N} r(\tilde{X})^\top \tilde{X}\Delta - \frac{\lambda}{n} r(X)^\top X\Delta \quad (3)$$

Combining the labeled and unlabeled terms on the right-hand side, we can rewrite the basic inequality as:

$$\frac{1-\lambda}{2n} \|X\Delta\|_2^2 + \frac{\lambda}{2N} \|\tilde{X}\Delta\|_2^2 \leq \underbrace{\frac{1}{n} (\varepsilon - \lambda r(X))^\top X\Delta}_{=: T_1(\lambda)} + \underbrace{\frac{\lambda}{N} r(\tilde{X})^\top \tilde{X}\Delta}_{=: T_2(\lambda)}. \quad (4)$$

**Step 2: lower bound on the left-hand side.** Let

$$\hat{\Sigma}_n := \frac{1}{n} X^\top X, \quad \hat{\Sigma}_N := \frac{1}{N} \tilde{X}^\top \tilde{X}, \quad \Sigma := \mathbb{E}[XX^\top],$$

and assume  $\lambda_{\min}(\Sigma) \geq \kappa_0 > 0$ . Under the sub-Gaussian design assumption  $\|X\|_{\psi_2} \leq K$ , Vershynin's covariance concentration theorem (e.g., [?, Theorem 4.6.1]) implies that there exist absolute constants  $c, C > 0$  such that, whenever

$$\min(n, N) \geq c K^4 \frac{d + \log(6/\delta)}{\kappa_0^2}, \quad (5)$$

the following event holds with probability at least  $1 - \delta/3$ :

$$\mathcal{E}_{\text{op}} := \left\{ \|\hat{\Sigma}_n - \Sigma\|_{\text{op}} \leq \frac{\kappa_0}{2}, \quad \|\hat{\Sigma}_N - \Sigma\|_{\text{op}} \leq \frac{\kappa_0}{2} \right\}. \quad (6)$$

On  $\mathcal{E}_{\text{op}}$ , Weyl's inequality yields

$$\lambda_{\min}(\hat{\Sigma}_n) \geq \lambda_{\min}(\Sigma) - \|\hat{\Sigma}_n - \Sigma\|_{\text{op}} \geq \frac{\kappa_0}{2}, \quad \lambda_{\min}(\hat{\Sigma}_N) \geq \frac{\kappa_0}{2}.$$

Therefore, for any  $\lambda \in [0, 1]$  and any  $\Delta \in \mathbb{R}^d$ ,

$$\begin{aligned} \frac{1-\lambda}{2n} \|X\Delta\|_2^2 + \frac{\lambda}{2N} \|\tilde{X}\Delta\|_2^2 &= \frac{1}{2} \Delta^\top \left( (1-\lambda)\hat{\Sigma}_n + \lambda\hat{\Sigma}_N \right) \Delta \\ &\geq \frac{1}{2} \lambda_{\min} \left( (1-\lambda)\hat{\Sigma}_n + \lambda\hat{\Sigma}_N \right) \|\Delta\|_2^2 \\ &\geq \frac{1}{2} \left( (1-\lambda)\lambda_{\min}(\hat{\Sigma}_n) + \lambda\lambda_{\min}(\hat{\Sigma}_N) \right) \|\Delta\|_2^2 \\ &\geq \frac{1}{2} \cdot \frac{\kappa_0}{2} \|\Delta\|_2^2 = \frac{\kappa_0}{4} \|\Delta\|_2^2. \end{aligned} \tag{7}$$

**(Optional refinement).** One may equivalently track the explicit concentration radii  $\epsilon_n(\delta), \epsilon_N(\delta)$  defined by

$$\|\hat{\Sigma}_n - \Sigma\|_{\text{op}} \leq \epsilon_n(\delta), \quad \|\hat{\Sigma}_N - \Sigma\|_{\text{op}} \leq \epsilon_N(\delta),$$

which implies

$$\frac{1}{2} \Delta^\top \left( (1-\lambda)\hat{\Sigma}_n + \lambda\hat{\Sigma}_N \right) \Delta \geq \frac{1}{2} \left( \kappa_0 - (1-\lambda)\epsilon_n(\delta) - \lambda\epsilon_N(\delta) \right) \|\Delta\|_2^2.$$

For readability, we use the uniform bound (7) in the main proof.

**Step 3: Decomposition of  $T_1(\lambda)$  and  $T_2(\lambda)$ .** We expose the bias–fluctuation structure of each term and then combine them to exploit the cancellation built into the PPI objective.

*Labeled term.* Define

$$A_n(\lambda) := \frac{1}{n} \sum_{i=1}^n X_i (\varepsilon_i - \lambda r(X_i)) = m_1(\lambda) + Z_1(\lambda),$$

where

$$m_1(\lambda) := \mathbb{E}[X\varepsilon] - \lambda \mathbb{E}[Xr(X)], \quad Z_1(\lambda) := A_n(\lambda) - m_1(\lambda).$$

Since  $\theta^*$  minimizes the population squared loss,  $\mathbb{E}[X\varepsilon] = 0$ , hence

$$m_1(\lambda) = -\lambda \mathbb{E}[Xr(X)]. \tag{8}$$

*Unlabeled term.* Define

$$B_N := \frac{1}{N} \sum_{j=1}^N \tilde{X}_j r(\tilde{X}_j) = m_2 + Z_2, \quad m_2 := \mathbb{E}[Xr(X)], \quad Z_2 := B_N - m_2.$$

**Step 4: Combined bound via bias cancellation.** Recall

$$T_1(\lambda) + T_2(\lambda) = \langle A_n(\lambda), \Delta \rangle + \lambda \langle B_N, \Delta \rangle.$$

Substituting the decompositions,

$$T_1(\lambda) + T_2(\lambda) = \langle m_1(\lambda) + \lambda m_2, \Delta \rangle + \langle Z_1(\lambda) + \lambda Z_2, \Delta \rangle.$$

The bias cancels exactly:

$$m_1(\lambda) + \lambda m_2 = -\lambda \mathbb{E}[Xr(X)] + \lambda \mathbb{E}[Xr(X)] = 0.$$

Therefore, by Cauchy–Schwarz and the triangle inequality,

$$T_1(\lambda) + T_2(\lambda) \leq (\|Z_1(\lambda)\|_2 + \lambda \|Z_2\|_2) \|\Delta\|_2. \quad (9)$$

It remains to control  $\|Z_1(\lambda)\|_2$  and  $\|Z_2\|_2$ .

*Bounding  $\|Z_1(\lambda)\|_2$ .* Assume  $X$  is mean-zero sub-Gaussian with

$$K_X := \sup_{\|u\|_2=1} \|u^\top X\|_{\psi_2} < \infty,$$

and assume  $\varepsilon$  and  $r(X)$  are sub-Gaussian with  $\|\varepsilon\|_{\psi_2} \leq K_\varepsilon$  and  $\|r(X)\|_{\psi_2} \leq K_r$ . Then each coordinate of  $X(\varepsilon - \lambda r(X))$  is sub-exponential. By a vector Bernstein inequality, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta/2$ ,

$$\|Z_1(\lambda)\|_2 \leq C K_X \sqrt{\frac{d + \log(2/\delta)}{n}} \sqrt{\text{Var}(\varepsilon - \lambda r(X))}, \quad (10)$$

where

$$\text{Var}(\varepsilon - \lambda r(X)) = \text{Var}(\varepsilon) + \lambda^2 \text{Var}(r(X)) - 2\lambda \text{Cov}(\varepsilon, r(X)). \quad (11)$$

*Bounding  $\lambda \|Z_2\|_2$ .* For  $u \in \mathbb{S}^{d-1}$ ,

$$\|u^\top (Xr(X))\|_{\psi_1} = \|(u^\top X)r(X)\|_{\psi_1} \leq C \|u^\top X\|_{\psi_2} \|r(X)\|_{\psi_2} \leq C K_X K_r,$$

so  $Xr(X)$  is a sub-exponential random vector with parameter  $CK_X K_r$ . By vector Bernstein, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta/2$ ,

$$\lambda \|Z_2\|_2 \leq \lambda C K_X K_r \sqrt{\frac{d + \log(2/\delta)}{N}}. \quad (12)$$

*Combined upper bound.* Applying a union bound to (10) and (12), we obtain that, with probability at least  $1 - \delta$ ,

$$T_1(\lambda) + T_2(\lambda) \leq \left[ C K_X \sqrt{\frac{d + \log(2/\delta)}{n}} \sqrt{\text{Var}(\varepsilon - \lambda r(X))} + \lambda C K_X K_r \sqrt{\frac{d + \log(2/\delta)}{N}} \right] \|\Delta\|_2. \quad (13)$$

**Step 5: Final bound and optimization over  $\lambda$ .** Let

$$A := C K_X \sqrt{\frac{d + \log(3/\delta)}{n}}, \quad B := C K_X K_r \sqrt{\frac{d + \log(3/\delta)}{N}}.$$

On the event  $\mathcal{E}_{\text{op}}$  in (6), we have the lower bound (7):

$$\frac{1-\lambda}{2n} \|X\Delta\|_2^2 + \frac{\lambda}{2N} \|\tilde{X}\Delta\|_2^2 \geq \frac{\kappa_0}{4} \|\Delta\|_2^2, \quad \forall \lambda \in [0, 1].$$

Moreover, by (9)–(13) (with  $\delta/3$  in each Bernstein bound), with probability at least  $1 - 2\delta/3$ ,

$$T_1(\lambda) + T_2(\lambda) \leq \left[ A \sqrt{\text{Var}(\varepsilon - \lambda r(X))} + B \lambda \right] \|\Delta\|_2.$$

By a union bound with  $\mathbb{P}(\mathcal{E}_{\text{op}}) \geq 1 - \delta/3$ , the following holds with probability at least  $1 - \delta$ :

$$\frac{\kappa_0}{4} \|\Delta\|_2^2 \leq \left[ A \sqrt{\text{Var}(\varepsilon - \lambda r(X))} + B \lambda \right] \|\Delta\|_2.$$

Hence, for any  $\lambda \in [0, 1]$ ,

$$\|\Delta\|_2 \leq \frac{4}{\kappa_0} \left[ A \sqrt{\text{Var}(\varepsilon - \lambda r(X))} + B \lambda \right]. \quad (14)$$

Using the variance expansion

$$\text{Var}(\varepsilon - \lambda r(X)) = \sigma_\varepsilon^2 + \lambda^2 \sigma_r^2 - 2\lambda\rho, \quad \sigma_\varepsilon^2 := \text{Var}(\varepsilon), \quad \sigma_r^2 := \text{Var}(r(X)), \quad \rho := \text{Cov}(\varepsilon, r(X)),$$

we may minimize the right-hand side of (14) over  $\lambda \in [0, 1]$ .

Combining these two displays with the basic inequality (4) yields

$$\|\Delta\|_2 \leq \frac{4}{\kappa_0} \left[ C K_X \sqrt{\frac{d + \log(1/\delta)}{n}} \sqrt{\sigma_\varepsilon^2 + \sigma_r^2 \lambda^2 - 2\rho\lambda} + \lambda C K_X K_r \sqrt{\frac{d + \log(1/\delta)}{N}} \right]. \quad (15)$$

Applying a union bound over the three events and adjusting constant.

**Closed-form minimizer.** Consider

$$g(\lambda) = A \sqrt{\sigma_\varepsilon^2 + \sigma_r^2 \lambda^2 - 2\rho\lambda} + B\lambda, \quad \lambda \in [0, 1].$$

If  $A^2 \sigma_r^2 \neq B^2$ , any interior stationary point satisfies

$$A \frac{\sigma_r^2 \lambda - \rho}{\sqrt{\sigma_\varepsilon^2 + \sigma_r^2 \lambda^2 - 2\rho\lambda}} = -B,$$

which yields

$$\lambda_{\text{unc}} = \frac{\rho(A^2 \sigma_r^2 - B^2) + B \sqrt{(\sigma_\varepsilon^2 \sigma_r^2 - \rho^2)(A^2 \sigma_r^2 - B^2)}}{\sigma_r^2 (A^2 \sigma_r^2 - B^2)}. \quad (16)$$

The bound-optimal choice is the projection onto  $[0, 1]$ :

$$\lambda^* = \min\{1, \max\{0, \lambda_{\text{unc}}\}\}. \quad (17)$$

### 3 Finite-Sample Bound for Stratified PPI

The data are partitioned into  $K$  strata indexed by a random variable  $S \in \{1, \dots, K\}$  with  $\mathbb{P}(S = k) = w_k$ . Let  $n_k$  and  $N_k$  denote the numbers of labeled and unlabeled samples in stratum  $k$ , and define the stratified PPI estimator

$$\hat{\theta} = \arg \min_{\theta} \sum_{k=1}^K w_k L_{k,\lambda_k}^{\text{PP}}(\theta),$$

where

$$L_{k,\lambda_k}^{\text{PP}}(\theta) = \frac{1}{2n_k} \|Y_{(k)} - X_{(k)}\theta\|_2^2 + \lambda_k \left( \frac{1}{2N_k} \|f(\tilde{X}_{(k)}) - \tilde{X}_{(k)}\theta\|_2^2 - \frac{1}{2n_k} \|f(X_{(k)}) - X_{(k)}\theta\|_2^2 \right).$$

#### 3.1 Assumptions

We impose the following standard conditions.

- (A1)  $X \in \mathbb{R}^d$  is mean-zero sub-Gaussian with covariance  $\Sigma = \mathbb{E}[XX^\top]$  and  $\lambda_{\min}(\Sigma) \geq \kappa_0 > 0$ .
- (A2) The noise  $\varepsilon$  is conditionally mean-zero and sub-Gaussian, with  $\mathbb{E}[\varepsilon^2 | S = k] \leq \sigma_{\varepsilon,k}^2$ .
- (A3) The residual  $r(X)$  is deterministic given  $f$ , and  $\mathbb{E}[r(X)^2 | S = k] \leq \sigma_{r,k}^2$ .

#### 3.2 Main Proof

**Step 1: Basic inequality.** Let  $\Delta = \hat{\theta} - \theta^*$ . By optimality of  $\hat{\theta}$ ,

$$\sum_{k=1}^K w_k (L_{k,\lambda_k}^{\text{PP}}(\hat{\theta}) - L_{k,\lambda_k}^{\text{PP}}(\theta^*)) \leq 0.$$

Substituting  $Y = X\theta^* + \varepsilon$  and  $f(X) = X\theta^* + r(X)$  yields

$$\begin{aligned} \sum_{k=1}^K w_k \left( \frac{1 - \lambda_k}{2n_k} \|X_{(k)}\Delta\|_2^2 + \frac{\lambda_k}{2N_k} \|\tilde{X}_{(k)}\Delta\|_2^2 \right) &\leq \sum_{k=1}^K w_k \left[ \frac{1}{n_k} (\varepsilon_{(k)} - \lambda_k r(X_{(k)}))^\top X_{(k)}\Delta \right. \\ &\quad \left. + \lambda_k \frac{1}{N_k} r(\tilde{X}_{(k)})^\top \tilde{X}_{(k)}\Delta \right]. \end{aligned} \quad (18)$$

Under the basic inequality, we have

$$\frac{1}{2} \Delta^\top \hat{\Sigma} \Delta \leq T_1 + T_2,$$

where

$$T_1 := \sum_{k=1}^K \frac{w_k}{n_k} (\varepsilon_{(k)} - \lambda_k r(X_{(k)}))^\top X_{(k)}\Delta, \quad T_2 := \sum_{k=1}^K w_k \frac{\lambda_k}{N_k} \tilde{X}_{(k)} r(\tilde{X}_{(k)})^\top \Delta.$$

$$\hat{\Sigma} := \sum_{k=1}^K w_k \left( \frac{1 - \lambda_k}{n_k} X_{(k)}^\top X_{(k)} + \frac{\lambda_k}{N_k} \tilde{X}_{(k)}^\top \tilde{X}_{(k)} \right).$$

With these definitions, (18) can be written compactly as

$$\frac{1}{2} \Delta^\top \hat{\Sigma} \Delta \leq T_1 + T_2, \quad (19)$$

where

$$T_1 := \sum_{k=1}^K \frac{w_k}{n_k} (\varepsilon_{(k)} - \lambda_k r(X_{(k)}))^\top X_{(k)} \Delta, \quad T_2 := \sum_{k=1}^K w_k \frac{\lambda_k}{N_k} r(\tilde{X}_{(k)})^\top \tilde{X}_{(k)} \Delta.$$

**Step 2: Lower bound on the left-hand side.** For each stratum, define the sample Gram matrices

$$\hat{\Sigma}_{n,k} := \frac{1}{n_k} X_{(k)}^\top X_{(k)}, \quad \hat{\Sigma}_{N,k} := \frac{1}{N_k} \tilde{X}_{(k)}^\top \tilde{X}_{(k)}, \quad \Sigma_k := \mathbb{E}[XX^\top \mid S = k].$$

Under (A1), we use the same sub-Gaussian parameter  $K$  for all strata and assume  $\lambda_{\min}(\Sigma_k) \geq \kappa_0$  for all  $k$ .<sup>1</sup>

By covariance concentration (e.g., Vershynin), there exist absolute constants  $c, C > 0$  such that if

$$\min_k(n_k, N_k) \geq c K^4 \frac{d + \log(6K/\delta)}{\kappa_0^2}, \quad (20)$$

then with probability at least  $1 - \delta/3$  the event

$$\mathcal{E}_{\text{op}} := \left\{ \max_{k \in [K]} \|\hat{\Sigma}_{n,k} - \Sigma_k\|_{\text{op}} \leq \frac{\kappa_0}{2}, \quad \max_{k \in [K]} \|\hat{\Sigma}_{N,k} - \Sigma_k\|_{\text{op}} \leq \frac{\kappa_0}{2} \right\}$$

holds. On  $\mathcal{E}_{\text{op}}$ , Weyl's inequality gives

$$\lambda_{\min}(\hat{\Sigma}_{n,k}) \geq \kappa_0/2, \quad \lambda_{\min}(\hat{\Sigma}_{N,k}) \geq \kappa_0/2, \quad \forall k \in [K].$$

Therefore,

$$\begin{aligned} \frac{1}{2} \Delta^\top \hat{\Sigma} \Delta &= \frac{1}{2} \Delta^\top \sum_{k=1}^K w_k \left( (1 - \lambda_k) \hat{\Sigma}_{n,k} + \lambda_k \hat{\Sigma}_{N,k} \right) \Delta \\ &\geq \frac{1}{2} \sum_{k=1}^K w_k \left( (1 - \lambda_k) \lambda_{\min}(\hat{\Sigma}_{n,k}) + \lambda_k \lambda_{\min}(\hat{\Sigma}_{N,k}) \right) \|\Delta\|_2^2 \\ &\geq \frac{1}{2} \sum_{k=1}^K w_k \cdot \frac{\kappa_0}{2} \|\Delta\|_2^2 = \frac{\kappa_0}{4} \|\Delta\|_2^2. \end{aligned} \quad (21)$$

**Step 3: Decomposition of  $T_1$  and  $T_2$  (bias cancellation).** Define, for each stratum  $k$ ,

$$A_{n,k}(\lambda_k) := \frac{w_k}{n_k} X_{(k)}^\top (\varepsilon_{(k)} - \lambda_k r(X_{(k)})), \quad B_{N,k} := \frac{w_k}{N_k} \tilde{X}_{(k)}^\top r(\tilde{X}_{(k)}).$$

Then

$$T_1 + T_2 = \left\langle \sum_{k=1}^K A_{n,k}(\lambda_k) + \sum_{k=1}^K \lambda_k B_{N,k}, \Delta \right\rangle. \quad (22)$$

---

<sup>1</sup>Equivalently, one may strengthen (A1) to hold conditionally within each stratum.

Let

$$m_k := \mathbb{E}[Xr(X) \mid S = k] \in \mathbb{R}^d.$$

Under (A2) and  $\mathbb{E}[X\varepsilon \mid S = k] = 0$ , we have

$$\mathbb{E}[A_{n,k}(\lambda_k)] = -w_k \lambda_k m_k, \quad \mathbb{E}[B_{N,k}] = w_k m_k.$$

Hence the bias cancels exactly at the level of the weighted sum:

$$\mathbb{E}\left[\sum_{k=1}^K A_{n,k}(\lambda_k) + \sum_{k=1}^K \lambda_k B_{N,k}\right] = 0. \quad (23)$$

Define the centered fluctuations

$$Z_{1,k}(\lambda_k) := A_{n,k}(\lambda_k) - \mathbb{E}[A_{n,k}(\lambda_k)], \quad Z_{2,k} := B_{N,k} - \mathbb{E}[B_{N,k}].$$

Then by (22)–(23),

$$T_1 + T_2 = \left\langle \sum_{k=1}^K Z_{1,k}(\lambda_k) + \sum_{k=1}^K \lambda_k Z_{2,k}, \Delta \right\rangle \leq \left( \left\| \sum_{k=1}^K Z_{1,k}(\lambda_k) \right\|_2 + \left\| \sum_{k=1}^K \lambda_k Z_{2,k} \right\|_2 \right) \|\Delta\|_2.$$

**Step 4: Fluctuation bounds (vector Bernstein).** Fix  $\delta \in (0, 1)$ . Under (A1)–(A3), each coordinate of  $X(\varepsilon - \lambda_k r(X)) \mid (S = k)$  and  $Xr(X) \mid (S = k)$  is sub-exponential. A vector Bernstein inequality (applied within each stratum and then combined via Cauchy–Schwarz) yields that with probability at least  $1 - \delta/3$ ,

$$\left\| \sum_{k=1}^K Z_{1,k}(\lambda_k) \right\|_2 \leq CK \sqrt{d + \log(6/\delta)} \left( \sum_{k=1}^K \frac{w_k^2}{n_k} \text{Var}(\varepsilon - \lambda_k r(X) \mid S = k) \right)^{1/2}, \quad (24)$$

and similarly,

$$\left\| \sum_{k=1}^K \lambda_k Z_{2,k} \right\|_2 \leq CK \sqrt{d + \log(6/\delta)} \left( \sum_{k=1}^K \frac{w_k^2 \lambda_k^2}{N_k} \mathbb{E}[r(X)^2 \mid S = k] \right)^{1/2}. \quad (25)$$

Using (A3),  $\mathbb{E}[r(X)^2 \mid S = k] \leq \sigma_{r,k}^2$ .

Combining (3.2)–(25), we obtain that with probability at least  $1 - \delta/3$ ,

$$T_1 + T_2 \leq CK \sqrt{d + \log(6/\delta)} \left[ \left( \sum_{k=1}^K \frac{w_k^2}{n_k} \text{Var}(\varepsilon - \lambda_k r(X) \mid S = k) \right)^{1/2} + \left( \sum_{k=1}^K \frac{w_k^2 \lambda_k^2}{N_k} \sigma_{r,k}^2 \right)^{1/2} \right] \|\Delta\|_2. \quad (26)$$

**Step 5: Final bound.** On the intersection of  $\mathcal{E}_{\text{op}}$  and the event in (26), which holds with probability at least  $1 - \delta$  by a union bound, the basic inequality (19), the curvature bound (21), and (26) imply

$$\frac{\kappa_0}{4} \|\Delta\|_2^2 \leq T_1 + T_2 \leq \Gamma(\lambda_{1:K}) \|\Delta\|_2,$$

where

$$\Gamma(\lambda_{1:K}) := CK \sqrt{d + \log(6/\delta)} \left[ \left( \sum_{k=1}^K \frac{w_k^2}{n_k} \text{Var}(\varepsilon - \lambda_k r(X) \mid S = k) \right)^{1/2} + \left( \sum_{k=1}^K \frac{w_k^2 \lambda_k^2}{N_k} \sigma_{r,k}^2 \right)^{1/2} \right].$$

Therefore,

$$\|\hat{\theta} - \theta^*\|_2 = \|\Delta\|_2 \leq \frac{4}{\kappa_0} \Gamma(\lambda_{1:K}). \quad (27)$$

**(Explicit variance expansion).** For each stratum,

$$\text{Var}(\varepsilon - \lambda_k r(X) | S = k) = \sigma_{\varepsilon,k}^2 + \lambda_k^2 \sigma_{r,k}^2 - 2\lambda_k \rho_k, \quad \rho_k := \text{Cov}(\varepsilon, r(X) | S = k),$$

so (27) can be optimized over  $\lambda_k \in [0, 1]$  separately for each  $k$ .

### 3.3 Why stratification can improve over unstratified PPI

The bound (27) depends on stratum-wise variance proxies

$$V_k(\lambda) := \text{Var}(\varepsilon - \lambda r(X) | S = k),$$

and on the residual scales  $\sigma_{r,k}^2 = \mathbb{E}[r(X)^2 | S = k]$ .

**Variance reduction from stratum-specific tuning.** An unstratified PPI estimator constrained to use a single global  $\lambda$  is governed by

$$V_{\text{glob}}^* := \min_{\lambda \in [0,1]} \sum_{k=1}^K \frac{w_k^2}{n_k} V_k(\lambda),$$

whereas stratification allows  $\lambda_k$  to be tuned separately, yielding

$$V_{\text{strat}}^* := \sum_{k=1}^K \frac{w_k^2}{n_k} \min_{\lambda \in [0,1]} V_k(\lambda).$$

Since pointwise minimization dominates joint minimization,  $V_{\text{strat}}^* \leq V_{\text{glob}}^*$ , with strict improvement whenever the stratum-wise optimal  $\lambda_k^*$  differ.

**Better use of unlabeled data under heterogeneity.** The unlabeled term in (27) is

$$U(\lambda_{1:K}, N_{1:K}) = \left( \sum_{k=1}^K \frac{w_k^2 \lambda_k^2}{N_k} \sigma_{r,k}^2 \right)^{1/2}.$$

Fix  $\lambda_{1:K}$  and impose  $\sum_{k=1}^K N_k = N$ . Minimizing

$$\sum_{k=1}^K \frac{w_k^2 \lambda_k^2}{N_k} \sigma_{r,k}^2$$

over  $N_k > 0$  yields (by Lagrange multipliers)

$$N_k^* = \frac{N w_k \lambda_k \sigma_{r,k}}{\sum_{j=1}^K w_j \lambda_j \sigma_{r,j}},$$

and therefore

$$U(\lambda_{1:K}, N_{1:K}^*) = \frac{1}{\sqrt{N}} \sum_{k=1}^K w_k \lambda_k \sigma_{r,k}.$$

Hence the optimal unlabeled allocation scales as

$$N_k^* \propto w_k \lambda_k \sigma_{r,k}.$$

## 4 Asymptotic Results: PPI for Linear regression

Assume squared loss

$$\ell_\theta(x, y) = \frac{1}{2}(x^\top \theta - y)^2,$$

and define

$$\varepsilon := Y - X^\top \theta^*, \quad e_f := f(X) - X^\top \theta^*, \quad \Sigma_X := \mathbb{E}[XX^\top].$$

Then

$$\nabla \ell_{\theta^*}(X, Y) = -X\varepsilon, \quad \nabla \ell_{\theta^*}^f(X) = -Xe_f, \quad H_{\theta^*} = \Sigma_X.$$

**Asymptotic covariance.** The asymptotic covariance reduces to

$$\Sigma^\lambda = \Sigma_X^{-1} \left( r\lambda^2 \text{Cov}(Xe_f) + \text{Cov}(X(\lambda e_f - \varepsilon)) \right) \Sigma_X^{-1}$$

If the linear model is correctly specified,  $\text{Cov}(e_f, \varepsilon) = 0$ , then

$$\text{Cov}(Xe_f, X\varepsilon) = 0,$$

and therefore

$$\Sigma^\lambda = \Sigma_X^{-1} \left( (1+r)\lambda^2 \text{Cov}(Xe_f) + \text{Cov}(X\varepsilon) \right) \Sigma_X^{-1}.$$

**Optimal  $\lambda^*$ .** The trace-minimizing  $\lambda^*$  becomes

$$\lambda^* = \frac{\text{Tr}\left(\Sigma_X^{-1} \text{Cov}(X\varepsilon, Xe_f) \Sigma_X^{-1}\right)}{(1+r)\text{Tr}\left(\Sigma_X^{-1} \text{Cov}(Xe_f) \Sigma_X^{-1}\right)}.$$

Under correct linear specification (i.e.,  $\text{Cov}(X\varepsilon, Xe_f) = 0$ ),

$$\lambda^* = 0.$$