

## Appendix

### A.1 Proof of Theorem 1

**Theorem 1.** *If  $v_g$  is a linear function of the features, that is,  $v_g(x) = \theta_*^\top \phi(x)$ , then OIS-LS is an unbiased estimator, that is,  $\mathbb{E}_l[\tilde{\theta}_n] = \theta_*$ .*

*Proof.* The proof is given by the following derivation:

$$\begin{aligned}
\mathbb{E}_l[\tilde{\theta}_n] &= \mathbb{E}_l \left[ \left( \sum_{k=1}^n \phi_k \phi_k^\top \right)^{-1} \sum_{k=1}^n \rho_k Y_k \phi_k \right] \\
&= \mathbb{E}_{l_X} \left[ \left( \sum_{k=1}^n \phi_k \phi_k^\top \right)^{-1} \sum_{k=1}^n \mathbb{E}_{l_{Y|X}} [\rho_k Y_k | X_k] \phi_k \right] \\
&= \mathbb{E}_{l_X} \left[ \left( \sum_{k=1}^n \phi_k \phi_k^\top \right)^{-1} \sum_{k=1}^n \mathbb{E}_{g_{Y|X}} [Y_k | X_k] \phi_k \right] = \mathbb{E}_{l_X} \left[ \left( \sum_{k=1}^n \phi_k \phi_k^\top \right)^{-1} \sum_{k=1}^n v_g(X_k) \phi_k \right] \\
&= \mathbb{E}_{l_X} \left[ \left( \sum_{k=1}^n \phi_k \phi_k^\top \right)^{-1} \sum_{k=1}^n \phi_k \phi_k^\top \theta_* \right] = \mathbb{E}_{l_X} \left[ \left( \sum_{k=1}^n \phi_k \phi_k^\top \right)^{-1} \left( \sum_{k=1}^n \phi_k \phi_k^\top \right) \right] \theta_* = \theta_*. \quad \square
\end{aligned}$$

### A.2 Proof of Theorem 2

**Theorem 2.** *Even if  $v_g$  is a linear function of the features, that is,  $v_g(x) = \theta_*^\top \phi(x)$ , the WIS-LS estimator defined in (6) is a biased estimator, that is,  $\mathbb{E}_l[\hat{\theta}_n] \neq \theta_*$ .*

*Proof.* : We prove it by providing a counterexample to the claim that  $\mathbb{E}_l[\hat{\theta}_n] = \theta_*$ . Consider  $\mathcal{X} = \{x\}$  and  $\phi(x) = 1$ . It is easy to see that in this case  $v_g = \theta_* = \mathbb{E}_g[Y_k]$ . Then the WIS-LS estimator  $\hat{\theta}_n$  reduces to the WIS estimator:

$$\hat{\theta}_n = \left( \sum_{k=1}^n \rho_k \right)^{-1} \sum_{k=1}^n \rho_k Y_k = \hat{v}_g,$$

which is a biased estimator, that is,  $\mathbb{E}_l[\hat{v}_g] \neq v_g$ . Hence, in general,  $\mathbb{E}_l[\hat{\theta}_n] \neq \theta_*$ .  $\square$

### A.3 Proof of Theorem 3

**Theorem 3.** *The OIS-LS estimator  $\tilde{\theta}_n$  is a consistent estimator of the MSE solution  $\theta_*$  given in (4).*

*Proof.* Due to the strong law of large numbers

$$\frac{1}{n} \sum_{k=1}^n \phi_k \phi_k^\top \xrightarrow{w.p.1} \mathbb{E}_{l_X} [\phi_k \phi_k^\top]; \quad \frac{1}{n} \sum_{k=1}^n \rho_k Y_k \phi_k \xrightarrow{w.p.1} \mathbb{E}_l [\rho_k Y_k \phi_k] = \mathbb{E}_{l_X} [\mathbb{E}_{g_{Y|X}} [Y_k | X_k] \phi_k].$$

Then it follows that  $\tilde{\theta}_n \xrightarrow{w.p.1} \theta_*$ .  $\square$

### A.4 Proof of Theorem 4

**Theorem 4.** *The WIS-LS estimator  $\hat{\theta}_n$  is a consistent estimator of the MSE solution  $\theta_*$  given in (4).*

*Proof.* It is very similar to the above proof. The only difference is that here we have to show  $\frac{1}{n} \sum_{k=1}^n \rho_k \phi_k \phi_k^\top \xrightarrow{w.p.1} \mathbb{E}_{l_X} [\phi_k \phi_k^\top]$ . However, it again follows due to the strong law of large numbers noting that  $\mathbb{E}_{l_{XY}} [\rho_k \phi_k \phi_k^\top] = \mathbb{E}_{l_X} [\mathbb{E}_{l_{Y|X}} [\rho_k | X_k] \phi_k \phi_k^\top] = \mathbb{E}_{l_X} [\phi_k \phi_k^\top]$ .  $\square$

## A.5 Proof of Theorem 5

**Theorem 5.** *If the features form an orthonormal basis, then the OIS-LS estimate  $\tilde{\theta}_n^\top \phi(x)$  of input  $x$  is equivalent to the OIS estimate of the outputs corresponding to  $x$ .*

*Proof.* Let  $\Phi$  denote to be the feature matrix the rows of which contain the feature vectors of different unique inputs:  $\Phi = (\phi(x_1), \dots, \phi(x_{|\mathcal{X}|}))^\top$ , where  $x_1, \dots, x_{|\mathcal{X}|}$  are different unique inputs. Then the vector containing the estimated conditional expectation of outputs for each unique input according to the OIS-LS estimator can be written as

$$\Phi \tilde{\theta}_n = \Phi \left( \sum_{x \in \mathcal{X}} n_x \phi(x) \phi(x)^\top \right)^{-1} \sum_{x \in \mathcal{X}} \left( \sum_{i=1}^{n_x} \rho_{x,i} Y_{x,i} \right) \phi(x) = \Phi (\Phi^\top \mathbf{N} \Phi)^{-1} \Phi^\top \mathbf{y},$$

where  $n_x$  is the number of times input  $x$  is observed among  $n$  samples,  $Y_{x,i}$  is the output corresponding to the  $i$ th occurrence of input  $x$  and  $\rho_{x,i}$  is the corresponding importance-sampling ratio. Here,  $\mathbf{N}$  is a diagonal matrix where the  $i$ th diagonal element contains  $n_{x_i}$ :  $\mathbf{N} = \text{diag}(n_{x_1}, \dots, n_{x_{|\mathcal{X}|}})$  and  $\mathbf{y} = (\sum_{i=1}^{n_{x_1}} \rho_{x_1,i} Y_{x_1,i}, \dots, \sum_{i=1}^{n_{x_{|\mathcal{X}|}}} \rho_{x_{|\mathcal{X}|},i} Y_{x_{|\mathcal{X}|},i})^\top$ .

Note that, due to orthonormality of the features,  $\Phi$  is necessarily a square matrix and full rank. Therefore, it follows that the vector of the estimates can be written as

$$\Phi \tilde{\theta}_n = \Phi \Phi^{-1} \mathbf{N}^{-1} \Phi^\top \Phi^\top \mathbf{y} = \mathbf{N}^{-1} \mathbf{y}.$$

The element of this vector corresponding to any input  $x$  is the ordinary importance-sampling estimator of its corresponding outputs:  $n_x^{-1} \sum_{i=1}^{n_x} \rho_{x,i} Y_{x,i}$ .  $\square$

## A.6 Proof of Theorem 6

**Theorem 6.** *If the features form an orthonormal basis, then the WIS-LS estimate  $\hat{\theta}_n^\top \phi(x)$  of input  $x$  is equivalent to the WIS estimate of the outputs corresponding to  $x$ .*

*Proof.* The proof is similar to the proof of Theorem 5. First, we write the vector of the estimates according to the WIS-LS estimate as

$$\Phi \hat{\theta}_n = \Phi \left( \sum_{x \in \mathcal{X}} \left( \sum_{i=1}^{n_x} \rho_{x,i} \right) \phi(x) \phi(x)^\top \right)^{-1} \sum_{x \in \mathcal{X}} \left( \sum_{i=1}^{n_x} \rho_{x,i} Y_{x,i} \right) \phi(x) = \Phi (\Phi^\top \mathbf{R} \Phi)^{-1} \Phi^\top \mathbf{y},$$

where  $\mathbf{R}$  is a diagonal matrix with each diagonal element containing the total summation of the importance-sampling ratios corresponding to each input:  $\mathbf{R} = \text{diag}\left(\left(\sum_{i=1}^{n_{x_1}} \rho_{x_1,i}\right), \dots, \left(\sum_{i=1}^{n_{x_{|\mathcal{X}|}}} \rho_{x_{|\mathcal{X}|},i}\right)\right)$ . Hence, the vector of estimates can be written as

$$\Phi \hat{\theta}_n = \Phi \Phi^{-1} \mathbf{R}^{-1} \Phi^\top \Phi^\top \mathbf{y} = \mathbf{R}^{-1} \mathbf{y},$$

The element of this vector corresponding to any input  $x$  is the WIS estimate of its corresponding outputs:  $\left(\sum_{i=1}^{n_x} \rho_{x,i}\right)^{-1} \sum_{i=1}^{n_x} \rho_{x,i} Y_{x,i}$ .  $\square$

## A.7 Proof of Theorem 7

**Theorem 7.** *At termination, the algorithm defined by (7) is equivalent to the WIS-LS method in the sense that if  $\lambda_0 = \dots = \lambda_t = \gamma_0 = \dots = \gamma_{t-1} = 1$  and  $\gamma_t = 0$ , then  $\theta_t$  defined in (7) equals  $\hat{\theta}_t$  as defined in (6), with  $Y_k \doteq G_k^t$ .*

*Proof.* When  $\gamma_0 = \dots = \gamma_{t-1} = 1$ ,  $\gamma_t = 0$  and also  $\lambda_0 = \dots = \lambda_t = 1$ , then

$$\mathbf{b}_{k,t} = \prod_{j=k}^{t-1} \rho_j G_k^t \phi_k = \rho_k^t G_k^t \phi_k, \quad \mathbf{A}_{k,t} = \prod_{j=k}^{t-1} \rho_j \phi_k \phi_k^\top = \rho_k^t \phi_k \phi_k^\top.$$

Hence, the solution can be written as  $\theta_t = \mathbf{A}_t^{-1} \mathbf{b}_t = \left(\sum_{k=0}^{t-1} \mathbf{A}_{k,t}\right)^{-1} \sum_{k=0}^{t-1} \mathbf{b}_{k,t} = \left(\sum_{k=0}^{t-1} \rho_k^t \phi_k \phi_k^\top\right)^{-1} \sum_{k=0}^{t-1} \rho_k^t G_k^t \phi_k$ , which is the WIS-LS solution.  $\square$

### A.8 Derivations of the recursive updates of $\mathbf{b}_{k,t}$ and $\mathbf{A}_{k,t}$ in $t$

The derivations are given below:

$$\begin{aligned}
\mathbf{b}_{k,t+1} &= \rho_k \sum_{i=k+1}^t C_k^{i-1} (1 - \gamma_i \lambda_i) G_k^i \boldsymbol{\phi}_k + \rho_k C_k^t G_k^{t+1} \boldsymbol{\phi}_k \\
&= \rho_k \sum_{i=k+1}^{t-1} C_k^{i-1} (1 - \gamma_i \lambda_i) G_k^i \boldsymbol{\phi}_k + (1 - \gamma_t \lambda_t) \rho_k C_k^{t-1} G_k^t \boldsymbol{\phi}_k \\
&\quad + \rho_t \gamma_t \lambda_t \rho_k C_k^{t-1} (G_k^t + R_{t+1}) \boldsymbol{\phi}_k \\
&= \mathbf{b}_{k,t} + \rho_k C_k^t R_{t+1} \boldsymbol{\phi}_k + (\rho_t - 1) \gamma_t \lambda_t \rho_k C_k^{t-1} G_k^t \boldsymbol{\phi}_k, \\
\mathbf{A}_{k,t+1} &= \rho_k \sum_{i=k+1}^t C_k^{i-1} \boldsymbol{\phi}_k ((1 - \gamma_i \lambda_i) \boldsymbol{\phi}_k - \gamma_i (1 - \lambda_i) \boldsymbol{\phi}_i)^\top + \rho_k C_k^t \boldsymbol{\phi}_k (\boldsymbol{\phi}_k - \gamma_{t+1} \boldsymbol{\phi}_{t+1})^\top \\
&= \rho_k \sum_{i=k+1}^{t-1} C_k^{i-1} \boldsymbol{\phi}_k ((1 - \gamma_i \lambda_i) \boldsymbol{\phi}_k - \gamma_i (1 - \lambda_i) \boldsymbol{\phi}_i)^\top \\
&\quad + \rho_k C_k^{t-1} \boldsymbol{\phi}_k ((1 - \gamma_t \lambda_t) \boldsymbol{\phi}_k - \gamma_t (1 - \lambda_t) \boldsymbol{\phi}_t)^\top + \rho_k C_k^t \boldsymbol{\phi}_k (\boldsymbol{\phi}_k - \gamma_{t+1} \boldsymbol{\phi}_{t+1})^\top \\
&= \rho_k \sum_{i=k+1}^{t-1} C_k^{i-1} \boldsymbol{\phi}_k ((1 - \gamma_i \lambda_i) \boldsymbol{\phi}_k - \gamma_i (1 - \lambda_i) \boldsymbol{\phi}_i)^\top + \rho_k C_k^{t-1} \boldsymbol{\phi}_k (\boldsymbol{\phi}_k - \gamma_t \boldsymbol{\phi}_t)^\top \\
&\quad + \rho_k C_k^t \boldsymbol{\phi}_k (\boldsymbol{\phi}_k - \gamma_{t+1} \boldsymbol{\phi}_{t+1})^\top - \gamma_t \lambda_t \rho_k C_k^{t-1} \boldsymbol{\phi}_k (\boldsymbol{\phi}_k - \boldsymbol{\phi}_t)^\top \\
&= \mathbf{A}_{k,t} + \rho_k C_k^t \boldsymbol{\phi}_k (\boldsymbol{\phi}_k - \boldsymbol{\phi}_t + \boldsymbol{\phi}_t - \gamma_{t+1} \boldsymbol{\phi}_{t+1})^\top - \gamma_t \lambda_t \rho_k C_k^{t-1} \boldsymbol{\phi}_k (\boldsymbol{\phi}_k - \boldsymbol{\phi}_t)^\top \\
&= \mathbf{A}_{k,t} + \rho_k C_k^t \boldsymbol{\phi}_k (\boldsymbol{\phi}_t - \gamma_{t+1} \boldsymbol{\phi}_{t+1})^\top + (\rho_t - 1) \gamma_t \lambda_t \rho_k C_k^{t-1} \boldsymbol{\phi}_k (\boldsymbol{\phi}_k - \boldsymbol{\phi}_t)^\top.
\end{aligned}$$