

Patient Representation

Derek Modzelewski

November 16, 2017

1 Weighting Patients

Not all patients had the same tissues sampled. Furthermore, not all patients have the same *number* of tissues sampled. Thus, we are more confident about the learned representations for some patients (the ones with more samples) than for other patients (the ones with fewer samples). This factors into our prediction of the tissue transform matrices and tissue centers, which then affect the prediction of patient representations.

Instead of solving for the optimal weights directly, we use the weights generated from a related problem which is significantly easier to work with. It is our belief that the weighting scheme will be optimal for our patient representation problem as well.

Theorem 1.1. *Given many estimations $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_n$ of parameter $\theta_0 \in \mathbb{R}^d$, where*

$$\hat{\theta} \equiv (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_n) = \theta_0 \mathbf{1}^T + B + u + r$$

$$\mathbb{E}[u] = \mathbb{E}[r] = \mathbf{0}$$

where B is the bias, u is the "uncertainty" in the bias and has mean 0 and covariance matrix ϵ , and r is the "random error" and has mean 0 and covariance DVD where D is a diagonal matrix where $D_{i,i} = 1/\sqrt{n_i}$ $\forall i \in [n]$, and u and r are independent

Then the estimation $\hat{\theta}$ of θ_0 which minimizes squared loss is given by $w^T \hat{\theta}$, where

$$w = k(2B^T B + \epsilon + DVD)^{-1} \mathbf{1}$$

for suitable $k \in \mathbb{R}$ s.t. $w^T \mathbf{1} = 1$

Proof. any estimation $\delta(\hat{\theta})$ of θ_0 using can be expressed as $\delta(\hat{\theta}) = \hat{\theta} \cdot w + v$ for some $w \in \mathbb{R}^n$ s.t. $\mathbf{1}^T w = 1$, some $v \in \mathbb{R}^d$, for given $\hat{\theta}$.

$$\begin{aligned} \mathbb{E}[||\delta(\hat{\theta}) - \theta_0||^2] &= \mathbb{E}[||\hat{\theta} \cdot w + v - \theta_0||^2] \\ &= \mathbb{E}[||\theta_0 \mathbf{1}^T w + Bw + uw + rw + v - \theta_0||^2] \\ &= \mathbb{E}[||\theta_0 + Bw + uw + rw + v - \theta_0||^2] \\ &= \mathbb{E}[||Bw + uw + rw + v||^2] \\ &= \mathbb{E}[|(B + u + r)w + v|^2] \\ &= \mathbb{E}[w^T (B + u + r)^T (B + u + r) w + 2v^T (B + u + r) w + v^T v] \\ &= \mathbb{E}[w^T (B + u + r)^T (B + u + r) w] + 2v^T \mathbb{E}[B + u + r] w + v^T v \\ &= w^T \mathbb{E}[(B + u + r)^T (B + u + r)] w + 2v^T Bw + v^T v \\ &= w^T \mathbb{E}[B^T B + u^T u + r^T r] w + 2v^T Bw + v^T v \\ &= w^T \mathbb{E}(B^T B + \epsilon + DVD) w + 2v^T Bw + v^T v \end{aligned}$$

At minimum, the gradient with respect to w and with respect to v will be normal to the feasible spaces of w and u , respectively.

v has full rank feasible region (no constraints), so only the 0 vector is normal to its feasible region:

$$\begin{aligned}
\Rightarrow 0 &= \nabla_v \left(\mathbb{E} [||\delta(\hat{\theta}) - \theta_0||^2] \right) \\
&= \nabla_v \left(w^T (B^T B + \epsilon + DVD) w + 2v^T Bw + v^T v \right) \\
&= \nabla_v \left(2v^T Bw + v^T v \right) \\
&= 2Bw + 2v \\
&\Rightarrow v = Bw
\end{aligned}$$

For w , our only constraint is that $\mathbf{1}^T w = 1$ \therefore for all feasible w , all vectors d normal to the feasible region can be expressed $d = k\mathbf{1}$ for some $k \in \mathbb{R}$

$$\begin{aligned}
\Rightarrow k\mathbf{1} &= \nabla_w \left(\mathbb{E} [||\delta(\hat{\theta}) - \theta_0||^2] \right) \\
&= \nabla_w \left(w^T (B^T B + \epsilon + DVD) w + 2v^T Bw + v^T v \right) \\
&= \nabla_w \left(w^T (B^T B + \epsilon + DVD) w + 2v^T Bw \right) \\
&= 2(B^T B + \epsilon + DVD)w + (2v^T B)^T \\
&= 2(B^T B + \epsilon + DVD)w + 2B^T v \\
&= 2(B^T B + \epsilon + DVD)w + 2B^T Bw \\
&= 2(2B^T B + \epsilon + DVD)w \\
&\Rightarrow w = k(2B^T B + \epsilon + DVD)^{-1} \mathbf{1}
\end{aligned}$$

for appropriate k

2 Weighting Patients

Not all patients had the same tissues sampled. Furthermore, not all patients have the same *number* of tissues sampled. Thus, we are more confident about the learned representations for some patients (the ones with more samples) than for other patients (the ones with fewer samples). This factors into our prediction of the tissue transform matrices and tissue centers, which then affect the prediction of patient representations.

Instead of solving for the optimal weights directly, we use the weights generated from a related problem which is significantly easier to work with. It is our belief that the weighting scheme will be optimal for our patient representation problem as well.

Theorem 2.1. *Given many estimations $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_n$ of parameter $\theta_0 \in \mathbb{R}^d$, where*

$$\hat{\theta} \equiv (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_n) = \theta_0 \mathbf{1}^T + r$$

Then the estimation $\tilde{\theta}$ of θ_0 which minimizes squared loss is given by $\hat{\theta} \cdot w + v$ where

$$w = k(\mathbb{E}[r^T r] - \mathbb{E}[r]^T \mathbb{E}[r])^{-1} \mathbf{1}$$

for suitable $k \in \mathbb{R}$ s.t. $w^T \mathbf{1} = 1$, and

$$v = -\mathbb{E}[r]w$$

Proof. any estimation $\delta(\hat{\theta}_.)$ of θ_0 using can be expressed as $\delta(\hat{\theta}_.) = \hat{\theta}_.w + v$ for some $w \in \mathbb{R}^n$ s.t. $\mathbf{1}^T w = 1$, some $v \in \mathbb{R}^d$, for given $\hat{\theta}_.$

$$\begin{aligned}\mathbb{E}[||\delta(\hat{\theta}_.) - \theta_0||^2] &= \mathbb{E}[||\hat{\theta}_.w + v - \theta_0||^2] \\ &= \mathbb{E}[||\theta_0 \mathbf{1}^T w + rw + v - \theta_0||^2] \\ &= \mathbb{E}[||rw + v||^2] \\ &= \mathbb{E}[w^T r^T r w + 2v^T r w + v^T v] \\ &= w^T \mathbb{E}[r^T r] w + 2v^T \mathbb{E}[r] w + v^T v\end{aligned}$$

At minimum, the gradient with respect to w and with respect to v will be normal to the feasible spaces of w and u , respectively.

v has full rank feasible region (no constraints), so only the 0 vector is normal to its feasible region:

$$\begin{aligned}\Rightarrow 0 &= \nabla_v \left(\mathbb{E}[||\delta(\hat{\theta}_.) - \theta_0||^2] \right) \\ &= \nabla_v \left(w^T \mathbb{E}[r^T r] w + 2v^T \mathbb{E}[r] w + v^T v \right) \\ &= \nabla_v \left(2v^T \mathbb{E}[r] w + v^T v \right) \\ &= 2\mathbb{E}[r] w + 2v \\ &\Rightarrow v = -\mathbb{E}[r] w\end{aligned}$$

For w , our only constraint is that $\mathbf{1}^T w = 1$ \therefore for all feasible w , all vectors d normal to the feasible region can be expressed $d = k\mathbf{1}$ for some $k \in \mathbb{R}$

$$\begin{aligned}\Rightarrow k\mathbf{1} &= \nabla_w \left(\mathbb{E}[||\delta(\hat{\theta}_.) - \theta_0||^2] \right) \\ &= \nabla_w \left(w^T \mathbb{E}[r^T r] w + 2v^T \mathbb{E}[r] w + v^T v \right) \\ &= \nabla_w \left(w^T \mathbb{E}[r^T r] w + 2v^T \mathbb{E}[r] w \right) \\ &= 2\mathbb{E}[r^T r] w + (2v^T \mathbb{E}[r])^T \\ &= 2\mathbb{E}[r^T r] w + 2\mathbb{E}[r]^T v \\ &= 2\mathbb{E}[r^T r] w - 2\mathbb{E}[r]^T \mathbb{E}[r] w \\ &= 2(\mathbb{E}[r^T r] - \mathbb{E}[r]^T \mathbb{E}[r]) w \\ &\Rightarrow w = k(\mathbb{E}[r^T r] - \mathbb{E}[r]^T \mathbb{E}[r])^{-1} \mathbf{1}\end{aligned}$$

for appropriate k

Remark 1. The optimal v and w values are independent of the dimension of the parameter space, d

Corollary 2.1.1. If we restrict our estimation to be in the convex hull of $\{\hat{\theta}_1, \dots, \hat{\theta}_n\}$, the optimal solution is given by $\hat{\theta}_.w$ where

$$w = k\mathbb{E}[r^T r]^{-1} \mathbf{1}$$

for suitable k as before.

Proof. Any point in the convex hull of $\{\hat{\theta}_1, \dots, \hat{\theta}_n\}$ may be expressed as $\hat{\theta}_.w$, which is equivalent to setting v to $\mathbf{0}$. In the previous proof, we arrived at

$$\begin{aligned}\dots k\mathbf{1} &= 2\mathbb{E}[r^T r] w + 2\mathbb{E}[r]^T v \\ &= 2\mathbb{E}[r^T r] w \qquad \qquad \qquad \therefore v = \mathbf{0}\end{aligned}$$

$$\Rightarrow w = k\mathbb{E}[r^T r]^{-1} \mathbf{1}$$

Remark 2. It may often be convenient to separate $\mathbb{E}[r^T r] = \mathbb{E}[r]^T \mathbb{E}[r] + \text{Var}(r)$, especially for biased estimators.

Corollary 2.1.2. In the special case that $r = u + v$ where

$$\mathbb{E}[u] = \mathbb{E}[v] = \mathbf{0}$$

and $\mathbb{E}[u^T u] = \phi I$ for some $\phi \in \mathbb{R}$, $\mathbb{E}[u^T v] = \mathbf{0}$, and $\mathbb{E}[v^T v]$ is diagonal where $\mathbb{E}[v^T v]_{i,i} = \gamma/n_i \ \forall i \in [n]$, for some $\gamma \in \mathbb{R}$, and we are restricted to estimations in the convex hull of $\{\hat{\theta}_1, \dots, \hat{\theta}_n\}$, the optimal estimator, $\tilde{\theta}$, is

$$\tilde{\theta} = \frac{\sum_{i=1}^n \frac{n_i}{n_i + \gamma/\phi} \hat{\theta}_i}{\sum_{i=1}^n \frac{n_i}{n_i + \gamma/\phi}}$$

Remark 3. This is equivalent to placing a $n_i/(n_i + \gamma/\phi)$ weight to each estimator θ_i

Proof.

$$\begin{aligned} \mathbb{E}[r^T r] &= \mathbb{E}[u^T u] + \mathbb{E}[u^T v] + \mathbb{E}[v^T u] + \mathbb{E}[v^T v] \\ &= \phi I + \mathbb{E}[v^T v] \end{aligned}$$

$\therefore \mathbb{E}[r^T r]$ is diagonal where $\mathbb{E}[r^T r]_{i,i} = \phi + \gamma/n_i \ \forall i \in [n]$
 $\Rightarrow \mathbb{E}[r^T r]^{-1}$ is diagonal where

$$\begin{aligned} \mathbb{E}[r^T r]_{i,i}^{-1} &= \frac{1}{\phi + \gamma/n_i} \\ &\propto \frac{n_i}{n_i + \gamma/\phi} \end{aligned}$$

Using Corollary @ref?? the weight given to estimator θ_i is $w_i \propto \frac{n_i}{n_i + \gamma/\phi}$

One interpretation of this is that each of our estimators, θ_i , are biased (and inconsistent) with unknown bias, but each estimator has independent bias with equal variance (in whatever system generates the bias), and variance proportional to the inverse of the number of samples. This is extremely related to our problem of weighting patients, as each patient's samples help estimate our linear representation of patients and tissues, but has unknown bias from other sources (confounders), and it is very reasonable to expect the variance to be proportional to the inverse of the number of samples.