

Consider parametrized families $\lambda(k, \sigma, \omega)$ and $Q(i, u; \theta, \omega)$. Consider

$$\begin{aligned}
\theta_{n+1} &= \theta_n - a(n) \left(\nabla_{\theta} Q(X_{n+1}, v_n; \theta_n, \omega_n) - \nabla_{\theta} f(Q(\theta_n, \omega_n)) \right) \Big|_{\theta=\theta_n} \\
&\quad - \nabla_{\theta} Q(X_n, U_n; \theta_n, \omega_n) \Big) \times \\
&\quad \frac{\left((1 - U_n)(r(X_n, 0) + \lambda(X_n, \sigma_n, \omega_n)) + U_n r_n(X_n, 1) + \max_{v \in \{0,1\}} Q(X_{n+1}, v; \theta_n, \omega_n) \right.}{\left. - f(Q(\theta_n, \omega_n)) - Q(X_n, U_n; \theta_n, \omega_n) \right) + a(n) \xi_{n+1}}, \\
\sigma_{n+1} &= \sigma_n - b(n) \left(Q(X_n, 1; \theta_n, \omega_n) - r(X_n, 0) + f(Q(\theta_n, \omega_n)) \right. \\
&\quad \left. - \max_{v \in \{0,1\}} Q(X_{n+1}, v; \theta_n, \omega_n) - \lambda(X_n, \sigma_n, \omega_n) \right) \times \\
&\quad \left(- \nabla_{\sigma} \lambda(X_n, \sigma_n, \omega_n) \right), \\
\omega_{n+1} &= \omega_n - b(n) \left(\nabla_{\omega} Q(X_{n+1}, v_n; \theta_n, \omega_n) - \nabla_{\omega} f(Q(\theta_n, \omega_n)) \right) \Big|_{\omega=\omega_n} \\
&\quad - \nabla_{\omega} Q(X_n, U_n; \theta_n, \omega_n) \Big) \times \\
&\quad \frac{\left((1 - U_n)(r(X_n, 0) + \lambda(X_n, \sigma_n, \omega_n)) + U_n r_n(X_n, 1) + \max_{v \in \{0,1\}} Q(X_{n+1}, v; \theta_n, \omega_n) \right.}{\left. - f(Q(\theta_n, \omega_n)) - Q(X_n, U_n; \theta_n, \sigma_n, \omega_n) \right) + a(n) \xi_{n+1}} \\
&\quad - b(n) \left(Q(X_n, 1; \theta_n, \omega_n) - r(X_n, 0) + f(Q(\theta_n, \omega_n)) \right. \\
&\quad \left. - \max_{v \in \{0,1\}} Q(X_{n+1}, v; \theta_n, \omega_n) - \lambda(X_n, \sigma_n, \omega_n) \right) \times \\
&\quad \left(\nabla_{\omega} Q(X_{n+1}, v_n; \theta_n, \omega_n) - \nabla_{\omega} f(Q(\theta_n, \omega_n)) \right) \Big|_{\omega=\omega_n} \\
&\quad - \nabla_{\omega} Q(X_n, U_n; \theta_n, \omega_n) - \nabla_{\omega} \lambda(X_n, \sigma_n, \omega_n) \Big),
\end{aligned}$$

The θ_n iteration is the SGD for the mean square error

$$\begin{aligned}
\mathcal{E}_1 &:= E \left[\left\| (1 - U_n)(r(X_n, 0) + \lambda(X_n, \sigma_n, \omega_n)) + U_n r_n(X_n, 1) \right. \right. \\
&\quad \left. \left. + \max_{v \in \{0,1\}} Q(X_{n+1}, v; \theta_n, \omega_n) - f(Q(\theta_n, \omega_n)) - Q(X_n, U_n; \theta_n, \omega_n) \right\|^2 \right].
\end{aligned}$$

The σ_n iteration is the SGD to minimize the mean square error

$$\begin{aligned}
\mathcal{E}_2 &:= E \left[\left\| Q(X_n, 1; \theta_n, \omega_n) - r(X_n, 0) + f(Q(\theta_n, \omega_n)) \right. \right. \\
&\quad \left. \left. - \max_{v \in \{0,1\}} Q(X_{n+1}, v; \theta_n, \omega_n) - \lambda(X_n, \sigma_n, \omega_n) \right\|^2 \right].
\end{aligned}$$