**Algorithm 1** Preference-Based Learning of Reward Functions

---

1: **Input:** Features $\phi$, horizon $N$, dynamics $f$, *iter*
2: **Output:** Distribution of **w**: $p(\mathbf{w})$
3: Initialize $p(\mathbf{w}) \sim \text{Uniform}(B)$, for a unit ball $B$
4: **While** $t < iter$:
5: $\quad W \leftarrow M$ samples from AdaptiveMetropolis($p(\mathbf{w})$)
6: $\quad (x^0, \mathbf{u}_R, \mathbf{u}_H^A, \mathbf{u}_H^B) \leftarrow \text{SynthExps}(W, f)$
7: $\quad I_t \leftarrow \text{QueryHuman}(x^0, \mathbf{u}_R, \mathbf{u}_H^A, \mathbf{u}_H^B)$
8: $\quad \varphi = \Phi(x^0, \mathbf{u}_R, \mathbf{u}_H^A) - \Phi(x^0, \mathbf{u}_R, \mathbf{u}_H^B)$
9: $\quad f_\varphi(\mathbf{w}) = \min(1, I_t \exp(\mathbf{w}^\top \varphi))$
10: $\quad p(\mathbf{w}) \leftarrow p(\mathbf{w}) \cdot f_\varphi(\mathbf{w})$
11: $\quad t \leftarrow t + 1$
12: **End for**