
Supplementary Material for “Learning efficient contextual categorical sequence generation models with pseudo-action policy gradient”

Anonymous Author(s)

Affiliation

Address

email

¹ **A Qualitative results**

² **A.1 Pseudo sentences produced by ARS-K algorithm.**



Greedy sentence: a man riding a motorcycle with a dog.

Pseudo sentence 1: a man riding a motorcycle with hay on it.

Pseudo sentence 2: a man riding a motorcycle with pack of sheep.



Greedy sentence: a group of people flying kites in a field.

Pseudo sentence 1: a group of teenagers standing in a field flying kites.



4

Greedy sentence: a man and woman are standing in front of a table.

Pseudo sentence 1: a man and woman are standing in an market.

Pseudo sentence 2: a man and woman are standing in the street.

Pseudo sentence 3: a man and woman are standing in to a tent.

Pseudo sentence 4: a man and woman are standing in front of a table.



Greedy sentence: a city that has a large white building on it.

Pseudo sentence 1: a city with a red traffic and a large building.

Pseudo sentence 2: a city intersection with a traffic light and a street sign.

Pseudo sentence 3: a city bus is driving down the street.

Pseudo sentence 4: a city street with a city street with cars parked on it.

Pseudo sentence 5: a city road with a traffic light and a street sign.

5

6 B Algorithms

Algorithm 1: ARS-K/ARSM(K=V) policy gradient for fine-tuning a contextual categorical sequence generation model with a discrete-action space of V actions.

input : MLE pretrained policy parameter θ , K
output : Fine-tuned policy parameter θ ;

```

while not converged do
    Receive random sample  $x, y$ ;
     $\theta_{\text{update}} = 0$ ;
    for  $t = 1 : T$  do
         $\phi_t = \mathcal{T}_\theta(z_{1:t-1}, x)$ ;
         $\mathbf{g}_{\phi_t} = 0$ ;
        Sample  $\pi_t \sim \text{Dirichlet}(\mathbf{1}_V)$ ;
        Let  $z_t = \text{argmin}_{i \in \{1, \dots, V\}} (\ln \pi_{ti} - \phi_{ti})$ ;
        Let  $j_1, \dots, j_K$  be randomly selected  $K$  reference categories without replacement;
        for  $k = 1, \dots, K$  do
            for  $v = 1, \dots, V$  (in parallel) do
                Let  $z_t^{v \leftarrow j_k} := \text{argmin}_{i \in \{1, \dots, V\}} (\ln \pi_i^{v \leftarrow j_k} - \phi_{ti})$  as the  $(v, k)$ th pseudo action;
            end for
        end for
        Let  $S_t = \text{unique}(\{z_t^{v \leftarrow j}\}_{v,j})$  which means  $S_t$  is the set of all unique values in
         $\{z_t^{v \leftarrow j}\}_{v,j}$ , denote the cardinality of  $S_t$  as  $|S_t|$ ;
        for  $s \in \{1, \dots, |S_t|\}$  (in parallel) do
            Let  $\tilde{z}_{ts} = S_t(s)$  be the  $s$ th unique pseudo action at time  $t$ ;
            If  $t < T$ , sample  $z_{t+1:T}^s \sim p_\theta(z_{t+1:T} | z_{1:t-1}, \tilde{z}_{ts}, x)$ ;
            Let  $f(z_t^{v \leftarrow j_k}) = r(z_{1:t-1}, z_{t:T}^s | x, y)$  if  $z_t^{v \leftarrow j_k} = \tilde{z}_{ts}$ ;
        end for
        for  $k = 1, \dots, K$  do
            Let  $\bar{f}_t = \frac{1}{V} \sum_{v=1}^V f(z_t^{v \leftarrow j_k})$ ;
            Let  $g_{\phi_{tv}} = g_{\phi_{tv}} + \frac{1}{K} (f(z_t^{v \leftarrow j_k}) - \bar{f}_t) (1 - V \pi_{j_k})$ ;
             $\theta_{\text{update}} = \theta_{\text{update}} + \eta_\theta \nabla_\theta \phi_t \mathbf{g}_{\phi_t}$ , with step-size  $\eta_\theta$ 
        end for
    end for
     $\theta = \theta + \theta_{\text{update}}$ 
end while

```

Algorithm 2: Binary-tree-ARSM policy gradient for fine-tuning a binary-tree contextual categorical sequence generation model.

input : MLE pretrained policy parameter θ , K
output : Fine-tuned policy parameter θ ;

while not converged **do**

```

    Receive random sample  $x, y$ ;
     $\theta_{\text{update}} = 0$ ;
    for  $t = 1 : T$  do
         $\phi_t = \mathcal{T}_\theta(\mathbf{z}_{1:t-1}, x)$ ;
        for  $l = 1 : D$  do
            Sample  $\pi_{tl} \sim \text{Uniform}(0, 1)$ ;
            Let  $b_{tl}^{(1)} = \mathbf{1}_{[\pi_{tl} < \sigma(\phi_t, b_{t(1:l-1)})]}$ ;
            Let  $b_{tl} = b_{tl}^{(1)}$ ;
            Let  $b_{tl}^{(2)} = \mathbf{1}_{[\pi_{tl} > \sigma(-\phi_t, b_{t(1:l-1)})]}$ ;
            if  $b_{tl}^{(1)} \neq b_{tl}^{(2)}$  then
                If  $l < D$ , sample  $b_{t(l+1:D)}^{(j)} \sim p_\theta(b_{t(l+1:D)}^{(1)} | \mathbf{z}_{1:t-1}, b_{t,1:l-1}, b_{tl}^{(j)}, x)$ ,  $j = 1, 2$ ;
                Let  $\mathbf{z}_t^{(j)} = \nu(b_{t(1:D)}^{(j)})$ ,  $j = 1, 2$ ;
                If  $t < T$ , sample  $\mathbf{z}_{t+1:T}^{(j)} \sim p_\theta(\mathbf{z}_{t+1:T} | \mathbf{z}_{1:t-1}, \mathbf{z}_t^{(j)}, x)$ ,  $j = 1, 2$ ;
                Let  $f_{tl}^{(j)} = r(\mathbf{z}_{1:t-1}, \mathbf{z}_{t:T}^{(j)} | \mathbf{x}, \mathbf{y})$ ;
                Let  $g_{\phi_t, b_{t(1:l-1)}} = \frac{1}{2}(f_{tl}^{(1)} - f_{tl}^{(2)})(1 - 2\pi_{tl})$ ;
                 $\theta_{\text{update}} = \theta_{\text{update}} + \eta_\theta \nabla_\theta \phi_{t, b_{t(1:l-1)}} g_{\phi_t, b_{t(1:l-1)}}$ , with step-size  $\eta_\theta$ 
            end if
        end for
         $\theta = \theta + \theta_{\text{update}}$ 
    end for
end while

```
