Research Presentation

Active Learning with V-Learning

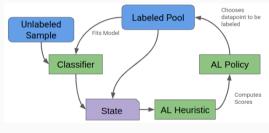
Thorben Werner March 31, 2022

Information Systems and Machine Learning Lab (ISMLL)
Institute for Computer Science
University of Hildesheim

Pool-Based Active Learning



 $\mathcal{L} \in \mathbb{R}^{\lambda imes m}$ Labeled Set $\mathcal{U} \in \mathbb{R}^{\mu imes m}$ Unlabeled Set $\phi_{ heta} := \mathbb{R}^m o \mathbb{R}^c$ Classifier $\mathcal{K} \in \mathbb{R}^{k imes m}$ Unlabeled Sample $\psi := \mathbb{R}^{k imes m} o \mathbb{R}^k$ Active Learning Heuristic $\pi_{\psi} := argmax \ \psi(\mathcal{S})$ Active Learning Policy



Active Learning Cycle

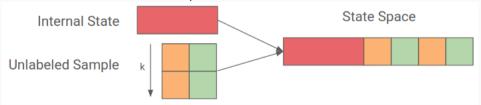
Active Learning with Q-Learning



State Space:

- 1. Internal state of the environment: content of labeled pool \mathcal{L} , state of the classifier θ , remaining budget, current F1-Score, etc. $\to \mathbb{R}^a$
- 2. Information about the unlabeled sample K: Output of the classifier $\phi_{\theta}(k_t)$, the datapoints themselves, other metrics, etc. $\to \mathbb{R}^{k \times b}$

Results in a flattened state space of $S \in \mathbb{R}^{a+k \times b}$



Active Learning with Q-Learning



 $\mathcal{A} \in [\mathsf{o}, \dots, k]$ Action Space $\mathcal{R} := \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ Reward Function $\tau := \{\mathcal{S}, \mathcal{A}, \mathcal{S}, \mathcal{R}, \mathbb{R}\}$ Transition

Problems:

- Fixed Sample size
- Expensive Transitions
- Same actions in different places

Active Learning with Q-Learning

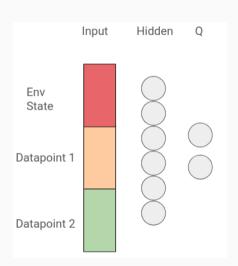


Sample Size k = 2

The same datapoint can appear in multiple places in the input

Both output nodes essentially learn the same function since the datapoints are sampled randomly and independently.

A permutation invariant network can fix the problem and create a ranking of actions



The Bellman Target



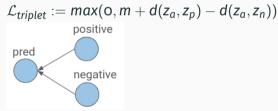
Q-Learning Target: What is the expected discounted improvement in F1-Score under the current policy when we choose a given datapoint?

The Bellman target does not rank actions, but tries to estimate their independent, global value under the current policy

Reinforcement Learning

$$\mathcal{L}_Q := \hat{\pi}(s_t, a_t) - r_t + \gamma \max_{a}(\pi(s_{t+1}))$$
pred target

Ranking (Triplet Loss)



V-Learning vs Q-Learning



TODO

Active Learning with V-Learning



Adaptations

- We use the batch dimension for representing the sample size k
- No action space needed

 $\mathcal{S} \in \mathbb{R}^\sigma$ State Space $\mathcal{R} := \mathcal{S} o \mathbb{R}$ Reward Function $V_{ heta} := \mathcal{S} o \mathbb{R}$ Agent Network

Active Learning with V-Learning



Example:

State is generated : $s \in \mathbb{R}^{b \times \sigma}$

Agents makes prediction : $v = V_{\theta}(s)$

Policy selects a point : $a = \underset{b}{\operatorname{argmax}} v_b$

Action is applied : s' = env(a)

Reward is observed : $r = \mathcal{R}(s')$

Transition is stored : $\tau = \{s_a, r, \bar{s}', done\}$

with :
$$\bar{s}' = \frac{1}{b} \sum_{i=0}^{b} s_i'$$

Storing Transitions in Memory



Q-Learning:

$$O(\tau) = \mathbf{s} \times \mathbf{a} \times \mathbf{s}' \times \mathbf{r} \times \mathbb{R}$$
$$O(\tau) = \mathbf{k}^2 \times \sigma^2 + 3$$

V-Learning:

$$O(\tau) = s_a \times a \times \overline{s}' \times r \times \mathbb{R}$$

$$O(\tau) = 2 \times \sigma^2 + 3$$