Lecture 8: Learning to Search – Imitation Learning in NLP

Kai-Wei Chang CS @ UCLA kw@kwchang.net

Couse webpage: https://uclanlp.github.io/CS269-17/



ML in NLP

Learning to search approaches Shift-Reduce parser

- Maintain a buffer and a stack
- Make predictions from left to right
- Three (four) types of actions: Shift, Reduce, Left, Right



Credit: Google research blog

Structured Prediction as a Search problem

◆ Decomposition of y often implies an ordering
⇒ a sequential decision making process





Notations



Goal: make joint prediction to minimize a joint loss

find $h \in H$ such that $h(x) \in Y(X)$ minimizing $E_{(x,y)\sim D}[loss(y,h(x))]$ based on Nsamples $(x_n, y_n) \sim D$



Credit Assignment Problem

When making a mistake, which local decision should be blamed?





An Analogy from Playing Mario From Mario Al competition 2009

Input:



Output: Jump in {0,1} Right in {0,1} Left in {0,1} Speed in {0,1}

Extrac from la (14 bir

nputer Science

High level goal: Watch an expert play and learn to mimic her behavior





Example of Search Space







A policy maps observations to actions





Labeled data \rightarrow Reference policy

✤ Given partial traj. a₁, a₂, ... a_{t-1} and true label y, the minimum achievable loss is $(a_t^*, a_{t+1}^*, ... a_T^*) = \arg\min_{(a_t, a_{t+1}, ..., a_T)} loss(y, \hat{y}(a))$





Labeled data \rightarrow Reference policy

- ✤ Given partial traj. a₁, a₂, ... a_{t-1} and true label y, the minimum achievable loss is (a_t^{*}, a_{t+1}^{*}, ... a_T^{*}) = arg min_(a_t, a_{t+1}, ... a_T) loss(y, ŷ(a))
- The optimal action is the corresponding a_t^*
- The optimal policy is the policy that always selects the optimal action
- Reference policy can be constructed by the gold label in the training phase



Ingredients for learning to search

Structured Learning

Input: $x \in X$ Truth: $y \in Y(x)$ Outputs: Y(x)Loss: $loss(y, \hat{y})$

Learning to Search

Search Space:

- state: $s \in S$
- action: $a \in A(a)$

- end state $e \in S$ Policies: $\pi(s) \rightarrow a$ Reference policy: π^{ref}



A Simple Approach

\diamond Collect trajectories from expert π^{ref}

Store as dataset $D = \{(o, \pi^{ref}(o)) | o \sim \pi^{ref} \}$

***** Train classifier π on D

• Let π play the game!



Learning a Policy[Chang+ 15, Ross+15]

At "?" state, we construct a cost-sensitive multi-class example (?, [0, .2, .8])



one-step deviations



Example: Sequence Labeling

Receive input:x =the monster ate the sandwichy =DtNnVbDtNn

Make a sequence of predictions:

x = the monster ate the sandwich $\hat{y} =$ DtDtDtDtDt

Pick a timestep and try all perturbations there:

x = the monster ate the sandwich

$\hat{y}_{Dt} = Dt$	Dt	Vb Dt	Nn	l=1
$\hat{y}_{Nn} = Dt$	Nn	Vb Dt	Nn	1=0
$\hat{y}_{Vb} = Dt$	Vb	Vb Dt	Nn	1=1

Compute losses and construct example:

 $(\{ w=monster, p=Dt, ... \}, [1,0,1])$

UCLA ENGINEERING Computer Science

Learning a Policy[Chang+ 15, Ross+15]

At "?" state, we construct a cost-sensitive multi-class example (?, [0, .2, .8])



one-step deviations



Analysis

roll-out → \downarrow roll-in	Reference	Mixture	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

Mixture: w.p. β use Ref, else use Learned.

[ICML 15]: Learning to search better than your teacher



Analysis

roll-out → \downarrow roll-in	Reference	Mixture	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

Mixture: w.p. β use Ref, else use Learned.

Roll-in with Ref:

unbounded structured regret





roll-out → \downarrow roll-in	Reference	Mixture	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

Mixture: w.p. β use Ref, else use Learned.

Roll-out with Ref: no local optimal if reference is sub-optimal



Analysis

roll-out → \downarrow roll-in	Reference	Mixture	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

Mixture: w.p. β use Ref, else use Learned.

Roll-in & Roll-out with current policy ignore Ref \Rightarrow reinforcement learning



Analysis

roll-out → \downarrow roll-in	Reference	Mixture	Learned
Reference	Inconsistent		
Learned	No local opt	Good	RL

Mixture: w.p. β use Ref, else use Learned.

Minimizes a combination of regret to Ref and regret to its own one-step deviations.

Competes with Ref when Ref is good.

Competes with local deviations to improve on suboptimal Ref



How to Program?

Sequential_RUN(examples)

- 1: for i = 1 to len(*examples*) do
- 2: $prediction \leftarrow predict(examples[i], examples[i], label)$
- 3: **loss**(prediction \neq examples[i].label)
- 4: end for

Decoder + loss + reference advice

Sample codes are available in VW

