Learning Cooperative Visual Dialog Agents with Deep Reinforcement Learning

Maria Lubeznov Weining Hu

Popularity on Virtual Assistant







Motivation – Trend in Combining Vision and Language

- Aiding visually impaired users in understanding their surroundings or social media content
- 2. Allowing medical personnel to better interpret medical scans
- Helping AR/VR applications where a user could chat in natural language and work with a virtual companion who is seeing what they are seeing based on a visual common ground



Motivation – Trend in Combining Vision and Language

- Aiding visually impaired users in understanding their surroundings or social media content
- 2. Allowing medical personnel to better interpret medical scans
- Helping AR/VR applications where a user could chat in natural language and work with a virtual companion who is seeing what they are seeing based on a visual common ground



Image captioning vs VQA vs Visual dialog

Inputs

- Dialog history
- Image + caption
- Question

Output

• Answer



VQA Q: How many people on wheelchairs ?

- A: Two
- Q: How many wheelchairs ? A: One

Captioning Two people are in a wheelchair and one is holding a racket.

Visual Dialog

- Q: How many people are on wheelchairs ?
- A: Two
- Q: What are their genders ?
- A: One male and one female
- Q: Which one is holding a racket ?
- A: The woman



Visual Dialog

- Q: What is the gender of the one in the white shirt ?
- A: She is a woman
- Q: What is she doing ?
- A: Playing a Wii game
- Q: Is that a man to her right
- A: No, it's a woman

Motivation2 - They are not perfect

Vision+language recipe -> DATASETS (Collected from Amazon mechanical turk)

Treat dialog as static supervised learning -> **Training (Ground truth answer provided)**

These lead to two bad consequence!

Visual Dialog: Problem #1



- Q: How many people on wheelchairs?
- A: Two
- Q: What gender are the people in the wheelchairs?
- A: One is female, one is male
- Q: Which one is holding the racket?
- A: The female



Model can't steer conversation and doesn't get to see the future consequences of its utterances during training

Visual Dialog: Problem #2



Q: How's the weather?



Ground truth: Sunny

Clear Can't see the sky Looks warm It's not raining I can't see the sky, but I see shadows, so probably sunny

Evaluation infeasible for utterances outside the dataset

Problem formulation

A-BOT

- Picks an image from dataset
- Answers question about the image

Q-BOT

- Asks questions
- Update its mental model of unseen image
- Makes prediction for image feature vector

RL reward for both agents

how close the prediction to true image



Visual dialog dataset - VisDial

VisDial

Training set (235M) 82,783 images Validation set (108M) 40,504 images

```
'questions': [
        'does it have a doorknob',
        'do you see a fence around the bear',
        . . .
     ],
      'answers': [
        'no, there is just green field in foreground',
        'countryside house',
        . . .
      1,
      'dialogs': [
          'image id': <COCO image id>,
          'caption': <image caption from COCO>,
          'dialog': [
              'question': <index of question in `data.questions` list>,
              'answer': <index of answer in `data.answers` list>,
              'answer options': <100 candidate answer indices from data.answers`>,
              'gt index': <index of `answer` in `answer options`>
            },
            ... (10 rounds of dialog)
```

Reinforcement Learning for Dialog Agents



Reinforcement Learning for Dialog Agents

Agent: Cooperative A-Bot and Q-bot

State:

$$s_t^A = [I, c, q_1, a_1, \dots, q_{t-1}, a_{t-1}, q_t] \qquad s_t^Q = [c, q_1, a_1, \dots, q_{t-1}, a_{t-1}]$$
Action: question and answer pair
$$(q_t, a_t)$$

Reward:

$$r_t\Big(\underbrace{s_t^Q}_{\text{state}}, \underbrace{(q_t, a_t, y_t)}_{\text{action}}\Big) = \underbrace{\ell\left(\hat{y}_{t-1}, y^{gt}\right)}_{\text{distance at }t-1} - \underbrace{\ell\left(\hat{y}_t, y^{gt}\right)}_{\text{distance at }t}$$

Policy networks for Q –Bot and A-Bot



Policy network - Q-BOT

Fact encoder - LSTM $(q_t, a_t) \rightarrow F_t^Q \in \mathbb{R}^{512}$

State encoder - LSTM $(F_1^Q, \dots, F_t^Q) \rightarrow S_t^Q$

Question decoder - LSTM which uses S_{t-1}^Q as initial hidden state and generates q_t by sequentially sampling words.

Feature regression network - 1 FC layer $\hat{y}_t = f(S_t^Q)$.

Parameters - θ_Q (LSTMs) and θ_f (FC layer)



Policy network - A-BOT

Question encoder - LSTM $q_t \rightarrow Q_t^A \in \mathbb{R}^{512}$

Fact encoder - same as in Qbot.

State encoder - LSTM $((y, Q_1^A, F_0^A), \dots, (y, Q_t^A, F_{t-1}^A)) \rightarrow S_t^A$

Answer decoder - LSTM which uses $S_t^{\widehat{A}}$ s initial hidden state and generates by sequentially sampling words a_t Parameters - (LSTMs)



Training

- 1. Supervised pre training on VisDial
 - a. Faster RL convergence
 - b. To prevent the bots from inventing their own uninterpretable language
- 2. Fine tuning with RL parameters are updated in response to experienced rewards.

Supervised training

BOTs are trained separately.

Q-BOT

• Is trained to generate the follow-up question by the questionnaire, given the caption and the QA history.

A-BOT

• Is trained to generate the response by the answerer, given the image, caption and the QA history.

CNN for image feature extraction - pre trained on ImageNet

Reinforcement fine tuning

Objective: maximize the expected reward, under agents' policies, for each dialog round (time t).

$$\max(J(\theta_A, \theta_Q, \theta_g)) \text{ where } J(\theta_A, \theta_Q, \theta_g) = \mathop{\mathbb{E}}_{\pi_Q, \pi_A} \left[r_t \left(s_t^Q, (q_t, a_t, y_t) \right) \right]$$

REINFORCE algorithm is applied to calculate the $\nabla_{\theta_Q} J$ and $\nabla_{\theta_A} J$.

Intuition:

- Informative QA round \rightarrow positive reward \rightarrow positive gradient(probability)
- Poor exchange \rightarrow negative reward \rightarrow negative gradient(probability)

$$\theta_f$$
 is updated in a "supervised" way.

Synthetic example



64 images.

 $V_A = \{1, 2, 3, 4\}$ $V_Q = \{X, Y, Z\} = \{\text{shape, color, style}\}$

RL iteration = 10K episodes

Results - supervised vs RL dialogues

Image + Caption	Human-Human dialog [4]	${\tt SL-pretrained}Q{\textbf{-}BOT}{\textbf{-}A}{\textbf{-}BOT}dialog$		RL-full-QAf Q-BOT-A-BOT dialog	
A couple of people in the snow	What are their genders? 1 man 1 woman Are they both adults? Yes Do they wear goggles? Looks like sunglasses Do they have hats on? Man does Are there any other people? No What color is man's hat? Black Is it snowing now? No What is woman wearing? Blue jacket and black pants Are they smiling? Yes Do you see trees? Yes on skis. Yes	How many people? Male or female? What are they wearing? What color are skis? Are they wearing goggles? Are they wearing goggles?	2 I can't tell I can't tell, they are far away I can't tell Yes Yes Yes Yes Yes Yes Yes	Are people male or female or both? What are they wearing? Is it snowing? Can you see any trees? Can you see any buildings in background? Can you see any buildings in background? Can you see any trees? Can you see any buildings in background? Can you see any buildings in background? Does it look like they are in park?	I think I is male, can't see other 1 They are all wearing snow pants and jackets It does not appear to be snowing at moment Yes, there are lot of trees in background No, I can not Yes, there are lot of trees in background No, I do not see any buildings Yes, there are lot of trees in background No, I do not see any buildings It does not appear to be
Ivripoonful com	Are there any people in picture? No, there are no people Is cheese white? Yes, cheese is white Are crackers white? No, cracker is brownish Are there lots of greens? No, only 1 sprig Was picture taken in kitchen? I cannot tell, it is too zoomed in How many crackers are there? There is 1 Is there design on plate? No, plate is plane white Is there design on plate? No, plate is plane white Is there good lighting? Yes, lighting is decent Is plate on table? I cannot tell I can only see plate as background Does cracker look delicious? Absolutely	Are there any people? What color is plate? Is plate on table? What color is plate? Is plate on table? Is there anything else on table?	No White I can't tell White I can't tell No No No No No	Is plate large or small? Is it on table? What kind of donut? What kind of chips? Any people? Is this inside or outside? What color is table? Any chairs? What color is table? Is it large or small?	Large I can't tell Plain with sugar on it Potato chips No I can't tell, it's close up White No I can't tell I can't see plate

A cracker with shredded cheese and leafy greens on a white plate.

Evaluation

Dataset: test split for VisDial v0.5 (~9.5k images) Input : Image + caption Rounds: 10

Output: Sort the entire test set in ascending distance to this prediction



Model	MRR	R@5	R@10	Mean Rank
SL-pretrain	0.436	53.41	60.09	21.83
Frozen-Q	0.428	53.12	60.19	21.52
Frozen-f	0.432	53.28	60.11	21.54
RL-full-QAf	0.428	53.08	60.22	21.54
Frozen-Q-multi	0.437	53.67	60.48	21.13

(b) Visual Dialog Answerer Evaluation.

Evaluation

 $\ell 2$ distance to ground truth image in fc7 space

Round 1: What kind of pizza is it?





Pizza slice sitting on top of white plate.

Round 5: Is there anything else on plate? Yes, there are 2 other plates in background.







Group of people standing on top of lush green field.





Man in light-colored suit and tie standing next to woman in short purple dress.

Round 2: What color are pants? Navy blue.

Round 4: Are they indoors or outdoors? Outdoors.



Round 1: Are flowers in vase? Yes.



Pink flowers in one of

many ceramic vases.

People staring at man

on fancy motorcycle.



Round 3: What color are flowers? Purple, yellow and white.





Round 1: How many people are there?

Lot, too many to count.



Round 3: Does it look old or new? It looks new.



Round 9: Is it sunny out? Yes.



(c) Qualitative Retrieval Results.

$\ell 2$ distance to ground truth image in fc7 space







Weakness of the evaluation

- 1. No evaluation on the response time
- 2. No evaluation on how the initial caption quality affect the later dialog



Strengths and weaknesses of the approach

Strengths

- Self talk \rightarrow unlimited data. No need to collect very large datasets
- Evaluation is image guessing → no need to evaluate natural language generation.
- Method is agent driven learning to deal with actions consequences

Weaknesses

- No explanation of why the agent becomes less sure after 2 dialog rounds.
- Authors mention that A-BOT responses are not enough "human like".

Potential extensions

Our thoughts

- Address repetitive questions
- It may be a good idea to combine this method with some sort of CGAN from paper presented last week, to better emulate human dialogs.
- Evaluate the impact of different image captioning quality on final reward achievement

Follow up papers

- Evaluate if and how this method can be applied on human-machine interaction ([1]).
- Investigate how similar method can be applied on other tasks that require negotiation ([2])



[1] Evaluating Visual Conversational Agents via Cooperative Human-AI Games

[2] <u>Deal or No Deal? End-to-End Learning for Negotiation Dialogues</u>

[3] <u>Demo</u>

[3] <u>https://research.fb.com/visual-reasoning-and-dialog-towards-natural-language-</u> <u>conversations-about-visual-data/</u>

Appendix - technical details

- Supervised training for the first K rounds of dialog and transition to RL for the remaining 10 – K rounds
- In each epoch, K gradually annealed to 0.
- 15 epochs
- Adam optimizer with learning rate 0.001
- Gradients are clamped to [-5,5] to avoid explosion.