

Instantly share code, notes, and snippets.

shagunsodhani / **GuessWhat.md**

Created 5 years ago

☆ Star

<> **Code**

🔗 Revisions 1

🍴 Forks 1

Summary of "GuessWhat?! Visual object discovery through multi-modal dialogue" paper

📄 **GuessWhat .md**

GuessWhat?! Visual object discovery through multi-modal dialogue

🔗 Introduction

- The paper introduces *GuessWhat* - a two-player guessing game where the goal is to locate an object in a rich image scene.
- The game is used to produce a large scale dataset of visual question-answer pairs on the image.
- The paper also describes three tasks based on the game and provides a neural architecture based baselines for each task.
- [Link to the paper](#)

GuessWhat?! Game

- One player, called as the **oracle**, is randomly assigned an object in the given image.
- The second player, called as the **questioner**, tries to locate the object, given just the image.
- The **questioner** can ask a series of questions about the object and the **oracle** can reply in "yes" or "no" or "not applicable".
- Once the **questioner** is confident of having identified the image, the **oracle** presents a list of objects to the **questioner** to choose from.

- A small penalty is added, every time a question is asked, so as to encourage informative questions only.

Dataset

- A filtered subset of images from MSCOCO is used as the image set.
- Two separate tasks create on Amazon Mechanical Turk (AMT) - for the role of **oracle** and **questioner**.
- Data was post processed -- both manually and using AMT -- to account for things like spelling mistakes and validation.
- Final dataset comprises of 150K thousand human game iterations with 800K question-answer pairs on 60K images.
- Dataset is available at <https://guesswhat.ai/download>

Interesting Observations

- an average number of questions, given the number of objects in the image, grows at a rate between logarithmic and linear possibly because:
 - **questioner** does not have access to the list of images while asking the question.
 - **questioner** might ask a few extra questions just to be sure.
- **questioner** uses abstract object properties like "human/object/furniture" early in the conversation and quickly switch over to spatial/visual aspects like "left/right or table/chair"/
- The paper also includes the analysis of how factors like size of the unknown object, its position, total number of objects etc affect the accuracy of humans (playing on AMT).

Model

- Given an image, a set of segments objects (along with their category and pixel-wise segmentation mask) and the object to be identified.
- The **questioner** generates a series of questions to ask from the **oracle**

Oracle

- Modelled as a single hidden layer MLP.
- **Input**: Concatenate embeddings for the
 - image - obtained using FC8 features from VGG Net - fixed during training
 - cropped target object - obtained using FC8 features from VGG Net - fixed during training

- spatial information about the target object (bounding box coordinates, normalised wrt coordinates of the centre) in form of a vector - fixed during training
- category of target object - dense categorical embedding - trained
- current question asked by the **questioner** - encoded by LSTM - trained
- **Output:** One of the three answers - "yes", "no", "not applicable"
- **Loss:** Negative log-likelihood of correct answer
- **Performance:** The best model achieves the test error of 21.5% and uses question embedding + category embedding + spatial information.

Questioner

- Question performs two sub tasks:
 - **Guesser**
 - **Input:** Concatenate embeddings for the
 - image - obtained using FC8 features from VGG Net - fixed during training
 - dialogue - the series of question-answer pair are embedded into fixed size vectors using an LSTM or HRED encoder
 - Objects are represented by:
 - concatenation of their category and spatial features
 - passing through an MLP (which shares parameters across objects)
 - Perform a dot product between input embedding and embedding for the objects, followed by softmax, to obtain the probability distribution of the objects.
 - **Performance:**
 - The best model achieves the test error of 38.7% and uses LSTM alone (no VGG features were used)
 - VGG features did not improve the performance, probably because the questions and the objects captured all the information already
 - Maybe using some different network for image features may help
 - **Generator**
 - HRED, conditioned over the VGG features (from the image) and questions asked so far (if any) is used to generate the natural language questions by maximising conditional log-likelihood.
 - The questions generated by the **generator** are answered by the **oracle** which is ground truth at train time and trained oracle at test time. The paper acknowledges this shortcoming.

- 5 questions are generated before triggering the **guesser**
- **Performance:** The best model achieves an error of 38.7% using human generated dialogues.
- the performance is also deteriorated by the fact that **oracle** and **guesser** are not perfect.

Comments

- The paper has provided a very detailed analysis of the dataset and have experimented with various combinations to find the best embedding model.
- The sample questions show in the paper indicates that the **generator** model produces the same question many times - indicating poor generalisation.