

research papers

Sam Jones



Show and Tell

Bleu-1 score

I input

$p(S|I)$ likelihood (S given I)
of producing a target sequence words

$$S = \{s_1, s_2, \dots\}$$

Encoder RNN \Rightarrow feature vector

Replace encoder RNN with CNN

CNN \rightarrow image encoder of last hidden layer

Learning descriptions:

Description close to an image in a vector space

RNN starts with image processed by CNN

Joint embedding

Encode variable length input into fixed dimensional vector

$$\theta^* = \arg \max_{\theta} \sum_{(I, S)} \log p(S|I; \theta)$$

θ = parameters of model

I = image

S = correct transcription

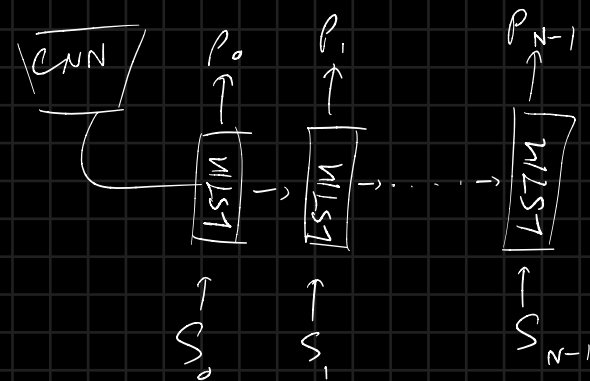
$$\textcircled{1} \log p(S|I) = \sum_{t=0}^T \log p(S_t | I, S_0, \dots, S_{t-1})$$

during training,

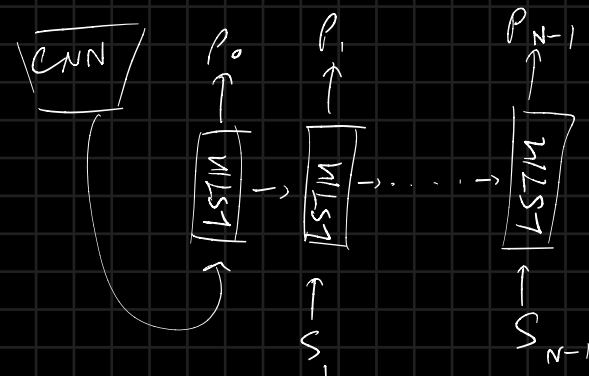
one example (S, I)

optimizer $\textcircled{1}$

LSTM net



or...



one hot vector S_t of vocab-size
start & stop word.

$t=1$, image is input

only feed once at start \rightarrow according
to results.

Loss:

$$L(I, s) = - \sum_{t=1}^N \log p_z(S_t)$$

Correct word at each time step

write all parameters

Sampling from model, until end-token

BLEU Score

Use a pre-trained model on ImageNet

Dropout ✓ for over fitting prevention

512 dimension embedding (probably hidden layer)

Real Time American Sign Language Recognition with Convolutional Neural Networks

Google Net

ILSVRC 2012 Dataset

Obtain video of user signing

Classify each frame to a letter

Reconstructing and displaying the most likely word

from classification scores

Considerations:

Environment (lighting & camera)

occlusion (out of view)

Sign Boundary Detection

Co-articulation

what heuristics to use?

Hand cropping

Joint detection? coordinates of joints

physical component

Image processing per frame

$$\text{Loss} = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_i y_i}}{\sum_{j=1}^c e^{f_i y_j}} \right)$$

$$f_i(z) = \frac{e^{z_i}}{\sum_{k=1}^c e^{z_k}}$$

Bottle neck: one frame per second

Let the user move onto the next letter

Surry university Finger Spelling
Massey University

Datasets

Make horizontal flips of image

confusion matrix:

True label

predicted label

Reinitialize classification layer (Pre-trained model)

Can easily get confused if letters look the same.

Contrast adjustment
Background subtraction
Cropping

use CNN to locate the hand