

Deep learning and chatbot

Le LI

iAdvize & Université d'Angers

Meetup à Pau, 03/05/2017

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 RNN introduction
 - RNN and derivatives
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

1 Context

- Introduction about iAdvize
- Demanding of chatbot

2 Overview of chatbot

- Definition
- Types of chatbot

3 General task and procedures

- Task
- Training and inference procedure

4 RNN introduction

- RNN and derivatives
- Seq2Seq

5 Simple work flow for constructing seq2seq based chatbot

- Pre-processing steps
- Training the seq2seq model in TF

A conversational commerce platform,



Figure: logo

- object: help clients to increase their conversion rate and customer's satisfaction.
- ways: detect valuable visitors on the internet, social media etc.
- platform: integrate client's website, social media, apps etc.
- channels: chat, call, video connecting visitors and clients(consultants or experts).
- clients: Airfrance, Voyage sncf etc.

1 Context

- Introduction about iAdvize
- Demanding of chatbot

2 Overview of chatbot

- Definition
- Types of chatbot

3 General task and procedures

- Task
- Training and inference procedure

4 RNN introduction

- RNN and derivatives
- Seq2Seq

5 Simple work flow for constructing seq2seq based chatbot

- Pre-processing steps
- Training the seq2seq model in TF

Demanding

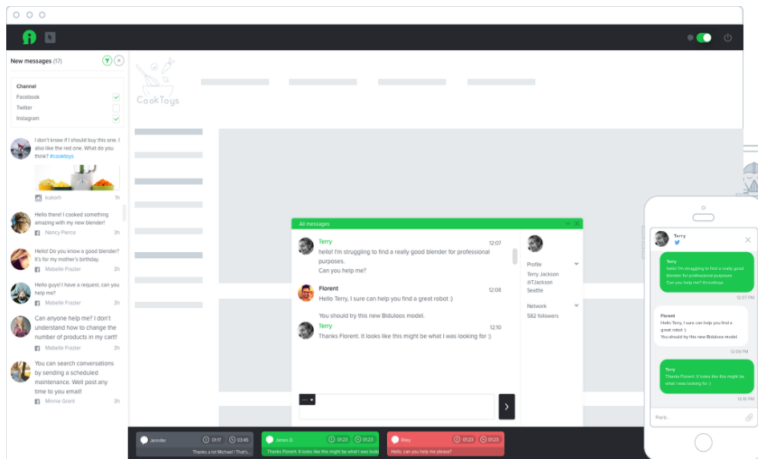


Figure: interface

why chatbot? a big volume of visitors (6000 visitors/hour) vs a small amount of consultants.

- answer simple questions in some certain scenario.
- reduce repeated work of consultants.

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 RNN introduction
 - RNN and derivatives
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

Definition

chatbot: a Conversational Agent or Dialog Systems that can interact with customers by having natural conversations indistinguishable from human.
eg: (Facebook (M), Apple (Siri), Google etc).

Outline

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 RNN introduction
 - RNN and derivatives
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

- **Retrieval-based model**: choose a response from a prefixed set of responses.
 - expression-match, ML techniques(word2vec) etc.
 - pros: No grammar mistakes.
 - cons: work bad on unseen cases for which no pre-defined response exist. disability to refer back.
- **Generative-model**: generate response from scratch.
 - Recurrent Neural Network(eg:seq2seq)
 - pros: ability to refer back and cope with new cases.
 - cons: grammar mistakes, huge amount of training samples.

- **Short text conversation:** one single question consecutively with one single answer.
- **Long text conversation:** multiple questions and responses changes(not necessarily consecutive).

- **More formal platform:** log of chat between a visitor and a consultant, forum.
- **Less formal platform:** twitter, facebook, messenger: abbreviation, emoji, non-existing words.

- **Closed domain**: the chatbot is trying to achieve a very specific goal in a certain scenario(eg:delivery).
- **Open-domain**: the chatbot is able to handle conversations with open subject.

Outline

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 RNN introduction
 - RNN and derivatives
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

A little math

Given an original message \mathbf{x} , find most likely response \mathbf{y}^*

$$\mathbf{y}^* = \underset{\mathbf{y} \in \mathcal{Y}}{\operatorname{argmin}} \mathbb{P}(\mathbf{y}|\mathbf{x}).$$

Construct a model that can score responses and then find highest scoring response.

- original message \mathbf{x} : "il n'y a donc pas d'autres solutions que par carte bancaire ?"
- response \mathbf{y} : "C'est la seule solution je suis désolé :/"

Since we score a sequence of tokens y_1, y_2, \dots, y_m in \mathbf{y} , conditional on x_1, x_2, \dots, x_n in \mathbf{x}

$$\mathbb{P}(y_1, y_2, \dots, y_m | x_1, x_2, \dots, x_n) = \prod_{i=1}^m \mathbb{P}(y_i | y_{i-1}, \dots, y_1, x_1, x_2, \dots, x_n) \quad (1)$$

Outline

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 RNN introduction
 - RNN and derivatives
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

Given a corpus of pair (\mathbf{x}, \mathbf{y}) ,

- training objective to maximize

$$\sum_{\mathbf{x}, \mathbf{y}} \ln \mathbb{P}(y_1, y_2, \dots, y_m | x_1, x_2, \dots, x_n).$$

- inference objectives
 - 1 random sample from (1).
 - 2 greedy search: taking most likely tokens at each time.
 - 3 determine the likelihood of a specific response candidate.

Outline

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 **RNN introduction**
 - **RNN and derivatives**
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

Review Neural Network

- A simple NN with 1 layer:

input: $\mathbf{x}_n \in \mathbb{R}^p$,

hidden layer: $\mathbf{h}_n = \sigma_1(W_i \mathbf{x}_n + B_i)$,

output: $\hat{\mathbf{y}}_n = \sigma_2(W_o \mathbf{h}_n + B_o)$.

where σ_1 : activation function and σ_2 : sigmoid function.

- Eg: \mathbf{x}_n : pixels of images; \mathbf{y}_n real character in the image, eg [cat,dog,car,house]; $\hat{\mathbf{y}}_n$: estimated character

[0.8,0.2,0.1,0] \rightarrow cat

[0.1,0.7,0.1,0.1] \rightarrow dog

[0.15,0.05,0.65,0.15] \rightarrow car

[0,0,0.5,0.95] \rightarrow house

A simple RNN

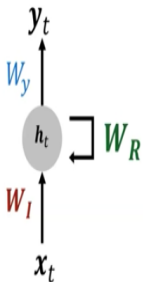
- objective:

$$\text{minimize: } \sum_{n=1}^N (\mathbf{y}_n - \hat{\mathbf{y}}_n)^2.$$

with respect to parameters: W_i, W_o, B_i, B_o .

- optimization: gradient descent, stochastic gradient descent etc.
- disadvantage: order of \mathbf{x} is important, translation, speech recognition etc. "I am interested in Neural network".

Recurrent neuron



$$h^{(t)} = g_h(W_I x^{(t)} + W_R h^{(t-1)} + b_h)$$

$$y^{(t)} = g_y(W_y h^{(t)} + b_y)$$

Figure: structure of RNN

Unrolling of RNN

Unrolling a recurrent network into a feed-forward network

$$\mathbf{h}^{(t)} = g_h(W_I \mathbf{x}^{(t)} + W_R \mathbf{h}^{(t-1)} + \mathbf{b}_h)$$

$$\mathbf{y}^{(t)} = g_y(W_y \mathbf{h}^{(t)} + \mathbf{b}_y)$$

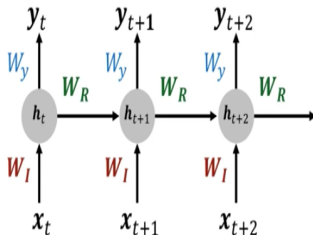


Figure: Unrolling of RNN

Long Short Term Memory

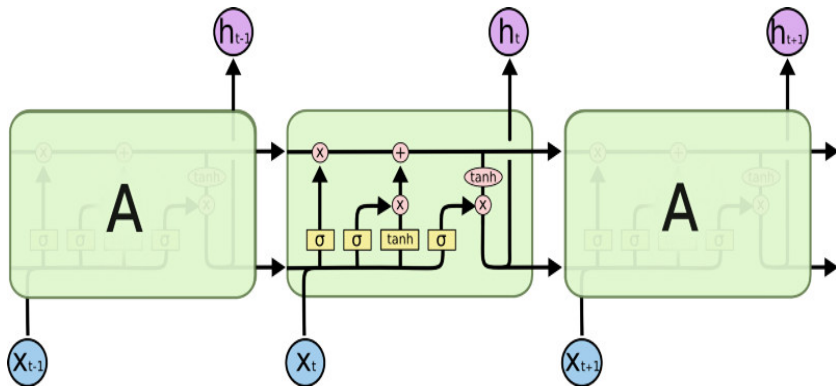


Figure: LSTM

- Long short term memory(LSTM):

Forget gate layer : $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$,

Input gate layer : $i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$,

Input : $\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$,

Cell state : $c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$,

Output gate layer : $o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$,

Output : $h_t = o_t * \tanh(c_t)$.

Variants of LSTM

- Gated recurrent unit(GRU): combines the forget and input gates into a single “update gate.” It also merges the cell state and hidden state.
- Coupled LSTM.
- Depth Gated RNNs().

Outline

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 RNN introduction
 - RNN and derivatives
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

Sequence to Sequence

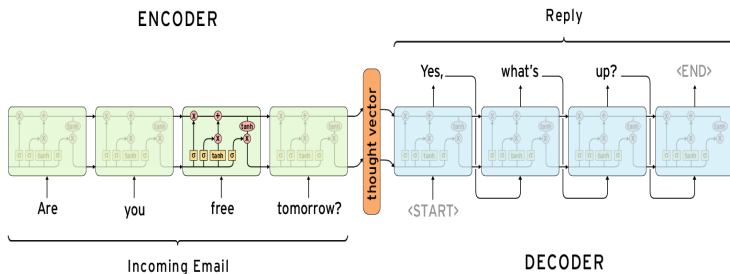


Figure: seq2seq

- Encoder: original message $\mathbf{x} \rightarrow$ representation of message c_n (last cell state). No output h_n in encoder step.
- Decoder: $h_{t-1}, c_{t-1} +$ response tokens $y_t \rightarrow$ output h_t, c_t .

Outline

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 RNN introduction
 - RNN and derivatives
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

- language detection.
- remove accents: j'ai aimé \rightarrow j'ai aime
- normalization: '*www.abc.com*',113 \rightarrow url, number_digit
- separate contraction: j'ai \rightarrow je ai
- tokenization: Pau est beau \rightarrow [*Pau,est,beau*]
- padding: [*pau,est,beau,pad,pad*]
- bucketing: grouper des questions et reponses.

Outline

- 1 Context
 - Introduction about iAdvize
 - Demanding of chatbot
- 2 Overview of chatbot
 - Definition
 - Types of chatbot
- 3 General task and procedures
 - Task
 - Training and inference procedure
- 4 RNN introduction
 - RNN and derivatives
 - Seq2Seq
- 5 Simple work flow for constructing seq2seq based chatbot
 - Pre-processing steps
 - Training the seq2seq model in TF

Tensorflow and other framework

- Tensorflow
- Keras, TFlearn etc.

For Further Reading I



A. Kannan.

smart reply: automatied response suggestion for email.

<https://arxiv.org/abs/1606.04870>.



L. Sutskever.

Sequence to Sequence Learning with Neural Networks.

<https://arxiv.org/abs/1409.3215>



Q. Le.

A Neural Conversational Model. <https://arxiv.org/abs/1506.05869>