Sequence to Sequence 27.5.2019

Sequence-to-Sequence



The quick brown fox

Handwriting recognition: The input can be (x, y) coordinates of a pen stroke or pixels in an image.



Speech recognition: The input can be a spectrogram or some other frequency based feature extractor.

Sequence-to-Sequence

- Speech recognition
- Handwriting recognition
- Machine translation
- Speech synthesis (text to speech)
- Summarization
- Text generation

Speech Recognition: Hybrid HMM-DNN



- Use feed-forward net to predict phones
- Use existing HMM/graph structure to form words and sentences
- Problem: requires alignment for frame-wise training!
- (Did not work for translation and generation)

Connectionist Temporal Classification (CTC)

- Map from $X = [x_1, x_2, \dots, x_T]$ to $Y = [y_1, y_2, \dots, y_U]$, where
 - Both X and Y can vary in length.
 - The ratio of the lengths of X and Y can vary.
 - We don't have an accurate alignment of X and Y
- Key ideas:
 - Use epsilon padding to ensure equal lengths
 - Marginalize over all possible alignments

CTC – Blank (eps) Symbol



First, merge repeat characters.

Then, remove any ϵ tokens.

The remaining characters are the output.

Example

input (X)	x_6	<i>x</i> ₅	X_4	<i>x</i> ₃	<i>x</i> ₂	\boldsymbol{x}_{l}
alignment	t	а	а	а	С	С
output (Y)	t		а		2	C

Valid Alignments	Invalid Alignments	
ϵ C C ϵ a t	c ϵ c ϵ a t	corresponds to Y = [c, c, a, t]
c c a a t t	c c a a t	has length 5
C a ϵ ϵ ϵ t	C ϵ ϵ ϵ t t	missing the 'a'

CTC Loss Function – illustrated





We start with an input sequence, like a spectrogram of audio.

The input is fed into an RNN, for example.

The network gives $p_t (a \mid X)$, a distribution over the outputs {h, e, l, o, ϵ } for each input step.

h	е	ϵ		ε		0	0
h	h	е		ε	ε	ε	0
ϵ	е	ϵ		ϵ	ε	0	0

With the per time-step output distribution, we compute the probability of different sequences

h e I I o e I I o h e I o

By marginalizing over alignments, we get a distribution over outputs.

CTC Loss Function

To be precise, the CTC objective for a single (X, Y) pair is:

$$p(Y \mid X) \hspace{.1in} = \hspace{.1in} \sum$$

The CTC conditional **probability** $A \in \mathcal{A}_{X,Y}$ marginalizes over the set of valid alignments $\prod_{t=1}^{T} p_t(a_t \mid X)$ computing the **probability** for a single alignment step-by-step.

Marginalize over all Alignments?

- Valid alignments
 - Retain the original sequence of tokens
 - Optionally insert epsilons between tokens
- Enumerating all alignments is way too expensive \rightarrow DP!





Summing over all alignments can be very expensive.

Dynamic programming merges alignments, so it's much faster.

CTC Alignment Graph



score of the subsequence $Z_{1:s}$ after t input steps.

CTC Training

- Loss can be analytically computed \rightarrow Gradient!
- Maximum Likelihood estimate, including sequence info

$$\sum_{X,Y)\in\mathcal{D}} -\log p(Y \mid X) \quad
ightarrow min.$$

CTC Inference

- How to handle epsilons and repeats?
 - "naive" way (regexp)
 - Beam search



alphabet $\{\epsilon, a, b\}$ and a beam size of three.

CTC for Speech Recognition



"spikey" activations for non-blanks

"traditional" HMM-GMM

"hybrid" HMM-DNN

Table 1. Label Error Rate (LER) on TIMIT. CTC and hybrid results are means over 5 runs, \pm standard error. All differences were significant (p < 0.01), except between weighted error BLSTM/HMM and CTC (best path).

System	LER
Context-independent HMM	38.85%
✓ Context-dependent HMM	35.21%
BLSTM/HMM	$33.84 \pm 0.06\%$
Weighted error BLSTM/HMM	$31.57\pm0.06\%$
CTC (best path)	$31.47 \pm 0.21\%$
CTC (prefix search)	$30.51\pm0.19\%$

First "end-to-end" in 2006!

CTC References

- "Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks." A. Graves et al., ICML2006
- <u>https://distill.pub/2017/ctc/</u>
- TensorFlow
 - <u>https://www.tensorflow.org/api_docs/python/tf/nn/ctc_loss</u>
 - <u>https://www.tensorflow.org/api_docs/python/tf/nn/ctc_beam_search_decod_er</u>
- cuDNN
 - <u>https://docs.nvidia.com/deeplearning/sdk/cudnn-developer-guide/index.html#cudnnCTCLossAlgo_t</u>