Speech Processing with Deep Learning

Leo Ebert Tuananh Vu Moritz Kirchner



Universidad del País Vasco

Euskal Herriko Unibertsitatea





Agenda

Introduction
 Preprocessing
 Automatic Speech Recognition
 Postprocessing
 Demo
 Conclusion

Speech Processing - First Question WHY

- WHY Speech Processing



- WHY using NN

- 1. Language is complicated -> Dialects
- 2. nothing sounds similar -> pronunciation
- 3. depends on who is speaking -> Understand also people never heard speaking before
- 4. need of using context to "understand"

Speech Processing - an Introduction



heartbeat.fritz.ai

Preprocessing

First step: Feed sound waves into the computer

- 1. Turn Soundwave into numbers
- 7.5 22.5 Sampling 6000 4000 2000
- 2. Record the height of the wave at equally spaced points
 - \rightarrow Also knowns as sampling

23.2346566756,45.231453 4534, 12.123254 ...

Preprocessing



Transformation to a spectrogram



Automatic Speech Recognition

Acoustic Model

- Deals with the audio waveforms predicting what phoneme each waveform corresponds to
- Many approaches exists (Usually neural network, but N-gram possible)



Automatic Speech Recognition

Language Model

- Probability distribution over a word or word sequence
- Can base on probabilistic methods or neural networks

A fixed-window neural Language Model

$\begin{array}{l} \begin{array}{l} \text{output distribution} \\ \hat{y} = \operatorname{softmax}(Uh + b_2) \in \mathbb{R}^{|V|} \\ \end{array}$ $\begin{array}{l} \text{se I am learning} \\ \hline{y} = \operatorname{softmax}(Uh + b_2) \in \mathbb{R}^{|V|} \\ \hline{hidden layer} \\ h = f(We + b_1) \\ \hline{werthing} \\ e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}] \\ \hline{f} \\ e = [e^{(1)}; x^{(2)}; x^{(3)}; x^{(4)} \\ \hline{f} \\ x^{(1)} \\ x^{(2)}; x^{(3)}; x^{(4)} \\ \hline{f} \\ x^{(1)} \\ \end{array}$



Trigram Probability

Corpus: I am happy because I am learning

$$P(happy|Iam) = \frac{C(Iam happy)}{C(Iam)} = \frac{1}{2}$$

Probability of a trigram:

Post-processing - considerations

• punctuation restoration and capitalization

`Let's eat, kids. ´ vs. `Let's eat kids. ´

- Ability to process unknown words
- \rightarrow as always use preferably as little computational resource



Demo

import seaborn as sns import tensorflow as tf rows = 3 cols = 3 n = rows*cols fig, axes = plt.subplots(rows, cols, figsize=(10, 12)) for i, (audio, label) in enumerate(waveform_ds.take(n)): r = i // cols c = i % cols ax = axes[r][c] ax.plot(audio.numpy()) ax.set_yticks(np.arange(-1.2, 1.2, 0.2)) label = label.numpy().decode('utf-8') ax.set_title(label)

plt.show()



simple audio.ipynb - Colaboratory (google.com)

def plot_spectrogram(spectrogram, ax):

Convert to frequencies to log scale and transpose so that the time is # represented in the x-axis (columns). An epsilon is added to avoid log of zero. log_spec = np.log(spectrogram.T+np.finfo(float).eps) height = log_spec.shape[0] width = log_spec.shape[1] X = np.linspace(0, np.size(spectrogram), num=width, dtype=int) Y = range(height) ax.pcolormesh(X, Y, log_spec)

fig, axes = plt.subplots(2, figsize=(12, 8))
timescale = np.arange(waveform.shape[0])
axes[0].plot(timescale, waveform.numpy())
axes[0].set_title('Waveform')
axes[0].set_xlim([0, 16000])
plot_spectrogram(spectrogram.numpy(), axes[1])
axes[1].set_title('Spectrogram')
plt.show()







Conclusion



amazon alexa

Appendix

Post-processing - models

÷

n-gram based approach	fast results	 needs a lot of disk space handling of out-of-vocabulary words
bidirectional RNN	high accuracy	size of input ↑ => time per token ↑

► TruncBiRNN

- forwards normal RNN
- backwards only a specified window is considered

Post-processing - models

