



通通学

Audio processing in TensorFlow

📖 转载IT好文 👤 Chris ⌚ 4个月前 (07-24) 👁 72次浏览 🚩 已收录 💬 0条评论



An implementation of the Short Time Fourier Transform

We found audio processing in TensorFlow hard, here is our fix

There are countless ways to perform audio processing. The usual flow for running experiments with Artificial Neural Networks in TensorFlow with audio inputs is to first preprocess the audio, then feed it to the Neural Net.

What happens though when one wants to perform audio processing somewhere in the middle of the computation graph?

TensorFlow comes with an implementation of the [Fast Fourier Transform](#) , but it is not enough.

In this post we will explain how we implemented it and provide the code so that the Short Time Fourier Transform can be used anywhere in the computation graph.

Audio preprocessing: the usual approach

When developing a Speech Recognition engine using Deep Neural Networks we need to feed the audio to our Neural Network, but... what is the right way to preprocess this input?

There are 2 common ways to represent sound:

- ⦿ *Time domain*: each sample represents the variation in air pressure.
- ⦿ *Frequency domain* : at each time stamp we indicate the amplitude for each frequency.

Despite the fact that Deep Neural Networks are extremely good at learning features automagically, it is always a good idea to rely on known features that carry the information needed for the task that we are trying to solve.

For most application, a Speech Recognition Engine included, the features we are interested in are encoded in the frequency domain representation of the sound.

The spectrogram and the Short Time Fourier Transform



A [spectrogram](#) shows how the frequency content of a signal changes over time and can be calculated from the time domain signal.

The operation, or transformation, used to do that is known as the [Short Time Fourier Transform](#) .

We could let the Neural Network figure out how to learn this operation, but it turns out to be quite complex to learn with 1 hidden layer. (refer to the [Universal approximation theorem](#))

We could add more layers, but we want to keep the complexity of the Neural Networks as small as possible and learn features only where it is most needed.

We have used the example of developing an Automatic Speech Recognition engine, but the use of the spectrogram as input to Deep Neural Nets is common also for similar tasks involving non-speech audio like noise reduction, music genre classification, whale call detection, etc.

A particular project that we want to mention is [Magenta](#) , from the [Google Brain team](#) , who's aim is to advance the state of the art in machine intelligence for music and art generation.

Why TensorFlow?

We mainly use TensorFlow when implementing Artificial Neural Networks and, because we haven't found an implementation of the Short Time Fourier Transform in TF, we decided to implement our own.

There can also be multiple reasons why a deep learning practitioner might want to include the Short Time Fourier Transform (STFT for my friends) in the computation graph, and not just as a separate preprocessing step.

Keep in mind that we haven't focused on making this efficient. It should (and will) be improved before being used in production.

What we need to know

In order to understand how the STFT is calculated, we need to understand how to compute the Discrete Fourier Transform.

Discrete Fourier Transform — DFT

This part can appear quite technical for those who are not familiar with these concepts, but we think it is important to go through some maths in order to give a complete understanding of the code.



Theory

Fourier analysis is fundamentally a method for expressing a function as a sum of periodic components, and for recovering the function from those components. When both the function and its Fourier transform are replaced with discretized counterparts, it is called the discrete Fourier transform (DFT).

Given a vector \mathbf{x} of n input amplitudes such as:

```
{x[0], x[1], x[2], x[3], ..., x[N-1]}
```

the Discrete Fourier Transform yields a set of n frequency magnitudes.

The DFT is defined by this equation:

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{\frac{-j2\pi kn}{N}}$$

DFT equation

- ⦿ k is used to denote the frequency domain ordinal
- ⦿ n is used to represent the time-domain ordinal
- ⦿ N is the length of the sequence to be transformed.

Fast Fourier Transform

The

Fast Fourier Transform is an efficient implementation of the DFT equation. The signal must be restricted to be of size of a power of 2.

This explains why we want N (the size of the signal in input to the DFT function) to be power of 2 and why it must be zero-padded otherwise.

We can detect whether x is a power of 2 very simply in python:

We only need half of it

Real sine waves can be expressed as the sum of complex sine waves using

Euler's identity

$$x[n] = A \cos(2\pi \frac{kn}{N}) = \frac{A}{2} e^{j \cdot 2\pi \cdot kn/N} + \frac{A}{2} e^{-j \cdot 2\pi \cdot kn/N}$$

Because the DFT is a linear function, the DFT of a sum of sine waves is the sum of the DFT of each sine wave. So for our spectral case, we get 2 DFTs, one for the positive frequencies and one for the negative frequencies, which are symmetric.

This symmetry occurs for real signals that can be viewed as an infinite (or finite in our case) sum of sine waves.

Windowing

Truncating a signal in the time domain will lead to ripples appearing in the frequency domain.

This can be understood if we think of truncating the signal as if we applied a rectangular window. Applying a window in the time domain results in a convolution in the frequency domain.

The ripples are caused when we convolve the 2 frequency domain representations together.

Find out more about [spectral_leakage](#) if you're interested.

Here is an example of an implementation of windowing in Python:

Zero-phase padding

In order to use the FFT, we need to have the input signal to have a power of 2 length. If the input signal does not have the right length, we have to append zeros to the signal itself both at the beginning and at the end.

Because the zero sample is originally at the center of the input signal, we have to split the padded signal through the middle and swap the order of these 2 parts.

The next code snippet shows how to do this in TensorFlow for a batch of inputs:

FFT, Magnitude and Phase

We now have everything we need to calculate the magnitude of the spectrogram in decibels and the phase of the signal:

Short Time Fourier Transform

We now know how to compute the DFT to evaluate the frequency content of a signal.

The STFT is used to analyze the frequency content of signals when that frequency content varies with time.

We can do this by:

- ⦿ Taking segments of the signal.
- ⦿ Windowing those out from the rest of the signal, and applying a DFT to each segment.
- ⦿ Sliding this window along each segment.

We get DFT coefficients as a function of both time and frequency.

The complete code is divided in 2 parts: *helpers.py* and *stft.py*.

Conclusion

The possibility of doing the STFT in TensorFlow allows Machine Learning practitioners to perform the transformation of a signal, from time-domain to frequency domain, anywhere in the computation graph.



New tools always bring new ideas and we hope this post will be the source of new ideas for developing new Deep Learning solutions.

About Cisco Emerge

At Cisco Emerge, we are using the latest machine learning technologies to advance the future of work.

Find out more on [our website](#).

来源：<https://medium.com/towards-data-science/audio-processing-in-tensorflow-208f1a4103aa>

📢 通通学，版权所有 | 如未注明，均为原创 | 本网站采用BY-NC-SA协议进行授权，转载请注明本文链接

<http://www.tongtongxue.com/archives/10202.html>

♡ 喜欢 (1)

🏷️ TensorFlow

« [RGB566与RGB888的区别](#)

[Learning Xamarin.Forms – Part 3: Navigatic](#)





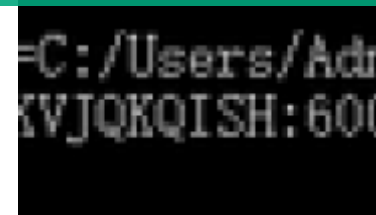
手把手 | 亲测好用！Google发布了一个新的Tensorflow物体识别API



深度学习服务器环境配置:
Ubuntu17.04+Nvidia GTX



Tensorflow学习之实现简单的神经网络



Tensorflow学习之基本概念

- 手把手 | 亲测好用！Google发布了一个新的Tensorflow物体识别API
- Tensorflow学习之实现简单的神经网络
- Tensorflow学习之调优
- Tensorflow 基础

- 深度学习服务器环境配置: Ubuntu17.04+Nvidia GTX 1080+CUDA
- Tensorflow学习之基本概念
- Tensorflow之MNIST解析
- 使用Tensorflow训练循环神经网络语言模型

您必须 [登录](#) 才能发表评论！如果您还没有注册，请点击[这里](#)，快速免费完成注册吧！

版权声明

本站的文章和资源来自互联网或者站长的原创，按照 CC BY -NC -SA 3.0 CN 协议发布和共享，转载或引用本站文章应遵循相同协议。如果有侵犯版权的资源请尽快联系站长，我们会在24h内删除有争议的资源。

意见反馈

邮箱:support@tongtongxue.com



友情链接

新浪微博 RSS 网站地图 ITeye

CSDN 博客园 阿里云 七牛

微信公众号



