# ADL x MLDS 2017 Fall
# HW1 - Sequence Labeling

2017/10/02
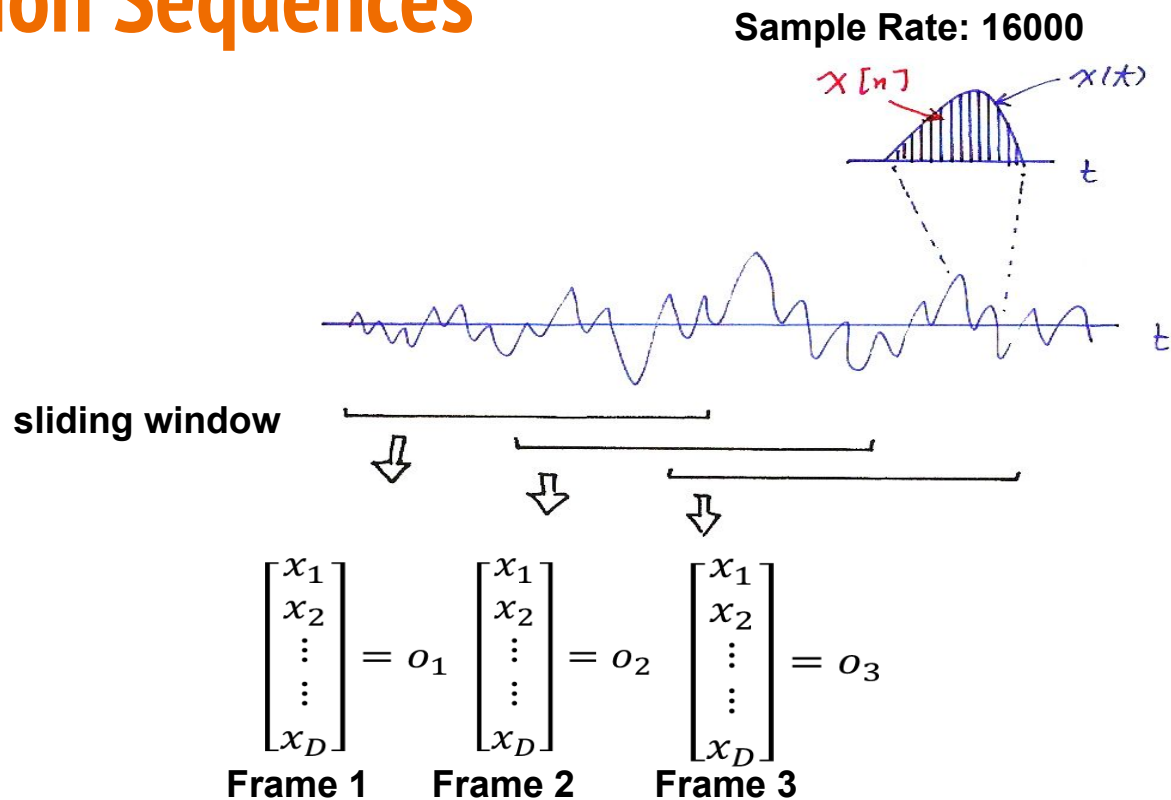adlxmlds@gmail.com

# Outline

- Task Description
    - Phone sequence labeling
    - TIMIT Dataset and Data Format
- Recurrent Neural Networks & Convolutional Neural Networks
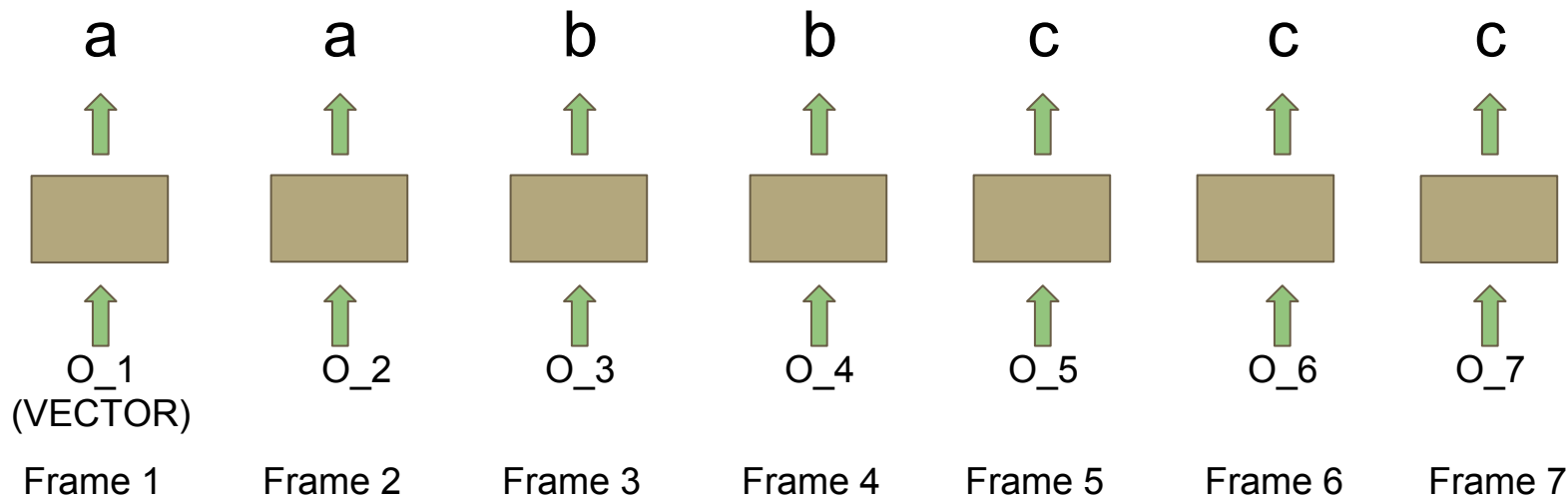- Kaggle
- Grading
- Format and Submission Rules

# Speech Recognition

- In speech processing
  - Each word consist of syllables
  - Each syllables consist of phone
  - "青色" → "青(ㄑㄧㄥ)色(ㄙㄜˋ)" → "ㄑ" (syllables)
    青:TSI --I –N (phone)
    色:S--@ (phone)

- Each time frame, with an observance (vector) mapped to a phone.

# Observation Sequences

**Sample Rate: 16000**

$x[n]$   $x(t)$

**sliding window**

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = o_1 \quad \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = o_2 \quad \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = o_3$$

**Frame 1**   **Frame 2**   **Frame 3**

ref: DSP lect 2.0

# Framewise Prediction

a     a     b     b     c     c     c

$O\_1$
(VECTOR)    $O\_2$    $O\_3$    $O\_4$    $O\_5$    $O\_6$    $O\_7$

Frame 1    Frame 2    Frame 3    Frame 4    Frame 5    Frame 6    Frame 7

# Phone Prediction

- What really matters in speech recognition is the **final phone sequence**, not the framewise alignment.
- That is, the final evaluation in this homework is based on the **phone sequence**.
- **You have to trim the frame-level sequence into phone sequence.**

# Trimmimg on Framewise Sequence (1/2)

- **Remove <u>consecutive duplicate</u> labels**

  Framewise prediction: {a, a, b, b, c, c, c}

  Phone prediction     : {a, b, c}

**You need to report result in <u style="color:red">phone sequence</u>**

# Trimmimg on Framewise Sequence (2/2)

- **Remove <u>only leading and tailing</u> silence**

  Framewise prediction: {<sil>, <sil>, a, a, b, <sil>, c, c, <sil>}

  Phone prediction        : {a, b, <sil>, c}

**You need to report result in <span style="color:red">phone sequence</span>**

# Dataset (1/2)

- **TIMIT**(**T**exas **I**nstrument and **M**assachusetts **I**nstitute of **T**echnology)
- Well-transcribed speech of American English speakers of different sexes and dialects.
- Designed for the development and evaluation of ASR systems.
- Features
  - MFCC: 39 dim
  - FBank: 69 dim

# Dataset (2/2)

Each instance consist of 3 parts: Speaker ID, Sentence ID, Frame ID

Ex:

Speaker   ID: faem0

Sentence ID: si1392

Frame     ID: 37

Instance ID

f a e m 0 _ s i 1 3 9 2 _ 3 7

Speak ID
Start with f: female
Start with m: male

Sentence ID

Frame ID

# Data Format (1/3)

- WAV file: Speaker-Sentence_ID + .wav → Check by your ears
- ARK file: Instance ID + features

# Data Format (2/3)

- LAB file: Instance ID + , + label
- 48 phones
- Map them to 39 phones by **yourselves**

```
1  maeb0_si1411_1,sil
2  maeb0_si1411_2,sil
3  maeb0_si1411_3,sil
4  maeb0_si1411_4,sil
5  maeb0_si1411_5,sil
6  maeb0_si1411_6,sil
7  maeb0_si1411_7,sil
8  maeb0_si1411_8,sil
9  maeb0_si1411_9,sil
10 maeb0_si1411_10,sil
11 maeb0_si1411_11,r
12 maeb0_si1411_12,r
13 maeb0_si1411_13,r
14 maeb0_si1411_14,r
15 maeb0_si1411_15,r
16 maeb0_si1411_16,r
17 maeb0_si1411_17,r
18 maeb0_si1411_18,r
19 maeb0_si1411_19,ix
20 maeb0_si1411_20,ix
21 maeb0_si1411_21,ix
22 maeb0_si1411_22,ix
```

# Data Format (3/3)

- MAP file: 2 mapping

  (1) 48 phones - 39 phones

  (2) 48 phones - 48 English characters

Delimiter: '\t'



MAP (1)



MAP (2)

# Evaluation

- **Average Phone Sequence Edit Distance**
  - Compare your trimmed phone sequence with correct ones
- **Edit Distance = Insertion + Deletion + Substitution**
- Consider the following case, edit distance = I + D + S = 2 + 1 + 2 = 5

# Recurrent Neural Networks

# RNN - Unfolded View

# RNN

- With Hidden Layer (Memory Layer), RNN can learn more long-term information.
    - **Sequential Information.**
- With **LSTM** gated-extension, the RNN can learn longer and longer.

# Long Short-term Memory (LSTM)



Other part of the network

Signal control the output gate

(Other part of the network)

Output Gate

Special Neuron: 4 inputs, 1 output

Memory Cell

Forget Gate

Signal control the forget gate

(Other part of the network)

Signal control the input gate

(Other part of the network)

Input Gate

LSTM

Other part of the network

http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Convolutional Neural Network

# Jointly train RNN with CNN

Output_k        Output_k+1        Output_k+2

RNN         RNN         RNN

n-dim vector       n-dim vector       n-dim vector

CNN         CNN         CNN

[V_k-1, V_k, V_k+1]     [V_k, V_k+1, V_k]     [V_k+1, V_k+2, V_k+3]

Frame_k            Frame_k+1          Frame_k+2

# Jointly train RNN with CNN

# CNN on acoustic features

Take feature MFCC for example:

39 dim

V_K-1    V_K    V_K+1

# CNN on acoustic features

Take feature MFCC for example:



CNN

39 dim

V_K-1    V_K    V_K+1

# CNN on acoustic features

Take feature MFCC for example:

# CNN on acoustic features

Take feature MFCC for example:



39 dim

CNN

V_K-1    V_K    V_K+1

Try different experiment settings and write down your observation in the report !

# CNN Lectures

- Machine Learning 2016 Fall

https://www.youtube.com/watch?v=FrKWiRv254g

- Tóth, László. "Convolutional Deep Maxout Networks for Phone Recognition", Interspeech, 2014

# HW Rules

# HW Rules

- Please write shell script to run your code.
- There should be hw1_rnn.sh, hw1_cnn.sh, hw1_best.sh
- Please follow the script usage below:
    - ./hw1_rnn.sh $1 $2
    - ./hw1_cnn.sh $1 $2
    - ./hw1_best.sh $1 $2
    - $1: the data directory, $2: output filename
- Ex: ./hw1_best.sh myData/ best.csv

# HW Rules

- Please implement RNN-based to fulfill the task
- Please also implement CNN+RNN-based to fulfill the task
- Please use python with version >= 3.5
- Please do not use extra dataset
- Allowed package includes:
    - PyTorch v0.2.0
    - tensorflow r1.3
    - Keras 2.0.7

# Kaggle

# Kaggle (1/3)

- Kaggle: https://www.kaggle.com/t/0d61e84f89594f998b12d999fa4b4d5f
- competition will started at **2017/10/5 12:00 (GMT+8).**
- Please create **ONE** account using your school mail (Ex: NTU)
- For students taking this class, your title on leaderboard should start with **your student ID**
    - Ex: b03xxxxxx_SamIsTheBest
- At most **5** submissions per day
- Indivisual task, do not team up!
- No score counted if
    - you create more accounts to get more submissions == cheating!
    - your title does not conform to the naming rules

# Kaggle (2/3)

- Testing set is divided into two sets: **public** and **private**
- Your performance on leaderboard during the competition is based on the public set
- After deadline, the private set will be evaluated
- Remember to choose **2 submissions** for the final evaluation before deadline, otherwise Kaggle will select for you
- Please do not attempt to fit the public set

# Kaggle (3/3)

- Submission format: a **.csv** file with then content as below
- Remember to map the framewise output to 39 phones
- Remember to map phones to English letter
- Remember to trim <sil>
- With header row: "id,phone_sequence"
- Instance ID + , + predicted phone sequence

```
id,phone_sequence
fadg0_si1279,HrLAJarDeBLMrDcLMwU
fadg0_si1909,vbLAFKnLhyUwJmrBJLAwLSyLAwKr
fadg0_si649,lwLJctryJvrCaBgLHwDLKyDwLHJywDLHrLHwJnDryJLAbLMtrBwLABsmrQ
fadg0_sx109,SJKyJBnLMwDJLHIyDyCrFLABSaDwDJLMyLAJBc
fadg0_sx19,vnBFDwDnDyUJFQSrLABwDatwDLhyJBcDLJwByD
fadg0_sx199,lynDyKnBLkrwDyKwmwLhaIymyLAJLhcIKrLMwLAwLksLzwJLAJwUyJwJ
```

# Grading

# Grading Policy

I. Baseline (6%)

II. Ranking (8%)

III. Report (4%)

IV. Bonus (2%).

V. Notice

# Grading Policy -- Baseline&Ranking

- Pass the public baseline (3%)
- Pass the private baseline (3%)
- Ranking (8%) For those passing the private baseline, your score will be linearly grade, rounded to the 2nd decimal place
  - Ex: if 100 people pass the baseline, you will get 6 points if you're at 25th place.
- We will run your code to make sure your leaderboard performance is aligned with your submission

# Grading Policy -- Report(4%)

- Do not exceed **4** pages and <span style="color:red">written in Chinese</span>
- Model description (2%)
    - RNN (1%)
    - RNN+CNN (1%)
- How to improve your performance (1%)
    - Write down the method that makes you outstanding
    - Describe the model or technique (0.5%)
    - Why do you use it (0.5%)
- Experimental results and settings (1%)
    - Compare and analyze the results between RNN and CNN (0.5%)
    - Compare and analyze the results with other models (0.5%)
        - other models can be variant of basic RNN, like LSTM, or some novel ideas you use

# Grading Policy -- Bonus(2%)

- TAs will select about 5 persons, according to both **creativity** and **performance** (top 10%) for introducing your model during the class
- If you are chosen, you have to present in order to get the bonus

# Grading Policy -- Notice

- Please fill the [late submission form](#) first only if you will submit HW late
- Please push your code before you fill the form
- There will be 25% penalty per day for late submission, so you get 0% after four days
- You can still upload your result on Kaggle, although it won't be counted in your grade

- You get 0% if the required script has bug.

    - If the error is due to the format issue, please come to fix the bug at the announced time, or you will get 10% penalty afterwards

# Submission Rules

# Submission Rules

- Please refer to this link **first**.
- Create hw1 directory under ADLxMLDS2017
- Under hw1, there should be:
    - report.pdf
    - hw1_rnn.sh  // should run your RNN model
    - hw1_cnn.sh  // should run your CNN+RNN model
    - hw1_best.sh // should run your best-performed model
    - model_rnn.py, model_cnn.py, model_best.py and other necessary files
    - *In model_rnn.py, model_cnn.py and model_best.py should include your training codes.
- Please do not upload TIMIT dataset to Github
- If your model are too big for github, upload to a cloud space and write it in your script to download the model
- Your script should be done within 10 mins excluding model donwloading

# Deadline

1. Kaggle deadline: **2017/10/28 12:00 (GMT+8)**
2. Github code & report deadline: **2017/10/28 23:59 (GMT+8)**

# FAQ

# Q1: 使用的lib 限制

A:

除了拿來training的lib有限制以外，其他lib在使用的時候只要沒有使用外部的dataset都是可以的。並且記得在report中註明使用的lib名稱以及版本。

Ex: sklearn的train, test, split沒有用到助教的其他data，所以可以使用。

# Q2: 請問助教會跑training的程式嗎？

A:

不會。我們所規定的十分鐘只包含testing。除非我們認為有必要就會請你們來跑training的code。

# Q3: Dataset在哪裡下載？

A:

Dataset可以從Kaggle上下載。

# Q4: 執行的時候助教要怎麼知道我是使用哪一種feature?

A:

在助教的電腦上，data directory結構如右所示。

而助教在測試的時候，我們argv只會輸入"data/"。所以

同學必須要自己設定好你們需要的檔案路徑讓助教output

正確的答案。

```
data/
----fbank/
--------test.ark
--------train.ark
----label/
--------train.lab
----mfcc/
--------test.ark
--------train.ark
-----phones/
--------48_39.map
----48phone_char.map
```

# Q5: Training label和feature的instance_ID順序不一樣，是要自己去對齊嗎？

A:

是的！這部分要麻煩同學自己去對齊！

# FAQ

- If you have other questions,
  - please contact TAs via adlxmlds@gmail.com
  - post your questions on facebook group
  - go to TA office hours
    - 王昱翔 Mon 16:00-17:30 (電二531)
    - 樊恩宇 Fri 10:30-12:00 (明達526)
    - 古志文 Fri 14:30-16:00 (德田524)