# Automatic Speech Recognition: Introduction

Peter Bell

#### Automatic Speech Recognition— ASR Lecture 1 11 January 2020

• • = • • = •

# Automatic Speech Recognition — ASR

#### Course details

- Lectures: About 18 lectures, delivered live on Teams for now
- Labs: Weekly lab sessions using Python, OpenFst (openfst.org) and later Kaldi (kaldi-asr.org)
  - Lab sessions will start in Week 3 exact format TBA.
- Assessment:
  - First five lab sessions worth 10%
  - $\bullet\,$  Coursework, building on the lab sessions, worth 40%
  - Open book exam in April or May worth 50%
- People:
  - Course organiser: Peter Bell
  - Guest lecturers: Hiroshi Shimodaira and Yumnah Mohammied
  - TA: Andrea Carmantini
  - Demonstrators: Chau Luu and Electra Wallington

http://www.inf.ed.ac.uk/teaching/courses/asr/

- If you have taken:
  - Speech Processing and either of (MLPR or MLP)
    - Perfect!
  - either of (MLPR or MLP) *but not* Speech Processing (probably you are from Informatics)
    - You'll require some speech background:
      - A couple of the lectures will cover material that was in Speech Processing
      - Some additional background study (including material from Speech Processing)
  - Speech Processing *but neither of* (MLPR or MLP) (probably you are from SLP)
    - You'll require some machine learning background (especially neural networks)
      - A couple of introductory lectures on neural networks provided for SLP students
      - Some additional background study

イロト イポト イヨト イヨト

- Series of weekly labs using Python, OpenFst and Kaldi
- They count towards 10% of the course credit
- Labs start week 3 exact arrangements TBA
- You will need to work in pairs
- Labs 1-5 will give you hands-on experience of using HMM algorithms to build your own ASR system
  - These labs are an important pre-requisite for the coursework take advantage of the demonstrator support!
- Later optional labs will introduce you to Kaldi recipes for training acoustic models – useful if you will be doing an ASR-related research project

# What is speech recognition?

・日・・ モ・・ モ・

æ

# What is speech recognition?







・ロト ・日ト ・ヨト ・ヨト

æ

#### Speech-to-text transcription

- Transform recorded audio into a sequence of words
- Just the words, no meaning.... But do need to deal with acoustic ambiguity: "Recognise speech?" or "Wreck a nice beach?"
- Speaker diarization: Who spoke when?
- Speech recognition: what did they say?
- Paralinguistic aspects: how did they say it? (timing, intonation, voice quality)
- Speech understanding: what does it mean?

A (1) < A (2) < A (2) </p>

# Why is speech recognition difficult?

白 ト イ ヨ ト イ ヨ ト

크

Many sources of variation

Speaker Tuned for a particular speaker, or speaker-independent? Adaptation to speaker characteristics

白 ト イ ヨ ト イ ヨ ト

Many sources of variation

Speaker Tuned for a particular speaker, or speaker-independent? Adaptation to speaker characteristics

Environment Noise, competing speakers, channel conditions (microphone, phone line, room acoustics)

• • = • • = •

Many sources of variation

Speaker Tuned for a particular speaker, or speaker-independent? Adaptation to speaker characteristics

Environment Noise, competing speakers, channel conditions (microphone, phone line, room acoustics)

Style Continuously spoken or isolated? Planned monologue or spontaneous conversation?

白 ト イヨト イヨト

Many sources of variation

Speaker Tuned for a particular speaker, or speaker-independent? Adaptation to speaker characteristics

Environment Noise, competing speakers, channel conditions (microphone, phone line, room acoustics)

- Style Continuously spoken or isolated? Planned monologue or spontaneous conversation?
- Vocabulary Machine-directed commands, scientific language, colloquial expressions

白 ト イヨト イヨト

Many sources of variation

Speaker Tuned for a particular speaker, or speaker-independent? Adaptation to speaker characteristics

Environment Noise, competing speakers, channel conditions (microphone, phone line, room acoustics)

- Style Continuously spoken or isolated? Planned monologue or spontaneous conversation?
- Vocabulary Machine-directed commands, scientific language, colloquial expressions

Accent/dialect Recognise the speech of all speakers who speak a particular language

・ 回 ト ・ ヨ ト ・ ヨ ト

Many sources of variation

Speaker Tuned for a particular speaker, or speaker-independent? Adaptation to speaker characteristics

Environment Noise, competing speakers, channel conditions (microphone, phone line, room acoustics)

- Style Continuously spoken or isolated? Planned monologue or spontaneous conversation?
- Vocabulary Machine-directed commands, scientific language, colloquial expressions

Accent/dialect Recognise the speech of all speakers who speak a particular language

Other paralinguistics Emotional state, social class, ...

・日・ ・ヨ・ ・ヨ・

Many sources of variation

Speaker Tuned for a particular speaker, or speaker-independent? Adaptation to speaker characteristics

Environment Noise, competing speakers, channel conditions (microphone, phone line, room acoustics)

- Style Continuously spoken or isolated? Planned monologue or spontaneous conversation?
- Vocabulary Machine-directed commands, scientific language, colloquial expressions

Accent/dialect Recognise the speech of all speakers who speak a particular language

Other paralinguistics Emotional state, social class, ....

Language spoken Estimated 7,000 languages, most with limited training resources; code-switching; language change

• As a classification problem: very high dimensional output space

白 ト イヨト イヨト

- As a classification problem: very high dimensional output space
- As a sequence-to-sequence problem: very long input sequence (although limited re-ordering between acoustic and word sequences)

白 ト イ ヨ ト イ ヨ ト

- As a classification problem: very high dimensional output space
- As a sequence-to-sequence problem: very long input sequence (although limited re-ordering between acoustic and word sequences)
- Data is often noisy, with many "nuisance" factors of variation in the data

白 ト イ ヨ ト イ ヨ ト

- As a classification problem: very high dimensional output space
- As a sequence-to-sequence problem: very long input sequence (although limited re-ordering between acoustic and word sequences)
- Data is often noisy, with many "nuisance" factors of variation in the data
- Very limited quantities of training data available (in terms of words) compared to text-based NLP
  - Manual speech transcription is very expensive (10x real time)

・ 回 ト ・ ヨ ト ・ ヨ ト …

- As a classification problem: very high dimensional output space
- As a sequence-to-sequence problem: very long input sequence (although limited re-ordering between acoustic and word sequences)
- Data is often noisy, with many "nuisance" factors of variation in the data
- Very limited quantities of training data available (in terms of words) compared to text-based NLP
  - Manual speech transcription is very expensive (10x real time)
- Hierachical and compositional nature of speech production and comprehension makes it difficult to handle with a single model

(日本) (日本) (日本)

• We generally represent recorded speech as a sequence of acoustic feature vectors (observations), X and the output word sequence as W

- We generally represent recorded speech as a sequence of acoustic feature vectors (observations), X and the output word sequence as W
- At recognition time, our aim is to find the most likely W, given X

- We generally represent recorded speech as a sequence of acoustic feature vectors (observations), X and the output word sequence as W
- At recognition time, our aim is to find the most likely W, given X
- To achieve this, statistical models are trained using a corpus of labelled training utterances (X<sup>n</sup>, W<sup>n</sup>)

# Representing recorded speech (X)



Represent a recorded utterance as a sequence of *feature vectors* 

Reading: Jurafsky & Martin section 9.3

(日)

# Labelling speech (W)



Labels may be at different levels: words, phones, etc. Labels may be *time-aligned* – i.e. the start and end times of an acoustic segment corresponding to a label are known

Reading: Jurafsky & Martin chapter 7 (especially sections 7.4, 7.5)

Aligning the sequences  $X^n$  and  $W^n$  for each training utterance

• • = • • = •

э

Aligning the sequences  $X^n$  and  $W^n$  for each training utterance



高 ト イ ヨ ト イ ヨ ト

э

Aligning the sequences  $X^n$  and  $W^n$  for each training utterance



向下 イヨト イヨト

э

Aligning the sequences  $X^n$  and  $W^n$  for each training utterance



向下 イヨト イヨト

크

Aligning the sequences  $X^n$  and  $W^n$  for each training utterance



向下 イヨト イヨト

э

Aligning the sequences  $X^n$  and  $W^n$  for each training utterance



#### In performing recognition:

Searching over all possible output sequences W

to find the most likely one

Aligning the sequences  $X^n$  and  $W^n$  for each training utterance



#### In performing recognition:

Searching over all possible output sequences W to find the most likely one

The **hidden Markov model** (HMM) provides a good solution to both problems

周 ト イ ヨ ト イ ヨ ト

## The Hidden Markov Model



- A simple but powerful model for mapping a sequence of continuous observations to a sequence of discrete outputs
- It is a generative model for the observation sequence
- Algorithms for training (forward-backward) and recognition-time decoding (Viterbi)

## The Hidden Markov Model



- A simple but powerful model for mapping a sequence of continuous observations to a sequence of discrete outputs
- It is a generative model for the observation sequence
- Algorithms for training (forward-backward) and recognition-time decoding (Viterbi)
- Later in the course we will also look at newer all-neural, fully-differentiable "end-to-end" models

# Hierarchical modelling of speech



・ロト ・日ト ・ヨト ・ヨト

æ

# "Fundamental Equation of Statistical Speech Recognition"

If X is the sequence of acoustic feature vectors (observations) and W denotes a word sequence, the most likely word sequence  $W^*$  is given by

$$\mathsf{W}^* = rg\max_{\mathsf{W}} \mathsf{P}(\mathsf{W} \mid \mathsf{X})$$

# "Fundamental Equation of Statistical Speech Recognition"

If X is the sequence of acoustic feature vectors (observations) and W denotes a word sequence, the most likely word sequence  $W^*$  is given by

$$W^* = \arg \max_{W} P(W \mid X)$$

Applying Bayes' Theorem:

$$P(W \mid X) = \frac{p(X \mid W)P(W)}{p(X)}$$

$$\propto p(X \mid W)P(W)$$

$$W^* = \arg \max_{W} \underbrace{p(X \mid W)}_{\text{Acoustic}} \underbrace{P(W)}_{\text{Language}}$$

$$model$$

$$\mathsf{W}^* = \arg\max_{\mathsf{W}} p(\mathsf{X} \mid \mathsf{W}) P(\mathsf{W})$$

Use an acoustic model, language model, and lexicon to obtain the most probable word sequence  $W^*$  given the observed acoustics X



э

지마지 지금 에 문 에 문 에

#### Phonemes

- abstract unit defined by linguists based on contrastive role in word meanings (eg "cat" vs "bat")
- 40-50 phonemes in English
- Phones
  - speech sounds defined by the acoustics
  - many allophones of the same phoneme (eg /p/ in "pit" and "spit")
  - limitless in number
- Phones are usually used in speech recognition but no conclusive evidence that they are the basic units in speech recognition
- Possible alternatives: syllables, automatically derived units, ...

・ロト ・日ト ・ヨト

## **Evaluation**

- How accurate is a speech recognizer?
- String edit distance
  - Use dynamic programming to align the ASR output with a reference transcription
  - Three type of error: insertion, deletion, substitutions
- Word error rate (WER) sums the three types of error. If there are *N* words in the reference transcript, and the ASR output has *S* substitutions, *D* deletions and *I* insertions, then:

$$WER = 100 \cdot \frac{S + D + I}{N}\% \qquad Accuracy = 100 - WER\%$$

• Speech recognition evaluations: common training and development data, release of new test sets on which different systems may be evaluated using word error rate

高 ト イ ヨ ト イ ヨ ト



ヘロア 人間 アメヨア 人間 アー

æ

# Example: recognising TV broadcasts







イロト イヨト イヨト イヨト

# Reading

- Jurafsky and Martin (2008). *Speech and Language Processing* (2nd ed.): Chapter 7 (esp 7.4, 7.5) and Section 9.3.
- General interest:
  - The Economist Technology Quarterly, "Language: Finding a Voice", Jan 2017. http://www.economist.com/technology-quarterly/2017-05-01/language
  - The State of Automatic Speech Recognition: Q&A with Kaldi's Dan Povey, Jul 2018. https://medium.com/descript/the-state-of-automaticspeech-recognition-q-a-with-kaldis-dan-poveyc860aada9b85

・ロト ・回ト ・ヨト ・ヨト