# Data Mining and Machine Learning

# HMMs for peech Recogniti

# Word and Sub-Wo el HMMs

Peter Jančovič

UNIVERSITY OF BIRMINGHAM

# Content

- Word level HMMs

- Sub-word HMMs
  - Phoneme-level HMMs

- Context-sens https://eduassistpro.github.io/
  - Biphone HMMs
  - Triphone HMMs

- Triphone HMM training issues

- Phoneme Decision Trees (PDTs)

Assignment Project Exam Help

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Word Level HMMs

- Early systems (1980s) used <u>word</u> level HMMs

- I.e. each word modelled by a single, dedicated HMM (c.f. "zero" picture)

  - Advantag

  - Good perf                                                    t modelling
    of word-dependent vari

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 6 state HMM of the digit 'zero'

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Word Level HMMs

- Disadvantages:
  - Many examples of each word needed for training
  - Fails to ex _____ oken language
- Word-level systems typica _____ ed to well-defined, demanding, small _____ y applications

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Sub-Word Level HMMs

- Build HMMs for a complete set of sub-word 'building blocks'

- Construct word-level HMMs by concatenation of sub-word HMMs

- E.g.  $slide = /$

**/ s /**          **/ l /**          **/ aɪ /**          / d /

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Sub-Word Level HMMs

- Advantages

  - Able to exploit regularities in speech patterns

  - More efficient use of training data - e.g. in phoneme                                    f aI v /) and "nine" (/n                                    to /aI/ model.

  - Flexibility - acoustic m                    e built **immediately** for words which did not occur in the training data

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Phoneme-Level HMMs

- Why choose phonemes rather than any other sub-word unit?

- Disadvantages

  - Phonemes the contrastive properties o _____ unds within a language - not their c _____ with HMM assumptions!

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Advantages of Phoneme-HMMs

- Completeness & compactness – approx. 50 phonemes required to describe English

- Well studied                                      ation of 'speech kno                                      tion differences due to accent...

- Availability of extensive phoneme-based pronunciation dictionaries

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Context-Sensitivity

- Problem
  - Acoustic realization of a phoneme depends on the c                curs
  - Think of                 he "k" sound in the words "book sho                hick"

Slide 10

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Biphones and Triphones

- Solution
  - **Context-sensitive** phoneme-level HMMs
  - E.g.
    - 'biphon <span style="color:red">https://eduassistpro.github.io/</span> p"
    - 'triphones' : (b-u-S) i <span style="color:red">edu_assist_pro</span> p"
- Almost all systems use triphone HMMs

<span style="color:red">Assignment Project Exam Help</span>

<span style="color:red">Add WeChat edu_assist_pro</span>

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Triphones - problems

- Increased number of model parameters

  – Need more (well-chosen) training data

- Which trip

  – If a word                                    ains a triphone
  which was not in the tra                    which triphone
  HMM should we use?

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Number of parameters

- If there are 50 phones, the maximum number of triphone HMMs is $50^3 = 125,000$

- Most ruled out by **phonological** constraints – most phone triples

- But many are legal

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example: Model Parameters

- Each model has 3 emitting states

- Each state modelled as, say, a 10 component Gaussian mixture

- Each feature                                                al

- Hence numb                                          odel is:

$$3 \times (10 \times (40 + 40 \qquad \qquad 7$$

| Number of states | Number of mixture components | Mean vector | Variance vector | Mixture weight | Transition probs |

UNIVERSITY OF BIRMINGHAM

# Acoustic model parameters

- So, even if we only have 1,000 acoustic models (instead of 125,000), total acoustic model parameters will be 2,457,000

- Too many to quantity of data

- Most common solution is **ameter tying**

- **Different** HMMs share **sa** ters

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Tied variance

- Variances are more costly to estimate than means

- Simple solution – divide set of all HMMs into classes, so that within a class all HMM state PDFs have same v <span style="color:red">https://eduassistpro.github.io/</span>

- This is **tied variance**

- If **all** HMM state PDFs sh          e variance, the variance is referred to as **grand variance**

<span style="color:red">Assignment Project Exam Help</span>

<span style="color:red">Add WeChat edu_assist_pro</span>

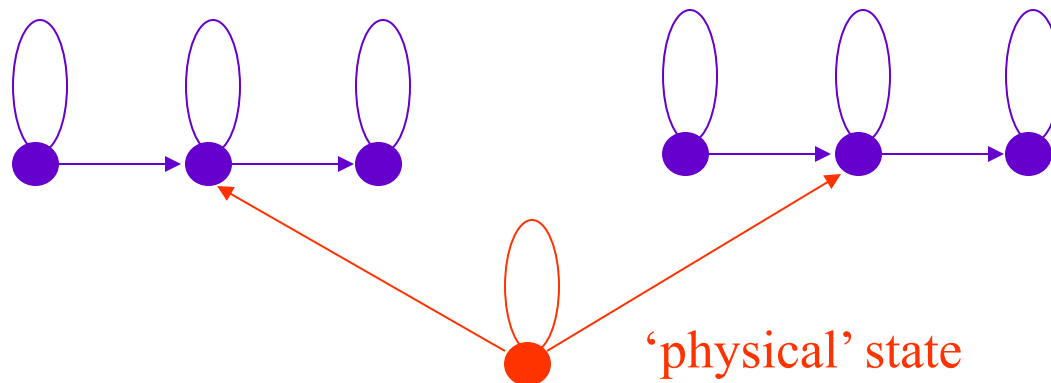UNIVERSITY OF BIRMINGHAM

# Phone decision trees

- Most common approach to general HMM tying is **decision tree clustering**

- Decision tree clustering can be applied to individual states or to w                                  nsider states

- Basic idea is t                                         t           nes are likely to in                     ffe  ts

'Logical'
models

'physical' state

Data Mining and Machine Learning
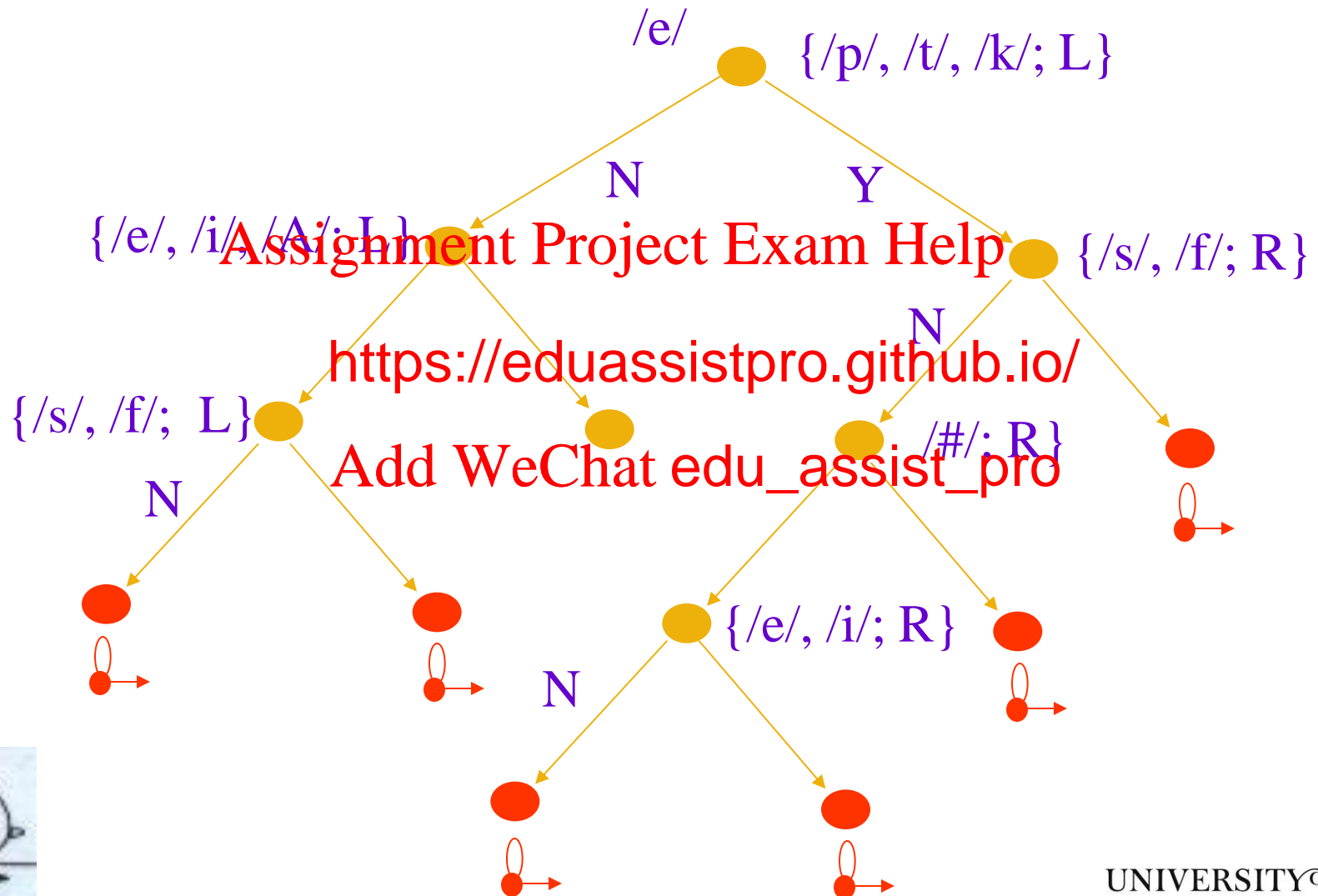
UNIVERSITY OF
BIRMINGHAM

# Phonetic knowledge

- For example, we know that /ʃ/ and /s/ are both unvoiced fricatives, produced in a similar manner

- Therefore we might **hypothesise** that, for example, an utterance ~~ed by~~ /ʃ/ might be similar to

- This is the basic idea behi                    tree clustering

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Phone Decision Tree

/e/ {/p/, /t/, /k/; L}

N Y

{/e/, /i/, /A:; L}

{/s/, /f/; R}

N

{/s/, /f; L}

N

{/#/; R}

{/e/, /i/; R}

N

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Summary

- Word-level and Sub-Word HMMs

- Phoneme-level HMMs

Assignment Project Exam Help

- Context-sens

  – Biphones & https://eduassistpro.github.io/

- Triphone decision trees

  Add WeChat edu_assist_pro

Data Mining and Machine Learning

**UNIVERSITY** OF
**BIRMINGHAM**