

# Data Mining and Machine Learning

Assignment Project Exam Help

Language Automatic  
Speech R <https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

Peter Jančovič



# Objectives

- Understand role of language model in speech recognition
- Approaches
  - Rule-Based
  - Statistical Language Models
- N-gram Language Models

Assignment Project Exam Help

g:

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# Speech Recognition: Statistical Methods

- Given an unknown utterance  $y$ , want to find the word sequence  $W$  such that  $P(W / y)$  is maximised

Assignment Project Exam Help

By Bayes' T

$$P(W | y) = \frac{p(y | W)P(W)}{p(y)}$$

Add WeChat <https://eduassistpro.github.io/>  
edu\_assist\_pro

- $P(W)$  - probability that the word sequence  $W$  is in application language - **language model probability**



# Language Modelling

- **Language Model (Grammar)** used to compute the probability  $P(W)$  that the sequence of words  $W$  'belongs to' the language

Assignment Project Exam Help

- Constrains recognition possible interpretations

<https://eduassistpro.github.io/>

- Basically there are two types of :  
Add WeChat edu\_assist\_pro
  - Rule-based (traditional) lang
  - Probabilistic language model



# Rule-Based Language Models

- Language models in linguistics and natural language processing typically **rule-based**
- A rule-based language model consists of:
  - A set of **non-terminal units** (e.g. sentence, noun-phrase, verb-phras
  - A set of **te**
  - A set of **rules**, defining how non-terminal units can be expanded into sequences of non-terminal and terminal units
- Corresponds to formal notion of grammar – like in school



# Rule-Based Language Models

- Let  $S$  denote the non-terminal root node corresponding to ‘sentence’
- A sequence of words is **grammatical** if it can be derived from  $S$  by a sequence of rules
- Example: Consider the sentence  
“The cat devoured the tiny mouse”

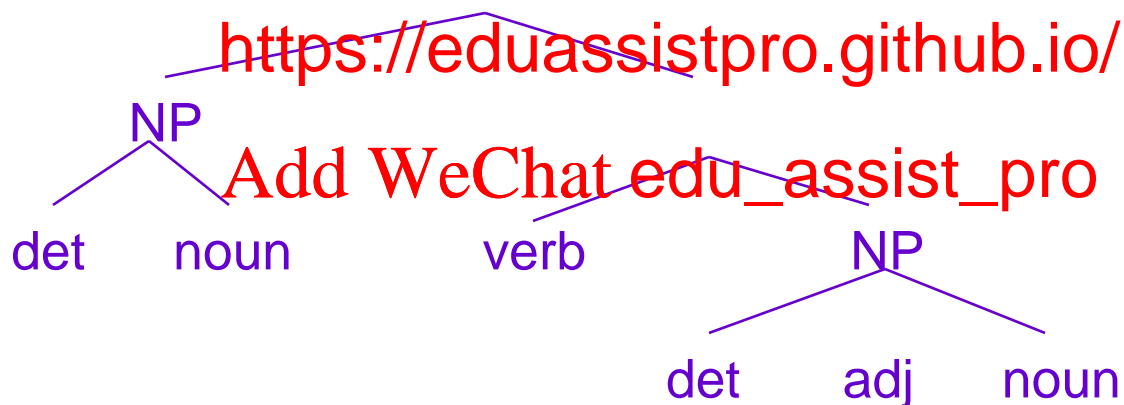


(From Geoffrey Finch, “How to study linguistics”, MacMillan, 1998)

# Rule-Based Language Models

## ■ Example rules:

- S :- NP + VP
- NP:- det + noun
- VP:- verb + NP
- NP:- det + adj + noun
- det:- “the”
- noun:- “cat”
- verb:- “devoured”
- adj:- “tiny”
- noun:- “mouse”



The cat devoured the tiny mouse



# Rule-Based Language Models

- Disadvantages

- Normally applied to **written** language
- A deterministic model of this type may not be able to acc **f spoken** language  
<https://eduassistpro.github.io/>
- Cannot easily handle un **Add WeChat edu\_assist\_pro**
- Cannot be derived automatically from example data and is - based on human knowledge





# Rule-Based Language Models

## ■ Advantages

- Can model complex structure, e.g. non-local dependencies
- Significant knowledge already exists
- Much effort has already been put into the construction of large language models of this type

“She ran, waving enthusiastically, across the bridge”



# Finite State Language Models

- Describe all possible sentences as routes through a finite state network
- Typically 'hand-crafted' using graphical design tools
- Not normally used for vocabularies greater than ~1,000 words

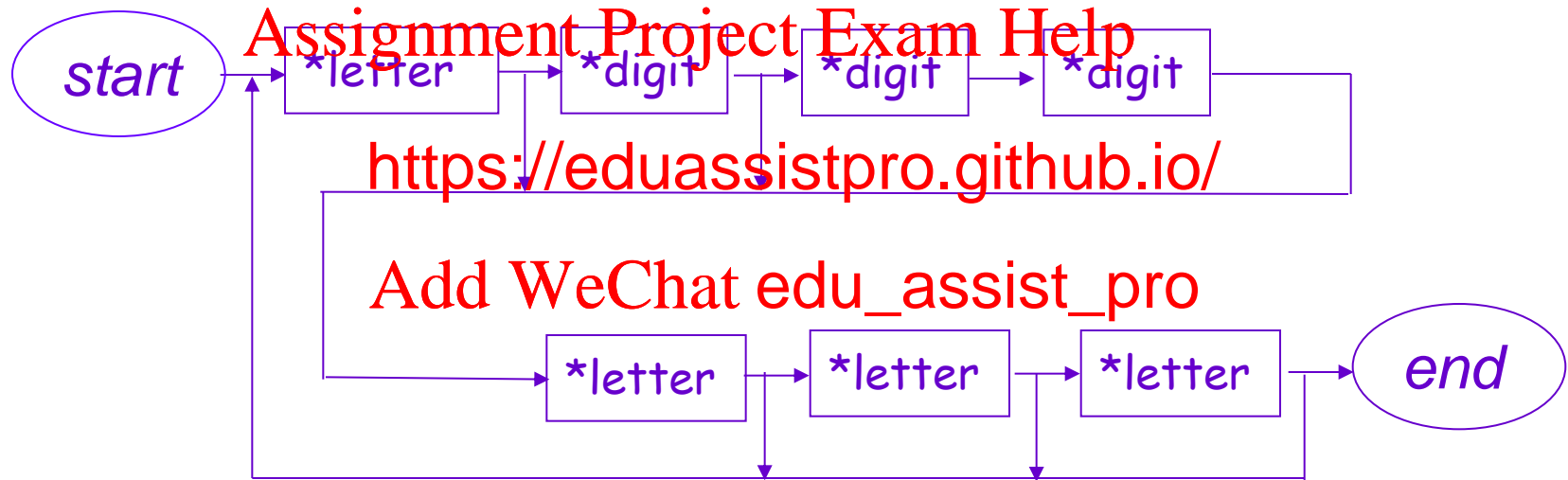
Assignment Project Exam Help

<https://eduassistpro.github.io/>

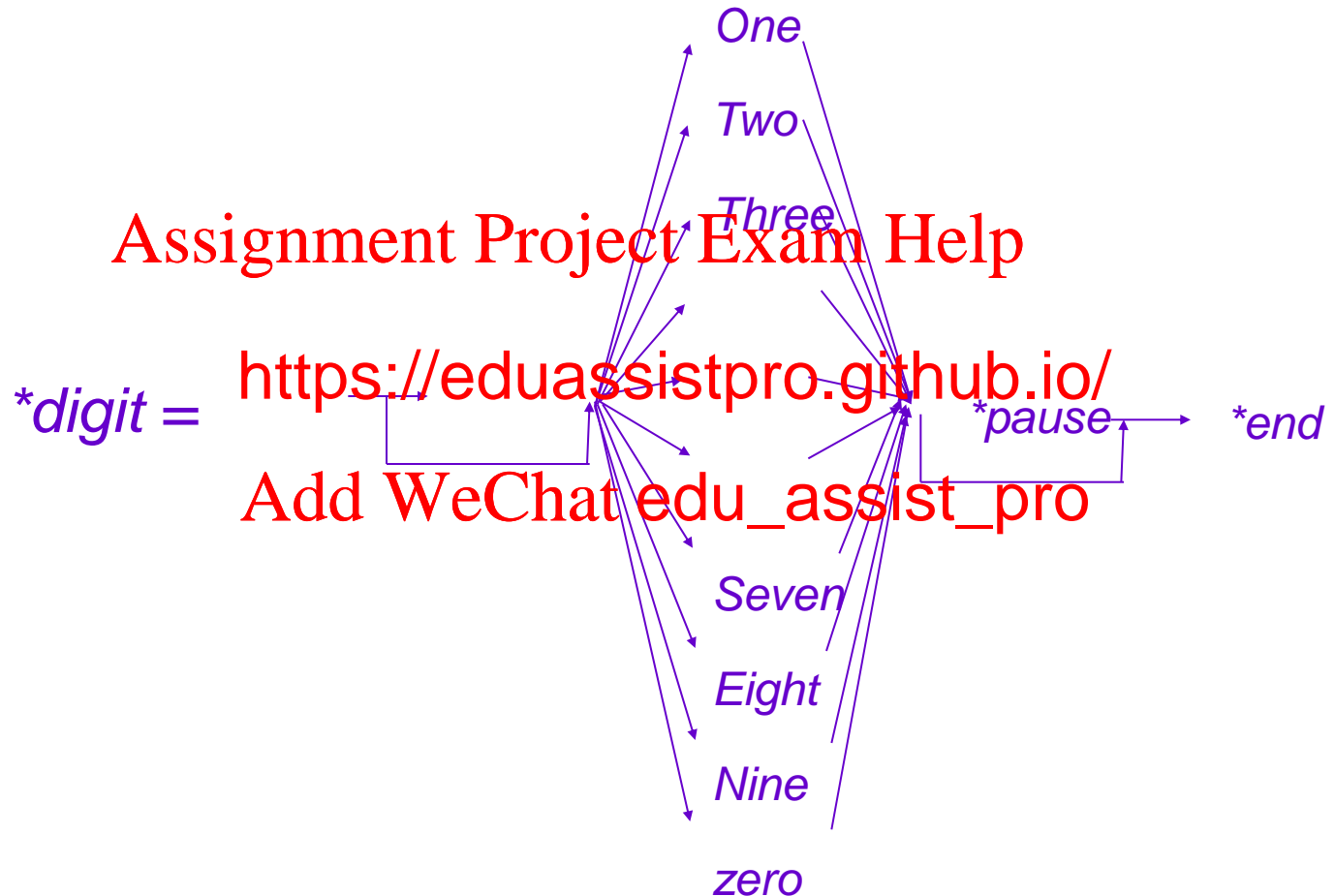
Add WeChat edu\_assist\_pro



# Finite-State Syntax



# Expansion of Macros



# Statistical Language Models

- With a rule based language model, a sequence of words  $W$  is either
  - in the language (**grammatical**) or
  - outside the (**cal**)
- With a statistical language model, the probability of a sequence of words  $W$  is in the language (**grammatical**) with probability  $P(W)$
- The most common statistical language model is known as the  **$N$ -gram model**



# N-gram Language Models

- Let  $W = W_1, W_2, \dots, W_K$  be a sequence of words
- In general:

$$P(W) = P(W_1)P(W_2/W_1) \dots P(W_k/W_1, \dots, W_{k-1}) \dots P(W_K/W_1, \dots, W_{K-1})$$

- In an  $N$ -gram model:

$$P(W_k/W_{k-1}, W_{k-2}, \dots, W_1)$$

i.e. the probability of the  $k^{th}$  word in the sequence depends only on identities of the previous  $N-1$  words

- The most commonly used  $N$ -gram models are 2-gram (**bigram**) and 3-gram (**trigram**) models



# Bigram and Trigram Models

- In a **Bigram Language Model**, we assume:

$$P(W_k/W_{k-1}, W_{k-2}, \dots, W_1) = P(W_k/W_{k-1})$$

Assignment Project Exam Help

- Similarly, in **Model**, we assume: <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

$$P(W_k/W_{k-1}, W_{k-2}, \dots, W_1)$$

- These probabilities can be **estimated from data**



# Estimation of Bigram Probabilities

- For example, given a training text, an estimate of the bigram probability  $P(W_2/W_1)$  is given by:

$$P(W_2/W_1) = \frac{N(W_1, W_2)}{N(W_1)}$$

where:

<https://eduassistpro.github.io/>

- $N(W_1, W_2)$  = number of time pair  $W_1, W_2$  occurs in the training text
- and  $N(W_1)$  = number of times the word  $W_1$  occurs in the training text





# Bigram Probabilities - Example

- Consider the training text:

*“John sat on the old chair. John read the old book. John was interesting. The book was interesting”*

- Suppose this is used to train a bigram grammar.
- ‘the’ occurs 3 times in the training text. The bigrams ‘the old’ and ‘the book’ occur 2 and 1 times respectively. Hence

$$P(\text{'old'} | \text{'the'}) = 2/3, \text{ and } P(\text{'book'} | \text{'the'}) = 1/3.$$

- Similarly, if the symbol # denotes start of sentence, then

$$P(\text{'john'} | \#) = 3/4, \text{ and } P(\text{'the'} | \#) = 1/4$$



# Example Continued

- The probability of the sentence  $S$   
*“John sat on the old chair”* is given by:

$$P(S)$$

$$= P(\text{john}/\epsilon) \cdot P(\text{sat}/\text{john}) \cdot P(\text{on}/\text{sat}) \cdot P(\text{the}/\text{on}) \cdot P(\text{old}/\text{the}) \\ \cdot P(\text{chair}/$$

$$= 3/4 \cdot 1/3 \cdot$$

- Similarly **Add WeChat edu\_assist\_pro**

$$P(\text{“The old chair”}) = 1/12$$

$$P(\text{“John read the old chair”}) = 1/12$$

- But  $P(\text{“John read the interesting book”}) = 0$



# Bigram & Trigram Estimation

- Most practical systems use a trigram language model
- In reality, there is never enough text to estimate trigram probabilities in this simple way
- E.g. experiments with trigram language models for a 1,000 word vocabulary
  - using 1.5 million and 300,000 words to test the models
  - 23% of the trigrams in the test set were absent from the **training** corpus
- Hence much more sophisticated training procedures are needed



# Estimation of $N$ -gram statistics

- In general, there will not be enough data to estimate  $N$ -gram statistics reliably.
- Possible solutions:
  - Robust estimation
  - Deleted interpolation
  - ‘Back-off’

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# Deleted interpolation

- ‘Interpolate’ trigram probability from estimated trigram, bigram and unigram probabilities:

Assignment Project Exam Help

$$\hat{P}(w_3 | w_2 w_1) \approx \lambda_1 P(w_3 | w_2 w_1) + \lambda_2 P(w_3 | w_2) + \lambda_3 P(w_3)$$

Add WeChat edu\_assist\_pro

- Estimate  $\lambda_1, \lambda_2, \lambda_3$  through recognition experiments



# ‘Backoff’

- Decide how many examples  $T$  are needed for robust estimation.
- Then:

$$\hat{P}(w_3 | w_2 w_1) = \begin{cases} P(w_3 | w_1 w_2 w_3) & |w_1 w_2 w_3| \geq T \\ P(w_3 | w_1 w_2) & |w_1 w_2| \geq T \\ P(w_3) & \text{otherwise} \end{cases}$$



# $N$ -gram Language Models - Summary

- Advantages

- Can be trained automatically from data
- Probabilistic model
- Consistent
- Mathematically sound

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# $N$ -gram Language Models - Summary

- Disadvantages

- Large amounts of training data needed
- Difficult to gain knowledge
- Cannot measure accuracy.

*“She walked al. hand in p ickly  
across the bridge”*





# Summary

- Role of language modelling in speech recognition
- Rule based language models
- Finite state l
- *N*-gram lang
- Difficulty of estimating *N*-grams

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat: edu\_assist\_pro

