

# COMP4650 / COMP6490

## Document Analysis 2018

Assignment Project Exam Help

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

Add WeChat: edu\_assist\_pro

Gabriela F



Australian  
National  
University

# Sequence labeling II

Weakness of Markov approaches is that it limits the context from which the current token is tracked

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

Anything outside the context window has no impact on the decision being made...

# Sequence labeling II

CRFs are indeed basically the sequential version of logistic regression

Assignment Project Exam Help

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

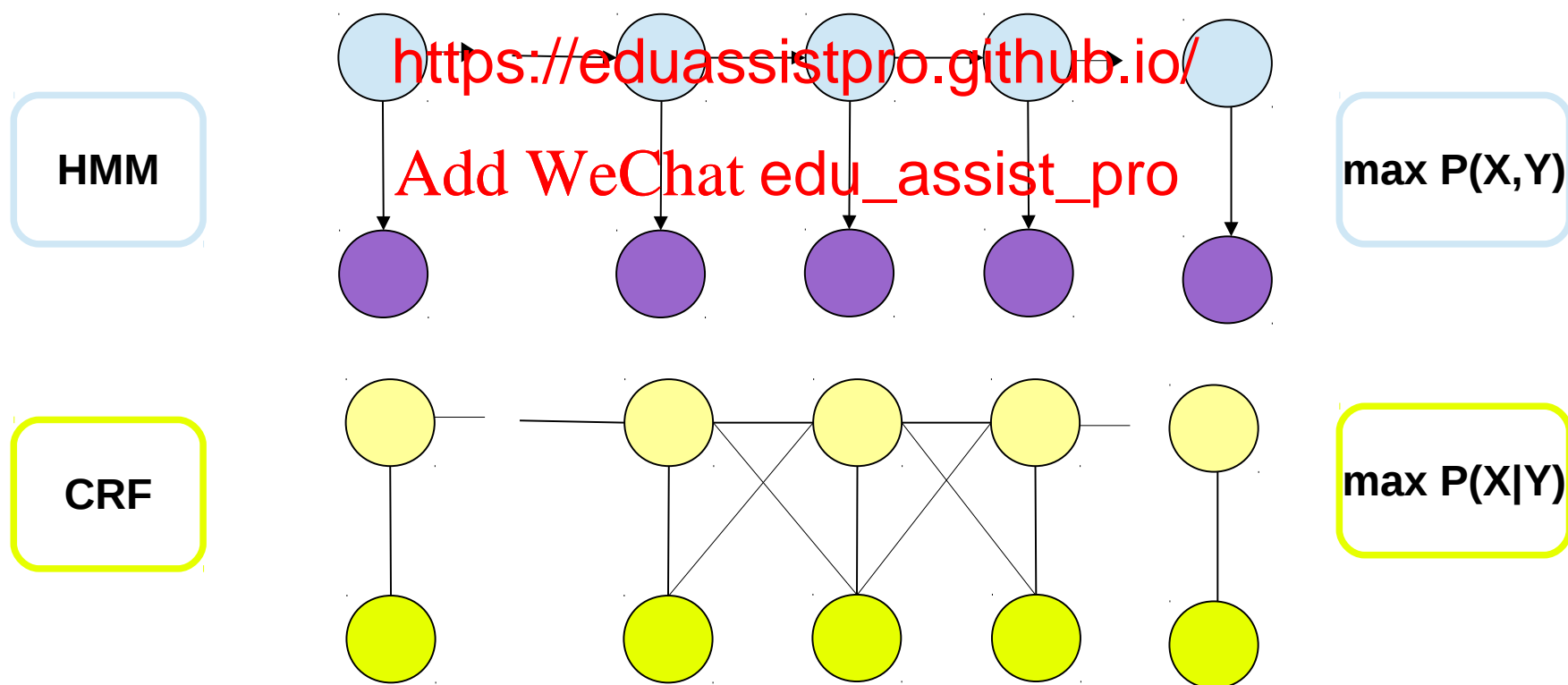
... whereas logistic regression is a linear classification, **CRFs are a log-linear model for sequential labels.**

# Sequence labeling II

CRFs can define a much larger set of features

HMMs are necessarily local in nature because they're constrained to binary transition and emission feature functions

which force each word to depend only on the current label and each label to depend only on the previous label



# Sequence labeling II

## How to label a sentence using CRF?

The naive way is to calculate  $p(\text{labels}|\text{sequence})$  for every possible labeling  $l$ , and then choose the label that maximizes this probability.

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

However, this is intractable.

Add WeChat edu\_assist\_pro

A better way is to realize that (linear-chain) CRFs satisfy an **optimal substructure** property that allows us to use a dynamic programming algorithm to find the optimal label, e.g., the **Viterbi algorithm** for HMMs.

# Information Extraction

BIO encoding

*(B) beginning*

Assignment Project Exam Help

*(I) inside*

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

*(O) other*

Add WeChat edu\_assist\_pro

$2n + 1$  tags, where  $n$  is the number  
of entity types

# BIO encoding

Without the B tag IO tagging is unable to distinguish between two entities of the same type that are right next to each other.

Assignment Project Exam Help

Since this situation does not exist in the real world, there is at least some punctuation or other device to distinguish them.

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

Add WeChat edu\_assist\_pro

+ IO tagging may be sufficient

+ advantage of using only  $n + 1$  tags

# Word-by-word feature encoding

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro



# Named Entity Recognition as sequence labeling

The features available to the classifier during training and classification are those in the boxed area

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Neuronal algorithm for NER

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

- Use a CRF layer on top of the bi-LSTM
- Use Viterbi for decoding for selecting the most likely tag sequence

# Information Extraction

- Extracting time and dates
  - Question answering
  - Calendar assistance
  - Personal as <https://eduassistpro.github.io/>
  - ... Add WeChat edu\_assist\_pro

Needs normalization!

So we can reason about them...

# Extracting time and dates

- Absolute → map to calendar dates
- Relative → some other r  
– A week from last Tuesday.  
– seconds, minutes, days, weeks, centuries, etc
- Duration → spans of time with different granularities

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Extracting time and dates

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Extracting time and dates

## Lexical triggers

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

A fare increase initiated <TIMEX3>last w by UALCorp's United Airlines was matched by competitors over <TIMEX3>the weekend</TIMEX3>, marking the second successful fare increase in<TIMEX3>two weeks</TIMEX3>.

(Pustejovsky et al. 2005, Ferro et al. 2005)

# Extracting time and dates

Sequence labeling with BIO encoding

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Extracting time and dates

Normalization: *VALUE attribute*

from the ISO 8601 standard for encoding temporal values

(ISO8601, 2004)

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro



# Events Extraction

- Identify mentions of events in texts

events can be assigned to point (or interval) in time

- sequence labeling

Assignment Project Exam Help

- BIO encoding

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

- usually applied s                      ning methods

Add WeChat edu\_assist\_pro

Feature	Explanation
Character affixes	Character-level prefixes and suffixes of target word
Nominalization suffix	Character level suffixes for nominalizations (e.g., <i>-tion</i> )
Part of speech	Part of speech of the target word
Light verb	Binary feature indicating that the target is governed by a light verb
Subject syntactic category	Syntactic category of the subject of the sentence
Morphological stem	Stemmed version of the target word
Verb root	Root form of the verb basis for a nominalization
WordNet hypernyms	Hypernym set for the target

# Events Extraction

Events + temporal expressions → Temporal ordering of the events

- Timeline

Classify events according to temporal relations

- Similar to Relation Extraction: relations between entities, relations are between events
- Finite set of temporal relations (Allen, 1984)

Useful for Q&A and summarization

TimeBank corpus

Allen  
relations  
between  
temporal  
events

**Assignment Project Exam Help**

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

**Add WeChat edu\_assist\_pro**

# Benchmarking

How do you know your method is working?

How good it is in respect to other methods?

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

Have a B !!!

# What is a baseline?

- Information that is used as a starting point by which to compare

Assignment Project Exam Help

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

- Benchmark Add WeChat edu\_assist\_pro

- Something you want to beat

# How a baseline looks like?

- Random assignment
- Majority class voting
- Simple heuristics
- Simple Machine Learning techniques
- Simple feature sets
- The system/method you want to beat!

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

# Summary

- **Named entities:** who and who's class (type)
- **Relation extraction:** who is doing what
- **Temporal ex** facilitate reasoning
- **Events:** facts

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro