

COMP90042

SUBJECT EXAM REVIEW

PREPROCESSING

- ▶ Sentence segmentation
- ▶ Tokenization
- ▶ Word normalization
 - ▶ Derivational vs. inflectional morphology
 - ▶ Lemmatisation vs. stemming
- ▶ Stop words

TEXT CLASSIFICATION

- ▶ Building a classification system
- ▶ Evaluation metrics
- ▶ Algorithms
- ▶ Text classification tasks

LEXICAL SEMANTICS

- ▶ Lexical relationships (*-nyms*)
- ▶ Structure of WordNet
- ▶ Similarity metrics
- ▶ Approaches to Word Sense Disambiguation

DISTRIBUTIONAL SEMANTICS

- ▶ Matrices for distributional semantics
- ▶ Association measures
 - ▶ Calculating (P)PMI from a co-occurrence matrix
- ▶ Count-based models
 - ▶ Basics of singular value decomposition (SVD)
- ▶ Predict-based models
 - ▶ Skip-gram, CBOW
- ▶ Cosine similarity

PART OF SPEECH TAGGING

- ▶ English parts-of-speech
- ▶ Tagsets
 - ▶ **not:** fine-grained tags of any particular tagset
- ▶ Approaches

SEQUENCE MODELS FOR TAGGING

- ▶ Markov Models vs Hidden Markov Model
 - ▶ mathematical formulation of HMM, assumptions
- ▶ Training on fully observed data, e.g., tagging
- ▶ Viterbi algorithm
- ▶ Unsupervised HMMs: hard EM / soft EM
 - ▶ Forward-backward

INFORMATION EXTRACTION

- ▶ Named entity recognition
 - ▶ Models
 - ▶ Tagging formalisms (BIO)
- ▶ Relation extraction:
 - ▶ How to frame the problem using binary and multi-class classifiers
- ▶ Differences between supervised models and OpenIE.

TOPIC MODELS

- ▶ Differences between text classification and topic modelling
- ▶ Differences between LDA and HMMs
- ▶ Applications

CONTEXT-FREE GRAMMARS

- ▶ Basic syntax of English
 - ▶ **not:** detailed grammar (see Q9 from 2017)
- ▶ The context-free grammar formalism
- ▶ Parsing
 - ▶ CYK algorithm

PROB. CFGS

- ▶ Ambiguity in grammars
- ▶ Probabilistic context free grammars: rules, generative process, probability of a tree
- ▶ PCYK algorithm for parsing
- ▶ Comparing to Viterbi and other ‘decoding’ methods

DEPENDENCY GRAMMAR

- ▶ Notion of dependency between words
- ▶ Dependency grammars and dependency parse trees
 - ▶ Projectivity vs non-projectivity
 - ▶ Transition based parsing algorithm
- ▶ **not:** graph based parsing
- ▶ **not:** detailed dependency edge inventory

N-GRAM LANGUAGE MODELS

- ▶ Derivation
- ▶ Smoothing techniques
 - ▶ Add- k
 - ▶ Interpolation vs. backoff
 - ▶ Absolute discounting
 - ▶ **not:** Kneser-Ney, continuation counts etc.
- ▶ Perplexity

RNN LANGUAGE MODELS

- ▶ Basics of neural network structure
- ▶ How to frame LM as a word-by-word classification task
 - ▶ feed-forward classifiers vs recurrent neural networks
- ▶ Links to seq2seq as used in MT, and classifiers used for other NLP tasks
- ▶ **not:** mathematical details of formulation

QUESTION ANSWERING

- ▶ Major approaches
- ▶ Information Retrieval QA pipeline
 - ▶ Passage retrieval
 - ▶ Answer extraction

INFORMATION RETRIEVAL FOUNDATIONS

- ▶ “Information need”
- ▶ TF*IDF weighting, components
 - ▶ Cosine similarity
- ▶ Efficient indexing
- ▶ Querying algorithm

IR INDEXING AND QUERYING

- ▶ Posting list compression
 - ▶ Use of gaps between document ids
 - ▶ vbyte encoding
 - ▶ opt-pfor-delta encoding
- ▶ WAND algorithm
- ▶ Index construction: static vs incremental
- ▶ Phrase search
 - ▶ positional index (intersection, extra information etc.)
 - ▶ **NOT** suffix array

IR QUERYING, EVALUATION AND L2R

- ▶ Query completion
 - ▶ trie+RMQ algorithm
 - ▶ Motivation, Data sources
- ▶ Relevance feedback (why, types)
- ▶ Evaluation methods
 - ▶ precision @ k, (Mean)AveragePrecision, RBP
 - ▶ research test collections
- ▶ Reranking IR system outputs using learned classifier

MACHINE TRANSLATION

- ▶ Motivation
- ▶ Word alignment with IBM model 1
 - ▶ **not:** mathematical derivation of alignment posterior
- ▶ Phrase based model; stack decoding algorithm
- ▶ Sequence to sequence model
 - ▶ **not:** mathematical formulation
- ▶ Evaluation
 - ▶ manual evaluation
 - ▶ automatic evaluation with BLEU

EXAM STRUCTURE

- ▶ Worth 50 marks
- ▶ Parts:
 - ▶ A: short answer [15]
 - ▶ B: method questions [17]
 - ▶ C: algorithm questions [10]
 - ▶ D: short essay [8]
- ▶ 2 hours in duration
 - ... 2 minutes 24 seconds / mark

SHORT ANSWER

- ▶ Several short questions
 - ▶ 1-2 sentence answers for each
 - ▶ 1 mark per question
- ▶ Often
 - ▶ definitional, e.g., *what is X?*
 - ▶ conceptual, e.g., *relate X and Y? What is the purpose of Z?*
 - ▶ may call for an example illustrating a technique/problem

METHOD QUESTIONS

- ▶ Longer answer
 - ▶ larger questions 5-7 marks each
 - ▶ broken down into parts
- ▶ Focus on analysis and understanding, e.g.,
 - ▶ contrast different methods
 - ▶ outline or analyze an algorithm
 - ▶ motivate a modelling technique
 - ▶ explain or derive mathematical equation

ALGORITHMIC QUESTIONS

- ▶ Perform algorithmic computations
 - ▶ numerical computations for algorithm on some given example data
 - ▶ present an outline of an algorithm on your own example
- ▶ 2 questions, each worth 4-6 marks.
- ▶ You won't be required to simplify maths, i.e., you can leave things as fractions

ESSAY QUESTION (8 MARKS)

- ▶ Expect to write 1 page
- ▶ Several broad topics in WSTA given, you should select **one**
 - ▶ no marks given for attempting many
- ▶ Provide
 - ▶ Definition and motivation
 - ▶ Relation to multiple tasks discussed in the class
 - ▶ Compare/contrast use across these tasks

WHAT TO EXPECT

- ▶ Even coverage of topic from the semester
- ▶ Be prepared for concepts that have not yet been assessed by homework / project
- ▶ Guest lectures are *fair game*
- ▶ Prescribed reading is *fair game*