

Machine learning

Tokenization

NLP
language
Model

Chat-GPT

RegEx

Second Year Project /

Lecture 1

Welcome ❤

OSEMN

Data Science is Osemn (Awesome) •
Obtain, Scrub, Explore, Model & Interpret
A lot of data is unstructured! (free text)

HPC

High Performance Computer at ITU is good for intense tasks, which might be beneficial for this course

Project

Hand-in a week early for feedback!

One of:

1. Named entity recognition NER the easier
2. Relation extraction

NER

RE

Determine entities in text: Person, Organisation & Relations between 2 entities (Ex "used for")

Introduction

Semantics

The study of meaning in a word/sentence
Can often be very ambiguous

Ambiguity

Lexical a word has several meanings
Syntactic relational; who does the action?
Semantic ambiguity of meaning!
Pragmatics contextual ambiguity

Tokenization

We divide a string into its parts!
i.e. words etc.

RegEx in python

import re

re.compile('<pattern>').split('<text>') splits on pattern

.findall('<text>') finds all matching patterns

Rob's Opinion: find words is risky because input disappears!

Lecture 2

Architecture

Command Line

<Command> -<option> <value> <input>

Useful

cat print contents

head/tail (-n 2) prints first/last 2 lines

wc counts words (Line Words Chars)

less interactive reading (<Search> for search)

| do several actions in a row

> write to file (

grep <regexp>

-i ignore capitalisation -o match only regexp (not line)

-n return line numbers -c count lines that match

-v invert

sort

-n numeric -r reverse

uniq remove adjacent duplicates

-c count number of words

cut/paste extract columns/paste columns

-d -f

tr replaces characters

sed "s;<FIND>;<REPLACE>;g"

-i in-place edit

Experimental Setup

Data Splits



Train → Dev (Same as Train → Val)

We might want to avoid f.ex. same user in both data sets, since we want to optimize for new users! Time, Document, Domain, writer

We overfit model to training data

We overfit design decisions to validation data

NOT k-fold or train-val (These are standard splits)

instead we (1) shuffle (2) split ... This can be beneficial, BUT NOT BEST PRACTICE (stratified is better)

why? due to stratified ensuring small correlation between train & test data

Baselines

Metrics

Choose a metric that reflects the goal

Common Metrics: Accuracy, Recall, Precision, F1

Balanced Classes
 ↓
 classes
 ↓
 classes
 ↓
 classes
 ↓
 classes
 can be macro-averaged

Reporting

Report only answering of research question

- hyperparameters in appendix

Results should be similar with similar setups

Results should be the same with same code/data

Statistical Significance

$$H_0: M(A) - M(B) = 0$$

$$H_A: M(A) - M(B) > 0$$

we want to prove $M(A) > M(B)$ by one-sided hypothesis test for $H_{A, \text{one-sided}}$

Bootstrap ; how surprising is it, that model A

outperforms model B?

(1) bootstrap $n(10,000)$ random datasets

(2) empirically create distribution for metric (diff. in metric A-B)

(3) calculate p-value for observed metric for A

ASO ; Almost Stochastic Order (see implementation in slides)

Ground Truth

Collection

Top-down: what data do we want?

Bottom-up: what data can we get?

Sample should be representative!!!

Be aware of bias & limitations

Sampling bias

Non-response bias

Annotation

→ Guideline → annotate → evaluate →

Gold standard → ground truth annotation
how well do humans mimic gold standard?

There is a tradeoff between having guidelines that produce similar results for a variety of annotators and correctness!

Guidelines goal, labels, what is annotated
& difficult cases!

P_o : observed agreement

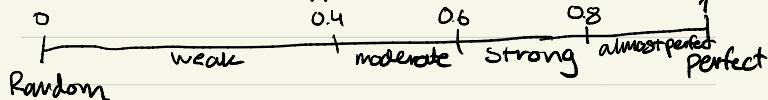
P_e : expected agreement

$$K_b = \frac{P_o - P_e}{1 - P_e}$$

Inter-annotator agreement

Cohen's Kappa: chance-corrected accuracy

Fleiss Kappa: more than 2 annotators



Dataset Statement

9 Q's to answer!

Sequence Labelling

POS Tagging

A type of sequence labelling
We use wordnet to find syntactic types for words... we want to disambiguate ambiguous cases.

Penn Treebank 52 tags

Universal POS tagging 17 tags → works for most languages!

Classes

Closed class classes that are relatively fixed → no additions

Open Class frequently changing / new words are often added!

English Web Treebank weblogs, newsgroups

Noun

simple test; "the" can be put in front!

Verb

refers to someone / something

Pronouns

word used to describe relation

Adposition

NLP Basic Classifiers

Predictor

the function that maps $X \rightarrow Y$
 $Y = f(X) + \epsilon$ $\hat{Y} = f(X)$

Simple NLP tasks

- * Spam or not spam?
- * Sequence Labelling (POS)

Feature types

"High-level" engineered features
"Low-level" directly derived from raw data

Bag-of-Words

We only care what is present, and ignore context!

Reduce Vocabulary

lemmatisation reduce words to base form
stemming approx. lemma by removing endings
sequence of n words in a text.

Generative Bayes Naïve

optimises 0-1 Loss (lower bound error)
Independence between features is assumed, thus $p(x, y) \approx p(y) \prod_n p(x_n | y)$

Conditional humm... Sigmoid

we model $p(y|x)$ directly!
Model p of $y=1$

Sigmoid linear model: $O(wx+b)$ is also Logistic regression

Regularisation

$\lambda \sum_i w_i^2$ (Lasso!) add penalty to loss

Language Models

Language Model

We want to predict the next word in a sentence! $P(\text{word} \mid \text{history})$
useful for spelling corrections / suggestions

Markov Assumption

We assume independence, thus being able to do the following simplification:

N-grams

We consider the last $n-1$ words when predicting the next word in a sentence!

$$P(w \mid h)$$

Baseline

$$= \text{count}(h, w) / \text{count}(h)$$

$P(w) = \frac{1}{\text{vocab}}$ (all words have the same P)

Intrinsic

We "win" if we predict words in the actual corpus!

Extrinsic

Perplexity

$$= P(w_1, \dots, w_N)^{\frac{1}{N}} \quad (\text{inverse } P \text{ of sequence})$$

the lower the better! $\approx \log \# \text{choices}$

Zipf's Law



Words not in vocabulary get 0 probability
few words are frequent, most words are rare!

Problem: Underestimating the P of unseen words

+ Replace unseen words with <OOV> token

Feature engineering → + maybe even categories of OOV tokens?

Out Of Vocabulary

+ Smoothing: move probability mass to OOV tokens

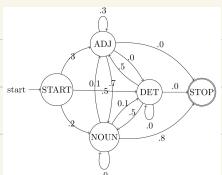
$$+ \text{Add-1 smoothing } P(w \mid h) = \frac{\text{count}(h, w) + 1}{\text{count}(h) + \text{vocab}}$$

Interpolation

Back-off smoothing: we combine our n-gram probabilities with (n-1)-grams...
Linear combinations!

$$P(w|w_1, w_2, w_3) = \beta_1 (P(w|w_1, w_2, w_3)) + \\ \beta_2 (P(w|w_2, w_3)) + \\ \beta_3 (P(w|w_3)) + \beta_4 (P(w))$$

Markov Chains



Hidden Markov Model

Assumes dependency between labels of concurrent words.

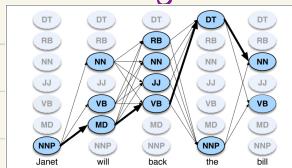
Transition Probability: the P of a label x_2 following the label x_1 . It is a graph!

In real life we do not know the ground truth transition probabilities!

We look at $p(w|t)$
 \uparrow
word tag

Evolution of Decoding

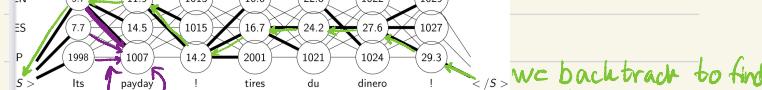
- (1) all labels get the same label
- (2) all tokens individually get most likely label
- (3) Viterbi Algorithm find most likely path!



Follow example
from slides :-)

Probabilities will be:

$$\underbrace{P(\text{history})}_{\text{Previous}} \cdot \underbrace{p(\text{word}|\text{label})}_{\text{emission}} \cdot \underbrace{P(\text{label}|\text{label}_m)}_{\text{transition}}$$



this number is the minimum negative log-probabilities of these edges... the THICK line is the chosen path!

we backtrack to find
most likely sequence

Words in Neural Nets

Deep Learning

The future! (Describes NNs that are deep)

Word Representation

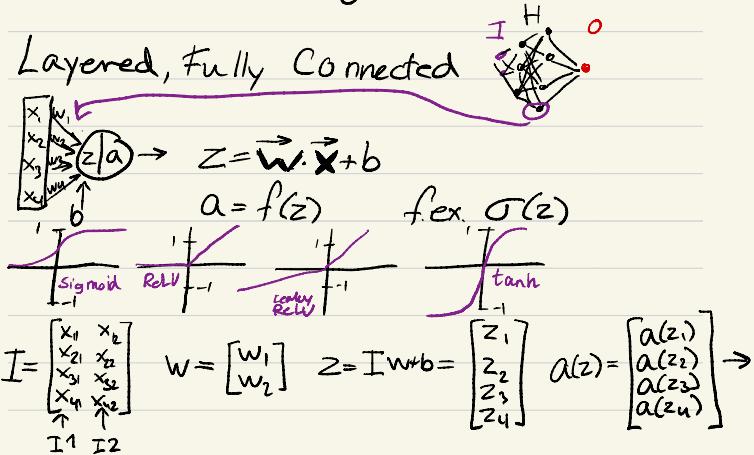
Instead of Binary (Bag-Of-Words) with CountVectorizer, we can represent words as vectors consisting of d floats

FFNN
A neuron

Activations

With words

Layered, Fully Connected



Word Embedding

How to Standardize

Save vector of floats for each word in matrix E . Dense Continuous features. NNs need same shape of input, but we might have different lengths...

We aggregate! *f.ex average!*

How to find w -vectors?

Right now,
next week; new
method!

We initialize E at random, and update the floats like all other weights in the network; values will shift depending on the task.

Is like putting Viterbi on top of a neural network, BUT! without emission... or something :)

Emission probabilities come from neural net!

CRF Layer

Vector Semantics

Representing Meaning

In NLP we use things like WordNet where we have synsets;

Collections of synonymous in meaning

We find the higher-order word describing the meaning of the word.

BUT! Thesaurus is subjective and human maintained

Sparse

One-hot-Encoding of features
Individual word-vectors will be orthogonal even when the meaning is similar!

Cosine Similarity

$\cos_{uv} = \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|}$ we map the cosine similarity to the meaning; But how in practice?

Dense

we encode features of a word, instead of the word itself
The dimensionality d is a hyperparameter

Co-occurrence

Matrix: $M \times M$
symmetric

What other words does word X appear with
hyperparameter: window size around word
Normalize PMI tfidf mutual information

Word2Vec

Methods

Framework for learning context-based word probabilities... adjust word vectors to max. P

CBOW Sum context \rightarrow output word

Skipgram word \rightarrow project context

We do not care about order, just about words being close to each other! (within window)

RNNs & neural LM

Language is discrete; There is no continuous path between 2 entities

N-gram LM

$$\begin{aligned} p(w_1, w_2, \dots, w_N) &= p(w_1) \cdot p(w_2|w_1) \cdot p(w_3|w_2, w_1) \cdot \\ &\quad p(w_4|w_3, w_2, w_1) \cdots p(w_N|w_{N-1}, \dots, w_2, w_1) \\ &= p(w_1) \prod_{i=1}^N p(w_i|w_1, \dots, w_{i-1}) \end{aligned}$$

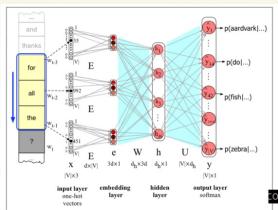
probability of word sequence

Markov assumption: we assume independence from everything but last $(N-1)$ words \approx P

Smoothing due to sparsity (Zipf's law)

Kneser-Ney, Witten-Bell, Laplace

Neural LM

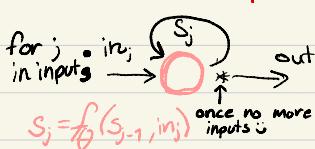


FFNN \Rightarrow continuous-space word embeddings
 \Rightarrow similar idea to N-gram, but NN

$$\begin{aligned} \rightarrow & \text{Input: } x_t \in \{0, 1\}^{|V|} \text{ (one-hot encoding)} \\ & \text{Output: } \hat{y}_t \in (0, 1)^{|V|} \text{ (probability distribution)} \\ & e = [Ex_{t-3}; Ex_{t-2}; Ex_{t-1}] \\ & h = \sigma(We + b) \\ & z = Uh \\ & \hat{y}_t = \text{softmax}(z) \end{aligned}$$

Recurrent Neural Nets

2 components



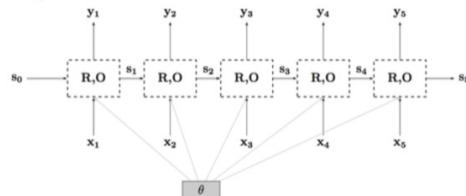
Can process unbounded input, and retains internal state across steps \circlearrowleft

memory s Memory of state

function f update function for state

The cell that takes input, and is run as a cycle until input is processed

Unrolling over time



Advanced RNNs

RNN math

Back to RNNs... gimme dimensions ❤

$$x_t = \mathbf{E} w_t$$

$$h_t = g(\mathbf{V} x_t + \mathbf{U} h_{t-1})$$

$$\hat{y}_t = \mathbf{W} h_t$$

$$\mathbf{E} \in \mathbb{R}^{|\mathcal{V}| \times V}$$

$$\mathbf{V} \in \mathbb{R}^{h \times i}$$

$$\mathbf{W} \in \mathbb{R}^{|V| \times h}$$

$$\mathbf{U} \in \mathbb{R}^{h \times h}$$

Observation:
Matrix $\in \mathbb{R}^{M \times N}$

$|V|$: vocabulary size

e : embedding size

h : hidden layer size

can map a vector
from dim. $\mathbb{R}^N \rightarrow \mathbb{R}^M$

Use cases

Sequence tagging/transducer, Classification/acceptor
Sequence generator (auto-regressive model /LM)

Regularization

How do we avoid over fitting a RNN?

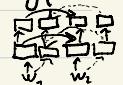
Drop-out by some probability p we drop a feature completely (set it to 0), ONLY during training (at test time we rescale by $1-p$)

Deepening

It can be beneficial to know the future!

Bi-directional RNN's we combine 

Will improve performance ☺

Stacked RNN's the hidden states of 1 layer is used in the next; 

Gradients

During backpropagation we might encounter vanishing or exploding gradients... a very real problem for deep models ☹

GRU

Gated recurrent unit (keep track of memory and hidden states separately!) at every step we consider overwriting memory &

LSTM

Long Short-Term Memory

Contextualized LMs

Contextualized Embeddings

Attention is a model that uses a weighted sum of recurrent hidden states
But! we are missing context... a word might have several meanings

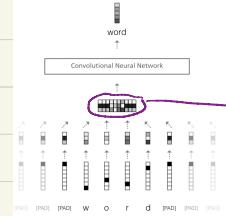
ELMo

Embeddings from Language Models

💡 each token representation is dependent on the entire input sentence

ELMo's Diet: characters 🍎 (in this paper at least)

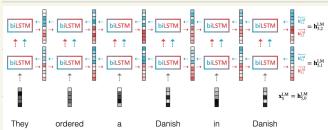
- (1) assign each character an embedding
- (2) CNN (filtering (several different size) pooling)
- (3) Feed Forward with convolutional output
- (4) Output new dense word-representation



How to capture context?

ELMo uses Long Short-Term Memory

We use a bi-LSTM ^{twice!} and combine each direction's hidden states for a token as the new contextualized embedding



Training ELMo!

Goal: predict word on basis of context

Assumption: tokens always appear in sequence

BERT

4 steps

Normalization

Tokenization

Subword Segmentation

subword examples:

{a, b, c, ab, ac, bc, abc}

Bidirectional Encoder Representations from Transformers

uses NFD to normalize input - irreversible

as defined earlier in course

look at how likely characters are to appear together, merge most likely characters together into "new" character (f.ex. 'th')

The goal might be to locate morphemes

Special Tokens

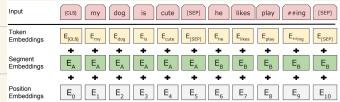
[SEP] - end of sentence/segment

[CLS] - beginning of document

Input

we sum embeddings on 3 levels;

- Subword/Token
- Segment
- Position



BERT Encoder

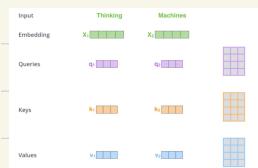
Self-attention

12 transformers (heads) each has;

Query: Learn what to focus on from current word

Key: How important as context

Value: Representation of subword (used as output)



$$\begin{aligned} q_i &= W^Q x_i \\ k_i &= W^K x_i \\ v_i &= W^V x_i \end{aligned}$$

Multi-head

Instead of a single self-attention layer, we give the input to n heads, then flatten the outputs \rightarrow hopefully learning different things

ROBERTA

Larger batches & learning rates + more data

Bias in NLP

Bias

is often un-intended but not necessarily un-beneficial... For some tasks it's good!

Disparity

Outcome outcomes are biased e.g. women in kids
men on skin's

Error

error occurrence is biased

Sources

Data who/where does data come from?

Annotation who annotates/validates data?

Embeddings are their embeddings encoding bias?

Models is the model prone to a specific error?

Designer what design decisions were made?

Data

Include demographic information in encoding embeddings or resample to be more representative!

Annotation

Train and evaluate annotators. Realize and support that some things are ambiguous!
We can try to do interventions but there is no magic/consistently good fix!!

Embeddings

Models often amplify problematic training data...

Models

Rob's writing tips!

Anatomy

Abstract, Introduction

Related work, Data, Analysis, Method, Results
Experiments, Results, Discussion, Related work

Conclusion, References

Structure

Funnel-like structure 

Good ideas

Abstract write as the very last thing!
Introduction task description → end on main idea
start with an example!

Related Work describe relation and relevant methodology!

Methodology whole pipeline + motivation for choices

Experiments Data + Setup + Results

Results Compare to 3 baselines!

Analysis what did we improve on? why? how?

Project

27/3 Baseline predictions

29/3 or 3/4 project presentation

12/4 project proposal

19/5 draft hand-in

26/5 final hand-in