

DEEP LEARNING For Natural Language Processing

Lecture 1: introduction, deep learning, language modeling, CNN Caio Corro

ABOUT THE COURSE

Other (deep learning for) NLP courses

There are many very good courses available online, see the resource page on my website.

- ► Good news for you: plenty of resources to learn! :)
- ► Bad news for me: how can my course be useful for you? :(

How I build this course

- I often observe many misconception in the NLP community (for example about attention, named entity recognition and BIO-tagging, CRFs, ...)
- 2. I built the course around a selection of these topics so you don't make the same mistakes

ABOUT THE COURSE

Other (deep learning for) NLP courses

There are many very good courses available online, see the resource page on my website.

- ► Good news for you: plenty of resources to learn! :)
- ► Bad news for me: how can my course be useful for you? :(

How I build this course

- I often observe many misconception in the NLP community (for example about attention, named entity recognition and BIO-tagging, CRFs, ...)
- 2. I built the course around a selection of these topics so you don't make the same mistakes

Outline

- Study a few important neural architectures in details in order to understand the building blocks of most of neural networks used in NLP
- Introduction to structured prediction
- ► Learn « advanced » PyTorch tricks that are useful in NLP (and other fields)

ABOUT THE COURSE

Grading scheme

► 100%: lab exercises, 2-3 students per group

Lab exercises

- ► Lab 1: convolution for text classification (not graded)
- ► Lab 2: language modeling
- ► Lab 3: part-of-speech tagging

How will I grad projects?

- Experimental results are important, but the grade will be unrelated to how good your model works
- ► Quality of data analysis and result analysis
- ► Description of the model and what you explored (i.e. math, figures, etc)
- Explanation of your implementation (in the report!)

NATURAL LANGUAGE PROCESSING

▶ ...

Applications

- ► Machine Translation
- Information Retrieval
- Sentiment Analysis
- Automatic Summarization
- ► Human-Computer Interaction

- Data type
- Structured: language is sequential
- Discrete: text processing
- Continuous: speech processing

Difficulties

- ► Many languages (> 6000)
- ► Ambiguity

We will focus on text, even if many languages are unwritten!



- Constantly evolving
- ► Noisy

► Programming languages: prescriptive grammar \Rightarrow easy to build a rule-based system

Natural languages: comes in many shapes and forms!

► Programming languages: prescriptive grammar \Rightarrow easy to build a rule-based system

Natural languages: comes in many shapes and forms!

Gender-inclusive language

Les étudiants sont en grève. Les étudiant(e)s sont en grève. Les étudiant.e.s sont en grève. Les étudiant.e.s sont en grève.

Ils sont en grève.

Iels sont en grève.

- ► Programming languages: prescriptive grammar \Rightarrow easy to build a rule-based system
- ► Natural languages: comes in many shapes and forms!



Boomers and « l'Académie Française » are apparently not capable to adapt to variation and diversity in languages, but you want your software to be smarter than them!

Problem: Named Entity Recognition (NER)

Task: given a large corpus a text files, extract all named entities:

- ► Person names
- ► Place names
- ► Dates

. . .

► Locations

After Sherborne, Turing studied as an undergraduate from 1931 to 1934 at King's College, Cambridge, where he was awarded first-class honours in mathematics. In 1935, at the age of 22, he was elected a fellow of King's on the strength of a dissertation in which he proved the central limit theorem. Unknown to the committee, the theorem had already been proven, in 1922, by Jarl Waldemar Lindeberg. A blue plaque at the college was unveiled on the centenary of his birth on 23 June 2012 and is now installed at the college's Keynes Building on King's Parade

Potential tags: LOCATION ORGANIZATION DATE MONEY PERSON PERCENT TIME

Extract of <u>https://en.wikipedia.org/wiki/Alan_Turing</u> tagged with <u>http://nlp.stanford.edu:8080/ner/</u>

Easy solution

Use a gazetteer:

Class	Gazetteer excerpt
Locations	sunderland, sundsvall, sundsvl, sunnyvale, sunrise, sunset beach, sunshine coast, suourland, suouroy
Organizations	xyratex, y e data, ya, yadkin, yahoo, yahoo maps, yahoo personal email, yahoo personals, yahoo sbc
Person Names	jerrine, jerrod, jerrol duane, jerrold, jerrome, jerry, jerry ben, jerry brown, jerry claude, jerry coleman
Other	january, february, sunday, monday, toyota prius, ford focus, apple ipod, microsoft zune, tylenol, advil

Gazetteer from [Carlson et al., 2009]

Downsides

- ► Language is ambiguous, e.g. « Tim Cook »
- ► Language is evolving, e.g. new people, new places, ...
- ► Typos, « non-standard » writing (e.g. tweets)

In general for natural language processing

- ► More or less intuitive for a human who speaks a given language
- Fuzzy decision, lots of « weak » contradictory signals

Even for syntactic analysis!

« The spokesperson, <u>Μιχάλης Χατζόπουλος</u>, has declared that... »

SENTENCE CLASSIFICATION

Sentiment analysis

- ► Input: sentence
- > Output: Positive? Neutral? Negative?

Inputs are of different length!

This movie is great! I saw that movie some years ago. The food was disgusting...

SENTENCE CLASSIFICATION

Sentiment analysis

- ► Input: sentence
- > Output: Positive? Neutral? Negative?

Inputs are of different length!

This movie is great! I saw that movie some years ago. The food was disgusting...

Natural language inference

- ► Input: premise and hypothesis (2 sentences)
- Output: Entailment? Neutral? Contradiction?

John likes Baltimore a lot. Johne likes Baltimore.

TAGGING

Part-of-speech tagging

- ► Input: sentence
- Output: grammatical category for each word of the sentence

PRP	VB	DET	NN
They	walk	the	dog

TAGGING

Part-of-speech tagging

- ► Input: sentence
- Output: grammatical category for each word of the sentence

PRP	VB	DET	NN
They	walk	the	dog

Named entity recognition with BIO tags

- ► Input: sentence
- ► Output: BIO tags + 6 classes

(person, location, group, creative work, product, corporation)

B-Per	I-Per	0	0	B-Loc
Neil	Armstrong	visited	the	moon

CHUNKING

Named entity recognition

- ► Input: sentence
- Output: chunks (person, location, group, creative work, product, corporation)



CHUNKING

Named entity recognition

- ► Input: sentence
- Output: chunks (person, location, group, creative work, product, corporation)



Nested named entity recognition

- ► Input: sentence
- Output: nested chunks



SEQUENCE GENERATION

Machine translation

- ► Input: sentence in a source language
- Output: sentence in a target language

They walk the dog



SEQUENCE GENERATION

Machine translation

- ► Input: sentence in a source language
- Output: sentence in a target language

They walk the dog

Image captioning

- ► Input: image
- Output: sentence describing the image





 \Rightarrow

 \Rightarrow

Spongebob is cooking burgers

SYNTACTIC PARSING

Syntactic dependency parsing

- ► Input: sentence
- Output: bi-lexical dependencies between words



SYNTACTIC PARSING

Syntactic dependency parsing

- ► Input: sentence
- Output: bi-lexical dependencies between words



Constituency parsing

- ► Input: sentence
- Output: hierarchical phrase-structure



SEMANTIC PARSING

SQL parsing

- ► Input: sentence
- ► Output: SQL query

I want to book a flight from Paris to Rome.

 \downarrow

SELECT * FROM flight WHERE from = "paris" AND to = "rome"

SEMANTIC PARSING

SQL parsing

- ► Input: sentence
- ► Output: SQL query

I want to book a flight from Paris to Rome.

 $\downarrow \downarrow$

SELECT * FROM flight WHERE from = "paris" AND to = "rome"

Abstract Meaning Representation (AMR) parsing

- ► Input: sentence
- ► Output: graph

The boy want to go.



DEEP LEARNING

Neural networks

- ► Input: « raw » features, e.g. an image, a sentence
- Hidden layers: sequence of non-linear transformations
- Output: linear classification/regression

Parameter estimation cookbook

- ► Loss function minimization
- Stochastic gradient descent
- ► Regularization (weight decay, dropout, ...)



DEEP LEARNING FOR NLP

Neural architectures

- ► How to efficiently represent text data?
- ► How to design the layers of the network?
- ► How to output structured representations?

This is what what makes deep learning for NLP interesting!

Training

- Annotation in NLP is expensive: unsupervised pre-training is essential
- ► Structured outputs are difficult to learn: exponential search space

Problematic even for sequence generation

BACKGROUND: MACHINE LEARNING

RECALL

The « old school » machine learning pipeline



Example of classifiers

- ► Decision Tree:
 - ► Make a decision considering a limited number of features
 - ► Use conjunction of features to make a prediction
- ► K-nearest neighbors:
 - ► All features are used and considered equals
- ► Linear models:
 - ► Weight features so they are more or less important to make a decision

BINARY CLASSIFICATION

Data

- ► Input: $\mathbf{x} \in \mathbb{R}^d$
- ► Output: $y \in \{0,1\}$ or $y \in \{-1,1\}$

Parameters

- ► Projection vector: $\mathbf{a} \in \mathbb{R}^d$
- ▶ Intercept: $b \in \mathbb{R}$

BINARY CLASSIFICATION

Data

- ► Input: $\mathbf{x} \in \mathbb{R}^d$
- ► Output: $y \in \{0,1\}$ or $y \in \{-1,1\}$

Parameters

- ► Projection vector: $\mathbf{a} \in \mathbb{R}^d$
- ▶ Intercept: $b \in \mathbb{R}$

Decision

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{a}^{\mathsf{T}} \mathbf{x} + b \ge 0, \\ 0 & \text{otherwise.} \end{cases} \quad \text{or} \quad f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{a}^{\mathsf{T}} \mathbf{x} + b \ge 0, \\ -1 & \text{otherwise.} \end{cases}$$

BINARY CLASSIFICATION

Data

- ► Input: $\mathbf{x} \in \mathbb{R}^d$
- ► Output: $y \in \{0,1\}$ or $y \in \{-1,1\}$

Parameters

- ► Projection vector: $\mathbf{a} \in \mathbb{R}^d$
- ▶ Intercept: $b \in \mathbb{R}$

Decision

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{a}^{\mathsf{T}} \mathbf{x} + b \ge 0, \\ 0 & \text{otherwise.} \end{cases} \quad \text{or} \quad f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{a}^{\mathsf{T}} \mathbf{x} + b \ge 0, \\ -1 & \text{otherwise.} \end{cases}$$

Probabilistic output

$$p(y = 1 | \mathbf{x}) = \sigma(\mathbf{a}^{\mathsf{T}}\mathbf{x} + b) = \frac{\exp(\mathbf{a}^{\mathsf{T}}\mathbf{x} + b)}{1 + \exp(\mathbf{a}^{\mathsf{T}}\mathbf{x} + b)}$$

 $p(y = 0 | \mathbf{x}) = 1 - \sigma(\mathbf{a}^{\mathsf{T}}\mathbf{x} + b)$

INTUITION



TRAINING 1/2

Training data

Assume access to n iid labeled datapoints: $D = \left\{ \langle \mathbf{x}^{(i)}, y^{(i)} \rangle \right\}_{i=1}^{n}$

Loss function

Function defined as $\ell : \{0,1\} \times \mathbb{R}^d \to \mathbb{R}_+$

- ► Perceptron loss:
 - ► Binary hinge loss:
 - ► Binary negative log-likelihood: $\ell(y, w) = -w \times y + \log(1 + \exp(w))$

Training objective

Regularization term

$$\min_{\mathbf{a} \in \mathbb{R}^{d}, \ b \in \mathbb{R}} \frac{1}{n} \sum_{\langle \mathbf{x}, y \rangle \in D} \ell(y, \mathbf{a}^{\mathsf{T}} \mathbf{x} + b) + \alpha \times r(\mathbf{a})$$

Actually not a set

$$\ell(y, w) = \max(0, 1 - w \times (2y - 1))$$

$$\ell(y, w) = \max(0, -w \times (2y - 1))$$

$$\ell(y, w) = \max(0, 1 - w \times (2y - 1))$$

$$\ell(y, w) = \max(0, -w \times (2y - 1))$$
$$\ell(y, w) = \max(0, 1 - w \times (2y - 1))$$

$$f(y, w) = \max(0, 1 - w \times (2y - 1))$$

$$f(y,w) = \max(0, 1 - w \times (2y - 1))$$

$$\ell(y, w) = \max(0, 1 - w \times (2y - 1))$$

$$\ell(y, w) = \max(0, -w \times (2y - 1))$$

$$\ell(y, w) = \max(0, 1 - w \times (2y - 1))$$

$$\ell(y, w) = \max(0, -w \times (2y - 1))$$

TRAINING 2/2

Training objective

Training algorithms

► Specialized solution

Simple solution: (sub-)gradient descent

 $\min_{\mathbf{a} \in \mathbb{R}^{d}, \ b \in \mathbb{R}} \frac{1}{n} \sum_{\langle \mathbf{x}, y \rangle \in D} \ell(y, \mathbf{a}^{\mathsf{T}} \mathbf{x} + b) + \alpha \times r(\mathbf{a})$

Regularization term

FIRST-ORDER METHODS IN OPTIMIZATION

Foundations and Trends[®] in Machine Learning

Optimization with Sparsity-Inducing Penalties

Francis Bach, Rodolphe Jenatton, Julien Mairal and Guillaume Obozinski **Amir Beck**

Can be way faster!

(depending on the choice of loss and regularization)

nou

(SUB-)GRADIENT DESCENT

$$\theta^{(t+1)} = \theta^{(t)} - \eta \times \nabla_{\theta} g(\theta)$$

How to choose the step-size?

- ► Fixed step size
- Diminishing step size
- Line-search (search for the best step-size at a given point)

Convergence

- ► Theoretical guarantee
- ► Fast and efficient in practice

For linear models! Neural networks optimization relies on different heuristics

MULTICLASS CLASSIFICATION WITH K CLASSES

Data

- ► Input: $\mathbf{x} \in \mathbb{R}^d$
- ► Output: $\mathbf{y} \in E(k)$ or $y \in \{1...k\}$

Set of one-hot vectors of size k

Parameters

- ► Projection vector: $\mathbf{A} \in \mathbb{R}^{k \times d}$
- ► Intercept: $\mathbf{b} \in \mathbb{R}^k$

MULTICLASS CLASSIFICATION WITH K CLASSES

Data

- ► Input: $\mathbf{x} \in \mathbb{R}^d$
- ► Output: $\mathbf{y} \in E(k)$ or $y \in \{1...k\}$

Set of one-hot vectors of size k

Parameters

► Projection vector: $\mathbf{A} \in \mathbb{R}^{k \times d}$

► Intercept:
$$\mathbf{b} \in \mathbb{R}^k$$

Decision

$$f(\mathbf{x}) = \underset{\mathbf{y} \in E(k)}{\operatorname{argmax}} \quad \mathbf{y}^{\mathsf{T}}(\mathbf{A}\mathbf{x} + \mathbf{b}) \qquad or$$

$$f(\mathbf{x}) = \underset{y \in \{1...k\}}{\operatorname{argmax}} [\mathbf{A}\mathbf{x} + \mathbf{b}]_{y}$$

MULTICLASS CLASSIFICATION WITH K CLASSES

Set of one-hot vectors of size k

Data

- ► Input: $\mathbf{x} \in \mathbb{R}^d$
- ► Output: $\mathbf{y} \in E(k)$ or $y \in \{1...k\}$

Parameters

► Projection vector: $\mathbf{A} \in \mathbb{R}^{k \times d}$

► Intercept:
$$\mathbf{b} \in \mathbb{R}^k$$

Decision

$$f(\mathbf{x}) = \underset{\mathbf{y} \in E(k)}{\operatorname{argmax}} \quad \mathbf{y}^{\mathsf{T}}(\mathbf{A}\mathbf{x} + \mathbf{b}) \qquad or \qquad f(\mathbf{x}) = \underset{y \in \{1...k\}}{\operatorname{argmax}} \quad [\mathbf{A}\mathbf{x} + \mathbf{b}]_{y}$$

Probabilistic output

$$p(\mathbf{y} | \mathbf{x}) = \frac{\exp(\mathbf{y}^{\mathsf{T}}(\mathbf{A}\mathbf{x} + \mathbf{b}))}{\sum_{\mathbf{y}'} \exp(\mathbf{y}^{\mathsf{T}}(\mathbf{A}\mathbf{x} + \mathbf{b}))}$$
ILLUSTRATION

Problem

- ► Input: features
- Output: 1-in-k prediction
 (e.g. select the part-of-speech tag associated with a given word)

Linear classifier w = Ax + b

- ► Input dim: 3
- ► Output dim: k=4 classes



TRAINING

Training data

Assume access to n iid labeled datapoints: $D = \left\{ \langle \mathbf{x}^{(i)}, \mathbf{y}^{(i)} \rangle \right\}_{i=1}^{n}$

Loss function

Function defined as $\ell : E(k) \times \mathbb{R}^d \to \mathbb{R}_+$

Hinge loss:
$$\ell(\mathbf{y}, \mathbf{w}) = \max(0, -\mathbf{w}^{\mathsf{T}}\mathbf{y} + 1 + \max_{\mathbf{y}' \neq \mathbf{y}} \mathbf{w}^{\mathsf{T}}\mathbf{y}')$$
Negative log-likelihood:
$$\ell(\mathbf{y}, \mathbf{w}) = -\mathbf{w}^{\mathsf{T}}\mathbf{y} + \log \sum_{\mathbf{y}'} \exp \mathbf{w}^{\mathsf{T}}\mathbf{y}'$$

Training objective

A

$$\min_{\boldsymbol{\epsilon} \in \mathbb{R}^{k \times d}, \quad \frac{1}{n} \sum_{\langle \mathbf{x}, \mathbf{y} \rangle \in D} \ell(\mathbf{y}, \mathbf{A}^{\mathsf{T}} \mathbf{x} + \mathbf{b}) \quad + \quad \alpha \times r(\mathbf{A})$$
$$\mathbf{b} \in \mathbb{R}^{k}$$

Actually not a set

Regularization term

LIMITS OF LINEAR CLASSIFIERS

TEXT CLASSIFICATION

Task

Predict if a movie review is good or bad.

« This movie is great! »

« Worst movie ever... »

Features

- 1. Uni-gram features:
 - ► One feature per word in the vocabulary
 - ➤ The feature is equal to one if the word is in the review
- 2. Bi-gram features:
 - ► One feature per couple of words in the vocabulary
 - ► The feature is equal to one if the words appear consecutively in the review

FEATURE SET

 Let V be the vocabulary. Unigram features Number of features: V We can use a dictionary to map a word to an index in a vector 	<i>x</i> =	0 1 1 0 0	bad movie great worst house
 Bigram features ➤ Number of features: V x V ➤ Similarly, use a dict! 	x =	$\begin{bmatrix} 1\\ 0\\ \cdots\\ 1 \end{bmatrix}$	This movie The movie is great

- > Number of features: |V|x|V|
- Similarly, use a dict!

Problem

- ► The feature vector is sparse
- ► The dot product can be unnecessarily expensive!
- ➤ There is no inductive bias: a sentence is not bag-of-words, it is a sequence!

28

- ► Can we always find a hyperplane that separate classes? <u>NO</u>
- ► Can we characterize formally in which cases we can? <u>YES</u>



- ► Can we always find a hyperplane that separate classes? <u>NO</u>
- ► Can we characterize formally in which cases we can? <u>YES</u>



- ► Can we always find a hyperplane that separate classes? <u>NO</u>
- ► Can we characterize formally in which cases we can? <u>YES</u>



- ► Can we always find a hyperplane that separate classes? <u>NO</u>
- ► Can we characterize formally in which cases we can? <u>YES</u>



- ► Can we always find a hyperplane that separate classes? <u>NO</u>
- ► Can we characterize formally in which cases we can? <u>YES</u>



- ► Can we always find a hyperplane that separate classes? <u>NO</u>
- ► Can we characterize formally in which cases we can? <u>YES</u>



A VERY SHORT INTRODUCTION TO DEEP LEARNING

MAIN IDEA

« Latent » hidden representation

- ► Compute a sequence of hidden representations so that classes are linear separable
- Train everything end-to-end via gradient descent!



Neural Network

« Complicated » non convex parameterized function $\ f_{\theta}: \mathcal{X} \to \mathcal{Y}$

MULTILAYER PERCEPTRON 1/2



x : input features
 z⁽ⁱ⁾ : hidden representations
 W : output logits
 θ = {A⁽¹⁾, b⁽¹⁾, ... } : trainable parameters
 σ : piecewise non-linear activation function



OPTIMIZATION

Backpropagation algorithm

- 1. Build computation graph
- 2. Back-propagate to compute the (sub-)gradient

Important: NLP often requires to allow to build dynamic computation graphs!

Beyond gradient descent (GD)

- Stochastic GD: approximate the objective using a random subset of datapoints (i.e. using a minibatch, or even a single datapoint)
- « improved » GD optimization algorithms:
 - ► Momentum

► Adam

▶ ...

In NLP, use Adam as default :)

REPRESENTATION LEARNING: COMPUTER VISION [LEE ET AL., 2009]



REPRESENTATION LEARNING: NATURAL LANGUAGE PROCESSING [VOITA ET AL., 2019]



GENERALIZATION

Overparameterized neural networks

- ➤ Networks where the number of parameters exceed the training dataset size.
- Can learn by heart the dataset,
 i.e. overfit the data -> does not generalize well to unseen data
- ► Are easier to optimize in practice

Monitoring the training process

- Loss should go down
- Training accuracy should go up
- Dev accuracy should go up

Regularization

Techniques to control parameters during learning and prevent overfitting

LEARNING WITH RANDOM INPUTS AND LABELS 1/2 [ZHANG ET AL., 2017]

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception		yes	yes	100.0	89.05
	1,649,402	yes	no	100.0	89.31
		no	yes	100.0	86.03
		no	no	100.0	85.75
(fitting random labels)		no	no	100.0	9.78
Inception w/o BatchNorm	1,649,402	no	yes	100.0	83.00
		no	no	100.0	82.00
(fitting random l	abels)	no	no	100.0	10.12
Alexnet	1,387,786	yes	yes	99.90	81.22
		yes	no	99.82	79.66
		no	yes	100.0	77.36
		no	no	100.0	76.07
(fitting random l	abels)	no	no	99.82	9.86
MLP 3x512	1,735,178	no	yes	100.0	53.35
		no	no	100.0	52.39
(fitting random labels)		no	no	100.0	10.48
MLP 1x512	1,209,866	no	yes	99.80	50.39
		no	no	100.0	50.51
(fitting random l	abels)	no	no	99.34	10.61

LEARNING WITH RANDOM INPUTS AND LABELS 2/2 [ZHANG ET AL., 2017]



DROPOUT 1/4 [HINTON ET AL., 2012, SRIVASTAVA ET AL., 2014]

How does dropout work?

- During training, we randomly "turn off" neurons, i.e. we randomly set elements of hidden layers to 0
- ► During test, we do use the full network



Intuition

- ► prevents co-adaptation between units
- equivalent to averaging different models

Units must be « independent »

DROPOUT 2/4 [HINTON ET AL., 2012]





DROPOUT 3/4 [HINTON ET AL., 2012]

Dropout layer

A dropout layer is parameterized by the probability of "turning off" a neuron

z' = Dropout(z, p = 0.5)

Implementation

 $oldsymbol{z} \in \mathbb{R}^n$: output of a hidden layer

 $p \in [0, 1]$: dropout probability

 $m{m} \in \{0,1\}^n$: mask vector

z': hidden values after dropout application

Forward pass:Backward pass: $m \sim \text{Bernoulli}(1-p)$ $\frac{\partial z'_i}{z_i} = \frac{m}{1-p}$ $z'_i = \frac{z_i * m_i}{1-p}$ \Rightarrow no gradient for
"turned off" neurons.

The mask **m** a vector of booleans stating if neuron is z_i kept ($m_i=1$) or « turned off' ($m_i=0$).

DROPOUT 4/4 [HINTON ET AL., 2012]

Where do you apply dropout?

- ➤ On the input of the neural network
- ► After activation function
- Do not apply dropout on the output logits

Which dropout probability should you use?

- Empirical question: you have to test!
- Dropout probability at different layers can be different (especially input vs. hidden layers)
- ► Usually between 0.1 and 0.5

Dropout variants

Dropout can be applied differently for special neural network architectures (e.g. convolutions, recurrent neural networks)

The language modeling task

Build a model that assigns a probability to any given sentence

 $p(y_1, y_2, \ldots, y_n)$

« Unsupervised » problem

The language modeling task

Build a model that assigns a probability to any given sentence

 $p(y_1, y_2, \ldots, y_n)$

Usefulness

► Text generation

 $y_1 \sim p(y_1)$ $y_2 \sim p(y_2 | y_1)$ $y_3 \sim p(y_3 | y_1, y_2)$...

« Unsupervised »

problem

The language modeling task

Build a model that assigns a probability to any given sentence

 $p(y_1, y_2, \ldots, y_n)$

Usefulness

► Text generation

 $y_1 \sim p(y_1)$ $y_2 \sim p(y_2 | y_1)$ $y_3 \sim p(y_3 | y_1, y_2)$...

➤ Test if a neural architecture is able to learn a good representation of language

p(«The dog is eating») > p(«Is eating dog the»)

« Unsupervised » problem

The language modeling task

Build a model that assigns a probability to any given sentence

 $p(y_1, y_2, \ldots, y_n)$

Usefulness

► Text generation

 $y_1 \sim p(y_1)$ $y_2 \sim p(y_2 | y_1)$ $y_3 \sim p(y_3 | y_1, y_2)$...

➤ Test if a neural architecture is able to learn a good representation of language

p(«The dog is eating») > p(«Is eating dog the»)

 Improve sequence generation problems by using unlabeled data, e.g. machine translation, speech recognition

$$\hat{y}_1...\hat{y}_n = \max_{y_1...y_n} \lambda \times p(y_1...y_n | x_1...x_m) + (1 - \lambda) \times p(y_1...y_n)$$

Translation accuracy

Target language adequacy

44

« Unsupervised »

problem

LANGUAGE MODEL EVALUATION



LANGUAGE MODEL EVALUATION



- Predict next word: suitable words should have a higher probability than incoherent ones
- Evaluate on data: better language model should give a higher probability to the test set

LANGUAGE MODEL EVALUATION



- Predict next word: suitable words should have a higher probability than incoherent ones
- Evaluate on data: better language model should give a higher probability to the test set

Perplexity

- Probability of test sentences normalized by the number of words
- ► The lower the better

$$PP(y_1, ..., y_n) = p(y_1, ..., y_n)^{-\frac{1}{n}}$$

$$= \sqrt[n]{\frac{1}{p(y_1,\ldots,y_n)}}$$

Intuition: longer sentences are less probable so we want to normalize to able to compare sentences of different sizes

N-GRAM MODEL 1/2

Auto-regressive model

The probability of a word is conditioned on previous words only, or, in other words, we assume a generative model that produce words from left to right.

 $p(y_1 \dots y_n) = p(y_1, \dots, y_{n-1})p(y_n | y_1, \dots, y_{n-1})$

N-GRAM MODEL 1/2

Auto-regressive model

The probability of a word is conditioned on previous words only, or, in other words, we assume a generative model that produce words from left to right.

$$p(y_1 \dots y_n) = p(y_1, \dots, y_{n-1})p(y_n | y_1, \dots, y_{n-1})$$

Markov chain

Memory-less auto-regressive model:

► 1st order Markov chain: $p(y_1, ..., y_n) = p(y_1) \prod_{i=2}^n p(y_i | y_{i-1})$

p(``The dog is eating'') = p(``The'')p(``dog'' | ``The'')p(``is'' | ``dog'')p(``eating'' | ``is'')

Bi-gram model

N-GRAM MODEL 1/2

Auto-regressive model

The probability of a word is conditioned on previous words only, or, in other words, we assume a generative model that produce words from left to right.

$$p(y_1 \dots y_n) = p(y_1, \dots, y_{n-1})p(y_n | y_1, \dots, y_{n-1})$$

Markov chain

Memory-less auto-regressive model:

► 1st order Markov chain: $p(y_1, ..., y_n) = p(y_1) \prod_{i=2}^n p(y_i | y_{i-1})$

 $p((The \ dog \ is \ eating \) = p((The \)p((dog \) \ | \ (The \)p((is \) \ | \ (dog \)p((eating \) \ | \ (is \)))$

► 2nd order Markov chain: $p(y_1, ..., y_n) = p(y_1) \quad p(y_2 | y_1) \quad \prod p(y_i | y_{i-1}, y_{i-2})$

 $\sum_{i=3}^{n} p(y_i | y_{i-1}, y_{i-2})$

Tri-gram model

Bi-gram model

N-GRAM MODEL 2/2

1st order Markov chain estimation

 $p((act) | (bcolor the bound to the color to the training data) = \frac{Number of times (bcolor to the training data)}{Number of times (bcolor to the training data)}$

Smoothing

Worse for high order models!

The distribution is really sparse, many sequences of words have a probability of zero.
N-GRAM MODEL 2/2

1st order Markov chain estimation

 $p((act) | (bcolor the bound to the color to the training data) = \frac{Number of times (bcolor to the training data)}{Number of times (bcolor to the training data)}$

Smoothing

Worse for high order models!

The distribution is really sparse, many sequences of words have a probability of zero.

► Introduce fake occurrences

 $p(\text{(a cat } \text{)} | \text{(b the } \text{)}) = \frac{1 + \text{Number of times (cat } \text{appears in the training data}}{1 + \text{Number of times (cat } \text{appears in the training data}}$

N-GRAM MODEL 2/2

1st order Markov chain estimation

 $p((act) | (bcolor the bound to the color to the training data) = \frac{Number of times (bcolor to the training data)}{Number of times (bcolor to the training data)}$

Smoothing

Worse for high order models!

The distribution is really sparse, many sequences of words have a probability of zero.

Introduce fake occurrences

 $p(\text{(cat } \text{)} | \text{(the } \text{)}) = \frac{1 + \text{Number of times (the cat } \text{) appears in the training data}}{1 + \text{Number of times (the } \text{) appears in the training data}}$

► Interpolation

$$\tilde{p}(\text{" cat "} | \text{" the big "}) = \lambda^{(1)} \times p(\text{" cat "} | \text{" the big "}) \qquad \text{Tri-gram} \\ + \lambda^{(2)} \times p(\text{" cat "} | \text{" big "}) \qquad \text{Bi-gram} \\ + \lambda^{(3)} \times p(\text{" cat "}) \qquad \text{Uni-gram} \\ \text{Uni-gram} \end{cases}$$

such that $\lambda^{(1)}, \lambda^{(2)}, \lambda^{(3)} \ge 0$ and $\lambda^{(1)} + \lambda^{(2)} + \lambda^{(3)} = 1$ 47

LIMITATIONS OF N-GRAMS MODELS

Pros

► Easy to estimate

Cons

Unobserved sequences will have small probabilities

The big cat The fat cat « big » and « fat » are semantically similar, therefore we would like both examples to have approx. the same probability even if « fat cat » doest not appear in training data.

Limited context, no long range dependencies



N-GRAM MODEL WITH A MULTILINEAR PERCEPTRON

Task

Predict the next word give a fixed size history.

► Bi-gram model: p((a cat * | a big *)) ► Tri-gram model: p((a cat * | a big *))

Task

Predict the next word give a fixed size history.

► Bi-gram model: p((a cat | a big)) ► Tri-gram model: p((a cat | a big))

Output

It's a classification problem, similar to MNIST digit classification!

Each output class represent a word, i.e. we map word to integers in $V = \{0, 1, ...\}$

~100.000 words

is standard

Task

Predict the next word give a fixed size history.

► Bi-gram model: p((a cat | a big)) ► Tri-gram model: p((a cat | a big))

~100.000 words

is standard

Output

It's a classification problem, similar to MNIST digit classification!

Each output class represent a word, i.e. we map word to integers in $V = \{0, 1, ...\}$



Input

- ► Bi-gram model: 1 word
- ➤ Tri-gram model: 2 words

One hot encoding of words

- ► Vocabulary mapped to a set of integers: $V = \{0, 1, ...\}$
- ► Input words are mapped to a one hot encoding vector

Input

- ► Bi-gram model: 1 word
- ➤ Tri-gram model: 2 words

One hot encoding of words

- ► Vocabulary mapped to a set of integers: $V = \{0, 1, ...\}$
- ► Input words are mapped to a one hot encoding vector



0 1 1 One-hot encoding of the word mapped to 2

Input

- ► Bi-gram model: 1 word
- ➤ Tri-gram model: 2 words

One hot encoding of words

- ► Vocabulary mapped to a set of integers: $V = \{0, 1, ...\}$
- ► Input words are mapped to a one hot encoding vector



Input

- ► Bi-gram model: 1 word
- ➤ Tri-gram model: 2 words

One hot encoding of words

- ► Vocabulary mapped to a set of integers: $V = \{0, 1, ...\}$
- ► Input words are mapped to a one hot encoding vector



LIMITATIONS OF NEURAL N-GRAM MODELS

Pros

► Generalize to unseen n-grams

Cons

- ► More expensive to estimate than the non-neural model
- ► The output layer is expensive to compute!
- Limited context, no long range dependencies

Recurrent networks solve this problem: next week!



WORD EMBEDDINGS

ISSUE WITH WORD AS INPUT

One hot encoding of words

- ► Vocabulary mapped to a set of integers: $V = \{0, 1, ...\}$
- Input words are mapped to a one hot encoding vector



WORD EMBEDDING

Dictionary

- Dictionary: map each word to an integer
- Embedding table: map an integers to a vector

These vectors are parameters that are trained with the rest of the network!

car:	0
red:	1
big:	2



Unknown word

At test time, you may see words that where not present in the training set :(

- 1. Add a special « unknown word » to the dictionary/embedding table
- 2. Learn it by randomly replacing other words with the « unknown word » during training

PYTORCH IMPLEMENTATION

```
emb_table = torch.nn.Embedding(
    num_embeddings = 10,
    embedding_dim = 100
)
```

Construct embedding table

```
words = torch.LongTensor([2, 5])
# x.shape = (2, 100)
x = emb_table(words)
```

Retrieve embeddings of the input

x.shape = (1, 200)
x = x.reshape(1, -1)

Reshape so that it looks like a mini batch of size one with 2 input words

CONVOLUTIONAL NEURAL NETWORK

SEQUENCE REPRESENTATION

Motivation

- ► Build token representation for out-of-vocabulary words (e.g. for tagging)
- Build sentence-level representation for sentence classification

This movie is great! I saw that movie some years ago. The food was disgusting...

To classify these sentences, we need to build a fixed size hidden representation to feed it to a MLP!

This movie is great !

.

.

This movie is great ! Word embeddings







These 2 sentences will have the same representation!

SEQUENCE REPRESENTATION

Motivation

- ► Build token representation for out-of-vocabulary words (e.g. for tagging)
- Build sentence-level representation for sentence classification

This movie is great! I saw that movie some years ago. The food was disgusting...

To classify these sentences, we need to build a fixed size hidden representation to feed it to a MLP!

Convolutional Neural Network

- 1. Convolution over words with a fixed size sliding window
- 2. Pooling operation



























hidden representation

SUMMARY

Convolution

- ► The size of sliding window is an hyper-parameter (2, 3, 4...)
- The hidden representation of convolutions with different sliding window sizes can be concatenated
- ► The sentence can be padded with begin/end of sentence tokens

Pooling and activation functions

- ► Pooling functions: max, min, mean
- ► Activation functions: tanh, relu, sigmoid

Word representation

Word representation can be constructed with a CNN on character embeddings (and possibly concatenated with a word embedding)