

## — Modeling and Transfer —

• How to deal with the scarcity of labeled data? What is the role of language models?

9c.1

1

## Language modeling motivated by problem solving (recap)

### q How to model natural language?

- 2 main approaches
- motivated by 2 types of problems with 2 types of resolutions

### q Grammar

- to answer the **problem of belonging**: • is a given sequence a sentence of language L?

### q Language model

- to answer the **continuation problem**: • what is the most likely continuation of a given sequence of the language L?

9c.2

2

## Probabilistic neural modeling

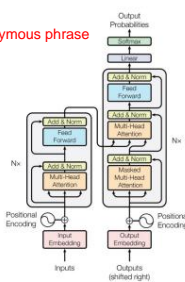
### q Automatic translation

- Transformer originally designed for translation automatic, from a phrase in a source language, obtain a synonymous phrase in a target language

• eg "Cork is extracted from the cork oak." -> "Cork is extracted from the cork oak."

### q Output generation

- Made in time steps
- "Cork is extracted from the cork oak." -> "Cork"
- "Cork is extracted from the cork oak tree. -> Cork" -> "is"
- "Cork is extracted from the cork oak tree. -> Cork is" -> "extracted"



9c.3

3

## Continuation problem

### q Text generation

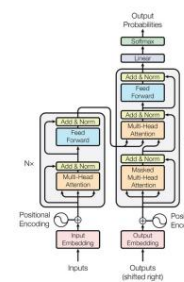
- Transformer can be reused to solve the continuation problem
- From the same language, go from an input segment to the segment that can follow it

• "Cork is extracted" -> "of"

• "Cork is extracted -> from" -> "the"

• "Cork is extracted -> from" -> "cork oak"

• "Cork is extracted -> from the cork oak" -> "."



9c.4

4

## Neural language model

### q Self-supervised learning

• to train a probabilistic neural language model of a given language, the training data are running ("raw") texts of that language

• Simply (automatically) remove the segment that comes next and ask for your prediction from the previous segment

### q Colossal training datasets

• There is no need to annotate the training set! Just collect mass of texts from the web...

• Much larger training datasets than (annotated) ones exist for any supervised task

9c.5

5

## Training (recap)

### q Forward passage

• **multi-class classification**: each item in the vocabulary is a class

• network output layer/vector  $\mathbf{y}^{\wedge}$  with vocabulary size

• predicted item corresponds to the component  $y^{\wedge}[i]$  with the highest value

### q Loss

• Correct output: one-hot vector  $\mathbf{y} = y[1] \dots y[m]$  with 1 component whose index is the index in the vocabulary of the next item, zero in the remaining

• loss function, e.g. cross entropy  $L(\mathbf{y}^{\wedge}, \mathbf{y}) = - \sum_i y[i] \log(y^{\wedge}[i])$

### q Backward passage

• gradient  $L'(\mathbf{y}^{\wedge}) \in L'(\mathbf{y}^{\wedge}[i])=1$  e  $L'(\mathbf{y}^{\wedge}[j])=0$  para  $q \neq j$

9c.6

6

## Language model performance

### q Loss in individual prediction (recap)

• correct output: a one-hot vector  $\mathbf{y} = y[1] \dots y[m]$  with 1 at index  $i$  of the next correct item (in the training data) and zero in the remaining indices

• prediction:  $\mathbf{y}^{\wedge} = y^{\wedge}[1] \dots y^{\wedge}[m]$  normalized by softmax: each component in  $[0,1]$  and add them all to 1

• being cross entropy:  $L(\mathbf{y}^{\wedge}, \mathbf{y}) = - \sum_i y[i] \log(y^{\wedge}[i])$

• individual loss:  $-(0 \times \log(y^{\wedge}[0]) + \dots + 1 \times \log(y^{\wedge}[i]) + \dots + 0 \times \log(y^{\wedge}[n]))$   
 $\log(y^{\wedge}[i])$

### q Model performance (training completed)

• for model tested on  $n+1$  items, average loss:

•  $\bar{y} - \bar{y}_j = 0 \log(y^{\wedge}) / n$

• lower value is better performance: lower loss/error

9c.7

7

## Using language models

### q Text generation

• Answer questions

• Fill out questionnaires

• Create fake news

• Dialogue

• ...

### q Leverage supervised tasks

• with transfer learning

9c.8

8

## Pre-train and tune

### q Pre- train

• First phase of training:

• train the neural model as a language model • with colossal volumes of training data, i.e. text

### q Fine -tuning

• Second training phase: •

continue training the neural model to solve the supervised task of interest • e.g. argument

extraction, emotion detection, inference,... • with annotated training data, i.e. of small volume

9c.9

9

## Linguistic regularities

q Pre-training: the (general) language model • First training phase: • learning "general" linguistic regularities

q Fine-tune: for the relevant (specific) task • Second training phase: • complement with learning of "specific" regularities of the task/application

• called downstream tasks

9c.10

10

## Tuning

### q Global

• all model weights are allowed to adjust during the period of tuning

### q Selective

• An extra layer is added to the model and only the weights of this extra layer are changed during tuning (in the encoder family – more on this later)

• Or in some designated layers the respective weights are left to adjust

9c.11

11

## Dealing with the scarcity of labeled data

### q Better performance

• very advantageous compared to just training for the task of interest, i.e. with only (less) annotated training data

• way to mitigate the scarcity of labeled data

9c.12

12

## Modeling linguistic knowledge

### q Symbolic approach

- Linguistic regularities are (manually) specified in a grammar
- Which solves the problem of belonging and provides representation of meaning based on inference
- Which will serve as a base to be embedded and enable applications

### q Neural approach

- Linguistic regularities are (automatically) captured in a language model
- Which solves the continuation problem and provides a representation of meaning based on vector space
- Which will serve as a basis for applications by complementary tuning of its weights

9c.13

13

## Continuation tasks

### q Next item ("next token")

- given a segment, predict which item should follow • called causal language modeling

9c.14

14

## Continuation tasks (with tokens)

### q Masked Language Modeling ("MLM") • predict the

- words that have been masked or (randomly) replaced in a sentence • e.g. for "I looked at my [MASK] and noticed that [MASK] was late.", predict which words should replace each occurrence of [MASK]

- eg for "I looked at my watch and realized she was delayed.", predict which words should replace those that were exchanged

### q Token deletion

- predict the words that were deleted (randomly) in a sentence

9c.15

15

## Continuation tasks (with sequences)

### q Prediction of the next sentence

- ("NSP – Next Sentence Prediction") • of two given sentences, indicate whether the second one can be the continuation of the first
- eg After "Cork is extracted from the cork oak." what follows? • "Represents 10% of Portuguese exports." • eg After "The birds fly in the sky.." what follows? • "Sailing is a water sport."

### q Deform sentences

- ("sentence warping") • start with a random word from the sentence/document, put the remaining tokens that did not appear at the beginning of the sentence, at the end of the sentence

9c.16

16

## Transfer learning / aprendizagem por transferência

q Tuning, after pre-training  $\tilde{y}$   
particular case of transfer learning

### q Transfer learning

("transfer learning")  $\tilde{y}$

knowledge learned in solving one problem is applied to solving another problem

9c.17

17

## Dealing with data scarcity / lidar com a escassez dos dados

### q Lack of annotated data

$\tilde{y}$  annotated data are scarce, difficult to obtain

$\tilde{y}$  transfer learning is one of the techniques that aims to mitigate this difficulty

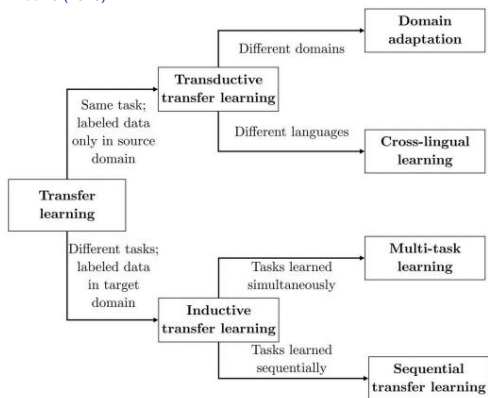
$\tilde{y}$  leveraging data from other tasks for the task of interest

9c.18

18

## Transfer taxonomy / taxonomia de transferência

Windows (2019):



9c.19

19

## Transductive learning by transfer (same task) / aprendizagem por transferência (a mesma tarefa)

### q Domain adaptation $\tilde{y}$ same task, eg

Named Entity Recognition  $\tilde{y}$  more annotated data in texts from one domain, eg news  $\tilde{y}$  but insufficient data for another domain, eg legal text  $\tilde{y}$  apply the model trained in the rich domain to the scarce domain

### q Aprendizagem multilingue ("cross-lingual learning")

$\tilde{y}$  the same task, e.g. POS tagging  $\tilde{y}$  more annotated data in one language, e.g. English  $\tilde{y}$  but insufficient data for another language, e.g. Portuguese  $\tilde{y}$  apply the model trained on English texts to Portuguese

9c.20

20

## Inductive transfer learning (different tasks)

(tarefas diferentes)

q Multi **-task learning** ÿ same language

and domain, eg English and

blogging ÿ but different tasks, eg sentiment analysis

and trend detection

emotions

ÿ merge the two training data sets and do a single model training

q Sequential transfer **learning** ÿ same language and domain, e.g. English and

blogging ÿ but different tasks, e.g. language

modeling and NER ÿ training in two phases: first on one task; then

continue training, focusing on a second task

BRUNNEN

9c.21

21

## - Conclusion -

q Index

ÿ Neural language models

ÿ Pre-train and tune

ÿ Tuning

ÿ Continuation tasks

ÿ Transfer learning

How to adopt progress benchmarks? What variants of

Do Transformers exist? How to assemble the input and output data in  
terms of different tasks?

BRUNNEN

9c.22

22