AI for Document Understanding

Final Presentation – 18.02.2019



TUM-DI-LAB Team



Wassim

Studies: Master Mathematics Studies: Master Mathematical Finance and Actuarial Science

Alican

Studies: Master Mathematics in Data Science

Cingis



Oliver

Studies: Master Mathematics

Hey, my name is John.

I work at a big accounting consulting company. Today is a bad day.

Investigate an audit of a big firm.

Done in 1995. No digital trace.





Invoices.



Thousands of them.

How it feels like?









Dataset

• RVL-CDIP dataset





First Steps



- + Unconventional approach
- Complicated preprocessing
- Hard to implement

- + Very good accuracy
- Complicated preprocessing



- + State-of-the-art
- + Simple preprocessing
- + Matches our project goals
- Training time

Classification Model



- Validation accuracy: 0.98
- Test accuracy: 0.98

Dataset

Preprocessing



Region-Based & Holistic CNNs





OCR Illustration

		41 EAST 42nd STI	REET, NEW YO	RK, N.Y. 10017, (212) 599-55
PROGRAM	NEWS	60 secs.	STATION	CBS
DATE	JUNE 11, 1968	3 7:21 PM	CITY	NEW YORK
	KENT C	IGARETTES 803	231	
(SFX	TRAIN WHISTLE)			
(MUSI	(C)			
MAN:	Here's the stor	y they still to	ell around	the railroad yard,
about one	night when Casey	Jones was driv	vin' extra l	hard. Another
train ahea	id, but watch out	. An accident!	But Case	y wouldn't jump
cause he w	vas haulin' Kent.	Roarin' towar	ds the blg	nest wreck in all
the univer	se. Casey grabb	ed the throttle	and pulled	in reverse. He
built up s	uch a head of st	eam, it was mag	nificent.	Backed all the
way to San	ta Fe, just to s	ave the Kent.		
VOICE	S: He saved the	train.		
MAN:	For the taste o	f Kent.		
VOICE	: He saved the	train.		
MAN:	For the taste o	f Kent.		
VOICE	: He saved the	train.		
MAN:	For the taste of	f Kent.		æ
VOICE	5: He saved the	train!		553
MAN	For the taste of	Kent.		00

OFFICES IN: NEW YORK • DETROIT • LOS ANGELES • WASHINGTON, D. C. • CHICAGO • AND OTHER PRINCIPAL CITIES wine Read To Magazine a detabase to assure the accuracy the mainteent and acquerate in caracteria to assore the accuracy of the accur

PROGRAM	NEWS	60 secs.	STATION	CBS
DATE	JUNE 11, 1968	7:21 PM	CITY	NEW YORK

KENT CIGARETTES 803231

(SEX TRAIN WHISTLE)

(MUSIC)

MAN: Here's the story they still tell around the railroad yard,

about one night when Casey Jones was driving extra hard. Another

't'rain ahead, but watch out. An accident! But Casey wouldn''t jump

OFFICES IN: NEW YORK + DETROIT + LOS ANGELES + WASHINGTON; D. C. + CHICAGO + AND OTHER PRINCIPAL CITIES While Radio TV Reports ind endeavors to assure the accuracy of material supplied by it, it cannot be responsible for mistakes or omissions Material supplied by Radio: TV Reports, Inc. may be used for the and reference purposes only. It may not be reproduced, sold or publicly demonstrated or exhibited.

OCR Illustration

PROGRAM NEWS 60 secs. STATION CBS DATE JUNE 11, 1968 7:21 PM CITY NEW YORK

PROGRAM	NEWS	60 secs.	STATION	C B S
DATE	JUNE 11, 1968	7:21 PM	CITY	NEW YORK

KENT CIGARETTES 803231

(srx TRAIN WHISTLE)

(MUSIC)

MAN: Here's the story they still tell around the railroad yard, about one nlqht when Casey Jones was drlvln' extra hard. Another train ahead, but watch out. An accident! But Casey wouldn't jump

KENT CIGARETTES 803231

(SFX TRAIN WHISTLE)

(MUSIC)

MAN: Here's the story they still tell around the railroad yard,

about one night when Casey Jones was driving extra hard. Another

'train ahead, but watch out. An accident! But Casey wouldn't jump

OFFICES IN: NEW YORK.DETROIT.LOS ANGELES.WASNINGTON D C.CHICAGO AND OTHER PRINCIPAL CITIES While Padro TV Reoriis imc endeevors to ...

OFFICES IN: NEW YORK + DETROIT + LOS ANGELES + WASHINGTON; D. C. + CHICAGO + AND OTHER PRINCIPAL CITIES While Radio TV Reports ind endeavors to assure the accuracy of material supplied by it, it cannot be responsible for mistanes or omissions Material supplied by Radio: TV Reports, Inc. may be used for the and reference purposes only. It may not be reproduced, sold or publicly demonstrated or exhibited:

OCR Pipeline



Preprocessing & Postprocessing

Pipeline

Preprocessing

• Upscaling

OCR

- Sharpness
- Contrast

Postprocessing

- Autocorrection
- Stopwords

Factors trained on

Usefulness of Preprocessing





28



Complete

Architecture

- No handcrafted features
- State-of-the-art results
- Mathematically interesting
- Curiosity about new networks

Decision for NNs

Named Entity

Research

NER

Read papers

Character

- 4 different papers
- Networks are comparable, yet different

• Embeddings and CNN

Casing

Information

Bilstm

- BiLSTM
- CRF

GloVe

• Output of NER (.json)

Implementation

Conditional

Random Field

Complete Character Architecture <u>Representati</u>

tation GloVe

Bilstm

Named Entity Recognition

Named Entities:

Organization Location Person

Complete

Bilstm

Named Entity Recognition

Named Entities:

Organization Location Person

Character

Facebook Inc. is endowing a new institute, led by

Christoph Luetge , devoted to the ethics of artificial

intelligence at the Technical University of Munich,

in Germany.

[10]

Complete

GloVe

Character

Bilstm

Named Entity Recognition

Named Entities:

Organization Location Person

Facebook Inc. is endowing a new institute, led by

Christoph Luetge, devoted to the ethics of artificial

intelligence at the Technical University of Munich,

in Germany.

[10]

BiLSTM-CNN-CRF: Complete Architecture



- Validation accuracy: 94.2%
- Test accuracy: 90.8%

Character

GloVe

Bilstm

CNN-Extracted Features

CNN: Transform character-level information into *character- representation*



Word-Embeddings using GloVe

Complete

Word Embeddings:

Research

NER

Pretrained word embeddings with GloVe ٠ 100 dimensions.

Named Entity

Recognition



Casing

Information

Bilstm

GloVe: Projection into 2D

GloVe

Character

Representation

Conditional

Random Field

Complete

Character

Bilstm

Casing Information

The casing of a word contains important information



Bi-Directional Long Short-Term Memory

Character

Casing

Information

Bilstm

GloVe

Complete

Capture context on both sides of a word of a sentence

Named Entity

Research

NER



Conditional

Random Field

Complete

Bilstm

Conditional Random Field

Which sequence of labels is most likely for the sentence?

[...] Technical University of Munich, in Germany.

Complete

Architecture

Which sequence of labels is most likely for the sentence?

Named Entity

[...] Technical University of Munich, in Germany.

Character

Representation

Casing

Information

Bilstm

GloVe

Not likely:

Research

NER



Conditional

Random Field

Complete

Which sequence of labels is most likely for the sentence?

Named Entity

[...] Technical University of Munich, in Germany.

Character

Representation

Casing

Information

Bilstm

GloVe

Not likely:

Research

NER



More likely:



Conditional

Random Field





Goals and Outcomes



Cloud Structure (Flask)



Thank you for your Attention



References

- Logo OpenNebula, <u>https://opennebula.org/referencing/</u>, 2019-02-08 [1]
- Logo Flask, <u>https://de.wikipedia.org/wiki/Datei:Flask_logo.svg</u>, 2019-02-08 [2]
- Logo Solr, <u>http://lucene.apache.org/solr/logos-and-assets.html</u>, 2019-02-08 [3]
- Classification by cre.ativo mustard from the Noun Project [4]
- Content by Jyoti Vyas from the Noun Project [5]
- Logo NER, <u>https://wordlift.io/blog/en/entity/named-entity-recognition/</u>, 2019-02-13 [6]
- Logo OpenNebula, <u>https://www.v3.co.uk/v3-uk/news/2356529/microsoft-gets-more-open-with-cross-platform-packerio-and-opennebula-cloud-tools</u>, 2019-02-13 [7]
- Logo Flask, <u>https://engineering.bitnami.com/articles/deploy-a-production-ready-mariadb-cluster-on-kubernetes-with-bitnami-and-helm.html</u>, 2019-02-14 [8]
- process by Gregor Cresnar from the Noun Project [9]
- Facebook Endows AI Ethics Institute at German University TUM, Jeremy Kahn, <u>https://www.bloomberg.com</u>, 2019-01-20 [10]
- Graph, https://thenounproject.com/mb.icons/collection/network/?i=1775952, 2019-02-14 [11]
- SVM-SVD, <u>https://thenounproject.com/search/?q=svm&i=1503831</u>, 2019-02-14 [12]
- Network, <u>https://thenounproject.com/smodgekar/collection/technology/?i=1714861</u>, 2019-02-14 [13]
- RVL-CDIP sample1, <u>http://www.cs.cmu.edu/~aharley/rvl-cdip/images/sample1.png</u>, 2019-02-16 [14]
- RVL-CDIP sample2, <u>http://www.cs.cmu.edu/~aharley/rvl-cdip/images/sample2.png</u>, 2019-02-16 [15]
- GloVe, <u>https://nlp.stanford.edu/projects/glove/</u>, 2019-02-16 [16]
- Paperwork, <u>https://de.depositphotos.com/62875115/stock-illustration-vector-paperwork-mood.html</u>, 2019-02-17 [17]
- Banana, https://www.alamy.com/banana-logo-template-vector-icon-illustration-design-image159202591.html, 2019-02-17 [19]
- Dashboard, <u>https://doc.lucidworks.com/fusion/2.0/Dashboards.html</u>, 2019-02-17 [20]
- Magnifier, Search by Mello from the Noun Project, 2019-02-17 [21]
- Data, database by Aiden Icons from the Noun Project, 2019-02-17 [22]
- Cloud-Structure, cloud storage by un delivered from the Noun Project, 2019-02-17 [23]
- Cloud, implementation by Tomas Knopp from the Noun Project, 20189-02-17 [24]