

METODE INTELIGENTE DE REZOLVARE A PROBLEMELOR REALE

Laura Dioşan
Text mining

Facultatea de Matematică şi Informatică
Universitatea Babeş-Bolyai

1

1

□ Natural Language Processing

- Text classification
- Language modelling
- Machine translation
- ...

2

From Languages to Information

□ Aim

- Automatically extracting meaning and structure from:
 - Human language text and speech (news, social media, etc.)
 - Social networks
 - Genome sequences
- Interacting with humans via language
 - Dialog systems/Chatbots
 - Question Answering
 - Recommendation Systems

3

From Languages to Information

□ How?

- Extracting information from language
 - Information Retrieval
 - Text Classification
 - Extracting Sentiment and Social Meaning
 - ...
- Interacting with humans via language
 - IBM's Watson
 - ...
 - Counseling conversations
 - <https://www.aclweb.org/anthology/Q16-1033.pdf>
 - https://web.stanford.edu/class/cs124/lec/counseling_slides.pdf
 - Understanding police officer respect
 - <https://www.pnas.org/content/114/25/6521>
 - https://nlp.stanford.edu/robvoigt/124_lecture/

4

Text classification

□ Definition

- Text categorization
 - Assign (predefined) labels (categories) to some documents
 - Documents
 - Technical reports, web pages, messages, books, etc.
 - Categories:
 - Topics (art, economy)
 - Attributes (spams, ham)

□ Data and various text classification problems

- https://nlpprogress.com/english/text_classification.html

	Cuvinte	Documente
Învăţare supervizată	Etichetarea părţilor de vorbire	Clasificarea textelor, Filtrarea, Detectarea subiectelor
Învăţare nesupervizată	Indexarea semantică, construcţia automată a tezaurilor, extragerea cuvintelor cheie	Clusterizarea documentelor, Detectarea subiectelor

5

Text classification

□ Automatic approaches

- Based on knowledges
 - From experts
 - Encoded as rules
- Based on learning
 - Experts label some training examples
 - The algorithm label the test examples
 - Learning
 - Supervised
 - Unsupervised

□ Evaluation

- Accuracy, Precision, Recall, AUC, etc.

6

Text classification

- Process
 - Data training analysis
 - **Document indexing**
 - Construct a document representation
 - Attributes / features
 - Weights of these features
 - Reduce the dimension of these representations
 - Feature selection
 - Feature extraction
 - Learning a classification model
 - Classification of test data
 - Text indexing
 - Apply the learnt model -> categories

7

Text classification

- Document indexing
 - Construct a document representation
 - Attributes / features
 - Weights of these features
 - 4 steps:
 1. Document linearization
 2. Filtering
 3. Canonical form
 4. Weighting

} Reduce the vocabulary size

8

Text classification – indexing

Step 1: linearization (segmentation)

- Reduce the document to a vector of terms (attributes)
 - *bag of words model*
 - A matrix
 - Documents -> lines
 - Terms -> columns
 - Each cell 1/0 – if the current term belongs to the current document
- 2-stage process for term identification
 - Format removal
 - E.g. Elimination of HTML tags
 - *Tokenization*
 - Parsing (segmentation)
 - Removal of punctuation mark
 - Conversion of capital letters

Initial	Linearizat
Interactive query expansion modifies queries using terms from a user. Automatic query expansion expands queries automatically.	Interactive query expansion modifies queries using terms from a user. Automatic query expansion expands queries automatically.

9

Text classification -> indexing

Paul 2: filtering

- Select some terms that are able to
 - Describe the content of the document
 - Make a difference between two documents
- Removal of stopwords
 - From a predefined list
 - Based on their frequencies (under a given threshold)

Segmentat	Filtrat
Interactive query expansion modifies queries using terms from a user. Automatic query expansion expands queries automatically.	Interactive query expansion modifies queries terms automatic query expansion expands queries automatically

10

Text classification – indexing

Step 3 – canonical form (normalization/standardisation)

- Lematisation
 - Morphological analysis
 - Based on context
 - Ex. "better" → "good"
- Stemming
 - Term-level
 - Ex. "computer", "computing", "compute" → "comput"
 - stemming algorithm
 - Martin Porter
 - WordNet

Filtrat	Reduc
Interactive query expansion modifies queries terms automatic query expansion expands queries automatically	Interact queri expan modifi queri term automat queri expan expand queri automat

11

Text classification – indexing

Step 4: weighting

- Pre-determined (manual)
 - Ex. Bag of Words
- Automatic (learnt)
 - Ex. embedded representations

12

Text classification - indexing

- Step 4: manual weighting
 - Usage of a particular model
 - Weights relative to
 - A single document
 - term frequency - TF
 - A collection of documents
 - inverse document frequency - IDF
 - hybrid TF and IDF
 - $TF \rightarrow$ cu cât un termen este mai frecvent într-un document, cu atât el este mai important pentru acel document
 - $IDF \rightarrow$ cu cât un termen apare în mai multe documente, cu atât el este mai puțin important în descrierea semantică a celui document
 - Range of frequencies
 - Binary \rightarrow presence / absence of a term (One-Hot encoding)
 - Real $[(0,1)] \rightarrow$ importance of that term
 - For a set of D documents and a set of T terms, the weight p_{ij} of term t_i in document d_j ($i=1,2,\dots,|T|$, $j=1,2,\dots,|D|$) can be:
 - binary: $p_{ij} = 1$, if t_i belongs to d_j
0, otherwise
 - TF : $p_{ij} = tf_{ij}$ (no of appearances of term t_i in document d_j)
 - $TF \cdot IDF$: $p_{ij} = tf_{ij} \cdot \log_{10}(|D|/df_i)$, where df_i = no of documents where term t_i is found

13

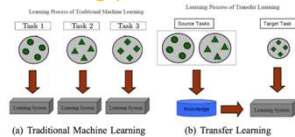
Text classification – indexing

- Step 4 – weight learning
 - Unsupervised learnt representation
 - Classical models
 - Probability of a target word based on k previous words
 - Document-level context
 - Semantic relation (boat – water)
 - Count-based models
 - Latent Semantic Analysis (LSA)
 - Latent Dirichlet Allocation (LDA)
 - Modern models
 - Word-level context
 - Semantic similarity (boat – sheep)

14

Text classification – indexing

- Step 4 – weight learning (word representations)
 - Methodologies
 - Traditional Machine Learning
 - Transfer learning
 - https://www.cse.ust.hk/~qyang/Docs/2009/tkde_transfer_learning.pdf



15

Text classification – indexing

Learning word representations (vectors)

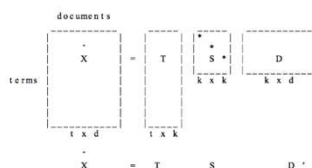
- Why?
 - Embeddings = parameters \rightarrow they can be learnt
 - Share representation across tasks
 - Lower dimensional space
- How?
 - Pre-NN
 - 1990 - Latent Semantic Analysis (LSA)
 - <http://lsa.colorado.edu/papers/JASIS.lsi.90.pdf>
 - 1992 - n-gram models <https://www.aclweb.org/anthology/J92-4003.pdf>
 - 2003 - Latent Dirichlet Allocation (LDA) - Documents are mixtures of topics and topics are mixtures of words <https://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>
 - NN-based
 - Word-level (2003 - ...)
 - Sentence (document) level (2014 - ...)
 - Contextual word-vectors (Word vectors compress all contexts into a single vector) (2016 - ...)

16

Text classification – indexing

Learning word representations (vectors)

- How?
 - Pre-NN
 - 1990 - Latent Semantic Analysis (LSA)
 - <http://lsa.colorado.edu/papers/JASIS.lsi.90.pdf>



17

Text classification – indexing

Learning word representations (vectors)

- How?
 - Pre-NN
 - 1992 - n-gram models
 - <https://www.aclweb.org/anthology/J92-4003.pdf>
 - <https://web.stanford.edu/~jura/sky/slp3/3.pdf>

Suppose we are learning a 4-gram Language Model.

~~the~~ ~~proctor~~ ~~started~~ ~~the~~ ~~class~~ ~~the~~ ~~students~~ ~~opened~~ ~~their~~ ~~books~~

discard condition on this

$$P(w|\text{students opened their}) = \frac{\text{count}(\text{students opened their } w)}{\text{count}(\text{students opened their})}$$

For example, suppose that in the corpus:

- "students opened their" occurred 1000 times
- "students opened their books" occurred 400 times
- $\rightarrow P(\text{books} | \text{students opened their}) = 0.4$
- "students opened their exams" occurred 100 times
- $\rightarrow P(\text{exams} | \text{students opened their}) = 0.1$

Should we have discarded the "proctor" context?

18

Text classification – indexing

Learning word representations (vectors)

- How? -> NN-based
 - Word-level (2003 - ...)
 - 2003 - N-gram Neural language model (Montreal - Bengio)
 - <https://www.jmlr.org/papers/volume3/bengio03a/bengio03a.pdf>
 - Probability of the target word based on previous k words
 - Estimation of the probability by an ANN -> prediction of the next word in a sequence
 - Embedding layer -> word embeddings
 - Intermediate layer(s) -> intermediate representation (non-linear); standard (HardTanh layer) or recurrent layers
 - softmax layer for producing probabilities over vocabulary -> main bottleneck
 - Cross-entropy loss (max the probability of the next word)

19

Text classification – indexing

Learning word representations (vectors)

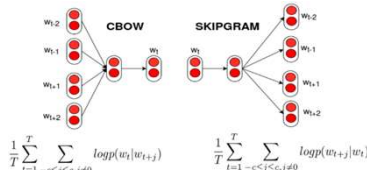
- How? -> NN-based
 - Word-level (2003 - ...)
 - 2008 - multi-task model (Princeton - Collobert)
 - https://ronan.collobert.com/pub/matos/2008_nlp_icml.pdf
 - Probability of MORE target words (a sequence of words)
 - Estimation of the probability by an ANN -similar to Bengio's model, but
 - a pairwise ranking criterion
 - outputs a higher score for a correct word sequence than for an incorrect one

20

Text classification – indexing

Learning word representations (vectors)

- How? -> NN-based
 - Word-level (2003 - ...)
 - 2013 - word2vec (Google - Mikolov)
 - <https://papers.nips.cc/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf>
 - Continuous BOW model -> try to predict a word based on its context (neighbours); input = neighbour words, output = target word
 - Skip-gram model -> try to predict context (neighbours) of a word; input = target word; output = neighbour words;
 - Visualisation of embeddings <https://ronxin.github.io/wevi/>
 - Trained vectors <https://radimrehurek.com/gensim/models/word2vec.html> or <https://github.com/S10p/word2vec-api#where-to-get-a-pretrained-models>



21

Text classification – indexing

Learning word representations (vectors)

- How? -> NN-based
 - Word-level (2003 - ...)
 - 2014 - GloVe (Stanford - Pennington, Manning)
 - <https://nlp.stanford.edu/projects/glove/>
 - 2017 - fastText (Facebook - Mikolov)
 - <https://arxiv.org/abs/1607.04606>
 - <https://radimrehurek.com/gensim/models/fasttext.html>

22

Text classification – indexing

Learning word representations (vectors)

- How? -> NN-based
 - Word-level (2003 - ...)
 - Sentence (document) level (2014 - ...)
 - 2014 - Paragraph embedding
 - <https://arxiv.org/abs/1405.4053>
 - 2015 - skip-thought vectors
 - Predict previous / next sentence with seq2seq model
 - <https://arxiv.org/abs/1506.06726>
 - 2015 - auto-encoders (semi or unsupervised)
 - <https://arxiv.org/abs/1511.01432>,
 - <https://arxiv.org/pdf/1511.06349.pdf>,
 - <https://arxiv.org/abs/1602.03483>

23

Text classification – indexing

Learning word representations (vectors)

- How? -> NN-based
 - Word-level (2003 - ...)
 - Sentence (document) level (2014 - ...)
 - Contextual word-vectors (Word vectors compress all contexts into a single vector) (2016 - ...)
 - 2016 context2vec <https://www.aclweb.org/anthology/K16-1006.pdf>
 - 2017 tagLM <https://arxiv.org/abs/1705.00108>
 - 2017 CoVe <https://papers.nips.cc/paper/2017/hash/20c86a628232a67e7bd46f76fba7ce12-Abstract.html>
 - 2018 ELMo <https://www.aclweb.org/anthology/N18-1202/>
 - 2018 ULMFIT <https://arxiv.org/abs/1801.06146>
 - 2019 BERT <https://arxiv.org/abs/1810.04805>
 - cross-lingual pre-training <https://arxiv.org/abs/1706.04901>
 - Code example <https://github.com/huggingface/naacl-transfer-learning-tutorial>
 - See also <https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/>

24

Text classification – indexing

Learning word representations (vectors)

□ How? -> NN-based

The diagrams show four different neural network architectures for learning word representations: word2vec (Mikolov et al. 2013), ELMO (Peters et al. 2018, ULMFiT (Howard & Ruder 2018), GPT (Radford et al. 2018)), Skip-Thought (Kiros et al. 2015), and BERT (Devlin et al. 2019). Each diagram shows how words are mapped to vectors in a neural network context.

25

Text classification

□ Process

- Data training analysis
 - Document indexing
 - Construct a document representation
 - Attributes / features
 - Weights of these features
 - Reduce the dimension of these representations
 - Feature selection
 - Feature extraction
 - Learning a classification model
- Classification of test data
 - Text indexing
 - Apply the learnt model -> categories

26

Text classification

□ Reducerea dimensiunii

- Are drept scop
 - Creșterea eficacității
 - Reducerea timpului de învățare a modelului de clasificare
 - Evitarea învățării pe de rost a modelului de clasificare
- Poate consta în
 - Selecția atributelor (*feature selection*)
 - o submulțime a atributelor inițiale (originale)
 - Extragerea atributelor
 - o mulțime de noi atribute determinate pe baza celor originale -> proiecția unui vector R -dimensional într-unul r -dimensional ($r < R$)
 - noile atribute (mai puține) reprezintă o transformare a atributelor originale

27

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor

□ Dându-se o mulțime de atribute $X_k = (x_{k1}, x_{k2}, \dots, x_{km})$ pentru un document $d_k \in D$, să se găsească o submulțime $X_{k,p} = (x_{k,1p}, x_{k,2p}, \dots, x_{k,ip}, \dots, x_{k,mp})$, cu $p < m$ care să optimizeze o funcție obiectiv $J(X_{k,p}^m)$

- Fc. obiectiv → eroarea de clasificare

□ Selecția implică

- O strategie de căutare pentru selecția submulțimilor candidat
 - căutare exhaustivă → toate submulțimile posibile → nefezabil
 - căutare strategică
 - prin ordonarea atributelor
 - pe baza unei metrici
 - și alegerea celor care depășesc un anumit prag
 - prin selectarea unei anumite submulțimi de atribute
 - se alege o submulțime optimă
- O funcție obiectiv pentru evaluarea acestor submulțimi candidat
 - măsură a calității unei submulțimi de atribute
 - ajută selecția unei noi submulțimi candidat

28

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor

□ Metode

- Nesupervizate
 - Clusterizare
 - Factorizarea matricilor
- Supervizate
 - Ordonarea atributelor
 - Selecția unei submulțimi de atribute

29

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Metode Nesupervizate → Clusterizare

□ Se grupează atributele în clusteri

- K-means
- Hierarchical clustering

The diagram shows four clusters of points labeled 'Feature group 1', 'Feature group 2', 'Feature group 3', and 'Feature group 4', illustrating the result of a clustering algorithm.

□ Se înlocuiesc (multe) atribute similare din același cluster cu centrul clusterului

30

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Metode Nesupervizate
→ factorizarea matricilor

- Analiza componentelor principale
- Descompunerea în valori singulare
- Factorizarea matricilor non-negative
- Isomap-uri
- Self-organized maps (SOMs)

31

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Prin ordonarea atributelor

- Pp. că avem n date $(\mathbf{x}_k, y_k), k=1,2,\dots,n$
 - $\mathbf{x}_k \in \mathbb{R}^n \rightarrow \mathbf{x}_k = (x_{k1}, x_{k2}, \dots, x_{kn})$
 - $y_k \in \mathbb{R}$
- Se calculează o funcție scor pentru fiecare pereche $S(i) = (x_{ki}, y_k)$
 - cu cât scorul este mai mare, cu atât variabila este mai importantă
- și se ordonează atributele în funcție de acest scor
- Notăție
 - $X_i \in \mathbb{R}^n \rightarrow X_i = (x_{i1}, x_{i2}, \dots, x_{in})$
 - $Y \in \mathbb{R}^n \rightarrow Y = (y_1, y_2, \dots, y_n)$

32

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Prin ordonarea atributelor

- Scoruri posibile
 - Coefficientul de corelație al lui Pearson
 - $R(i) = \text{cov}(X_i, Y) / (\text{var}(X_i) \text{var}(Y))^{1/2}$
 - $R(i) = \sum_{k=1}^n (x_{ki} - \bar{x}_i)(y_k - \bar{y}) / (\sum_{k=1}^n (x_{ki} - \bar{x}_i)^2 \sum_{k=1}^n (y_k - \bar{y})^2)^{1/2}$
 - $R^2(i) \rightarrow$ relație de dependență liniară între X_i și Y
 - Eroarea de clasificare
 - Mai mulți clasificatori cu o singură variabilă
 - $(x_{ki}, y_k), k=1,2,\dots,n$
 - Se stabilește eroarea de clasificare pt fiecare $i=1,2,\dots,n$
 - Se ordonează variabilele în funcție de eroare
 - Cu cât eroarea este mai mică cu atât variabila este mai importantă
 - Informația teoretică
 - Informația mutuală între densitatea variabilei X_i și densitatea variabilei Y
 - $I(i) = \int \int p(x,y) \log(p(x,y) / (p(x)p(y))) dx dy$
 - $p(x)$ – probabilitatea densității lui $x \rightarrow$ greu de estimat

33

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor
→ Prin ordonarea atributelor

- Critici
 - poate determina submulțimi de atribute redundante
 - nu ține cont de corelarea atributelor
 - un atribut nefolositor în izolație poate fi util în combinație cu alte atribute

34

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Prin alegerea unei submulțimi de atribute

- Căutarea
 - Căutare exhaustivă – toate submulțimile posibile \rightarrow nefezabilă
 - Căutare strategică – alegerea doar a unor submulțimi
- Funcția obiectiv – tipuri
 - Wrapper
 - Filter
 - Embedded

35

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Prin alegerea unei submulțimi de atribute

- Funcția obiectiv – tipuri
 - Wrapper
 - Funcția obiectiv este un clasificator care evaluează fiecare submulțime prin puterea ei predictivă
 - Alegerea atributelor este **dependentă** de performanța clasificatorului (algoritmului de învățare)
 - Algoritmul de învățare = cutie neagră pentru evaluarea submulțimii de atribute în funcție de puterea de învățare (clasificare) a acestora
 - Filter
 - Funcția obiectiv evaluează fiecare submulțime doar pe baza conținutului ei
 - Alegerea atributelor este **independentă** de performanța clasificatorului
 - Selecția atributelor este un pas anterior învățării
 - Embedded
 - Alegerea atributelor are loc **în timpul** învățării

36

UNIVERSITATEA BABEȘ-BOLYAI Facultatea de Matematică și Informatică		
Feat/output	Numeric	Categorical
Numerical	Pearson (linear) Spearman (non-lin)	LDA Rank-based meth
Categorical	ANOVA (linear) Kendall (non-lin)	Chi-Square Mutual Information

37

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Prin alegerea unei submulțimi de atribute → *Wrapper*

- Ideea de bază
 - Wrapper* → a înveli, a împacheta
 - Funcția obiectiv este un clasificator care evaluează fiecare submulțime prin puterea ei predictivă
 - Alegerea atributelor este **dependentă** de performanța clasificatorului (algoritmului de învățare)
 - Algoritmul de învățare = cutie neagră pentru evaluarea submulțimii de atribute în funcție de puterea de învățare (clasificare) a acestora
- Algoritm
 - Se alege o metodă de clasificare (învățare)
 - Se caută configurația optimă (submulțime de atribute și parametri ai clasificatorului)
 - Se alege o submulțime de atribute
 - Se repetă
 - Învățarea și optimizarea clasificatorului
 - cuantificarea performanței clasificatorului
 - alegerea unei noi submulțimi de atribute
 - până când se obține cea mai bună performanță în învățare

38

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Prin alegerea unei submulțimi de atribute → *Wrapper*

- Cum se alege o submulțime?
 - best-first*
 - branch-and-bound*
 - simulated annealing*
 - algoritmi genetici
 - greedy*
 - Forward selection*
 - Variablele sunt incorporate progresiv în submulțimi tot mai mari
 - Backward selection*
 - Variablele sunt eliminate progresiv din submulțime
- Cum se stabilește performanța algoritmului de învățare?
 - Validare
 - Validare-incrucșată
- Care algoritm de învățare să se folosească?
 - Arbori de decizie
 - Rețele neuronale
 - Mașini cu suport vectorial
 - Algoritmi evolutivi, etc

39

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor → Prin alegerea unei submulțimi de atribute → *Filter*

- Ideea de bază
 - Funcția obiectiv evaluează fiecare submulțime doar pe baza conținutului ei
 - Alegerea atributelor este **independentă** de performanța clasificatorului
 - Selecția atributelor este un pas anterior învățării
- Evaluare
 - Distanța sau măsura separabilității claselor
 - Ex. distanța (Euclideană, Hamming, etc) între clase
 - Corelația și măsuri de informație teoretică
 - Submulțimile bune conțin atribute
 - puternic corelate cu ieșirea
 - ne-corelate între ele
 - Măsuri liniare
 - Coefficientul de corelație
 - Măsuri neliniare
 - Informația mutuală

40

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Selecția atributelor
→ Prin alegerea unei submulțimi de atribute

- <http://jmlr.csail.mit.edu/papers/volume3/guyon03a/guyon03a.pdf>
- <http://jmlr.csail.mit.edu/proceedings/papers/v4/guerif08a/guerif08a.pdf>
- http://courses.cs.tamu.edu/rgutier/cs790_w02/I5.pdf

41

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

- Analiza documentelor de antrenament
 - Indexarea documentelor
 - Construirea unei reprezentări a documentelor → transformarea documentelor într-o formă interpretabilă de către clasificator
 - Obținerea unor concepte/termeni reprezentative(i) → atribute
 - Calcularea unor ponderi pt aceste atribute
 - Reducerea dimensiunii (a numărului de concepte/atribute/termeni reprezentative(i) pentru document)
 - Selecția atributelor
 - Extragerea atributelor
 - Învățarea unui model de clasificare
- Clasificarea noilor documente(de test)

42

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Reducerea dimensiunii → Extragerea atributelor

- Definire
 - Determinarea unei noi mulțimi de atribute determinate pe baza celor originale → proiecția unui vector R -dimensional într-unul r -dimensional ($r < R$)
 - Noile atribute (mai puține) reprezintă o transformare a atributelor originale
- Dându-se o mulțime de atribute $X_i = (x_{i1}, x_{i2}, \dots, x_{im})$, să se găsească o transformare $Z_i = g(X_i): R^m \rightarrow R^r$ cu $p < m$ astfel încât transformarea Z_i să păstreze (cea mai parte din) informația atributelor inițiale
 - Transformarea optimă – cea care nu determină creșterea probabilității de eroare
 - Transformarea poate fi
 - Lineară $y = Wx, W \in R^{r \times m}$
 - Ne-lineară – greu de determinat
 - Transformarea este ghidată de o funcție obiectiv care trebuie optimizată (min/max)
- Metode de extragere a atributelor în funcție de criteriul măsurat de funcția obiectiv:
 - Reprezentare a semnalului → transformarea are drept scop reprezentarea datelor cu o acuratețe cât mai bună într-un spațiu mai redus
 - Analiza componentelor principale (PCA)
 - Clasificare → transformarea are drept scop evidențierea discriminării între clase într-un spațiu mai mic
 - Analiza discriminantului linear (LDA)
- Metode – în funcție de relația dintre features
 - Relații liniare
 - PCA, LDA
 - Relații neliniare
 - Locally linear embedding → manifold learning
 - t-SNE
 - Auto-encoders

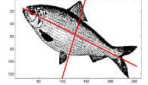
43

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Metode de reducere a dimensiunii → Extragerea atributelor → Analiza componentelor principale

- Scop
 - Transformarea unui set de variabile posibil corelate într-un set de variabile necorelate între ele (componente principale)
 - Prima componentă principală are cea mai mare varianță → cuantifică cea mai mare variabilitate posibilă a datelor
 - ACP determină axele care explică cel mai bine dispersia datelor (norul de puncte)
 - Descrierea datelor într-un spațiu dimensional mai mic
- Alte denumiri
 - Transformarea Karhunen-Loève (teoria comunicațiilor)



44

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Metode de reducere a dimensiunii → Extragerea atributelor → Analiza componentelor principale

- Tipologie
 - ACP liniară – date separabile liniar
 - ACP bazată pe kernele – date neseperabile liniar
- Algorithm
 - Pp că avem un set de date $X_i, i=1,2,\dots,n$ cu m atribute ($x_i \in R^m \Rightarrow x_i = (x_{i1}, x_{i2}, \dots, x_{im})$)
 - Scăderea mediei din fiecare dată (pe fiecare dimensiune) → centrarea datelor
 - $\bar{x}_j = x_{1j} - x_j, \text{ unde } x_j = (x_{1j} + x_{2j} + \dots + x_{nj})/n$
 - Calcularea matricii de covarianță C
 - $C = (c_{ij}), i, j = 1, 2, \dots, m, c_{ij} = \text{cov}(x_{ij}, x_{ij}), \text{ unde } x_{ij} = (x_{1ij}, x_{2ij}, \dots, x_{nij})$
 - $\text{cov}(X,Y) = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) / (n-1)$
 - Determinarea vectorilor proprii v_p și a valorilor proprii v_p (eigenvector, eigenvalue) corespunzătoare matricii de covarianță $A, v_p = V_p, v_p$
 - Alegerea componentelor și formarea vectorului de caracteristici (atribute)
 - Se ordonează vectorii proprii descrescător după valurile proprii → atributele în ordinea importanței
 - Formarea vectorului de caracteristici cu acei vectori proprii care se doresc a fi reținuți
 - Derivarea noilor date
 - Se înmulțește vectorul de caracteristici cu vectorul datelor centrate

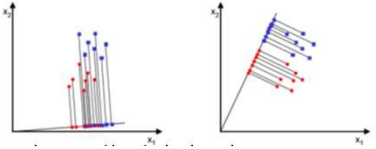
45

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Metode de reducere a dimensiunii → Extragerea atributelor → Analiza discriminantului liniar

- Scop
 - Determinarea separe datele
 - Modelarea dife
 - Proiectarea da
- observa o mai bună separabilitate a datelor → care este cea mai bună proiecție?
 - $y = w^T x$



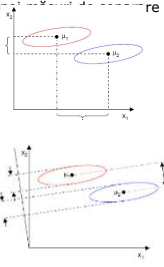
46

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Metode de reducere a dimensiunii → Extragerea atributelor → Analiza discriminantului liniar

- Găsirea celei mai bune proiecții necesită definirea unei funcții obiectiv între proiecțiile datelor
 - Distanța între proiecțiile mediilor corespunzătoare datelor din fiecare clasă
 - Nu este foarte bine pentru că nu se ține cont de dispersia datelor în interiorul claselor
 - Fisher → maximizarea raportului dintre diferența mediilor și împrăștierea în interiorul claselor
 - o proiecție astfel încât:
 - exemplele din aceeași clasă sunt proiectate foarte aproape unele de altele
 - proiecțiile mediilor fiecărei clase sunt cât mai departate unele de altele



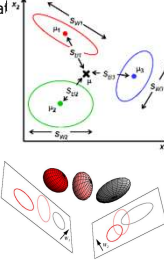
47

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

Metode de reducere a dimensiunii → Extragerea atributelor → Analiza discriminantului liniar

- Algorithm
 - Pp că:
 - există k clase,
 - μ_i – media instanțelor din clasa $i, i=1,2,\dots,k$
 - n – nr total de instanțe
 - n_i – nr de instanțe din clasa $i, i=1,2,\dots,k$
 - Se caută $k-1$ vectori de proiecție
 - Se calculează
 - Împrăștierea intra-clasă (scatter within class) S_w
 - $S_w = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \mu_i)(x_{ij} - \mu_i)^T$
 - Împrăștierea între clase (scatter between classes) S_b
 - $S_b = \sum_{i=1}^k n_i (\mu_i - \mu)(\mu_i - \mu)^T$, unde $\mu = 1/n \sum_{i=1}^k \sum_{j=1}^{n_i} x_{ij}$
 - Se maximizează
 - Raportul dintre
 - Pătratul diferenței mediilor (claselor) și
 - Împrăștierea intra-clasă
 - Soluție
 - $w = S_w^{-1}(\mu_1 - \mu_2)$



http://research.cs.tamu.edu/prism/lectures/pr/pr_110.pdf
<http://www.ditp.ro/ditp/ditp.htm>
http://www.music.mcgill.ca/~ich/classes/mumt611_05/classifiers/lda_theory.pdf

48

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

- Analiza documentelor de antrenament
 - Indexarea documentelor
 - Construirea unei reprezentări a documentelor → transformarea documentelor într-o formă interpretabilă de către clasificator
 - Obținerea unor concepte/termeni reprezentative(i) → atribut
 - Calcularea unor ponderi pt aceste atribute
 - Reducerea dimensiunii (a numărului de concepte/atribute/termeni reprezentative(i) pentru document)
 - Selecția atributelor
 - Extragerea atributelor
 - **Învățarea unui model de clasificare**
- Clasificarea noilor documente(de test)
 - Indexarea documentelor
 - Utilizarea modelului de clasificare pentru stabilirea categoriilor fiecărui document de test

49

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

- **Învățarea unui model de clasificare**
 - Alegerea unui algoritm de învățare
 - Arbori de decizie
 - Rețele neuronale artificiale
 - Mașini cu suport vectorial
 - Algoritmi evolutivi
 - Rețele Bayesiene
 - Fixarea/optimizarea parametrilor algoritmului
 - Cum se aleg parametrii?
 - Construirea modelului de clasificare și salvarea lui

50

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

- Analiza documentelor de antrenament
 - Indexarea documentelor
 - Construirea unei reprezentări a documentelor → transformarea documentelor într-o formă interpretabilă de către clasificator
 - Obținerea unor concepte/termeni reprezentative(i) → atribut
 - Calcularea unor ponderi pt aceste atribute
 - Reducerea dimensiunii (a numărului de concepte/atribute/termeni reprezentative(i) pentru document)
 - Selecția atributelor
 - Extragerea atributelor
 - Învățarea unui model de clasificare
- **Clasificarea noilor documente(de test)**
 - Indexarea documentelor
 - Utilizarea modelului de clasificare pentru stabilirea categoriilor fiecărui document de test

51

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Clasificarea automată a textelor – Învățare – proces

- Metode de reducere a dimensiunii
 - Extragerea atributelor
 - Analiza componentelor principale
 - Analiza componentelor independente
 - Scalare multidimensională
 - Hărți topografice

52

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

Natural Language Processing

- Text to numbers
 - <https://machinelearningmastery.com/prepare-text-data-deep-learning-keras/>
- Word embedding
 - <https://www.tensorflow.org/tutorials/word2vec>
 - <https://nlp.stanford.edu/projects/glove/>
 - <https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/>
 - <https://machinelearningmastery.com/develop-word-embeddings-python-gensim/>

53

UNIVERSITATEA BABEȘ-BOLYAI
Facultatea de Matematică și Informatică

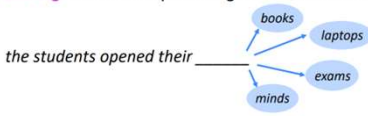
Natural Language Processing

- Text classification
- **Language modeling**
- Machine translation
- ...

54

Language modelling

- Language Modeling is the task of predicting what word comes next.

the students opened their 


- More formally: given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$:

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$$
 where $x^{(t+1)}$ can be any word in the vocabulary $V = \{w_1, \dots, w_{|V|}\}$
- A system that does this is called a **Language Model**.

55

Language modelling

- Practical tasks
 - Sentence or document classification (e.g. sentiment)
 - Sentence pair classification (e.g. NLI, paraphrase)
 - Word level (e.g. sequence labeling, extractive Q&A)
 - Structured prediction (e.g. parsing)
 - Generation (e.g. dialogue, summarization, machine translation)
 - See <https://medium.com/@sammyobama-rnn-machine-generated-political-speeches-c8abd18a2e90>



56

Language modelling

- Formal task
 - Assign a probability to a sentence
 - Machine Translation:
 - $P(\text{high winds tonite}) > P(\text{large winds tonite})$
 - Spell Correction
 - The office is about fifteen **minuets** from my house
 - $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$
 - Speech Recognition
 - $P(\text{I saw a van}) \gg P(\text{eyes awe of an})$
 - + Summarization, question-answering, etc., etc.!!

57

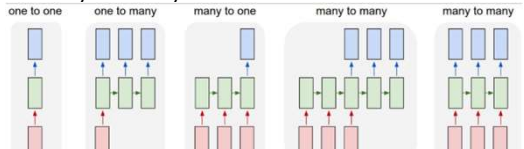
Language modelling

- Datasets
 - https://nlpprogress.com/english/language_modeling.html
- Evaluation
 - Extrinsic
 - Run two models on the same task and compare the accuracies
 - + : simple, relevant
 - : time consuming
 - Intrinsic
 - Perplexity
 - measure how well a probability distribution predicts a sample
 - Remember: Cross-entropy (H) -> classification performance (max)
 - Perplexity = 2^H (lower = better)
 - The Shannon Game:
 - How well can we predict the next word?
 - I always order pizza with cheese and **anchovies**
 - The 33rd President of the US was **fried rice**
 - I saw a **pepperoni**
 - Unigrams are terrible at this game. (Why?)
 - A better model of a text is one which assigns a higher probability to the word that actually occurs.

58

Language modelling

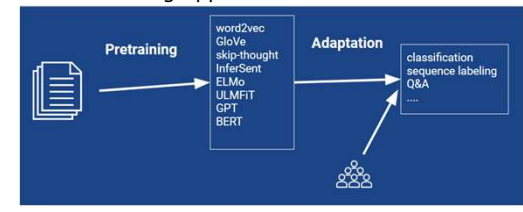
- How?
 - 1 to 1 – name entity recognition
 - 1 to many – image captioning
 - many to one – sentiment classification in sentences
 - many to many – machine translation



59

Language modelling

- How?
 - Traditional machine learning approaches
 - Transfer learning approaches



60

Language modelling

- How?
 - Count-based n-gram models
 - Approximate the history of observed words with just the previous n words
 - Neural n-gram models
 - Feedforward neural networks
 - Embed the n-gram history into a continuous space -> better capture the correlations between histories
 - Recurrent neural networks
 - Compress the entire history in a fixed length vector -> enable long range correlations

61

Language modelling

- How? -> Count-based n-gram models
 - Approximate the history of observed words with just the previous n words
 - Markov chains
 - Only previous history matters
 - Limited memory = previous k - 1 words
 - E.g. a 3-gram model:
 - a third order Markov model:
 - $p(w_1, w_2, \dots, w_n) = p(w_1) p(w_2 | w_1) p(w_3 | w_1, w_2) p(w_4 | w_2, w_3) \dots$
 - $p(w_n | w_{n-2}, w_{n-1})$
 - Probability estimation:
 - $p(w_3 | w_1, w_2) = \text{count}(w_1, w_2, w_3) / \text{count}(w_1, w_2)$
 - + scalable, able to be trained on trillions of words
 - + fast constant time evaluation of probabilities at test time
 - + smoothing methods for matching the empirical distribution of language (heaps' law)
 - large n-grams are sparse => hard to capture long correlations
 - symbolic similarity does not involve semantic similarity (dog vs. cat)

62

Language modelling

- How? -> Neural n-gram models -> Feedforward neural networks
 - Embed the n-gram history into a continuous space -> better capture the correlations between histories

63

Language modelling

- How?
 - NN-based -> standard NNs

A fixed-window neural Language Model

output distribution: $g = \text{softmax}(Uh + b_g) \in \mathbb{R}^{|V|}$

hidden layer: $h = f(We + b_h)$

concatenated word embeddings: $e = [e^{(1)}, e^{(2)}, e^{(3)}, e^{(4)}]$

words / one-hot vectors: $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$

Improvements over n-gram LM:

- No sparsity problem
- Don't need to store all observed n-grams

Remaining problems:

- Fixed window is too small
- Enlarging window enlarges W
- Window can never be large enough
- $x^{(1)}$ and $x^{(2)}$ are multiplied by completely different weights in W . No symmetry in how the inputs are processed.

64

Language modelling

- How?
 - NN-based -> recurrent NNs
 - Standard NN versus Recurrent NN
 - <https://www.deeplearningbook.org/contents/rnn.html>
 - <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
 - Char-level recurrent NN
 - <https://github.com/karpathy/char-rnn>

A Simple RNN Language Model

output distribution: $g^{(t)} = \text{softmax}(Ux^{(t)} + b_g) \in \mathbb{R}^{|V|}$

hidden states: $h^{(t)} = \sigma(Wx^{(t-1)} + Wh^{(t-1)} + b_h)$

word embeddings: $e^{(t)} = E x^{(t)}$

words / one-hot vectors: $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$

RNN Advantages:

- Can process any length input
- Computation for step t can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

RNN Disadvantages:

- Recurrent computation is slow
- In practice, difficult to access information from many steps back

65

Language modelling

- How?
 - NN-based -> recurrent NNs

Effect of vanishing gradient on RNN-LM

RNN Disadvantages:

- Recurrent computation is slow
- In practice, difficult to access information from many steps back

LM task: The writer of the books are planning a sequel

Correct answer: The writer of the books is planning a sequel

Syntactic recency: The writer of the books is (correct)

Sequential recency: The writer of the books are (incorrect)

LSTM

66

