

Text classification

Sergey Aksenov

Higher School of Economics

September 19, 2021

Sentiment analysis

1. **Task:** define expressed opinion of a text (negative, positive, neutral)
2. **Levels:**
 - ▶ classify the whole **document** (is a review positive or negative?)
 - ▶ does **a sentence** express negative or positive opinion? Does **a sentence** express an opinion?
 - ▶ identify what people like and dislike, extract specific aspects
3. **Challenges:** domain specific lexicons, sarcasm, negation, emoticons, abbreviations, slang and noisy user generated data

Intro

1. Know your data: explore the data, look for domain-specific features
2. Pretrain your models on noisy datasets
3. Augment labelled datasets if the classes are imbalanced
4. Use active learning if you want to annotated data

Deep averaging network **DAN**

Neural Bag-of-Words Model

1. **Task**: map an input sequence of tokens X to one of k labels
2. **Composition** function g averages word embeddings:

$$z = g(w \in X) = \frac{1}{|X|} \sum_{w \in X} v_w,$$

where v_w is a word embedding of word w

3. Estimate **probabilities** for each output label:
 $\hat{y} = \text{softmax}(W_s \times z + b)$ and **predict** the label with highest probability
4. **Training**: minimize cross-entropy error: $\sum_{p=1}^k y_p \log \hat{y}_p$

Deep averaging network **DAN**

The intuition is that each layer learns a more abstract representation of the input than the previous one. Add more layers:

$$z_i = g(z_{i-1}) = f(W_i \times z_{i-1} + b_i)$$

Word dropout: drop word tokens' entire word embeddings from the vector average

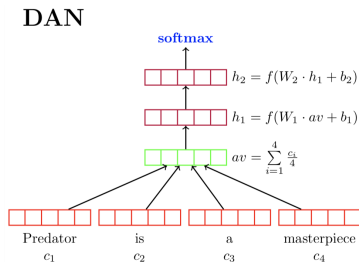


Figure: Deep Averaging Network

Deep averaging network DAN

Sentence	DAN	DRecNN	Ground Truth
a lousy movie that's not merely unwatchable, but also unlistenable	negative	negative	negative
if you're not a prepubescent girl, you'll be laughing at britney spears' movie-starring debut whenever it does n't have you impatiently squinting at your watch	negative	negative	negative
blessed with immense physical prowess he may well be, but ahola is simply not an actor	positive	neutral	negative
who knows what exactly godard is on about in this film, but his words and images do n't have to add up to mesmerize you.	positive	positive	positive
it's so good that its relentless, polished wit can withstand not only inept school productions, but even oliver parker's movie adaptation	negative	positive	positive
too bad, but thanks to some lovely comedic moments and several fine performances, it's not a total loss	negative	negative	positive
this movie was not good	negative	negative	negative
this movie was good	positive	positive	positive
this movie was bad	negative	negative	negative
the movie was not bad	negative	negative	positive

Figure: Predictions of DAN on real (top) and synthetic (bottom) sentences that contain negations and contrastive conjunctions

Deep averaging network **DAN**

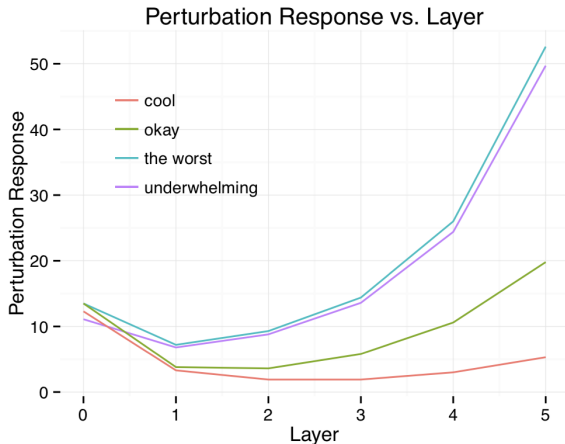


Figure: Perturbation analysis: in the template “the film’s performances were awesome” replace the final word with increasingly negative polarity words (cool, okay, underwhelming, the worst)

Convolution neural networks for text classification

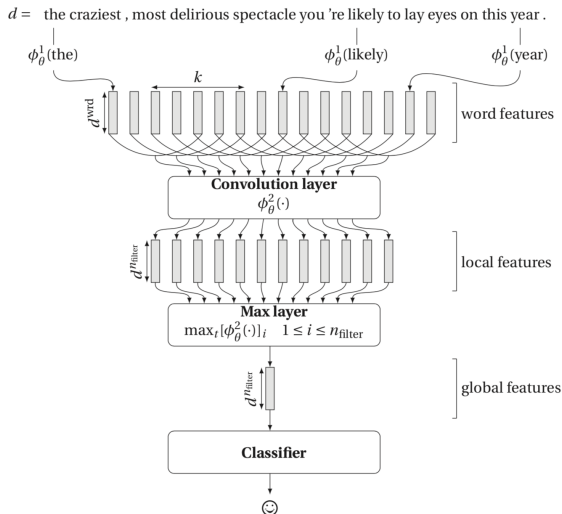


Figure: CNN for text classification

Convolution neural networks for text classification

1. Embedding layer:

$$\phi_{\theta}^1(w_1, w_2, \dots, w_T) = (E_{w_1} E_{w_2}, \dots, E_{w_T}) \in \mathbb{R}^{d_{wrd} \times T}$$

2. Convolutional layer:

$$\phi_{\theta}^2(w_t, \dots, w_{t+k}) = W_2(W_1 \oplus (E_{w_t}, \dots, E_{w_{t+k}}) + b_1) + b_2),$$

$\phi_{\theta}^2 \in \mathbb{R}^{n_{filter}}$ — filter, k — kernel (window) size,

$$W_1 \in \mathbb{R}^{n_h \times (k d_{wrd})}, W_2 \in \mathbb{R}^{n_{filter} \times n_h}$$

3. max pooling: feature-wise reduction

$$[\phi^3]_i = \max_t [\phi_{\theta}^2(\cdot)]_i$$

CNN for sentence classification **kim2014convolutional**

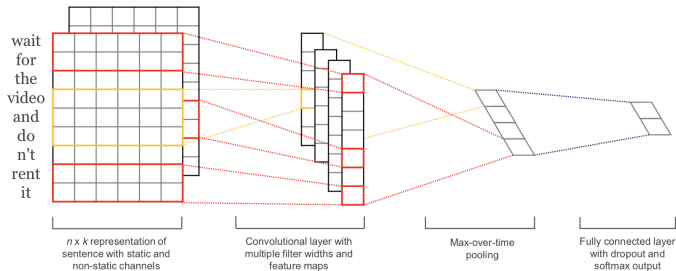


Figure: CNN with two channels for text classification

CNN for sentence classification zhang2015sensitivity

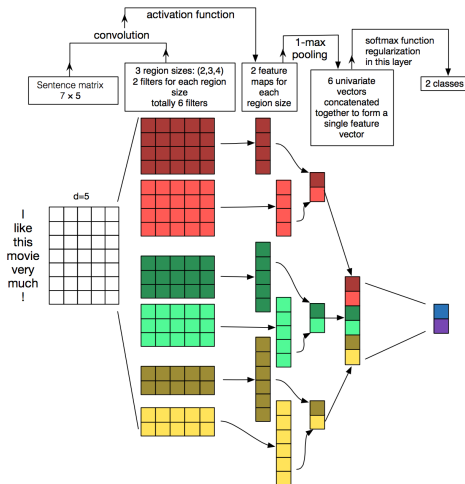


Figure: CNN with multiple channels for text classification

Today

How to improve your classifier?

SMOTE: Synthetic Minority Over-sampling Technique

The minority class is over-sampled by creating synthetic examples:

1. Take each minority class sample
2. Choose random samples k minority class nearest neighbors
3. take the difference between the feature vector (sample) under consideration and its nearest neighbor
4. multiply this difference by a random number between 0 and 1, and add it to the feature vector under consideration

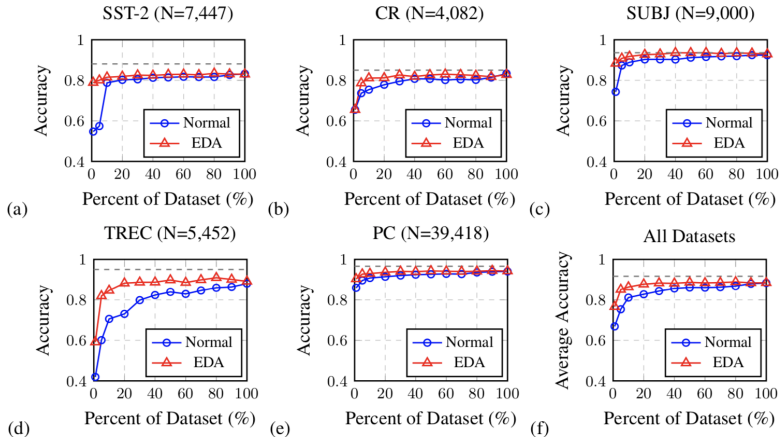
SMOTE is applied only to **feature vectors**, not raw texts!

Python: imbalanced-learn

EDA: Easy Data Augmentation Techniques **wei2019eda**

1. **Synonym Replacement (SR)**: Randomly choose n words from the sentence that are not stop words. Replace each of these words with one of its synonyms chosen at random.
2. **Random Insertion (RI)**: Find a random synonym of a random word in the sentence that is not a stop word. Insert that synonym into a random position in the sentence. Do this n times.
3. **Random Swap (RS)**: Randomly choose two words in the sentence and swap their positions. Do this n times.
4. **Random Deletion (RD)**: Randomly remove each word in the sentence with probability p .

EDA: Easy Data Augmentation Techniques **wei2019eda**



EDA: Easy Data Augmentation Techniques **wei2019eda**

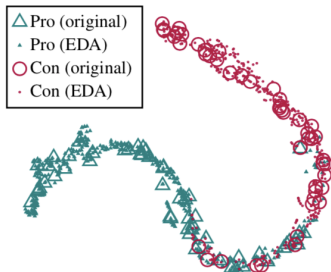


Figure: Latent space visualization of original and augmented sentences in the Pro-Con dataset

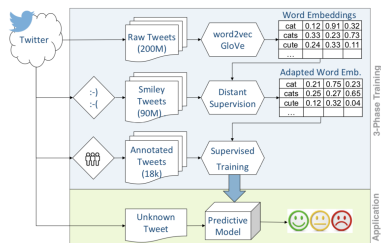
Deep model pretraining **severyn2015**unitn, **deriu2017**leveraging CNN for sentiment analysis

Preprocessing: URLs and usernames were substituted by a replacement token, the text was lowercased and finally tokenized

Creation of word embeddings: the word embeddings are learned on an unsupervised corpus containing 300M tweets

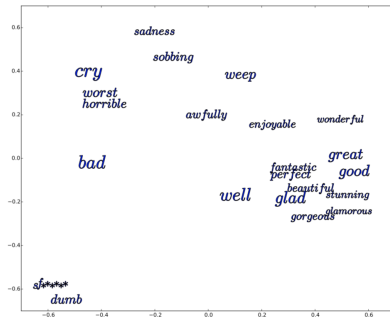
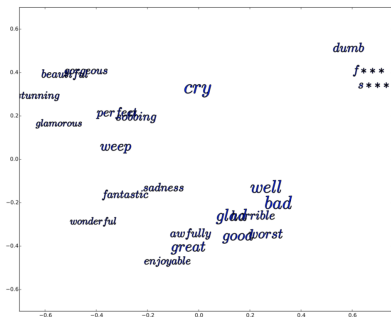
Distant-supervised phase: use emoticons to infer noisy labels on tweets in the training set

Supervised phase: the network is trained on the supervised training data.



Deep model pretraining **severyn2015**unitn, **deriu2017**leveraging

Word embeddings: before and after



AL for text classification with CNNs zhang2017active

Pool-based AL scenario:

1. L – labelled data, U – unlabeled data, $|L| \ll |U|$
2. Train on L , make queries to U to draw examples to be labeled
3. Query strategy: $x^* = \arg \max_{x_i \in U} \phi(x_i; \theta)$
4. Sampling strategy:
 - ▶ Random sampling
 - ▶ Uncertainty sampling:

$$- \sum_k P(y_i = k | x_i, \theta) \log P(y_i = k | x_i, \theta)$$

- ▶ Expected gradient length:

$$\max_{i \in x_i} P(y_i = k | x_i, \theta) \|\nabla J_{E(U)}(< x_i, y_i = k >; \theta)\|$$

AL for text classification with CNNs zhang2017active

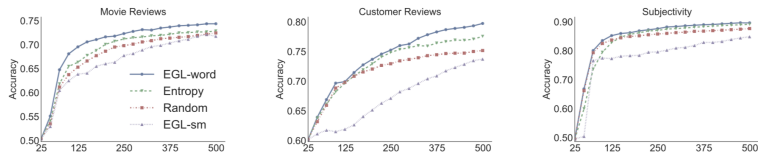


Figure: Number of labeled examples versus accuracy

Conclusions

1. Basic architectures can be improved with pretraining and distant supervision
2. Think of your data and the ways to augment it
3. We really need good representations and metrics for the datasets and tasks

Reading

1. Speech and Language Processing. Daniel Jurafsky, James H. Martin, Ch. 4 [url]
2. Natural Language Processing. Jacob Eisenstein, Ch. 2-4, [GitHub]
3. Neural Networks for NLP. Joav Goldberg, Ch. 15

Reference I