

slidenumbers: true

footer: A. Ferrara. **Language models. Part 3: Statistical language models.** [email](#), [course website](#), [slack](#), [github](#)

Information Retrieval

[fit]Language models

Part 3: Statistical Language Models. Prof. Alfio Ferrara

Master Degree in Computer Science

Master Degree in Data Science and Economics

Skip-Gram Models (I)

The data sparsity problem is mitigated but not solved in **n-gram models** by smoothing or back-off techniques. Indeed, small variations in similar texts can have a large effect on the probability estimation in a n-gram model.

Example (from Casablanca)

- The train for Marseilles leaves at five
- A train to Marseilles and Lyon leaves five o'clock

The two sentences provide basically the same information (there's a train going to Marseilles at five), but they do not have any 2-gram nor 3-gram in common.

This means that a n-gram model will assign very different probabilities to those sentences.

Skip-Gram Models (II)

The problem with the previous sentences is partially due to the presence of *noisy* words that are not relevant for the sentence meaning but change the chain of n-grams (such as *the, a, at, Lyon*)

To solve this, **Skip-Gram models** introduce a notion of *word context* that is not limited to $n - 1$ words as in n-grams, but that allow one to *skip* k words when indexing n-grams. The set S_{kn} of all the k skip n grams of a word sequence w_1, \dots, w_n is then

$$S_{kn} = \left\{ w_{i_1}, w_{i_2}, \dots, w_{i_n} : \sum_{j=2}^n (i_j - i_{j-1} - 1 \leq k, i_j > i_{j-1} \forall j) \right\}$$

Example

- The train for Marseilles leaves at five

```
1 def skip(sequence, n=2, s=2):
2     k_grams = []
3     for i in range(len(sequence)):
4         for z in range(s):
5             seq = [sequence[i]] + sequence[i+z+1:i+z+n]
6             if len(seq) == n and seq not in k_grams:
7                 k_grams.append(tuple(seq))
8     return k_grams
9
10 print(t1)
11 > ['the', 'train', 'for', 'marseilles', 'leaves', 'at', 'five']
12 skip(t1, n=2, s=3)
13 > [('the', 'train'), ('the', 'for'), ('the', 'marseilles'), ('train',
14     'for'),
15     ('train', 'marseilles'), ('train', 'leaves'), ('for', 'marseilles'),
16     ('for', 'leaves'), ('for', 'at'), ('marseilles', 'leaves'),
17     ('marseilles', 'at'), ('marseilles', 'five'), ('leaves', 'at'),
18     ('leaves', 'five'), ('at', 'five')]
```

Example

- The train for Marseilles leaves at five
- A train to Marseilles and Lyon leaves five o'clock

```
1 for e in [x for x in skip(t1, n=2, s=3)
2           if x in skip(t2, n=2, s=3)]:
3     print(e)
4
5 > ('train', 'marseilles')
6 > ('marseilles', 'leaves')
7 > ('leaves', 'five')
```

Continuous Bag of Words (CBOW) model

In the **Skip-gram model**, we aim at predicting the context of a word w given the word w , by collecting all the words *reachable* for w within a tolerance of k skip grams.

$train \rightarrow (train, for), (train, marseilles), (train, leaves)$

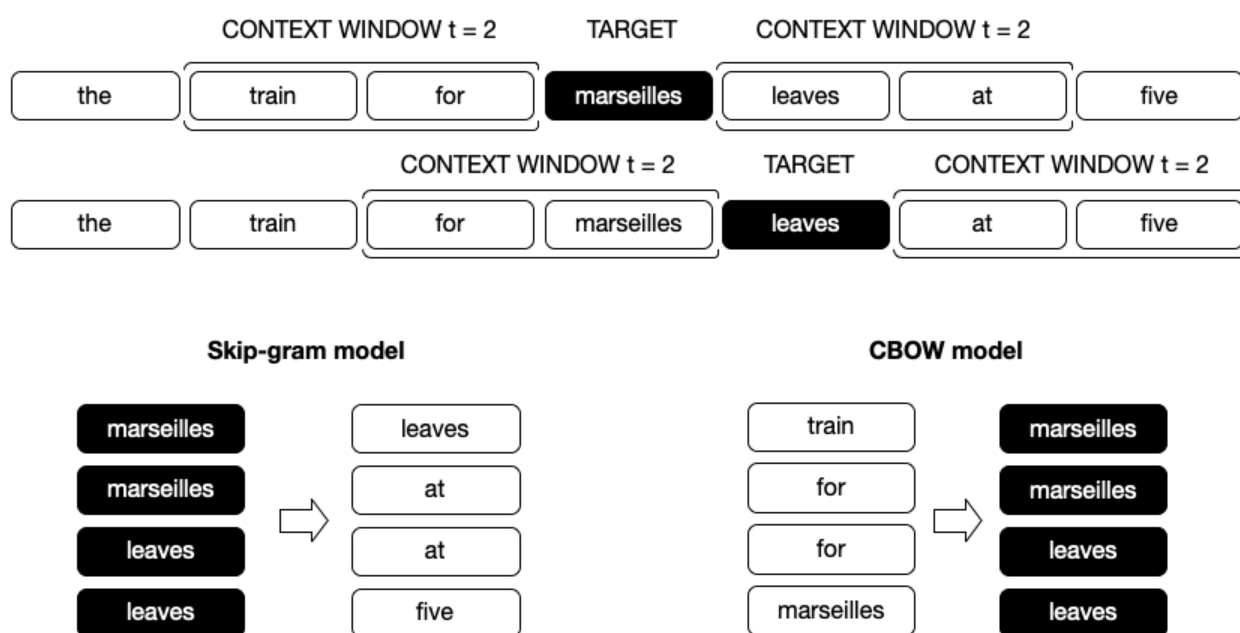
On the contrary, in the **Continuous Bag of Words model**, we aim at predicting a word w given the words appearing in its context

$(train, for), (train, marseilles), (train, leaves) \rightarrow train$

We will see more details about how to use **skip-gram** and **CBOW** models for learning in Part 4 of these lectures.

Word-Context models (I)

We can generalize the idea of exploiting the context of a word by computing a **word-context** matrix based on the counts of words appearing within t other words (before or after) the target word w .



Note that taking into account or not the word order (getting everything as a context or just words after) depends on the goals of the model:

- if we want to mimic the language (such as for generate text), order is relevant
- to focus on words embedding (and their semantics), order is less relevant

Word-Context models (II)

Having defined a notion of context, we can build a context matrix $C = [c_{ij}]$ by counting for each word w_i how many times each word w_j occurs in its context.

A row (or column) of C can be used then as a **vector representation** of w_i (or w_j column-wise) in the space of the other words.

Example (Casablanca)

```

1 from sklearn.metrics.pairwise import cosine_similarity
2
3 sigma = cosine_similarity(C, C)
4 j = list(C.index).index('train')
5 for i, x in sorted(enumerate(sigma[j]), key=lambda y: -y[1])[:5]):
6     print(C.index[i])
7
8 > train, soon, immediately, living, sick

```

Matrix factorization

The obtained matrix $C = [c_{ij}]$ can be factorized into rank- p matrices (using for example **singular value decomposition (SVD)** to obtain a more compact representation as ¹:

$$C \approx UV^T$$

where the rows of U are the embedding we are looking for.

However, the factorization is still dominated by the zero terms. One option is to use stochastic gradient descent in order to change the objective function of factorization to de-emphasize the zero entries by sampling zero entries at a lower rate. Such an approach results in **weighted matrix factorization**.

Weighted matrix factorization

Element-wise, the decomposition problem can be written as:

$$c_{ij} \approx \hat{c}_{ij} = u_i \cdot v_j,$$

where u_i and v_j are the i th and j th rows of matrices U and V , respectively. The goal of gradient descent is to estimate u_i , v_j , and bias parameters b_i and b_j by **minimizing**

$$\sum_i \sum_j (c_{ij} - b_i - b_j - u_i \cdot v_j)^2$$

where the error rate is $e_{ij} = c_{ij} - \hat{c}_{ij}$

Stochastic gradient descent

Initially, U and V are randomly initialized. In subsequent iterations U and V are updated for some randomly selected entries (i, j) in C according to

$$u_i \leftarrow u_i(1 - \alpha) + \alpha e_{ij} v_j$$

$$v_j \leftarrow v_j(1 - \alpha) + \alpha e_{ij} u_i$$

$$b_i \leftarrow b_i(1 - \alpha) + \alpha e_{ij}$$

$$b_j \leftarrow b_j(1 - \alpha) + \alpha e_{ij}$$

where α is the learning rate. Since the update is not executed for all the matrix entries, we can sample the entries in order to reduce the impact of zeros.

Negative sampling

A cycle of stochastic gradient descent samples through all the non-zero entries of C but only a random sample of zero entries.

The random sample size is always k times the number of non-zero entries.

The value of $k > 1$ is a user-driven parameter, and it implicitly controls the weight of positive and negative samples in the factorization.

GloVe embedding

In order to reduce the impact of the high variability of the words frequency, the **GloVe (Global Vectors for Word Representation)** method ² introduces two variations in the factorization process:

1. A dumping factor $\log(1 + c_{ij})$ is used instead of c_{ij} to reduce the impact of highly frequent words
2. The error associated with an entry (i, j) is weighted using a threshold M and a parameter α

$$weight(i, j) = \min \left\{ 1, \frac{c_{ij}}{M} \right\}^\alpha$$

Topic Modeling

Latent Dirichlet allocation

Latent Dirichlet allocation (LDA) is a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modeled as a finite mixture over an underlying set of **hidden variables**, called **topics**. Each topic is, in turn, modeled as an infinite mixture over an underlying set of topic probabilities.

Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM*, 55(4), 77-84.

Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan), 993-1022.

The generative model of LDA

Let $\phi^{(k)}$ be a discrete probability distribution over the vocabulary for the k_{th} topic

Let θ_d be a document distribution over the topics

Let z_i be the topic index for the word w_i

- For $k = 1, \dots, K \rightarrow \phi^{(k)} \sim \text{Dir}(\beta)$
 - For each $d \in D \rightarrow \theta_d \sim \text{Dir}(\alpha)$
 - Generate each word $w_i \in d: z_i \sim \text{Discrete}(\theta_d) w_i \sim \text{Discrete}(\phi^{(z_i)})$

Derivation of the joint distribution

$$p(w, z, \theta, \phi \mid \alpha, \beta) = p(\phi \mid \beta) p(\theta \mid \alpha) p(z \mid \theta) p(w \mid \phi, z)$$

Where, starting from data, we need to estimate z , ϕ , and θ .

- θ_d is a representation of d in the topic space
- z_i represents which topic generated the word w_i
- Each $\phi^{(k)}$ is a matrix $K \times W$ where $\phi_{ij} = p(w_i \mid z_j)$

In order to find the latent variables, we need to solve

$$p(\theta, \phi, z \mid w, \alpha, \beta) = \frac{p(\theta, \phi, z, w \mid \alpha, \beta)}{p(w \mid \alpha, \beta)}$$

Gibbs sampling

Gibbs Sampling is one member of a family of algorithms from the Markov Chain Monte Carlo (MCMC) framework

MCMC algorithms aim to construct a Markov chain that has the target posterior distribution as its stationary distribution. In other words, after a number of iterations of stepping through the chain, sampling from the distribution should converge to be close to sampling from the desired posterior.

Example

To sample from $p(X) = p(x_1, \dots, x_n)$ we can:

Randomly initialize X . Then, for each iterative step t

$$x_1^{t+1} \sim p(x_1 \mid x_2^t, x_3^t, \dots, x_n^t)$$

$$x_2^{t+1} \sim p(x_2 \mid x_1^{t+1}, x_3^t, \dots, x_n^t)$$

...

$$x_n^{t+1} \sim p(x_n \mid x_1^{t+1}, x_2^{t+1}, \dots, x_{n-1}^{t+1})$$

Gibbs sampling for LDA

We need to estimate $\phi^{(k)}$, θ_d and z_i . But note that

$$\theta_{d, z_i} = \frac{n(d, z_i) + \alpha}{\sum_{j \in Z} n(d, z_j) + \alpha},$$

where $n(d, z)$ is the number of times document d is assigned to topic z

$$\phi_{z,w_i} = \frac{n(z,w_i) + \beta}{\sum_{j \in W} n(z,w_j) + \beta},$$

where $n(z, w)$ is the number of times word w is assigned to topic z

Collapsed Gibbs sampling

The previous equations require to just estimate z_i

$$p(z_i \mid z_{-i}, \alpha, \beta, w),$$

where z_{-i} represents all topic allocations **except** for z_i . This results in

$$p(z_i \mid z_{-i}, \alpha, \beta, w) = \frac{p(z_i, z_{-i}, w \mid \alpha, \beta)}{p(z_{-i}, w \mid \alpha, \beta)} \propto p(z_i, z_{-i}, w \mid \alpha, \beta) = p(z, w \mid \alpha, \beta)$$

For a detailed discussion on how $p(z, w \mid \alpha, \beta)$ can be estimated, see this [introduction to topic modeling](#)

Example (I)

To have a meaningful example, we need some more documents

```
1 r = {'$match': {'character.movie.genres': 'war'}}
2 tm = MovieDialogCollection(db_name, collection,
3                             use_pos=False,
4                             drop_stopwords=True,
5                             pipeline=[r, p, s],
6                             pos_filter=['NOUN'],
7                             lemma=True)
```

Example (II)

```

1  from sklearn.decomposition import LatentDirichletAllocation
2
3  lda = LatentDirichletAllocation(n_components=50)
4  I = defaultdict(lambda: defaultdict(lambda: 0))
5  docs = []
6  for doc, tokens in tm.get_tokens():
7      docs.append(doc)
8      for token in tokens:
9          I[doc][token] += 1
10 X = pd.DataFrame(I).T
11 X.fillna(0, inplace=True)
12 X.shape
13 > (7073, 4244)

```

Example (III)

Document distribution over documents

```

1  theta = lda.fit_transform(X)
2  theta.shape
3  > (7073, 50)

```

Word distribution over topics

```

1  phi = lda.components_ / lda.components_.sum(axis=1)[:, np.newaxis]
2  phi.shape
3  > (50, 4244)

```

Example (IV)

Get most relevant documents per topic

```

1  documents = dict([(d, t) for d, t in tm])
2  topic = 1
3  for i, x in sorted(enumerate(theta[:,topic]), key=lambda y: -y[1])[:5]):
4      print(i, x)
5      print(documents[docs[i]], '\n')
6  >
7  4192 0.8599999999999999
8  What do you do in the joint besides pimp?
9
10 6874 0.8366666666666666
11 She'll never forgive me!
12

```



```
13 4222 0.8203732644422244
14 Let's get his clothes off quick.
15
16 4999 0.80399999999999722
17 I sent telegrams, I guess the military traffic held them up.
18
19 7042 0.80399999999999508
20 You will remember the name? Von Scherbach? VON SCHER-BACH!
21
```

Example (V)

Get most relevant words per topic

```
1 for i, x in sorted(enumerate(phi[topic,:]), key=lambda y: -y[1][:5]):
2     print(i, x)
3     print(X.columns[i], '\n')
4 >
5 37 0.347526097853741
6 sir
7
8 40 0.05967234120950544
9 lord
10
11 271 0.04346697422513307
12 thank
13
14 431 0.02282598214856176
15 lady
16
17 771 0.022410885193540186
18 tree
```

1. For matrix factorization and SVD have a look at Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to information retrieval*. Cambridge university press. [Chapter 18](#) ↩

2. Pennington, J., Socher, R., & Manning, C. D. (2014, October). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (pp. 1532-1543). ↩