

Learning Chords Embeddings

Sephora Madjiheurem

September 2, 2016

1 Skip gram: recap

The skip-gram allows to efficiently learn high-quality distributed vector representations that capture precise syntactic and semantic word relationships [1]. We give here a short reminder of how the skip-gram model works.

We define a text as a sequence of words drawn from a finite vocabulary of size W . A word can be described as a “one-hot” vector $w_t \in \{0, 1\}^W$, where exactly one entry is non-zero and the subscript t represent the position of the word in the text. Given a word w_t in a text, define the context of word w_t by $C(w_t) = \{w_{t+j}, -m \leq j \leq m, j \neq 0\}$, where m is the size of the context. We consider the conditional probability of a context given a word $p(w_{t+j}|w_t)$. The goal is to find word representations that are useful for predicting the surrounding words in a sentence. Formally, given a corpus of words of size T and the context of word w_t given by $C(w_t)$, the objective of the skip-gram model is to maximize the average log probability.

$$\frac{1}{T} \sum_{t=1}^T \sum_{w \in C(w_t)} \log p(w|w_t). \quad (1)$$

The parametrization of the skip-gram model uses the architecture depicted in Figure 1 (from Mikolov et al. [1]). In this model, each output is computed using softmax to obtain the posterior distribution of context words:

$$p(w_{t+j}|w_t) = \frac{\exp(v_{w_{t+j}}^T v_{w_t})}{\sum_{w=1}^W \exp(v_w^T v_{w_t})}, \quad (2)$$

where $-m \leq j \leq m, j \neq 0$, m is the size of the training context, v_w is the vector representation for w , and W is the number of words in the vocabulary.

Detailed derivations and explanations of the parameter learning for this original skip-gram model can be found in [3].

Because the computation cost of objective (1) is proportional to W which can be very large when working with text training data, Mikolov et al. [1] use instead an efficient approximation, known as negative sampling (see [2] for details).

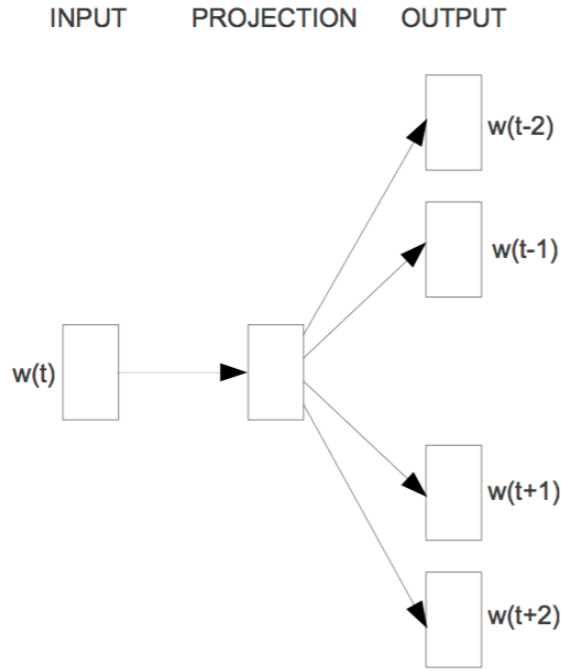


Figure 1: Skip-gram (Mikolov et al. [1])

1.1 Chord2Vec

Similarly to a text, we define a piece of music as a series of chords. A chord is a subset of notes drawn from a finite set of size N and can be represented

by a binary vector $\mathbf{c} = \{c_1, c_2, \dots, c_N\} \in \{0, 1\}^N$. We denote by \mathcal{N} the set of all notes and \mathcal{C} the set of all chords.

To adapt the skip-gram model to music data there are a few points that need to be considered:

1. A text can be represented as a sequence of words, where each word can be represented as a "one-hot" vector. In the case of music, we need a "many-hot" vector to represent a chord, as more than one note can be heard simultaneously.
2. The set of notes is smaller than the vocabulary considered when working with text data.

1.1.1 Model 1: bilinear

The first proposed linear model makes a conditional independence assumption between the notes in a context chord \mathbf{c} given a chord \mathbf{d} , i.e

$$p(\mathbf{c} = \mathbf{c}' | \mathbf{d}) = \prod_{i=1}^N p(c_i = c'_i | \mathbf{d}). \quad (3)$$

Using weights matrices $M \in \mathbb{R}^{D \times N}$ and $\tilde{M} \in \mathbb{R}^{N \times D}$, we can define a score function $h : \mathcal{N} \times \mathcal{C} \mapsto \mathbb{R}$:

$$h(i, \mathbf{d}) = \tilde{M}_{(:,i)} \frac{M\mathbf{d}}{|\mathbf{d}|_1} = \tilde{M}_{(:,i)} \mathbf{v}_d, \quad (4)$$

where $\tilde{M}_{(:,i)}$ denotes the i 'th row of \tilde{M} and \mathbf{v}_d is the D -dimensional vector representation of chord \mathbf{d} .

We can then use the sigmoid function to model the conditional probabilities:

$$p(c_i = c'_i | \mathbf{d}) = \begin{cases} \sigma(h(i, \mathbf{d})) = \frac{\exp(h(i, \mathbf{d}))}{1 + \exp(h(i, \mathbf{d}))} & c'_i = 1 \\ 1 - \sigma(h(i, \mathbf{d})) = \sigma(-h(i, \mathbf{d})) = \frac{1}{1 + \exp(h(i, \mathbf{d}))} & c'_i = 0 \end{cases} \quad (5)$$

The objective of this model is to maximize the log probability for each context chord \mathbf{d} :

$$\sum_{\mathbf{c}' \in C(\mathbf{d})} \log p(\mathbf{c} = \mathbf{c}' | \mathbf{d}) = \sum_{\mathbf{c}' \in C(\mathbf{d})} \sum_{i=1}^N \log p(c_i = c'_i | \mathbf{d}), \quad (6)$$

Where $C(\mathbf{d})$ denotes the neighborhood of chord \mathbf{d} .

Update equations

We define our loss function :

$$E := \log p(\mathbf{c} = \mathbf{c}' | \mathbf{d}), \quad (7)$$

and we also define

$$E_i := \log p(c_i = c'_i | \mathbf{d}), \quad (8)$$

where $\mathbf{c}' = \{c'_1, \dots, c'_N\} \in C(\mathbf{d})$.

Using stochastic gradient descent with learning rate $\eta > 0$, the weight update equations for the weights given by \tilde{M} and M respectively are:

$$\begin{aligned} \tilde{m}_{ij}^{(t+1)} &= \tilde{m}_{ij}^{(t)} - \eta \frac{\partial E_i}{\partial \tilde{m}_{ij}}, \\ m_{jk}^{(t+1)} &= m_{jk}^{(t)} - \eta \frac{\partial E}{\partial m_{jk}}, \end{aligned}$$

where \tilde{m}_{ij} and m_{jk} are elements of \tilde{M} and M respectively. We now compute the gradients. First observe that

$$\begin{aligned} E_i &= c'_i \log p(c_i = 1 | d) + (1 - c'_i) \log p(c_i = 0 | d) \\ &= c_i u_i - \log(1 + \exp(u_i)), \end{aligned}$$

where we write $u_i = h(i, \mathbf{d})$ to ease the reading. Then

$$\frac{\partial E_i}{\partial \tilde{m}_{ij}} = \frac{\partial E_i}{\partial u_i} \frac{\partial u_i}{\partial \tilde{m}_{ij}} = (c'_i - \frac{\exp(u_i)}{1 + \exp(u_i)}) v_{dj}$$

where v_{dj} denotes the j 'th element of \mathbf{v}_d .

$$\frac{\partial E}{\partial m_{jk}} = \frac{\partial E}{\partial v_{dj}} \frac{\partial v_{dj}}{\partial m_{jk}} = \left(\sum_{i=1}^N \frac{\partial E_i}{\partial u_i} \frac{\partial u_i}{\partial v_{dj}} \right) \frac{\partial v_{dj}}{\partial m_{jk}} = \sum_{i=1}^N (c'_i - \frac{\exp(u_i)}{1 + \exp(u_i)}) \tilde{m}_{ij} d_k$$

2 Model 2

3 Sequence autoencoding

3.1 Sequence-to-sequence: Recap

Sequence-to-sequence models allow to learn a mapping of input sequences of varying lengths to output sequences also of varying lengths [4]. It uses a neural network known as RNN Encoder-Decoder. Figure 2 depicts the model architecture. An LSTM encoder is used to map the input sequence

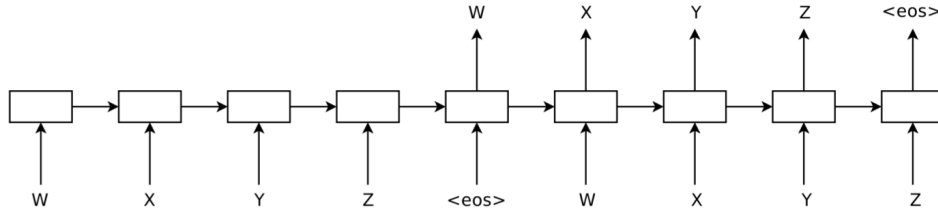


Figure 2: Sequence-to-sequence (Sutskever et al. [4])

to a fixed length vector, and another LSTM decoder is then used to extract the output sequence from this vector. The general goal is to estimate $p(y_1, \dots, y_{T'} | x_1, \dots, x_T)$, where x_1, \dots, x_T and $y_1, \dots, y_{T'}$ are the input and output sequences respectively, and T' and T need not to be equal.

The objective is given by:

$$\max_{\theta} \frac{1}{|\mathcal{T}|} \sum_{(X,Y) \in \mathcal{T}} \log p(Y|X, \theta), \quad (9)$$

where Y is a correct output given the input X and \mathcal{T} is the training set and θ is the set of the model parameters. The encoder and decoder are jointly trained to maximize the objective according to θ .

The model estimates the conditional probability $p(y_1, \dots, y_{T'} | x_1, \dots, x_T)$ by first obtaining the fixed-length vector representation v of the input sequence (given by the last state of the LSTM encoder) and then computing the probability of $y_1, \dots, y_{T'}$ with the LSTM decoder:

$$p(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1}) \quad (10)$$

3.2 Chord-to-chord

The sequence-to-sequence model can be used to learn embeddings for chords, by training the model to learn the context of a given chord.

In this setting, a chord is represented as a sequence of notes in some fixed ordering: $c \subseteq N$, where N is the ordered set of all possible notes. A chord can have an arbitrary size.

The goal is then to estimate $p(n_1^{(j)}, \dots, n_T^{(j)} | n_1, \dots, n_{T'})$, where the n 's are in N , $c_t = n_1, \dots, n_{T'}$ is an input chord and $c_{t+j} = n_1^{(j)}, \dots, n_T^{(j)}$ is the j^{th} neighbor of c .

If $C(c_t)$ denotes the set of chords that are in the neighborhood of the chord c_t , then the objective in (9) can be written as:

$$\max_{\theta} \frac{1}{W} \sum_{t=1}^W \sum_{c \in C(c_t)} \log p(c | c_t, \theta), \quad (11)$$

Where W is the size of the corpus of chords in the training data.

3.3 Chord-to-chords

An alternative to estimating the probability of single context chord is to estimate the probability of the entire neighborhood $p(C(c_t) | c_t)$, by combining all context chord in one longer output sequence.

The corresponding objective is then given by:

$$\max_{\theta} \frac{1}{W} \sum_{t=1}^W \log p(C(c_t) | c_t, \theta), \quad (12)$$

4 Results

Table 1: Addam, no regul, no attention

layers	units	batch size	train	validation	test	epochs
3	512	128	2.3684	3.0333	3.0416	1
3	512	128	2.3684	3.0333	3.0416	1
3	512	256	2.2201	3.0178	3.0035	3
3	1024	128	2.4899	3.0629	3.0821	6
3	1024	256	2.3052	3.0323	3.0383	1

Table 2: GD, no regul, no attention

layers	units	batch size	train	validation	test	epochs
3	512	128	2.2873	2.9922	2.9900	7
3	512	256	2.4412	2.9697	2.9716	9
3	1024	128	2.2295	3.0082	2.9739	5
3	1024	256	2.3736	2.9807	2.9636	8
3	512	128	2.2873	2.9922	2.9900	7
3	512	256	2.4412	2.9697	2.9716	9
3	1024	128	2.2295	3.0082	2.9739	5
3	1024	256	2.3736	2.9807	2.9636	8

Table 3: Seq2seq trained with gradient descent, perplexities

layers	units	batch size	train	validation	test	epochs
1	512	64	2.3081	2.9367	2.9488	5
1	512	128	2.3988	2.9258	2.9380	7
1	512	64	2.3770	2.9876	3.0109	8
1	512	128	2.4345	2.9802	2.9919	8
1	1024	64	2.2457	2.9483	2.9632	7
1	1024	128	2.3319	2.9297	2.9423	8
1	1024	64	2.3335	2.9873	3.0027	9
1	1024	128	2.3794	2.9867	3.0107	5
2	512	64	2.1998	2.9065	2.9234	6
2	512	128	2.3058	2.8968	2.9043	9
2	512	64	2.2248	2.9282	2.9479	5
2	512	128	2.3130	2.9241	2.9415	6
2	1024	64	2.1488	2.9222	2.9344	5
2	1024	128	2.2539	2.9090	2.9297	7
2	1024	64	2.1638	2.9274	2.9443	6
2	1024	128	2.2548	2.9302	2.9544	8

Table 4: Best seq 2 seq model, log likelihood

data set	train	validation	test	epochs
JSB Chorales	-0.8424	-1.1053	-1.1090	3
Piano-midi.de	-1.0026	-1.8653	-1.7795	1
MuseData	-1.3788	-1.4773	-1.5243	7
Nottingham	-0.4874	-0.5117	-0.4974	10
all data sets	-1.0739	-1.2012	-1.2207	7

Table 5: Linear models, log likelihood

D	train	validation	test
512	-9.6941	-9.5372	-9.7731
1024	-9.6934	-9.5364	-9.773

Table 6: Baseline models, log likelihood

	test set
marginal probability	-12.2349
random	-60.9969

References

- [1] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.
- [2] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 3111–3119. Curran Associates, Inc., 2013.
- [3] Xin Rong. word2vec parameter learning explained. *CoRR*, abs/1411.2738, 2014. URL <http://arxiv.org/abs/1411.2738>.
- [4] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. *CoRR*, abs/1409.3215, 2014. URL <http://arxiv.org/abs/1409.3215>.