

ON THE DIMENSIONALITY OF EMBEDDINGS FOR SPARSE FEATURES AND DATA

Maxim Naumov

Facebook, 1 Hacker Way, Menlo Park, CA, 94025

ABSTRACT

In this note we discuss a common misconception, namely that embeddings are always used to reduce the dimensionality of the item space. We show that when we measure dimensionality in terms of information entropy then the embedding of sparse probability distributions, that can be used to represent sparse features or data, may or not reduce the dimensionality of the item space. However, the embeddings do provide a different and often more meaningful representation of the items for a particular task at hand. Also, we give upper bounds and more precise guidelines for choosing the embedding dimension.

1 BACKGROUND

The mapping of concepts, objects or items into a vector space is called an embedding. The representation of n items in d dimensional vector space is used in many applications, which can broadly be split into two categories.

The first category is characterized by models that produce embeddings. In this case the embeddings represent the information obtained by the model about the item space. If the model represents a probability distribution over items in the dataset then we can interpret the output as an embedding of this probability distribution.

For instance, singular value decomposition Golub & Van Loan (2012) is used in latent semantic analysis/indexing (LSA/LSI) Dumais (2005) to produce low-rank approximations of the original data. In this case a low-rank approximation of the word-document matrix $C \in \mathbb{R}^{m \times n}$ for m words and n documents, where each row corresponds to a word and each column corresponds to a document, while non-zero entries represent the occurrence of words in documents, can be written as

$$C \approx WV^T \quad (1)$$

where the matrices $W \in \mathbb{R}^{m \times d}$ and $V \in \mathbb{R}^{n \times d}$ can be interpreted as embeddings of words and documents in a d dimensional space. A similar interpretation can be given to matrix factorization for collaborative filtering, with the caveat that empty matrix entries are unknown (rather than being 0.0) Koren et al. (2009); Frolov & Oseledets (2017).

Also, deep learning models based on auto-encoders Sedhain et al. (2015), multi-label classification Bengio et al. (2010) and neural machine translation (NMT) Neubig (2017) can be seen as generating a probability distribution over a set of classes/objects. Let this probability distribution be represented as a vector $\mathbf{p} \in \mathbb{R}^n$. Notice that the embeddings produced by the models, such as embeddings of target language words in NMT, represent an embedding of this probability distribution over n items into d dimensional space

$$\mathbf{v} = V^T \mathbf{p} \quad (2)$$

for some embedding matrix $V \in \mathbb{R}^{n \times d}$. Notice that during training we often start with an arbitrary probability distribution and therefore vector \mathbf{p} is dense. However, when training converges towards the end it is common for the probability distribution to concentrate more at particular items and therefore vector \mathbf{p} becomes sparse (and if needed it can be quantized to eliminate a long tail of small values). It will become clear in this note that the information the embedding V needs to represent will vary greatly based on the signature of vector \mathbf{p} . In particular, when \mathbf{p} is very sparse and its elements can be quantized to a few values the embedding dimension d needed to represent the information in it can be relatively small.

The second category is characterized by models that consume embeddings. In this case embeddings represent information obtained from the input features or data. It is common for these input features to be sparse, such as the history of user clicks on webpages or words present in a post. The sparse features are often represented by a list of integer indices that select items from a larger sequence/set (in contrast to dense features that are represented by a single floating point number). They can also be represented by a sparse vector \mathbf{p} as will be shown in the next section.

For instance, neural collaborative filtering He et al. (2017), wide & deep Cheng et al. (2016) and deep & cross recommendation systems Wang et al. (2017) all use embeddings to process sparse input features. The advantage of using embeddings instead of lists of sparse items is that we can measure distance between them in a more meaningful way. Also, notice that embedding elements represent the sparse features in some abstract space relevant to the model at hand, while integers simply represent an ordering of the input data.

The natural language processing models Kalchbrenner & Blunsom (2013); Sutskever et al. (2014) may fall somewhere in between these two categories. In particular, NMT models Neubig (2017) often use two embeddings one representing words in the source and another in the target language. On one hand, the source embedding can be seen as a sparse feature consumed by the model. It is characterized by a list of indices that selects words used in the input sentence. On the other hand, the target embedding can be seen as a representation of the probability distribution over words in the target language.

In this note we focus on the embedding of sparse vectors \mathbf{p} that can be used to represent sparse features and data that belong to the second category. We point out that the implications of choosing a particular value of a hyper-parameter d are not theoretically well understood. The choice is usually based on empirical experiments or resource limitations, such as compute and memory available on hardware platforms Park et al. (2018). A recent work attempted to explain the sizing of the embedding vectors based on pairwise inner product dissimilarity metric Yin & Shen (2018).

We propose an alternative approach based on the entropy information measure. We leverage the ideas of Tishby et al. (1999); Shwartz-Ziv & Tishby (2017) as well as Traub & Woziakowski (1980); Pinkus (1985); Donoho (2006), but in contrast to them we do not attempt to explain the behavior of neural networks or find ways to compress the parameters/data. We use the information measure to discuss the dimensionality and provide guidelines on sizing of the embedding vectors, i.e. we provide roofline and more precise models for selecting the dimensionality d .

2 INTRODUCTION TO EMBEDDINGS

Let n items be mapped into a d dimensional vector space. The vectors corresponding to n items are often organized into an embedding table, which can be seen as a tall matrix $V \in \mathbb{R}^{n \times d}$ with $n \gg d$ and that can be written as

$$V^T = [\mathbf{v}_1, \dots, \mathbf{v}_n] \quad (3)$$

where vector \mathbf{v}_i corresponds to i -th item.

A sparse feature is characterized by a list of integer indices, which can be represented as item lookups with different signature in the embedding table.

The item lookup with a single index is often encoded as a dense matrix-vector multiplication

$$\mathbf{v}_i = V^T \mathbf{e}_i \quad (4)$$

where vector $\mathbf{e}_i \in \mathbb{R}^n$ and

$$\mathbf{e}_i^T = [0, \dots, 1, \dots, 0] \quad (5)$$

has 1 in i -th position and 0 everywhere else. It is often referred to as a one-hot encoding vector.

Notice that we can select multiple items with some weights in a single lookup and express it as

$$\mathbf{u} = V^T \mathbf{a} = a_{i_1} \mathbf{v}_{i_1} + \dots + a_{i_k} \mathbf{v}_{i_k} \quad (6)$$

where vector $\mathbf{a} \in \mathbb{R}^n$ and

$$\mathbf{a}^T = [0, \dots, a_{i_1}, \dots, a_{i_k}, \dots, 0] \quad (7)$$

has weight $a_i \neq 0$ at $i = i_1, \dots, i_k$ and 0 everywhere else Jia et al. (2014); Paszke et al. (2017).

Further, we can generalize it to multiple lookups where each selects multiple items with some weight and encode it as a dense matrix-sparse matrix multiplication

$$U^T = V^T A \quad (8)$$

where sparse matrix $A \in \mathbb{R}^{n \times r}$ and

$$A = [\mathbf{a}_1, \dots, \mathbf{a}_r] \quad (9)$$

is composed of multiple vectors \mathbf{a}_j each corresponding to a single lookup¹ with non-zero elements corresponding to the items being selected. The output matrix $U \in \mathbb{R}^{r \times d}$ is the result of r lookups.

This setup is often used to state that the embedding vectors project n dimensional items space into d dimensional embedding vectors. In the next section we will examine this statement in more detail and show that it can be misleading and therefore lead to incorrect conclusions.

3 THE DIMENSIONALITY OF EMBEDDINGS

Notice that when we discussed space dimensions in the previous section, we never took into consideration the precise vector element type and how much information we can represent with it. Let us now incorporate it into our analysis by measuring the cardinality of the set it can describe and the information associated with it.

Recall that the entropy of an information source s is given by

$$H(s) = - \sum_{i=1}^n p_i \log_2 p_i \quad (10)$$

where p_i is the probability of i -th symbol to be communicated Shannon & Weaver (1949).

Notice that the embedding vector has d elements, each with s bits. Therefore, it could represent

$$g = 2^{ds} \quad (11)$$

values and if we interpret it as an information source $\mathbf{v}_.$ its entropy is

$$H(\mathbf{v}_.) = - \sum_{i=1}^g p_i \log_2 p_i \quad (12)$$

where p_i denotes the probability of i -th value being selected. Hence, if $p_i = 1/g$ is uniform then

$$H(\mathbf{v}_.) = ds \quad (13)$$

3.1 SINGLE LOOKUP OF SINGLE ITEM

Let a single lookup of a single item be done as shown in (4). In this case the i -th embedding vector represents the i -th item from the item space.

Notice that because we represent this lookup as a binary vector \mathbf{e}_i in (5) in n dimensional space, it can describe only n items and if we interpret it as an information source $\mathbf{e}_.$ its entropy is

$$H(\mathbf{e}_.) = - \sum_{i=1}^n p_i \log_2 p_i \quad (14)$$

where p_i denotes the probability of i -th item being selected. Note that if $p_i = 1/n$ is uniform then

$$H(\mathbf{e}_.) = \log_2 n \quad (15)$$

Therefore, the dimensionality of the item and embedding spaces, as measured by the information they can represent, can be compared by looking at (12) and (14). For instance, under the assumption of uniform probability of item/value selection, if $n = 20\text{M}$ and using (15) we have $H(\mathbf{e}_.) \approx 24.3$ then a single 32-bit element is enough to represent information in the item space.

¹Notice that the vector subscript in here has a different meaning than before. In (9) it denotes the j -th lookup, while in (5) it denoted the j -th item being selected.

3.2 SINGLE LOOKUP OF MULTIPLE ITEMS

Let a single lookup of multiple items be done as shown in (6). In this case the i -th embedding vector represents combinations of i -th and other items from the item space.

Let us first consider the case when this lookup is a binary vector in n dimensional space, i.e. $a_i \in \{0, 1\}$ for $i = 1, \dots, k$ in (7). In this case, the vector can describe h items, with

$$h = \binom{n}{k} = \frac{n!}{k!(n-k)!} \quad (16)$$

and if we interpret it as an information source its entropy is

$$H(\mathbf{a}_.) = - \sum_{i=1}^h p_i \log_2 p_i \quad (17)$$

where p_i is the probability of a combination being selected. Note that if $p_i = 1/h$ is uniform then

$$H(\mathbf{a}_.) = \log_2 h \quad (18)$$

$$\approx n \log_2 \frac{n}{n-k} + k \log_2 \frac{n-k}{k} + \quad (19)$$

$$+ \frac{1}{6} \log_2 \frac{n + 4n^2 + 8n^3}{((n-k) + 4(n-k)^2 + 8(n-k)^3)(k + 4k^2 + 8k^3)} - \frac{1}{2} \log_2 \pi$$

where we have used properties of logarithms and Ramanujan's approximation Ramanujan (1988)

$$\ln n! \approx n \ln n - n + \frac{1}{6} \ln(n + 4n^2 + 8n^3) + \frac{1}{2} \ln \pi \quad (20)$$

that are described in more detail in the appendix.

Therefore, the dimensionality of the item and embedding spaces, as measured by the information they can represent, can be compared by looking at (12) and (17). For instance, under the assumption of uniform probability of item/value selection, if $n = 20M$, $k = 100$ and using (19) we have $H(\mathbf{a}_.) \approx 1756.3$ then a vector with 64 elements with 32-bit per element is enough to represent information in the item space.

3.3 SINGLE LOOKUP OF MULTIPLE ITEMS WITH WEIGHTS

Let us now consider the case when this lookup is a vector in n dimensional space, with each of the a_i represented by t bits for $i = 1, \dots, k$ in (7). Notice that in this case the analysis of the previous section remains the same, but we can now select 2^t values for each position in the lookup. Therefore, the vector can describe h' items, with

$$h' = 2^{tk} h \quad (21)$$

Note that if we interpret it as an information source, under the assumption $p_i = 1/h'$ is uniform, then its entropy is

$$H(\mathbf{a}_.) = \log_2 h' = tk + \log_2 h \quad (22)$$

$$\approx tk + n \log_2 \frac{n}{n-k} + k \log_2 \frac{n-k}{k} + \quad (23)$$

$$+ \frac{1}{6} \log_2 \frac{n + 4n^2 + 8n^3}{((n-k) + 4(n-k)^2 + 8(n-k)^3)(k + 4k^2 + 8k^3)} - \frac{1}{2} \log_2 \pi$$

Therefore, the dimensionality of the item and embedding spaces, as measured by the information they can represent, can be compared by looking at (12) and (21). For instance, under the assumption of equal probability of item/value selection, if $n = 20M$, $k = 100$, $t = 16$ and using (23) we have $H(\mathbf{a}_.) \approx 3356.3$ then a vector with 128 elements with 32-bit per element is enough to represent information in the item space.

3.4 THE EFFECTS OF MINI-BATCH

We point out that the use of mini-batch of size r does not affect the dimensionality because each vector in the mini-batch is treated independently. This can be observed in (8) where r lookups from matrix A in embedding V generate r results in matrix U .

4 UPPER BOUNDS AND SIZING GUIDELINES FOR EMBEDDINGS

Notice that embeddings corresponding to sparse features do not necessarily reduce the dimensionality of data, as measured by the information they can represent. The dimensionality reduction depends on the size of the embedding vectors and the detailed signature of the input lookup vectors.

The upper bound (roofline) dimensionality of different types of lookups is provided in (15), (19) and (23) which can be compared with dimensionality of embeddings given in (13). We have tabulated a few sample lookup signatures and embedding dimensions for comparison in Tab. 1.

Input Lookup Signature			Embedding Dimensions			
n	k	H(.)	d (s=8)	d (s=16)	d (s=32)	H(.)
1M	1	19.9	4	2	1	32
10M	1	23.2	4	2	1	32
100M	1	26.5	4	2	1	32
1M	10	163.0	24	12	6	192
10M	10	196.3	28	14	7	224
100M	10	229.5	32	16	8	256
1M	100	1324.1	168	84	42	1344
10M	100	1656.3	208	104	52	1664
100M	100	1988.5	252	126	63	2016
1M	1000	9958.0	1248	624	312	9984
10M	1000	13281.2	1664	832	416	13312
100M	1000	16603.3	2076	1038	519	16608

Table 1: Entropy of a sample set of lookup signatures and embedding dimensions

Although, the function (18) achieves its maximum for larger n as shown in Fig. 1a, notice that for the sample lookup signatures in Tab. 1 the choice of k has a much larger effect on the entropy function, as shown on Fig. 1b. We point out that the function is plotted on a log scale in both plots.

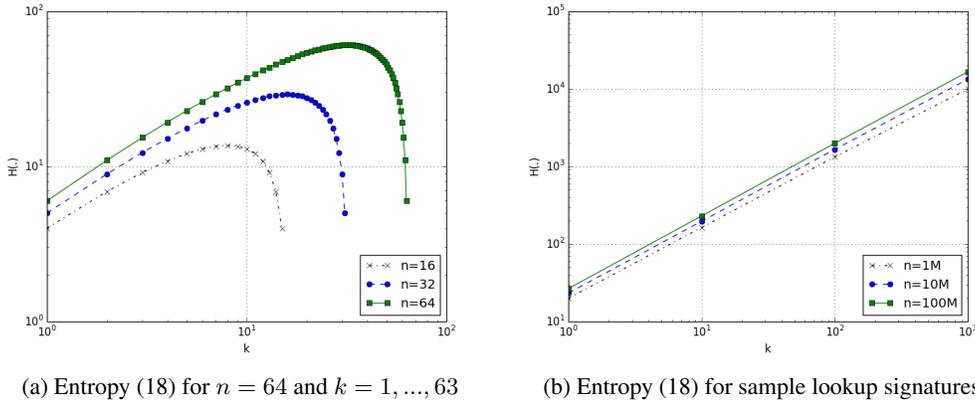


Figure 1: Entropy of a sample set of lookup signatures in log scale

In practice not all combinations will be exercised by the item lookups. We can make a pass over the input dataset to discover how many combinations are actually present and how often the same combinations are repeated. Then, we can use formulas (14) and (17) to estimate the dimensionality of the input dataset from the information perspective. Ultimately, the more precise measure of dimensionality of the input dataset can guide a better choice of the embedding vector size and element type to be used in the model.

Once again we point out that although embeddings do not necessarily reduce the dimensionality of lookups, they do provide a different and very useful representation of them. Note that the embedding values are learned during training. Therefore, the embeddings are found based on what the lookups represent/mean in some abstract space relevant to the model at hand.

5 CONCLUSION

We have discussed embeddings corresponding to dense and sparse probability distributions. We have analyzed the dimensionality of the item lookups and embeddings corresponding to sparse features and data. We have shown that using embeddings does not necessarily reduce the dimension of the lookups, as measured in terms of the information they can represent. We have also provided rooflines and more precise guidelines for choosing the embedding size for the dataset and model at hand.

ACKNOWLEDGEMENTS

The author would like to thank Aleksandr Ulanov, Dheevatsa Mudigere, Satish Nadathur and Misha Smelyanskiy for thoughtful comments as well as Mark Tygert and Juan Miguel Pino for insightful discussion on the use of embeddings in different applications and NMT models.

REFERENCES

- Samy Bengio, Jason Weston, and David Grangier. Label embedding trees for large multi-class tasks. In *Proc. Advances in Neural Information Processing Systems*, 2010.
- Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishu Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. Wide & deep learning for recommender systems. In *Proc. 1st Workshop on Deep Learning for Recommender Systems*, pp. 7–10, 2016.
- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. MIT Press, 3rd edition, 2009.
- David L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52:1289–1306, 2006.
- Susan T. Dumais. Latent semantic analysis. *Annual Review of Information Science and Technology*, 38:188 – 230, 2005.
- Evgeny Frolov and Ivan Oseledets. Tensor methods and recommender systems. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 7(3):e1201, 2017.
- Gene H. Golub and Charles F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, 4th edition, 2012.
- Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural collaborative filtering. In *Proc. 26th Int. Conf. World Wide Web*, pp. 173–182, 2017.
- Yanqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *CoRR*, 2014. URL <https://arxiv.org/abs/1408.5093>.
- Nal Kalchbrenner and Phil Blunsom. Recurrent continuous translation models. In *Proc. 2013 Conf. Empirical Methods in Natural Language Processing*, pp. 1700–1709, 2013.
- Ekaterina Karatsuba. On the asymptotic representation of the Euler gamma function by Ramanujan. *Journal of Computational and Applied Mathematics*, 135:225–240, 2001.
- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 8:30–37, 2009.
- Lucien Le Cam. The central limit theorem around 1935. *Statistical Science*, 1:78–96, 1986.
- Graham Neubig. Neural machine translation and sequence-to-sequence models: a tutorial. *CoRR*, 2017. URL <https://arxiv.org/abs/1703.01619>.

-
- Jongsoo Park, Maxim Naumov, Protonu Basu, Summer Deng, Aravind Kalaiah, Daya Khudia, James Law, Parth Malani, Andrey Malevich, Satish Nadathur, Juan Pino, Martin Schatz, Alexander Sidorov, Viswanath Sivakumar, Andrew Tulloch, Xiaodong Wang, Yiming Wu, Hector Yuen, Utku Diril, Dmytro Dzhulgakov, Kim Hazelwood, Bill Jia, Yangqing Jia, Lin Qiao, Vijay Rao, Nadav Rotem, Sungjoo Yoo, and Mikhail Smelyanskiy. Deep learning inference in facebook data centers: Characterization, performance optimizations and hardware implications. *CoRR*, 2018. URL <https://arxiv.org/abs/1811.09886>.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. *Proc. Advances in Neural Information Processing Systems*, 2017.
- Alan Pinkus. *n-widths in approximation theory*. Springer-Verlag, Berlin, 1985.
- Srinivasa Ramanujan. *The Lost Notebook and Other Unpublished Papers*. Springer, Berlin, 1988.
- Dan Romik. Stirlings approximation for $n!$: The ultimate short proof? *The American Mathematical Monthly*, 107:556–557, 2000.
- Suvash Sedhain, Aditya K. Menon, Scott Sanner, and Lexing Xie. Autorec: Autoencoders meet collaborative filtering. In *Proc. 24th Int. Conf. World Wide Web*, 2015.
- Claude Shannon and Warren Weaver. *The Mathematical Theory of Communication*. University of Illinois Press, 1949.
- Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. *CoRR*, 2017. URL <https://arxiv.org/abs/1703.00810>.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Proc. Advances in Neural Information Processing Systems*, pp. 3104–3112, 2014.
- Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method. In *Proc. 37th Allerton Conference on Communication, Control and Computing*, pp. 368–377, 1999.
- Joe F. Traub and Henryk Woziakowski. *A general theory of optimal algorithms*. Academic, New York, 1980.
- Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. Deep & cross network for ad click predictions. In *Proc. ADKDD*, pp. 12, 2017.
- Zi Yin and Yuanyuan Shen. On the dimensionality of word embeddings. In *Proc. Neural Information Processing Systems*, 2018.

6 APPENDIX

6.1 LOGARITHMIC IDENTITIES

We list a few handy logarithmic identities in this section Cormen et al. (2009). The multiplication and division operations are equivalent to

$$\log_b xy = \log_b x + \log_b y \quad (24)$$

and

$$\log_b x/y = \log_b x - \log_b y \quad (25)$$

respectively. Also, we can change logarithm from base b to c using the following

$$\log_b x = \log_c x / \log_c b \quad (26)$$

6.2 APPROXIMATION THEORY

There exist a number of approximations to the factorial function, including Stirling's approximation Le Cam (1986); Romik (2000)

$$n! \sim \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \quad (27)$$

and more accurate Ramanujan's approximation Ramanujan (1988); Karatsuba (2001)

$$n! \sim \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \left(1 + \frac{1}{2n} + \frac{1}{8n^2}\right)^{1/6} \quad (28)$$

Then, using (28) it follows that

$$\ln n! \sim n \ln n - n + \frac{1}{6} \ln(n + 4n^2 + 8n^3) + \frac{1}{2} \ln \pi \quad (29)$$

Finally, notice that using logarithmic properties and (29) we can write

$$\begin{aligned} \log_2 \frac{n!}{k!(n-k)!} &= \log_2 n! - \log_2(n-k)! - \log_2 k! \\ &\approx \frac{1}{\ln 2} \left(n \ln n - n + \frac{1}{6} \ln(n + 4n^2 + 8n^3) + \frac{1}{2} \ln \pi \right) \\ &\quad - \frac{1}{\ln 2} \left((n-k) \ln(n-k) - (n-k) + \frac{1}{6} \ln((n-k) + 4(n-k)^2 + 8(n-k)^3) + \frac{1}{2} \ln \pi \right) \\ &\quad - \frac{1}{\ln 2} \left(k \ln k - k + \frac{1}{6} \ln(k + 4k^2 + 8k^3) + \frac{1}{2} \ln \pi \right) \\ &= n \log_2 \frac{n}{n-k} + k \log_2 \frac{n-k}{k} + \\ &\quad + \frac{1}{6} \log_2 \frac{n + 4n^2 + 8n^3}{((n-k) + 4(n-k)^2 + 8(n-k)^3)(k + 4k^2 + 8k^3)} - \frac{1}{2} \log_2 \pi \end{aligned} \quad (30)$$

We will take advantage of this expression in this note.