

# An Investigation of Lie Groups in Machine Learning

Juncheng Wan, 120033910148, jorgenwan@gamil.com

June 7th, 2021

## 1 Introduction

As a graduate student of Computer Science and Engineering, I currently focus on machine learning and data mining<sup>1</sup>. In this area, I use machine learning algorithms to approximately solve non-convex problems, such as approximating complex conditional distribution.

One example is **Machine Translation**, which aiming at translating sentence from one language to another language. Specifically, given source sentence  $(x_1, \dots, x_m)$ , the machine needs to translate it into the target sentence  $(y_1, \dots, y_n)$ , where  $x_i, y_j, (1 \leq i \leq m, 1 \leq j \leq n)$  are all random variables. Thus, this problem is equal to approximating the distribution  $p(y_1, \dots, y_n | x_1, \dots, x_m)$  with parameterized distribution  $q_\theta(y_1, \dots, y_n | x_1, \dots, x_m)$ , where  $\theta$  indicates the model parameters. However, this problem is difficult for three reasons:

1. The maximal sentence length could be long. For example, in the common translation benchmarks, such as LDC Chinese-English task and WMT14 English-German task, there are sentences with length larger than 512.
2. The vocabulary size could also be large. The vocabulary size of English is from 50,000 to 10,000. Thus, we need to deal with at least  $50000^{512}$  possibilities.
3. The complex phrase composition, grammatical rules, and syntactic structure.

The above problems are not only for machine translation, but also for other natural language processing problems, such as information retrieval, sequence labeling, etc.

From my perspective, utilizing various symmetries of language have two advantages:

1. Reducing the size of the space to be modeled. For example, if there are sentences, the only difference between them is *the place of time adverbials*. Then, I think they are almost equal in meaning, due to the symmetry of the syntactic structure.
2. Providing interpretability of the model. Some machine learning algorithms are considered black-box models and lack interpretability for parameters of the submodule, such as neural networks.

As Lie group is a group of symmetries where the symmetries are continuous, I investigate Lie group in machine learning in this survey out of my interest. I will first give some definitions and concepts of

---

<sup>1</sup>The website of our lab is <http://apex.sjtu.edu.cn/>.

mathematics nouns and popular machine learning models. Then, I will give some basic results of Lie group in machine learning.

## 2 Definition

In this section, I clarify basic concepts and definitions.<sup>2</sup>

### 2.1 Feature Map Vector Space

**Definition 1.** Let  $S$  be some set, the **feature map vector space** is defined as:

$$V \triangleq \{f | f : S \rightarrow \mathbb{R}\} \quad (1)$$

First, I want to introduce the concept of feature map. This is a space of scalar-valued functions on the set  $S$ , representing the features of each point in  $S$ . For example, in computer vision,  $S$  could be  $\mathbb{R}^2$  which indicating the 2D coordinates and each  $f \in V$  could be a gray value function that assigning a gray value for each pixel between  $[0, 255]$ . In natural language processing,  $S$  could be  $\{1, 2, \dots, 512\}$  indicating the discrete positions of words in the sentence and  $\text{range}(f) \in [50000] \subseteq \mathbb{R}$  indicating the words assignment, where  $[50000]$  is the coded vocabulary.

### 2.2 $G$ -Equivariant/Invariant

**Definition 2.** Let  $G$  be a group,  $V_1, V_2$  be two feature map vector spaces. A map  $\Phi : V_1 \rightarrow V_2$  is  **$G$ -equivariant** with respect to actions  $\rho_1, \rho_2$  of  $G$  acting on  $V_1, V_2$  respectively if:  $\Phi[\rho_1(g)f] = \rho_2(g)\Phi[f]$  for any  $g \in G, f \in V_1$ .

**Definition 3.** A map  $\Phi : V_1 \rightarrow V_2$  is  **$G$ -invariant** with respect to actions  $\rho_1, \rho_2$  of  $G$  acting on  $V_1, V_2$  respectively if  $\Phi$  is  $G$ -equivariant and  $\rho_2$  is the identity map for any  $g \in G$ .

Then,  $G$ -equivariant, a concept from representation theory in undergraduate, is necessary. Because in machine learning I am interested in those operation  $\Phi$  that is equivariant under the group action. For example, there are rotations and transitions in images. I use  $\Phi$  to map the shallow color features contained in  $V_1$  to useful textural features  $V_2$  and require that they are the same despite the order of rotation or transition of the object.

### 2.3 Regular Representation

**Definition 4.** A **regular representation**  $\pi$  acting on feature map vector space  $V$  is defined as follows:

$$[\pi(g_\theta)(f)](g_\phi) \triangleq f(g_\theta^{-1}g_\phi) \quad (2)$$

---

<sup>2</sup>Though some concepts are fundamental for students of mathematics, they are fresh for me. Thus, I also write them down.

## 2.4 Lifting

Before introducing lifting, I need to define the homogeneous space:

**Definition 5.** Let the group  $G$  acts on a set  $\mathcal{X}$  via action  $\rho$ , if  $\forall x, x' \in \mathcal{X}, \exists g \in G : \rho(g)x = x'$ , then the action is **transitive** and  $\mathcal{X}$  is a **homogeneous space** with respect to  $G$ .

This means that all elements of  $\mathcal{X}$  are connected by the action.

**Definition 6.** Let  $\mathcal{X}$  be a homogeneous space with respect to some group  $G$ ,  $F : \mathcal{X} \rightarrow \mathbb{R}$  be a feature map vector space. The **lifting**  $\mathcal{L}$  maps  $f_{\mathcal{X}} \in F$  (supported on  $\bigcup_{i=1}^n \{x_i\} \subset \mathcal{X}$ ) to  $\mathcal{L}[f_{\mathcal{X}}]$  (supported on  $\bigcup_{i=1}^n s(x_i)H \subset G$ ) such that:

$$\mathcal{L}[f_{\mathcal{X}}](g) \triangleq f_{\mathcal{X}}(x_i) \text{ for } g \in s(x_i)H \quad (3)$$

## 2.5 Lie groups in Machine Learning

I introduce some common Lie groups used in machine learning and their parameterization. Let  $\exp : \mathfrak{g} \rightarrow G$  and  $\log : G \rightarrow \mathfrak{g}$  be the Exponential maps and Log maps in the courses. Let the map  $\nu : \mathfrak{g} \rightarrow \mathbb{R}^d$  represent a function to extracts the free parameters from the Lie algebra element. I give the Lie groups and their corresponding free parameters as follows:

- $G = T(n), t \in \mathbb{R}^n, \nu[\log(t)] = t.$
- $G = SO(2), R = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \in \mathbb{R}^{2 \times 2}, \nu[\log(R)] = \theta = \arctan(R_{10}/R_{01}).$
- $G = SE(2), R = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \in \mathbb{R}^{2 \times 2}, t \in \mathbb{R}^2, \nu[\log(tR)] = \begin{bmatrix} t' \\ \theta \end{bmatrix},$  where  $t' = V^{-1}t, V = \begin{bmatrix} a & -b \\ b & a \end{bmatrix}, a \triangleq \frac{\sin \theta}{\theta}, b \triangleq \frac{1 - \cos \theta}{\theta}.$
- $G = SO(3), R \in \mathbb{R}^{3 \times 3}, t \in \mathbb{R}^3, \nu[\log(R)] = \nu \left[ \frac{\theta}{2 \sin \theta} (R - R^{\top}) \right] = \frac{\theta}{2 \sin \theta} \begin{bmatrix} R_{21} - R_{12} \\ R_{02} - R_{20} \\ R_{10} - R_{01} \end{bmatrix},$  where  $\cos \theta = \frac{\text{Tr}(R) - 1}{2}.$
- $G = SE(3), R \in \mathbb{R}^{3 \times 3}, t \in \mathbb{R}^3, \nu[\log(tR)] = \begin{bmatrix} t' \\ r' \end{bmatrix},$  where  $t' = V^{-1}t, r' = \nu[\log(R)], V = I + \frac{1 - \cos \theta}{\theta^2} (R - R^{\top}) + \frac{\theta - \sin \theta}{\theta^3} (R - R^{\top})^2.$

## 2.6 Lie Convolution Model

First, I define the conventional convolution operation.

**Definition 7.** Let  $I \in \{f|f : \mathbb{R}^3 \rightarrow \mathbb{R}\} \triangleq V_1$  be the feature map of a 2D image with channels (the third dimation),  $K \in \{f|f : \mathbb{R}^3 \rightarrow \mathbb{R}\} \triangleq V_2$  be a kernel function,  $n_H, n_W, n_C$  be the height, width, channels of

image, feature map vector space  $V_3 \triangleq \{f|f : \mathbb{R}^2 \rightarrow \mathbb{R}\}$ . Then, the **convolution** operator  $\text{conv} : V_1 \times V_2 \rightarrow V_3$  is defined as follows:

$$[\text{conv}(I, K)](x, y) = \sum_{i=1}^{n_H} \sum_{j=1}^{n_W} \sum_{k=1}^{n_C} K(i, j, k) I(x+i-1, y+j-1, k) \quad (4)$$

This operator  $\text{conv}$  with backpropagation algorithm [5] achieves great improvement in image processing.

As a natural generalization of conventional convolution operation, group equivariant convolution is proposed first in paper [2], which enjoys a substantially higher degree of weight sharing than regular convolution layers.

**Definition 8.** The **group equivariant convolution**  $\Psi : \mathcal{I}_U \rightarrow \mathcal{I}_U$  is defined as :

$$[\Psi f](g) \triangleq \int_G \psi(g'^{-1}g) f(g') dg' \quad (5)$$

where  $\psi : G \rightarrow \mathbb{R}$  is the convolutional filter and the integral is defined with respect to the left Haar measure of  $G$ .

Recently, Lie convolution is proposed, such as papers [3]. The definition is rather similar to the group equivariant convolution. The discretized integral is defined as follows:

$$h(u_i) = (k \hat{*} f)(u_i) = \frac{1}{n_i} \sum_{j \in \text{nbhd}(i)} k(v_j^{-1} u_i) f(v_j) \quad (6)$$

where  $u_i$  is the group elements,  $k$  is the kernel function,  $\text{nbhd}(i) = \{v : d(u_i, v) \leq r\}$  is the local neighborhood of  $u_i$ . The distance  $d(u, v) \triangleq \|\log(u^{-1}v)\|_F$  is induced by Frobenius norm.

## 2.7 Lie Self Attention Model

The Self Attention model consists of Scaled Dot-Product Attention and Multi-Head Attention. Assume the input is a sequence of tokens. The Scaled Dot-Product Attention maps linearly those tokens to keys', queries', and values' spaces by matrices  $Q, K, V$ . Assume that queries and keys are of dimension  $d_k$ , and values of dimension  $d_v$ . It does the dot products of the query with all keys, divide each by  $\sqrt{d_k}$ , and applies a softmax function to obtain the weights on the values. The matrix of Scaled Dot-Product Attention is as:

$$\text{Scaled Dot-Product Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad (7)$$

As it is beneficial to linearly project the queries, keys and values  $h$  times with different, learned linear projections to  $d_k$ ,  $d_k$  and  $d_v$  dimensions, respectively, Multi-Head Attention is used. Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions.

$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O \\ \text{where head}_i &= \text{Scaled Dot-Product Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned} \quad (8)$$

where the projections are parameter matrices  $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$  and  $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$ .

Lie Self Attention model [6] is using the lifted group elements to substitute the operations between tokens. For example, I can define the followings:

- dot-product: for any two elements  $g, g' \in G$  as  $\frac{1}{\sqrt{d_v}} (W^Q f(g))^\top W^K f(g') \in \mathbb{R}$ .
- normalization of weights:  $\text{softmax}(\{\alpha_f(g, g')\} \mid g' \in G_f)$
- multihead equivariant self-attention:

$$V^m(g) = \int_{G_f} w_f(g, g') W^{V,m} f(g') dg' \in \mathbb{R}^{d_v/M} \quad (9)$$

From my view, the main difference is that the operation is on group elements, which represent some kind of symmetry.

## 2.8 Lie Auto Encoder Model

Auto-encoder is a generative neural network model whose encoder compresses the inputs into hidden vectors and the decoder restores the original input. Commonly, the hidden vector is modeled as Gaussian distribution.

Some paper [4] investigate the tangent space of a special Lie group manifold: upper triangular positive definite affine transform matrices (UTDATs). The application is that non-degenerate Gaussian distributions are isomorphic to UTDATs.

As UTDATs form a Lie group, one can work on the tangent spaces (Lie algebras) to make them suitable for machine learning models. Then project back to Lie group by exponential mapping.

## 3 Main Results

In this section, I collect some results of Lie group in machine learning. My work is more like a research porters.

**Theorem 1.** [1] *The function composition  $f \circ f_K \circ \dots \circ f_1$  of several equivariant functions  $f_k, k \in 1, 2, \dots, K$  followed by an invariant function  $f$ , is an invariant function.*

*Proof.* Consider group representations  $\pi_1, \dots, \pi_K$  that act on  $f_1, \dots, f_K$  respectively, and representation  $\pi_0$  that acts on the input space of  $f_1$ . If each  $f_k$  is equivariant with respect to  $\pi_k$ ;  $\pi_{k1}$  such that  $f_k \circ \pi_{k1} = \pi_k \circ f_k$ , and  $f$  is invariant such that  $f \circ \pi_k = f$ , then we have:

$$\begin{aligned} f \circ f_k \circ \dots \circ f_1 \circ \pi_0 &= f \circ f_k \circ \dots \circ \pi_1 \circ f_1 \\ &\vdots \\ &= f \circ \pi_k \circ f_k \circ \dots \circ f_1 \\ &= f \circ f_k \circ \dots \circ f_1 \end{aligned} \quad (10)$$

hence  $f \circ f_K \circ \dots \circ f_1$  is invariant.  $\square$

**Theorem 2.** [6] *The group equivariant convolution  $\Psi : \mathcal{I}_U \rightarrow \mathcal{I}_U$  defined as:  $[\Psi f](g) \triangleq \int_G \psi(g'^{-1}g) f(g') dg'$  is equivariant with respect to the regular representation  $\pi$  of  $G$  acting on  $\mathcal{I}_U$  as  $[\pi(u)(f)](g) \triangleq f(u^{-1}g)$ .*

*Proof.* Use the invariance of the left Haar measure.

$$\begin{aligned}
\Psi[\pi(u)f](g) &= \int_G \psi(g'^{-1}g) [\pi(u)f](g') dg' \\
&= \int_{uG} \psi(g'^{-1}g) f(u^{-1}g') dg' \\
&= \int_G \psi(g'^{-1}u^{-1}g) f(g') dg' \\
&= [\Psi f](u^{-1}g) \\
&= [\pi(u)[\Psi f]](g)
\end{aligned} \tag{11}$$

□

**Theorem 3.** *The lifting layer  $\mathcal{L}$  is equivariant with respect to the representation  $\pi$ .*

*Proof.* Note  $\mathcal{L}[\pi(u)f_{\mathcal{X}}](g) = \mathbf{f}_i$  for  $g \in s(ux_i)H$  and  $[\pi(u)\mathcal{L}[f_{\mathcal{X}}]](g) = \mathcal{L}[f_{\mathcal{X}}](u^{-1}g) = \mathbf{f}_i$  for  $g \in us(x_i)H$ . Hence  $\mathcal{L}[\pi(u)f_{\mathcal{X}}] = \pi(u)\mathcal{L}[f_{\mathcal{X}}]$  because the two cosets are equal:  $s(ux_i)H = us(x_i)H, \forall u \in G$ . □

In natural language processing, self attention [7] is a powerful model that achieves the best performance in amounts of real-life applications. For more details, please refer to Section 2 or the original paper.

**Theorem 4.** [6] *LieSelfAttention is equivariant with respect to the regular representation  $\pi$ .*

*Proof.* Let  $\mathcal{I}_U = \mathcal{L}(G, \mathbb{R}^D)$  be the space of unconstrained functions  $f : G \rightarrow \mathbb{R}^D$ . We can define the regular representation  $\pi$  of  $G$  acting on  $\mathcal{I}_U$  as follows:

$$[\pi(u)f](g) = f(u^{-1}g) \tag{12}$$

$f$  is defined on the set  $G_f = \bigcup_{i=1}^n s(x_i)H$  (i.e. union of cosets corresponding to each  $x_i$ ). Note  $G_{\pi(u)f} = uG_f$ , and  $G_f$  does not depend on the choice of section  $s$ .

Note that for all provided choices of  $k_c$  and  $k_l$ , we have:

$$\begin{aligned}
k_c([\pi(u)f](g), [\pi(u)f](g')) &= k_c(f(u^{-1}g), f(u^{-1}g')) \\
k_l(g^{-1}g') &= k_l((u^{-1}g)^{-1}(u^{-1}g'))
\end{aligned} \tag{13}$$

Hence for all choices of  $F$ , we have that

$$\begin{aligned}
\alpha_{\pi(u)f}(g, g') &= F(k_c([\pi(u)f](g), [\pi(u)f](g')), k_l(g^{-1}g')) \\
&= F(k_c(f(u^{-1}g), f(u^{-1}g')), k_l((u^{-1}g)^{-1}u^{-1}g')) \\
&= \alpha_f(u^{-1}g, u^{-1}g')
\end{aligned} \tag{14}$$

We thus prove equivariance for the below choice of LieSelfAttention  $\Phi : \mathcal{I}_U \rightarrow \mathcal{I}_U$  that uses softmax normalisation, but a similar proof holds for constant normalisation. Let  $A_f(g, g') \triangleq \exp(\alpha_f(g, g'))$ , hence Equation (10) also holds for  $A_f$ :

$$\begin{aligned}
[\Phi f](g) &= \int_{G_f} w_f(g, g') f(g') dg' \\
&= \int_{G_f} \frac{A_f(g, g')}{\int_{G_f} A_f(g, g'') dg''} f(g') dg'
\end{aligned} \tag{15}$$

Hence:

$$\begin{aligned}
w_{\pi(u)f}(g, g') &= \frac{A_{\pi(u)f}(g, g')}{\int_{G_{\pi(u)f}} A_{\pi(u)f}(g, g'') dg''} \\
&= \frac{A_f(u^{-1}g, u^{-1}g')}{\int_{uG_f} A_f(u^{-1}g, u^{-1}g'') dg''} \\
&= \frac{A_f(u^{-1}g, u^{-1}g')}{\int_{G_f} A_f(u^{-1}g, g'') dg''} \\
&= w_f(u^{-1}g, u^{-1}g')
\end{aligned} \tag{16}$$

Then we can show that  $\Phi$  is quivariant with respect to the representation  $\pi$  as follows:

$$\begin{aligned}
\Phi[\pi(u)f](g) &= \int_{G_{\pi(u)f}} w_{\pi(u)f}(g, g') [\pi(u)f](g') dg' \\
&= \int_{uG_f} w_f(u^{-1}g, u^{-1}g') f(u^{-1}g') dg' \\
&= \int_{G_f} w_f(u^{-1}g, g') f(g') dg' \\
&= [\Phi f](u^{-1}g) \\
&= [\pi(u)[\Phi f]](g)
\end{aligned} \tag{17}$$

□

**Theorem 5.** [3] Let operator  $\mathcal{K} : \mathbb{L}_2(X) \rightarrow \mathbb{L}_2(Y)$  be linear and bounded, let  $X, Y$  be homogeneous spaces on which Lie group  $G$  act transitively, and  $d\mu_X$  a Radon measure on  $X$ , then:

1.  $\mathcal{K}$  is a kernel operator, i.e.,  $\exists \tilde{k} \in \mathbb{L}_1(Y \times X) : (\mathcal{K}f)(y) = \int_X \tilde{k}(y, x) f(x) d\mu_X$
2. under the  $G$ -equivariance constraint of Eq. (3) the map is defined by a one-argument kernel:

$$\tilde{k}(y, x) = \frac{d\mu_X(g_y^{-1} \odot x)}{d\mu_X(x)} k(g_y^{-1} \odot x) \tag{18}$$

for any  $g_y \in G$  such that  $y = g_y \odot y_0$  for some fixed origin  $y_0 \in Y$

3. if  $Y \equiv G/H$  is the quotient of  $G$  with  $H = \text{Stab}_G(y_0) = \{g \in G \mid g \odot y_0 = y_0\}$  then the kernel is constrained via:

$$\forall_{h \in H}, \forall_{x \in X} : k(x) = \frac{d\mu_X(g_y^{-1} \odot x)}{d\mu_X(x)} k(h^{-1} \odot x) \tag{19}$$

*Proof.* 1. It follows from Dunford-Pettis Theorem, that if  $\mathcal{K}$  is linear and bounded it is an integral operator.

2. The left-equivariance constraint then imposes bi-left-invariance of the kernel  $\tilde{k}$  as follows, where  $\forall_{g \in G}$  and  $\forall_{f \in \mathbb{L}_2(X)}$ :

$$\begin{aligned}
(\mathcal{K} \circ \mathcal{L}_g^{G \rightarrow \mathbb{L}_2(X)})(f) &= (\mathcal{L}_g^{G \rightarrow \mathbb{L}_2(Y) \circ \mathcal{K}})(f) \Leftrightarrow \\
\int_X \tilde{k}(y, x) f(g^{-1}x) dx &= \int_X \tilde{k}(g^{-1}y, x) f(x) dx \stackrel{\text{in r.h.s. integral}}{\Leftrightarrow} x \leftarrow g^{-1}x \\
\int_X \tilde{k}(y, x) f(g^{-1}x) dx &= \int_X \tilde{k}(g^{-1}y, g^{-1}x) f(g^{-1}x) d(g^{-1}x) \Leftrightarrow \\
\int_X \tilde{k}(y, x) f(g^{-1}x) dx &= \int_X \tilde{k}(g^{-1}y, g^{-1}x) f(g^{-1}x) \frac{1}{|\det g|} dx
\end{aligned} \tag{20}$$

Since the final equation holds for all  $f \in \mathbb{L}_2(X)$  we obtain:

$$\forall_{g \in G} : \quad \tilde{k}(y, x) = \frac{1}{|\det g|} \tilde{k}(g^{-1}y, g^{-1}x) \tag{21}$$

Furthermore, since  $G$  acts transitively on  $Y$  we have that  $\forall_{y, y_0 \in Y} \exists_{g_y \in G}$  such that  $y = g_y y_0$  and thus

$$\tilde{k}(y, x) = \tilde{k}(g_y y_0, x) = \frac{1}{|\det g_y|} \tilde{k}(y_0, g_y^{-1}x) =: \frac{1}{|\det g_y|} k(g_y^{-1}x) \tag{22}$$

for every  $g_y \in G$  such that  $y = g_y y_0$  with arbitrary fixed origin  $y_0 \in Y$ .

3. Every homogeneous space  $Y$  of  $G$  can be identified with a quotient group  $G = H$ . Choose an origin  $y_0 \in Y$  s.t.  $\forall_{h \in H} : h y_0 = y_0$ , i.e.,  $H = \text{Stab}_G y_0$ , then

$$\tilde{k}(y_0, x) = \tilde{k}(h y_0, x) \Leftrightarrow k(x) = \frac{1}{|\det h|} k(h^{-1}x) \tag{23}$$

□

**Theorem 6.** [4] Let  $\mathbf{G}_i$  be the UTDAT and  $\mathbf{g}_i$  be the corresponding vector in its tangent Lie algebra at the standard Gaussian. Then:

$$\mathbf{G}_i = \begin{bmatrix} \sigma_{i1} & & & \mu_{i1} \\ & \sigma_{i2} & & \mu_{i2} \\ & & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} \tag{24}$$

$$\mathbf{g}_i = \begin{bmatrix} \phi_{i1} & & & \theta_{i1} \\ & \phi_{i2} & & \theta_{i2} \\ & & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}, \tag{25}$$

where:

$$\phi_{ik} = \log(\sigma_{ik}) \tag{26}$$

$$\theta_{ik} = \frac{\mu_{ik} \log(\sigma_{ik})}{\sigma_{ik} - 1} \tag{27}$$

*Proof.* By the definition of UTDAT, we can straightforwardly get the first equation.

Let  $\mathbf{H} = \mathbf{G}_i - \mathbf{I}$ , Using the series form of matrix logarithm, we have:

$$\begin{aligned} \mathbf{g}_i &= \log(\mathbf{G}_i) \\ &= \log(\mathbf{I} + \mathbf{H}) \\ &= \sum_{t=1}^{\infty} (-1)^{t-1} \frac{\mathbf{H}^t}{t}. \end{aligned} \tag{28}$$

By substituting  $\mathbf{H}$  into it, we get the second equation and the following:

$$\begin{aligned} \phi_{ik} &= \sum_{t=1}^{\infty} (-1)^{t-1} \frac{(\sigma_{ik} - 1)^t}{t} \\ &= \log(\sigma_{ik}) \end{aligned} \tag{29}$$

and:

$$\begin{aligned} \theta_{ik} &= \sum_{t=1}^{\infty} (-1)^{t-1} \frac{\mu_{ik} (\sigma_{ik} - 1)^{t-1}}{t} \\ &= \frac{\mu_{ik} \log(\sigma_{ik})}{\sigma_{ik} - 1} \end{aligned} \tag{30}$$

Alternatively, after we identify  $\mathbf{g}_i$  has the form as in the second equation, we can derive the exponential mapping by the definition of matrix exponential:

$$\begin{aligned} \mathbf{G}_i &= \exp(\mathbf{g}_i) = \sum_{t=0}^{\infty} \frac{\mathbf{g}_i^t}{t!} \\ &= \begin{bmatrix} \sum_{t=0}^{\infty} \frac{\phi_{i1}^t}{t!} & & \theta_{i1} \sum_{t=1}^{\infty} \frac{\phi_{i1}^{t-1}}{t!} \\ & \ddots & \vdots \\ 0 & \dots & 1 \end{bmatrix} \\ &= \begin{bmatrix} e^{\phi_{i1}} & \frac{\theta_{i1}}{\phi_{i1}} \left( \sum_{t=0}^{\infty} \frac{\phi_{i1}^t}{t!} - 1 \right) \\ & \ddots & \vdots \\ 0 & \dots & 1 \end{bmatrix} \\ &= \begin{bmatrix} e^{\phi_{i1}} & \frac{\theta_{ik} (e^{\phi_{i1}} - 1)}{\phi_{i1}} \\ & \ddots & \vdots \\ 0 & \dots & 1 \end{bmatrix}. \end{aligned} \tag{31}$$

□

## 4 Conclusion

To sum up, in this project, I investigate the applications of Lie group in machine learning. The applications are mainly in the improvement of convolution operator in computer vision, self attention operator in natural language processing, and auto-encoder. Most of the time, I am like a Theorem porter. In my eyes, I need to see the code for better understanding the detailed implementation.

## References

- [1] Benjamin Bloem-Reddy and Yee Whye Teh. “Probabilistic Symmetries and Invariant Neural Networks”. In: *J. Mach. Learn. Res.* 21 (2020), 90:1–90:61.
- [2] Taco Cohen and Max Welling. In: *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*. Vol. 48. 2016, pp. 2990–2999.
- [3] Marc Finzi et al. “Generalizing Convolutional Neural Networks for Equivariance to Lie Groups on Arbitrary Continuous Data”. In: *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*. Vol. 119. PMLR, 2020, pp. 3165–3176.
- [4] Liyu Gong and Qiang Cheng. “Lie Group Auto-Encoder”. In: *CoRR* abs/1901.09970 (2019).
- [5] Ian Goodfellow et al. *Deep learning*. Vol. 1. 2. MIT press Cambridge, 2016.
- [6] Michael Hutchinson et al. “LieTransformer: Equivariant self-attention for Lie Groups”. In: *CoRR* abs/2012.10885 (2020).
- [7] Ashish Vaswani et al. “Attention is All you Need”. In: *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*. 2017, pp. 5998–6008.