

---

# Learning Invariant Weights in Neural Networks (Supplementary material)

---

Tycho F.A. van der Ouderaa<sup>1</sup>

Mark van der Wilk<sup>1</sup>

<sup>1</sup>Imperial College London, UK

## APPENDIX A: DETAILED DERIVATION OF VARIATIONAL INVERENCE

Applying Variational Inference (VI) [?], we maximise the marginal likelihood w.r.t. parameters  $\theta = \text{vec}(\mathbf{W}_2)$  by minimizing the  $D_{\text{KL}}(\cdot||\cdot)$ -divergence between approximate posterior  $q(\mathbf{W}_2|\mu, \Sigma)$  and true posterior distribution of weights  $p(\mathbf{W}_2|\mathcal{D})$ , equivalent to maximizing the evidence lower bound (ELBO) denoted by  $\mathcal{L}$ :

$$\begin{aligned}
& \arg \min_{\mu, \Sigma} D_{\text{KL}}(q(\mathbf{W}_2|\mu, \Sigma)||p(\mathbf{W}_2|\mathcal{D})) \\
&= \arg \min_{\mu, \Sigma} \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} \left[ \log \frac{q(\mathbf{W}_2|\mu, \Sigma)}{p(\mathbf{W}_2|\mathcal{D})} \right] \\
&= \arg \min_{\mu, \Sigma} \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} \left[ \log \frac{q(\mathbf{W}_2|\mu, \Sigma)}{p(\mathbf{W}_2)p(\mathcal{D}|\mathbf{W}_2)} \right] + \log p(\mathcal{D}) \\
&= \arg \min_{\mu, \Sigma} \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} \left[ \log \frac{q(\mathbf{W}_2|\mu, \Sigma)}{p(\mathbf{W}_2)p(\mathcal{D}|\mathbf{W}_2)} \right] \\
&= \arg \min_{\mu, \Sigma} \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} [\log p(\mathbf{W}_2|\mu, \Sigma) - \log p(\mathbf{W}_2) - \log p(\mathcal{D}|\mathbf{W}_2)] \\
&= \arg \min_{\mu, \Sigma} \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} [\log p(\mathbf{W}_2|\mu, \Sigma) - \log p(\mathbf{W}_2)] - \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} [\log p(\mathcal{D}|\mathbf{W}_2)] \\
&= \arg \min_{\mu, \Sigma} D_{\text{KL}}(q(\mathbf{W}_2|\mu, \Sigma)||p(\mathbf{W}_2)) + \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} [-\log p(\mathcal{D}|\mathbf{W}_2)] \\
&= \arg \max_{\mu, \Sigma} \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} [\log p(\mathcal{D}|\mathbf{W}_2)] - D_{\text{KL}}(q(\mathbf{W}_2|\mu, \Sigma)||p(\mathbf{W}_2)) \\
&= \arg \max_{\mu, \Sigma} \mathcal{L}
\end{aligned}$$

We independently model the weight  $w_2^c$  for each class  $c$  with a full co-variance multivariate Gaussian distribution  $\mathcal{N}(w_2^c|\mu^c, \Sigma^c)$ , parameterised by mean vector  $\mu^c$  and lower-triangular (Cholesky) decomposition of the co-variance  $(\mathbf{L}^c)^T \mathbf{L}^c = \Sigma^c$  to avoid computational issues, following ?. We can view the variational posterior  $q(\mathbf{W}_2|\mu, \Sigma)$  as multivariate Gaussian over all classes with concatenated mean and block-diagonally stacked covariances from which we sample flattened matrix  $\mathbf{W}_2$  in one go, or -equivalently- sample row vectors  $w_2^c$  for each class and concatenate them to obtain matrix  $\mathbf{W}_2$ . By sampling  $L$  times from variational approximation  $\mathbf{W}_2^{(1)}, \mathbf{W}_2^{(2)} \dots \mathbf{W}_2^{(L)} \sim q(\mathbf{W}_2|\mu, \Sigma)$  we obtain a Monte Carlo estimate of  $\mathbb{E}_{\mathbf{W}} := \mathbb{E}_{\mathbf{W}_2 \sim q(\mathbf{W}_2|\mu, \Sigma)}$  required to compute the final ELBO or negative loss  $\mathcal{L}(\theta, \mathcal{D})$ :

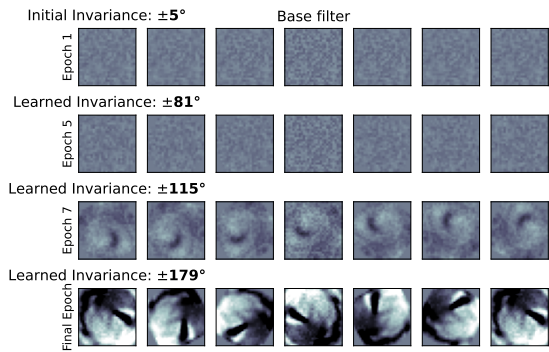
$$\begin{aligned}
\mathcal{L}(\theta, \mathcal{D}) &= \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} [\log p(\mathcal{D}|\mathbf{W}_2)] - D_{\text{KL}}(p(\mathbf{W}_2|\mu, \Sigma)||p(\mathbf{W}_2)) \\
&= \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma)} [\log p(\mathcal{D}|\mathbf{W}_2)] - \sum_c D_{\text{KL}}(\mathcal{N}(w_2^c|\mu^c, \Sigma^c)||p(w_2^c)) \\
&= \mathbb{E}_{q(\mathbf{W}_2|\mu, \Sigma^c)} [\log p(\mathcal{D}|\mathbf{W}_2)] - \sum_c D_{\text{KL}}(\mathcal{N}(w_2^c|\mu, \Sigma^c)||\mathcal{N}(\mathbf{0}; \Sigma_p)) \\
&= \underbrace{- \sum_l \sum_i^N -\log \sigma_{y_c^{(i)}} \left( \mathbb{E}_{T \sim p_{\eta}(T)} \left[ \mathbf{W}_2 \circ \phi \left( \mathbf{W}_1 \circ T \circ \mathbf{x}^{(i)} \right) \right] \right)}_{\text{Regular Average Cross-entropy}} - \underbrace{\sum_c \frac{1}{2} \left[ \log \frac{|\Sigma^c|}{|\Sigma_p|} - D + \text{tr} \{ \Sigma_p \Sigma^c \} + \mu^T \Sigma_p^{-1} \mu \right]}_{\text{Closed-form KL Regularizer}}
\end{aligned}$$

for every input  $\mathbf{x}^{(i)}$ , log soft-argmax output  $\sigma_{y_c}$  for class of corresponding label  $y_c^{(i)}$ , fixed first layer weights  $\mathbf{W}_1$ , prior weights  $\Sigma_p = \mathbf{I}\alpha$ , input dimensionality  $D$ , and trace  $\text{tr}(\cdot)$ . To allow for mini-batching, we use the Stochastic Variational Bayes Estimate (SGVB) from ? of the ELBO or negative loss  $\tilde{\mathcal{L}}(\theta, \mathcal{D})$ :

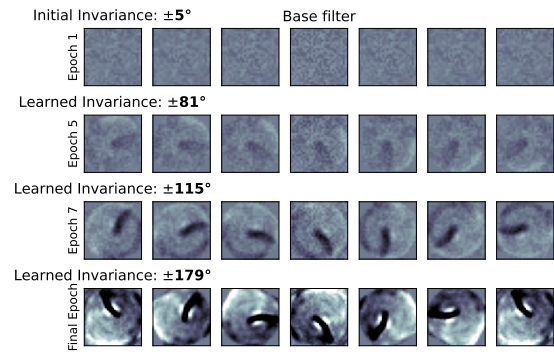
$$\tilde{\mathcal{L}}(\theta, \mathcal{D}) = -N \frac{1}{M} \sum_l \sum_i^M -\log \sigma_{y_c^{(i)}} \left( \mathbb{E}_{T \sim p_{\eta}(T)} \left[ \mathbf{W}_2 \circ \phi \left( \mathbf{W}_1 \circ T \circ \mathbf{x}^{(i)} \right) \right] \right) - \sum_c \frac{1}{2} \left[ \log \frac{|\Sigma^c|}{|\Sigma_p|} - D + \text{tr} \{ \Sigma_p \Sigma^c \} + \mu^T \Sigma_p^{-1} \mu \right]$$

where we can choose  $L = 1$  if we use a sufficiently large batch size.

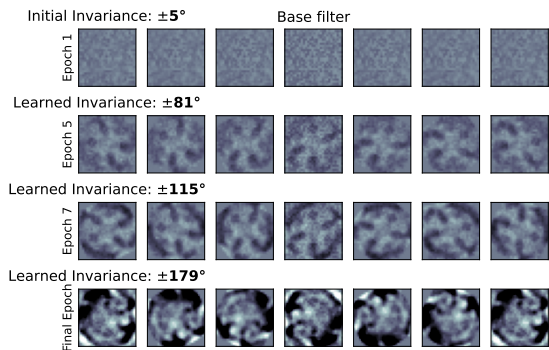
## APPENDIX B: WEIGHT VISUALIZATIONS OF LEARNED ROTATIONAL INVARIANCE



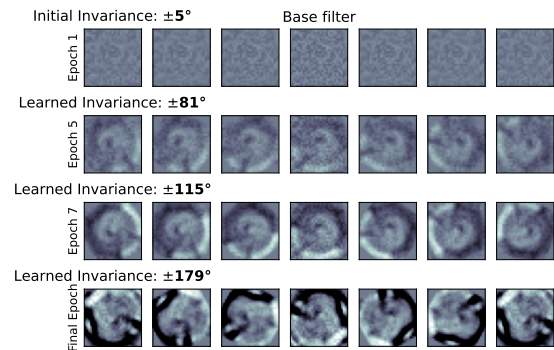
(a) Feature bank #1 over training iterations



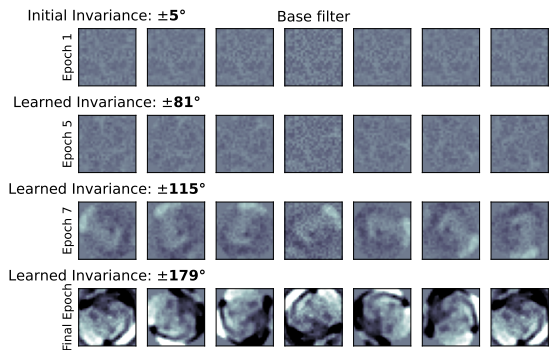
(b) Feature bank #2 over training iterations



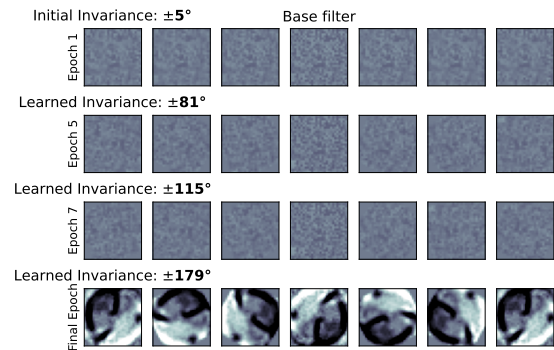
(c) Feature bank #3 over training iterations



(d) Feature bank #4 over training iterations



(e) Feature bank #5 over training iterations



(f) Feature bank #6 over training iterations

Figure 1: Illustration of the features banks over training iterations. Features are randomly initialised with almost no rotational invariance and converge to particular filters with full rotational invariance when trained on fully rotated MNIST data.

## APPENDIX C.1: ROTATIONAL INVARIANCE IN RFF NEURAL NETWORK

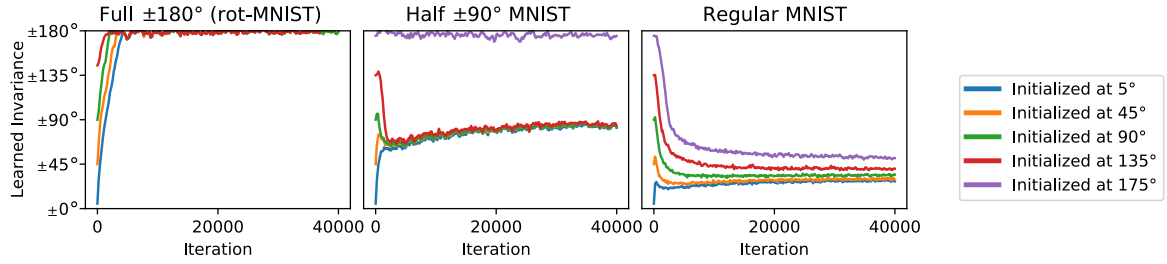


Figure 2: Predicted invariance over training iterations for different initial invariances for RFF neural network.

Model	Test Accuracy			ELBO		
	Fully rotated MNIST	Partially rotated MNIST	Regular MNIST	Fully rotated MNIST	Partially rotated MNIST	Regular MNIST
Fixed 5°	79.29	86.71	<b>96.00</b>	-1.07	-0.80	-0.36
Fixed 45°	87.35	91.13	95.93	-0.63	-0.49	<b>-0.26</b>
Fixed 90°	90.33	<b>91.69</b>	94.69	-0.52	<b>-0.44</b>	-0.30
Fixed 135°	91.19	91.04	92.13	-0.45	-0.45	-0.36
Fixed 175°	<b>91.57</b>	90.47	90.97	<b>-0.43</b>	-0.47	-0.45
Learned (5° Init)	<b>91.72</b>	<b>92.34</b>	<b>96.40</b>	<b>-0.43</b>	<b>-0.42</b>	<b>-0.26</b>
Learned (45° Init)	<b>91.65</b>	<b>92.31</b>	<b>96.42</b>	<b>-0.43</b>	<b>-0.42</b>	<b>-0.26</b>
Learned (90° Init)	<b>91.65</b>	<b>92.37</b>	<b>96.40</b>	<b>-0.43</b>	<b>-0.42</b>	<b>-0.26</b>
Learned (135° Init)	<b>91.66</b>	<b>92.37</b>	<b>96.10</b>	<b>-0.43</b>	<b>-0.42</b>	<b>-0.26</b>
Learned (175° Init)	<b>91.68</b>	<b>91.69</b>	95.64	<b>-0.43</b>	<b>-0.43</b>	<b>-0.26</b>

Table 1: Table containing Test Accuracy and ELBO scores after training for experiments with RFF network. In bold: the best scores for fixed invariance and, for learned invariances, all scores that surpass the best score using fixed invariance.

## APPENDIX C.2: ROTATIONAL INVARIANCE IN RELU NEURAL NETWORK

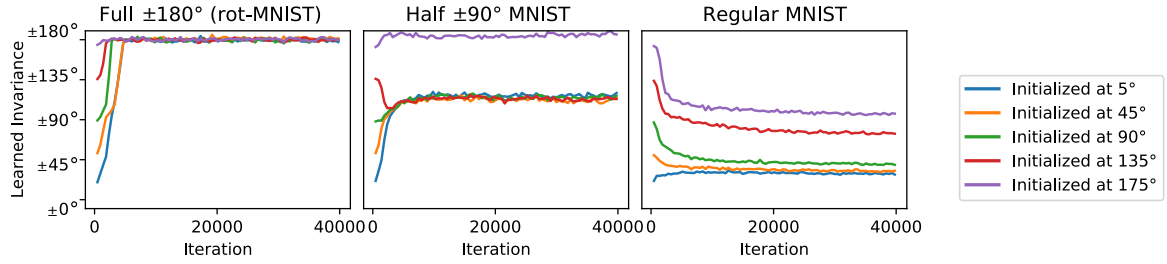


Figure 3: Predicted invariance over training iterations for different initial invariances of ReLU neural network with both input and output layer weights trained.

Model	Test Accuracy			ELBO		
	Fully rotated MNIST	Partially rotated MNIST	Regular MNIST	Fully rotated MNIST	Partially rotated MNIST	Regular MNIST
Fixed 5°	87.21	90.68	96.76	-0.28	-0.20	<b>-0.02</b>
Fixed 45°	95.24	96.46	98.13	-0.09	<b>-0.06</b>	-0.02
Fixed 90°	96.50	97.11	<b>98.14</b>	-0.07	-0.06	-0.03
Fixed 135°	97.15	<b>97.31</b>	97.79	<b>-0.06</b>	-0.06	-0.04
Fixed 175°	<b>97.53</b>	97.30	97.15	-0.07	-0.06	-0.06
Learned (0° Init)	97.34	97.13	<b>98.40</b>	-0.07	<b>-0.06</b>	<b>-0.02</b>
Learned (45° Init)	97.23	<b>97.36</b>	<b>98.27</b>	-0.07	<b>-0.05</b>	<b>-0.02</b>
Learned (90° Init)	97.28	<b>97.22</b>	<b>98.19</b>	-0.07	<b>-0.06</b>	<b>-0.02</b>
Learned (135° Init)	97.45	<b>97.29</b>	<b>98.33</b>	<b>-0.06</b>	<b>-0.05</b>	<b>-0.02</b>
Learned (175° Init)	97.23	<b>97.23</b>	<b>98.03</b>	<b>-0.06</b>	<b>-0.06</b>	<b>-0.03</b>

Table 2: Table containing Test Accuracy and ELBO scores after training for experiments of ReLU neural network with both input and output layer weights trained. In bold: the best scores for fixed invariance and, for learned invariances, all scores that surpass the best score using fixed invariance.

### APPENDIX C.3: DIFFERENT TRANSFORMATIONS IN RFF NETWORK

Model	Test Accuracy				ELBO			
	Fully rotated MNIST	Translated MNIST	Scaled MNIST	Regular MNIST	Fully rotated MNIST	Translated MNIST	Scaled MNIST	Regular MNIST
Regular MLP	79.29	66.07	89.25	95.16	-1.14	-1.49	-0.69	-0.39
+ Rotation	<b>92.59</b>	75.06	88.66	96.59	<b>-0.43</b>	-1.08	-0.62	-0.26
+ Translation	83.66	<b>87.81</b>	86.15	96.78	-0.82	<b>-0.64</b>	-0.72	-0.24
+ Scale	82.77	75.48	<b>91.31</b>	96.52	-0.84	-1.08	<b>-0.49</b>	-0.26
+ Affine	<b>92.64</b>	<b>87.77</b>	<b>90.58</b>	<b>97.38</b>	<b>-0.43</b>	<b>-0.64</b>	<b>-0.54</b>	<b>-0.21</b>

Table 3: Test Accuracy and ELBO for learned invariance using different transformations in a shallow RFF neural network.

### APPENDIX C.4: DIFFERENT TRANSFORMATION IN RELU NETWORK

Model	Test Accuracy				ELBO			
	Fully rotated MNIST	Translated MNIST	Scaled MNIST	Regular MNIST	Fully rotated MNIST	Translated MNIST	Scaled MNIST	Regular MNIST
Regular MLP	90.35	89.34	96.61	98.10	-0.06	-0.06	-0.03	-0.02
+ Rotation	<b>98.05</b>	94.08	97.62	98.64	<b>-0.05</b>	-0.06	-0.03	-0.02
+ Translation	93.59	<b>97.87</b>	97.98	98.76	-0.09	-0.06	-0.03	-0.02
+ Scale	93.80	94.30	<b>98.06</b>	98.35	-0.06	-0.06	-0.03	-0.02
+ Affine	<b>98.14</b>	<b>97.66</b>	<b>98.31</b>	<b>98.93</b>	<b>-0.05</b>	-0.06	-0.03	-0.02

Table 4: Test Accuracy and ELBO for learned invariance using different transformations in a shallow ReLU neural network.

### APPENDIX C.4: DIFFERENT TRANSFORMATION IN RELU NETWORK ON DATASETS WITH COMBINATIONS OF TWO INVARIANCES.

Model	Test Accuracy				ELBO			
	Fully rotated + Translated MNIST	Fully rotated + Scaled MNIST	Translated + Scaled MNIST	Regular MNIST	Fully rotated + Translated MNIST	Fully rotated + Scaled MNIST	Translated + Scaled MNIST	Regular MNIST
Regular MLP	53.36	80.71	75.50	98.10	<b>-0.26</b>	<b>-0.10</b>	<b>-0.12</b>	<b>-0.02</b>
+ Rotation	<b>85.35</b>	<b>95.66</b>	85.42	98.64	-0.31	-0.10	-0.27	-0.02
+ Translation	<b>83.84</b>	83.40	<b>91.77</b>	98.76	-0.42	-0.16	-0.19	-0.02
+ Scale	55.63	<b>89.81</b>	<b>86.04</b>	98.35	-0.39	-0.12	-0.17	-0.02
+ Affine	<b>89.37</b>	<b>95.88</b>	<b>91.95</b>	<b>98.93</b>	-0.37	-0.09	-0.18	-0.02

Table 5: Test Accuracy and ELBO for learned invariance using different transformations in a shallow ReLU neural network on datasets augmented by two subsequent transformations (rotation+translation, rotation+scaling and translation+scaling). Surprisingly, the regular MLP ends up with the best ELBO in this experiment. We did not consistently observe the best ELBO for the regular MLP throughout optimization, and find that we can still use our method and the ELBO to learn invariances in this case. Again, we observe that models with learned invariances achieve the highest test accuracy.

## APPENDIX D: DATASET DETAILS

All datasets have 60000 training examples and 10000 test examples and are created by taking regular MNIST or CIFAR-10 and applying random transformations:

**Regular MNIST Dataset:** MNIST handwritten digit database [?].

**Regular CIFAR-10 Dataset:** CIFAR-10 dataset with 10 classes [?].

**Partially rotated dataset:** Every sample rotated by radian angle  $\theta$ , sampled from  $\theta \sim U[-\frac{\pi}{2}, \frac{\pi}{2}]$ .

**Fully rotated dataset:** Every sample rotated by radian angle  $\theta$ , sampled from  $\theta \sim U[-\pi, \pi]$ .

**Translated dataset:** Translated samples relatively by  $dx$  and  $dy$  pixels, sampled from  $dx, dy \sim U[-8, 8]$ .

**Scaled dataset:** Every sample scaled around center with  $\exp(s)$ , sampled from  $s \sim U[-\log(2), \log(2)]$ .

## APPENDIX E: LIE GROUP GENERATORS

We follow ? and, similarly, utilise six matrix generators:

$$\begin{aligned} \mathbf{G}_{\text{transx}} = \mathbf{G}_1 &= \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, & \mathbf{G}_{\text{transy}} = \mathbf{G}_2 &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}, & \mathbf{G}_{\text{rot}} = \mathbf{G}_3 &= \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \\ \mathbf{G}_{\text{scalex}} = \mathbf{G}_4 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, & \mathbf{G}_{\text{scaley}} = \mathbf{G}_5 &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, & \mathbf{G}_{\text{shear}} = \mathbf{G}_6 &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \end{aligned}$$

To parameterise affine transformations we compute the following matrix exponential [?]:

$$T_{\epsilon} = \exp \left( \sum_i \epsilon_i \eta_i \mathbf{G}_i \right), \quad \epsilon \sim U[-1, 1]^k \quad (1)$$

Optionally, the values of  $\eta$  can be constrained to a positive range by passing them through a ‘softplus’-function, or in case of  $\eta_3 = \eta_{\text{rot}}$  to  $[-\pi, \pi]$  using a scaled ‘tanh’ function, preventing double coverage on the unit circle. In practice, however, we did not find such constraints necessary as long as  $\eta_{\text{rot}}$  is reasonably initialised (e.g.  $\eta = \mathbf{0}$ ).

By fixing certain  $\eta_i$  at 0, subsets of the generator matrices parameterise rotation, translation and scaling:

<p>For rotation only: Learn <math>\eta_3</math>. Fix <math>\eta_i = 0</math> for all <math>i \neq 3</math>.</p>	<p>For translation only: Learn <math>\eta_1</math> and <math>\eta_2</math>. Fix <math>\eta_i = 0</math> for all <math>i &gt; 2</math>.</p>	<p>For scaling only: Learn <math>\eta_4</math> and <math>\eta_5</math>. Fix <math>\eta_i = 0</math> for all <math>i \notin \{4, 5\}</math>.</p>
$\begin{aligned} T_{\epsilon}^{(\text{rot})} &= \exp \left( \sum_i \epsilon_i \eta_i \mathbf{G}_i \right) \\ &= \exp(\epsilon_3 \eta_3 \mathbf{G}_3) \\ &= \exp \left( \begin{bmatrix} 0 & -\epsilon_3 \eta_3 & 0 \\ \epsilon_3 \eta_3 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \right) \\ &= \begin{bmatrix} \cos(\epsilon_3 \eta_3) & -\sin(\epsilon_3 \eta_3) & 0 \\ \sin(\epsilon_3 \eta_3) & \cos(\epsilon_3 \eta_3) & 0 \\ 0 & 0 & 1 \end{bmatrix} \end{aligned}$	$\begin{aligned} T_{\epsilon}^{(\text{trans})} &= \exp \left( \sum_i \epsilon_i \eta_i \mathbf{G}_i \right) \\ &= \exp(\epsilon_1 \eta_1 \mathbf{G}_1 + \epsilon_2 \eta_2 \mathbf{G}_2) \\ &= \exp \left( \begin{bmatrix} 0 & 0 & \eta_1 \\ 0 & 0 & \eta_2 \\ 0 & 0 & 0 \end{bmatrix} \right) \\ &= \begin{bmatrix} 1 & 0 & \epsilon_1 \eta_1 \\ 0 & 1 & \epsilon_2 \eta_2 \\ 0 & 0 & 1 \end{bmatrix} \end{aligned}$	$\begin{aligned} T_{\epsilon}^{(\text{scale})} &= \exp \left( \sum_i \epsilon_i \eta_i \mathbf{G}_i \right) \\ &= \exp(\epsilon_4 \eta_4 \mathbf{G}_4 + \epsilon_5 \eta_5 \mathbf{G}_5) \\ &= \exp \left( \begin{bmatrix} \eta_4 & 0 & 0 \\ 0 & \eta_5 & 0 \\ 0 & 0 & 0 \end{bmatrix} \right) \\ &= \begin{bmatrix} \exp(\epsilon_4 \eta_4) & 0 & 0 \\ 0 & \exp(\epsilon_5 \eta_5) & 0 \\ 0 & 0 & 1 \end{bmatrix} \end{aligned}$