

# NEML CW

Ben James, Johns Noble

March 2025

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Constructing Metrics</b>	<b>1</b>
<b>3</b>	<b>Parametric Inference</b>	<b>3</b>
<b>4</b>	<b>Extending to Mixture Models</b>	<b>3</b>
<b>5</b>	<b>Conclusion</b>	<b>3</b>

## 1 Introduction

It is often very natural to attempt to model data using gaussian distribution. Gaussian distributions are well defined for data points which lie in euclidean space for any number of dimensions. We however know that most data should be modelled as to be taken from a manifold. There are methods to do manifold learning but once we have a manifold the notion of probability distributions need to be altered such that we are defining a probability measure across a manifold. It does not make sense for a region outside the manifold to possess non-zero measure. The Locally Adaptive Normal Distribution (LAND) aims to attempt to do this. Once we have a notion of normal distribution on a manifold given some mean and covariance, We can try and fit distribution to the data using MLE approaches. Finally, we will be able to use the fact that we can fit gaussians to data to be able to perform EM algorithms over data to fit a mixture of LANDS over data. TODO: Write something about how manifold learning is handled i.e. Do we know the manifold which data will be distributed on or is this something that needs to be learnt using graphs?

## 2 Constructing Metrics

In order to define the idea of distances over manifolds and in order to attempt to capture local behaviour of data, we need to a Riemannian Metric which acts

on tangent vectors. The Metric  $\mathbf{M}(\mathbf{x})$  gives an inner product  $\langle \mathbf{u}, \mathbf{M}(\mathbf{x})\mathbf{v} \rangle$ . This allows us to define geodesics on our manifold as the path which minimises the length:

$$\hat{\gamma} = \underset{\gamma}{\operatorname{argmin}} \int_0^1 \sqrt{\langle \gamma'(t), \mathbf{M}(\gamma(t))\gamma'(t) \rangle} dt, \quad \gamma(0) = \mathbf{x}, \quad \gamma(1) = \mathbf{y}$$

Once we have  $\hat{\gamma}$  it is natural to define our Exponential map given a  $\mathbf{v} \in T_{\mathbf{x}}M$  such that  $\mathbf{y} = \operatorname{Exp}_{\mathbf{x}}(\mathbf{v}) \in M$  where  $\hat{\gamma}(t) = \operatorname{Exp}_{\mathbf{x}}(t \cdot \mathbf{v})$ . We can then begin to define the *Log* mapping as:  $\operatorname{Log}_{\mathbf{x}}(\mathbf{y}) = \mathbf{v} \in T_{\mathbf{x}}M$ . In order to calculate each geodesic, we must express solve the equation relating christoffel symbols. This works out to be a second order differential ODE that in practice can be solved numerically.

A generalised view of thinking about the normal distribution is the distribution over a given space such that given a known mean and covariance for which we have maximum entropy. This way of defining the normal distribution allows us to extend to non euclidean spaces. It can be shown that given a Log map, the distribution that satisfies this is as follows:

$$p_M(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \exp\left(-\frac{1}{2}\langle \operatorname{Log}_{\boldsymbol{\mu}}(\mathbf{x}), \boldsymbol{\Sigma}^{-1}\operatorname{Log}_{\boldsymbol{\mu}}(\mathbf{x}) \rangle\right)$$

We also need to make considerations on how we can learn a metric tensor so that we can define a riemannian metric. The most obvious thing we can do is to use a single global metric tensor such that  $\operatorname{dist}^2(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M}(\mathbf{x}_i - \mathbf{x}_j)$  where  $\mathbf{M}$  is symmetric and positive definite. The issue with this is that a single global metric is never enough to capture the non linearity of manifolds. We can attempt to also learn several metric tensors for different sections of data. This involves picking a few centres  $(\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_r)$ , defining a fixed metric  $\mathbf{M}_r$  for each data point picking the closest metric. We can generalise this sort of approach by letting:

$$\mathbf{M}(\mathbf{x}) = \sum_{r=1}^R \frac{w_r(\mathbf{x})}{\sum_j w_j(\mathbf{x})}$$

Intuitively, we can think of the ratio of  $w_i$  to be the responsibility of each centre for  $\mathbf{x}$ . For example in the case of picking the nearest tensor, we can say that:

$$w_r(\mathbf{x}) = \begin{cases} 1, & \|\mathbf{x} - \mathbf{x}_r\|_{\mathbf{M}_r}^2 \leq \|\mathbf{x} - \mathbf{x}_j\|_{\mathbf{M}_j}^2, \forall j \\ 0, & \text{otherwise} \end{cases}$$

We can experiment with using different  $w$  functions and by using various number of centres. To summarise, a metric tensor is constructed using the inverse covariance of training data in a local space. The authors of LAND decided that when constructing metric tensors, that we should only consider the diagonals of the covariance to avoid issues with overfitting and to speed up computation. They also used for function  $w_n(\mathbf{x}) = \exp(-\frac{\|\mathbf{x}_n - \mathbf{x}\|_2^2}{2\sigma^2})$ . Defining for the  $d$  value in the diagonal of  $\mathbf{M}$  to be

$$\mathbf{M}(\mathbf{x}) = (\sum_{n=1}^N w_n(\mathbf{x})(x_{nd} - x_d)^2 + \rho)^{-1}$$

The  $\rho$  is here for numerical stability and  $\sigma$  can be tweaked depending on the density and dimensions of the data. In our experiments we begin by assuming that the data lies on a sphere so that we could use already existing Exponential and Log mappings and the metric tensor over a sphere. We were able to write an implementation that learnt a riemannian metric using the formula above.

### 3 Parametric Inference

We are tasked with finding the mean and covariance matrix of the distribution and we have implemented this using numerical methods. The objective function we are trying to maximise is as follows:

$$\underset{\mu \in M, \Sigma \in S_{++}^D}{\operatorname{argmin}} (\phi(\mu, \Sigma) = \frac{1}{2N} \sum_{n=1}^N \langle \operatorname{Log}_{\mu}(\mathbf{x}_n), \Sigma^{-1} \operatorname{Log}_{\mu}(\mathbf{x}_n) \rangle + \log(C(\mu, \Sigma)))$$

Where  $C$  here is the normalisation constant of normalisation to ensure its a probability distribution. To compute this value we require to be able to compute  $E_{N(0, \Sigma)}[\sqrt{|\operatorname{MExp}_{\mu}(\mathbf{v})|}]$ , intuitively can be thought of as the the expectation of the volume of the manifold. Since this is value is intractable we use monte carlo simulation techniques to be able to compute this by sampling over  $N(0, \Sigma)$ . We are able to work out two grad functions for the objective, one for the mean, one for the covariance. The grad function for the mean is in th form  $-\Sigma^{-1}(\dots)$ . Simply applying gradient descent with this grad function however results in unstable convergence due to the condition number of  $\Sigma$ . We end up choosing a direction which is independent from  $\Sigma$  by taking the direction to be  $-\Sigma \nabla \phi$  which can be proved to give a direction of objective function decrease. We also must be careful with applying gradient descent over the covariance matrix. This is because these covariance matrices must neccesarily be symmetric positive definit. In order to ensure that this always hold we can decompose  $\Sigma = PDP^T$  and therefore  $\Sigma^{-1} = PD^{-1}P^T = (PD^{-\frac{1}{2}})(PD^{-\frac{1}{2}})^T = A^T A$  Representing the objectives with  $A$  allows us to optimise over  $A$  and ensure our covariance matrices are valid.

### 4 Extending to Mixture Models

### 5 Conclusion

We were able to show that we could perform gradient descent in the space of parameters to find an approximation for the MLE of a normal distribution on a given manifold. We were also able given data points to learn a manifold and infer the geodesics to produce an exponential and log mapping by solving ODEs.

Interesting cases we found were ones where the variance of the distribution is high comparitive to the size of a spherical manifold. We found in these cases that generally convergence is slower and often inaccurate. This should make sense since we no longer have an injective exponential map and

the tails of the normal distribution are able to "wrap" back around. This effect causes the overall distribution to look less like a normal distribution and more like a distribution with heavy tails. We found that setting the start variance to 0.7 was a good value. We can come up with these values by establishing bounds on the probability that a point appears over  $\pi$  distance away in the tangent space.