

Neural Networks in the Emulation of Probability Distributions for Massive Computational Acceleration with Applications in Mathematical Biology

Alexander John Pinches

October 15, 2020

Abstract

The generation of the probability distributions of stochastic models such as those in Biology are computationally expensive time consuming relying on methods such as the Gillespie Algorithm for stochastically correct trajectories or slightly faster, but approximate, methods such as the Tau-Leaping Algorithm. To explore the large multi-dimensional parameter spaces of stochastic reaction networks present in biology we employ neural networks to emulate the probability distributions of the sampled trajectories much faster than using the Gillespie Algorithm. We explore the application of Long Short Term Memory (LSTM) arrays and another new approach in the form of an approximating Mixture Model whose parameters are generated by a feed forward neural network. Our source code is publicly available at <https://github.com/Parply/PDFEmulators>.

I Introduction

Biological systems are inherently stochastic in nature from the production of proteins from a regulating gene to the interaction of enzymes, or the growth of a population of cells and their observed heterogeneity. The stochasticity of these models create a distribution of members of each species in a model that is determined by the reaction rates. Distributions of these species at the steady state can vary greatly or remain similar with the change of reaction rates exhibiting multimodal distributions. Many models of interest cannot be solved analytically except for all but the simplest models and they rely on computationally expensive algorithms such as the Gillespie Algorithm *Székely and Burrage* [33]. This motivates the need for a computationally inexpensive method to emulate the solutions to such models to accelerate the exploration of their parameter space.

Simulated data from a mixture of 3 Gaussians are used to explore the predictive capability of the model and choose appropriate hyperparameters. In addition to providing a setting where the model will have to learn to emulate distributions of varying modality between 1 and 3. We demonstrate the model in a real world, simplistic and well explored model of a biological process we using a regulating gene producing mRNA and then a Protein. Such a reaction network exhibits heterogeneity due to the stochasticity of the process shown in *Shahrezaei and Swain* [30]. Protein production forms the basis of many biological systems and is stochastic in nature with the decomposition and production of mRNA and proteins and the switching of the promoter from active to inactive states all being stochastic. Seeking to emulate the distribution of proteins given some reaction rates which is not obtainable analytically we then expand this to the joint distribution of mRNA and Proteins with the Mixture Model to explore its capability in higher dimensions.

In pursuit of a faster method to emulate these systems we seek to apply neural networks which have become a popular method for emulating computationally complex problems *Cranmer et al.* [11]. Firstly in the use of LSTM arrays which have been used effectively in the prediction of time series *Siami-Namini and Namin* [31], text prediction *Guo et al.* [18] and even in complex video games with Google's DeepMind AI *Vinyals et al.* [35]. These problems rely on the problem to be represented as a series of related points in time. In our case representing the values of the bins as a series. They were used by *Wang et al.* [36] in the solutions of differential equations motivating our exploration of such neural networks in our problem. Secondly, we also apply a new method in the use of feed forward neural networks in the prediction of the parameters of a Mixture Model. Models of this type have seen applications for some problems but are not similar to the problem here as in *Bishop* [7] in estimating highly non-linear problems. Examples given were the kinematics of a simple 2-link robot arm or in the mapping $x = t + 0.3\sin(2\pi t) + \epsilon$ predicting x from a given t with ϵ white noise where a Multivariate Gaussian

Mixture Model is predicted for each x . Methods not utilising neural networks exist in Bayesian methods shown by *Beaumont* [5]. The complexity of such methods can get prohibitively large and rely on the computation of many random variables.

An issue unique to this problem compared to the previously mentioned is common loss functions such as the Mean Square Error $MSE = \frac{1}{T+1} \sum_{i=0}^T (y_i - x_i)^2$ have small gradients for $|y - x| \in (0, 1)$. Which we would expect when emulating probability distributions by their definition. To counteract this we propose the use of divergence measures as loss functions to prevent vanishing gradients caused by some loss functions.

II Mathematical Methods

Here we explain the mathematical methods employed in our models and in the generation of our data.

II.A Multilayer Perceptron

A multilayer perceptron or feedforward neural network is built up of layers of neurons *Murtagh* [23], *Trevor Hastie and Friedman* [34]. A neuron has inputs $X \in \mathbb{R}^n$ which can be the data or from another neuron, it has weights $W \in \mathbb{R}^n$ and a bias $\in \mathbb{R}$ It computes the activation as a function to the inputs, weights and bias $W^T X + b$. Then it applies an activation function $\sigma(\cdot)$ to the activation to give the output of the neuron as

$$y = \sigma(W^T X + b)$$

We can build multiple layers of multiple neurons that feed the output of the previous layer to the next with differing activations and sizes until we reach the output layer. We can define this neural network with $L + 1$ layers $z^{(0)}, \dots, z^{(L)}$ where $z^{(i)} \in \mathbb{R}^{n_i}$ with n_i the number of neurons in layer i. Input $x = z^{(0)} \in \mathbb{R}^d$ and output $\hat{y} = z^{(L)}$ We can then define for $l = 0, \dots, L - 1$ the forward propagation

$$z^{(l+1)} = \sigma\left(W^{(l)} z^{(l)} + b^{(l)}\right)$$

Where $W^{(l)} \in \mathbb{R}^{n_{l+1} \times n_l}$, $b^{(l)} \in \mathbb{R}^{n_{l+1}}$ and the activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ applied elementwise. We want to find the weights and biases such that we maximise the log-likelihood. To do this we must define a loss function for example we select the mean squared error with p the dimensions of the output.

$$MSE = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^p (y_j^{(i)} - \hat{y}_j^{(i)})^2$$

A feed forward neural network with the Identity Function $f(x) = x$ as activation simplifies to a Linear Model which follows from the matrix form of the layers of the neural network by using other activations however we can capture non-linear effects.

II.B Long Short Term Memory (LSTM)

LSTM cells have the ability to learn patterns in series based on an input and the previous cell state. This has lead to them being used extensively in areas such as text prediction and in modelling time series data where the autocovariance function demonstrates the need for a model that can remember long patterns. The ability to learn patterns has been proposed as a way for predicting functions in *Gers et al.* [15] as we believe that for a function with support \mathbb{X} and points along its support $\{x_1, \dots, x_n : x_i \in \mathbb{X} \forall i = 1, \dots, n\}$ it's value at $f(x_i)$ contains information about its value at $f(x_{i+1})$ and similarly $f(x_i)$ contains information about its value at $f(x_{i-1})$. The difference between points needs to be sufficiently small for this information to exist in our case the bins must be sufficiently small to capture the information. Using LSTM's we can stack arrays on top of one another and/or introduce Bidirectionality *Schuster and Paliwal* [29] to improve the models predictive ability. We achieve this by having one set of LSTM arrays going in one direction prediction forwards in time $1 \rightarrow T$ and another going in the opposite direction $T \rightarrow 1$. Then concatenating or summing their outputs. This has been proposed as a way for the network to learn more context about each point as it has predictions from both directions in text prediction and in our case, we expect it to improve the model's ability to predict in cases with long tails.

The update of the weights in the cells follow from the case of the feed forward neural network as the gates can be viewed as layers in feed forward neural network when we unfold and back propagate through time *Gers et al.* [15].

II.C Standard LSTM

A standard LSTM cell with forget gates as proposed in *Gers et al.* [15] of a given encoder size s and input dimension d consists of a the Cell Memory C_{t-1} this is the long term memory of the cell and the short term memory H_{t-1} which is the output of the cell. These are then fed through the cell's gates. The Forget Gate f_t that controls the forgetting of information stored in the cell memory C_{t-1} , the cell gate g_t which proposes new candidates to be added to the cell memory and the input gate i_t which controls the weighting of these inputs and finally the output gate o_t which controls the output of the cell. We also define $H_0 = 0$ and $C_0 = 0$. A graphical representation of this cell is shown in figure 1.

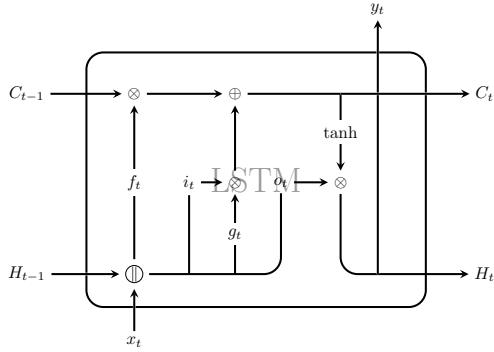


Figure 1: Standard LSTM Cell

Variables:

- $x_t \in \mathbb{R}^d$: Input Vector
- $C_t \in \mathbb{R}^s$: Cell Memory
- $H_t \in \mathbb{R}^s$: Cell Output
- $f_t \in \mathbb{R}^s$: Forget Gate
- $i_t \in \mathbb{R}^s$: Input Gate
- $g_t \in \mathbb{R}^s$: Cell Gate
- $o_t \in \mathbb{R}^s$: Output Gate
- $W_i \in \mathbb{R}^{s \times d}$: Weight Matrices
- $W_h \in \mathbb{R}^{s \times s}$: Weight Matrices
- $b_i \in \mathbb{R}^s$: Bias Vectors

Operators and functions:

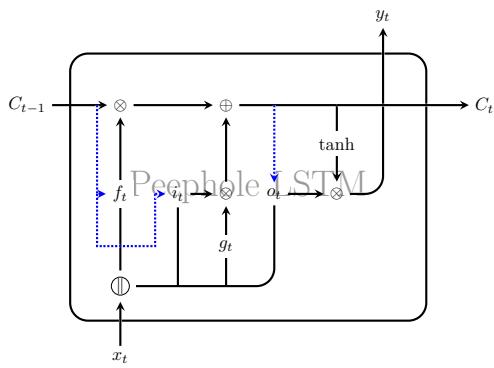
- \otimes = Hadamard Product
- \oplus = Element-wise Addition
- $\textcircled{\text{P}}$ = Concatenate
- $\sigma(x) = \frac{1}{1 + \exp(-x)}$: Sigmoid Function
- $\tanh(x) = \frac{\exp(2x) - 1}{\exp(2x) + 1}$: Hyperbolic Tangent Function

Gate functions:

$$\begin{aligned} i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}H_{t-1} + b_{hi}) \\ f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}H_{t-1} + b_{hf}) \\ g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}H_{t-1} + b_{hg}) \\ o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}H_{t-1} + b_{ho}) \\ C_t &= f_t \otimes C_{t-1} + i_t \otimes g_t \\ H_t &= o_t \otimes \tanh(C_t) \end{aligned}$$

II.D Peephole LSTM

The Peephole LSTM adds extra Peephole connections to feed the cell memory into the gates as proposed in *Gers and Schmidhuber* [14] and the gate functions are modified to use the C_i instead of H_i . We represent the additional connections by the blue dotted line in the graphical representation shown in figure 2. It's use in learning relationships over longer periods was shown in *Gers and Schmidhuber* [14] which could be useful in our model especially for a large number of bins.



Gate functions:

$$\begin{aligned} i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}C_{t-1} + b_{hi}) \\ f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}C_{t-1} + b_{hf}) \\ g_t &= \tanh(W_{ig}x_t + b_{ig}) \\ o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}C_{t-1} + b_{ho}) \\ C_t &= f_t \otimes C_{t-1} + i_t \otimes g_t \\ H_t &= o_t \otimes \tanh(C_t) \end{aligned}$$

Figure 2: Peephole LSTM Cell

II.E Activation functions

Example activation functions we could use include the identity function as previously mentioned. We ideally want easily differentiable activation functions and for activations on the final output layers to map to the space of what we are predicting.

$$\text{identity}(x) = x$$

A commonly used activation function that was used in *Wang et al.* [36] with their LSTM model was the rectified linear unit (ReLU) function. Its primary weakness being is as it is zero and has zero gradient for negative values so it cannot backpropagate on these values.

$$\text{ReLU}(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

This leads to variations of ReLU that overcome this problem. One such alternative is exponential linear units (ELU) which also has an extra parameter to optimise from *Clevert et al.* [10]. Where it was shown to obtain higher classification accuracy than ReLU.

$$\text{ELU}(x, \alpha) = \begin{cases} x & x \geq 0 \\ \alpha(\exp x - 1)x & x < 0 \end{cases}$$

In our application a PMF must have non negative values less than or equal to one and sum to one. An activation function that is commonly used for this is the standard unit Softmax function.

$$\sigma(x)_i = \frac{\exp(x_i)}{\sum_{j=1}^K \exp(x_j)} \text{ for } i = 1, \dots, K$$

The non-negativity comes from the exponential function and the summing to one is obviously true.

$$\sum_{i=1}^K \sigma(x)_i = \frac{\sum_{i=1}^K \exp(x_i)}{\sum_{j=1}^K \exp(x_j)} = 1$$

As it is non-negative and sums to unity then all $\sigma(x)_i \in [0, 1]$. We can obviously select from a large variety of activation functions both those available to us already in packages and we can create new activation functions. We here however limit ourselves to the activation functions available in Pytorch as of 1.5.0 in the exploration of improvements they can offer.

II.F Optimisers

In any neural network given a loss we need an algorithm to determine how we change the weights of the layers in a manner that will lead to convergence to a solution. Our choice of algorithm and its hyperparameters can determine if and how fast it converges to an optimal solution. Knowing which algorithm and the hyperparameters to use rely on testing and using performance on similar models as an indicator to determine which to use. The size of each update, the learning rate we denote by η .

II.F.1 Batching

Showed universally across different models is the increase in performance resulting from updating across a batch of data rather than individually for each sample or across the whole training dataset *LeCun et al.* [21]. The intuition behind this is that by using a batch we can help the model learn general features rather than fitting noise in the data and also train a model faster. Too large a batch however can lead to over generalisation a problem which is shown later in the LSTM model.

II.F.2 Gradient Descent

A commonly applied optimiser is that of Gradient Descent *LeCun et al.* [21]. Intuitively we would like to adjust the weights by moving in the direction that reduces the loss function. As we assume the samples from the data

are independent and identically distributed the loss is a summed across each sample for the model parameters θ for N samples and loss function L.

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N L(x_i, y_i, \theta)$$

$$\implies \nabla_{\theta} J(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} L(x_i, y_i, \theta)$$

This can be expensive for large datasets and lead to over fitting so instead we use Stochastic Gradient Descent *LeCun et al.* [21]. We approximate the expected gradient by using a random minibatch of the data $(x^{n_i}, y^{n_i})_{i=1}^{N'}, N' \ll N$. This means the estimated gradient is

$$g = \frac{1}{N'} \sum_{i=1}^{N'} \nabla_{\theta} L(x^i, y_i, \theta)$$

We adjust the parameter with learning rate $0 < \eta \leq 1$

$$\theta \leftarrow \theta - \eta g$$

This does introduce two training parameters the size of the minibatch and the learning rate. A general rule is to lower the learning rate as you increase the size of the batch. For example with the loss function as the MSE and $\sigma_{(l)}$ representing the activation function of layer l.

We can define the loss

$$L = \frac{1}{N'} \sum_{i=1}^{N'} (z_i^{(L)} - y_i)^2$$

We then take the derivative with respect to the output of layer L $z^{(L)}$

$$\frac{\partial L}{\partial z^{(L)}} = \frac{2}{N'} \sum_{i=1}^{N'} (y_i - z_i^{(L-1)}) \sigma'_L(z_i^{(L)})$$

We can then apply the chain rule to get

$$\frac{\partial L}{\partial z^{(l)}} = \prod_{i=l}^{L-1} \frac{\partial L}{\partial z^{(i+1)}} \frac{\partial z^{i+1}}{\partial z^{(i)}}^T$$

Then we can calculate the gradients

$$g(w^{(l)}) = \frac{\partial L}{\partial z^{(l)}} \frac{\partial z^{(l)}}{\partial w^{(l)}}^T$$

$$= \frac{\partial L}{\partial z^{(l)}} (z^{(l-1)})^T$$

$$g(b^{(l)}) = \frac{\partial L}{\partial z^{(l)}} \frac{\partial z^{(l)}}{\partial b^{(l)}}^T$$

$$= \frac{\partial L}{\partial z^{(l)}}$$

Using that $\frac{\partial z^{(l)}}{\partial b^{(l)}}$ is the identity and $\frac{\partial z^{(l)}}{\partial w^{(l)}} = z^{(l-1)}$. We then use the gradients to update each of the parameters as explained above. Repeating this for each batch. We then loop over new random samples of batches of the data for each training epoch updating the weights.

SGD guarantees convergence to a minimum but this can be a local minimum and the speed at which it converges can be slow *LeCun et al.* [21]. Extensions and variations of the SGD exist in the form of adding momentum *LeCun et al.* [21] and in adaptive learning rate gradient decent (Adagrad) *Duchi et al.* [12] but they were not applied to our models.

II.F.3 Adam

A now very popular optimisation algorithm Adam has seen wide spread use originally proposed in *Kingma and Ba* [19]. Using the first and second moments to control the learning rate of each parameter adaptively. Controlled by 3 parameters α the stepsize, β_1 controls the decay of the moving average of the first moment and β_2 the second moment. Some variations upon this algorithm can be seen in *Reddi et al.* [27] but we have not applied them to our models. We use as parameters the Pytorch defaults $\alpha = \eta$, $\beta_1 = 0.9$, $\beta_2 = 0.999$. In the paper α is the learning rate η , $0 < \epsilon$ is a small term to prevent division by zero usually 10^{-8} and g is as defined before for Gradient Descent $g = \nabla_{\theta} J(\theta)$. This gives us the simple update rule

$$\theta_{t+1} = \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

where we have

$$\begin{aligned}\hat{m}_t &= \frac{m_t}{1 - \beta_1^t} & \hat{v}_t &= \frac{v_t}{1 - \beta_2^t} \\ m_t &= (1 - \beta_1)g_t + \beta_1 m_{t-1} & v_t &= (1 - \beta_2)g_t^2 + \beta_2 v_{t-1} \\ m_0 &= 0 & v_0 &= 0\end{aligned}$$

m_t and v_t are exponential moving averages that are biased estimates for the first and second moments and \hat{m}_t, \hat{v}_t are the bias corrected estimates. We can see this by rewriting the formula of each and taking expectations.

$$\begin{aligned}m_t &= (1 - \beta_1)g_t + \beta_1 m_{t-1} & v_t &= (1 - \beta_2)g_t^2 + \beta_2 m_{t-1} \\ &= (1 - \beta_1) \sum_{i=0}^t \beta_1 g_i & &= (1 - \beta_2) \sum_{i=0}^t \beta_2 g_i^2 \\ E[m_t] &= (1 - \beta_1) \sum_{i=0}^t \beta_1 E[g_i] & E[v_t] &= (1 - \beta_2) \sum_{i=0}^t \beta_2 E[g_i^2] \\ &= E[g_t](1 - \beta_1^t) & &= E[g_t^2](1 - \beta_2^t) \\ \implies E[\hat{m}_t] &= E[g_t] & \implies E[\hat{v}_t] &= E[g_t^2]\end{aligned}$$

This leads to a key property of Adam in that the learning rate is scaled for each parameter. The Adam optimiser was also used in the source code provided with the LSTM model in *Wang et al.* [36].

II.G Mixture Models

A commonly used technique to represent distributions is a Mixture Model *Trevor Hastie and Friedman* [34]. Mixture models use the sum of a chosen number of individually weighted independent distributions.

Definition II.1. Mixture Model A mixture model with $n \geq 1$ distributions ϕ defined by a parameter vector $\boldsymbol{\theta}$ and weights $\alpha_i \in [0, 1]$ such that $\sum_{i=1}^n \alpha_i$ has probability mass function

$$f(x; \boldsymbol{\alpha}, \boldsymbol{\theta}, n) = \sum_{i=1}^n \alpha_i \phi(x; \boldsymbol{\theta}_i)$$

It has been shown that mixtures can get arbitrarily close to a given distribution for a finite number of distributions in the continuous case *Nguyen et al.* [26], *Nguyen and McLachlan* [25]. It is immediately obvious for the bounds on the number of distributions is at least as many as the modality of the output divided by the modality of our chosen distribution and at most the number of bins in the discrete case. The lower bound comes from the fact

that our mixture can have at modality of at most the number of mixtures with the used distributions being of modality 1. Our upper bound comes from the used distributions having location and scale parameters we can change independently. This leads to the obvious exact model consisting distributions with means inside each bin and weight α_i the probability of that bin and a variance of almost zero.

$$\text{var}(X) = 0 = \mathbb{E}((X - \mathbb{E}(X))^2) \quad (1)$$

$$(X - \mathbb{E}(X))^2 \geq 0 \implies \mathbb{P}((X - \mathbb{E}(X))^2 \neq 0) = 0 \quad (2)$$

$$\implies X = \mathbb{E}(X) \quad (3)$$

This outlines our ability to get the exact distribution if we have variance 0 in the discrete case. A model like this however, will likely overfit the data so in practice you wouldn't want to use this many distributions in the Mixture Model.

II.G.1 Gaussian Mixture Model

Where Mixture Models are applied the most commonly used distribution is the Normal distribution $\mathcal{N}(\mu, \sigma)$ *Trevor Hastie and Friedman* [34]. The attractiveness of using a Gaussian mixture are its use on many problems shows it to be effective and it's easily interpreted. We do however, as the Gaussian distribution is continuous have to discretise by using the value of the CDF across each bin. As the support of this distribution is $x \in \mathbb{R}$ we ignore the probability outside the range covered by the bins. This can cause problems in the faith we have in divergence metrics that aren't symmetrical such as the Kullback-Lieber Divergence as the model doesn't necessarily sum to 1 across the range covered by the bins. However, as the model converges towards its solutions and distance metrics such as the L1 Norm, MSE and Hellinger Distance reduce we can have more faith in the values of non-symmetric divergence measures. This also motivates the use of discrete distributions that can also have their mean and variance varied as alternative distributions in our Mixture Model. We show below an example of this Mixture Model with 3 distributions 3.

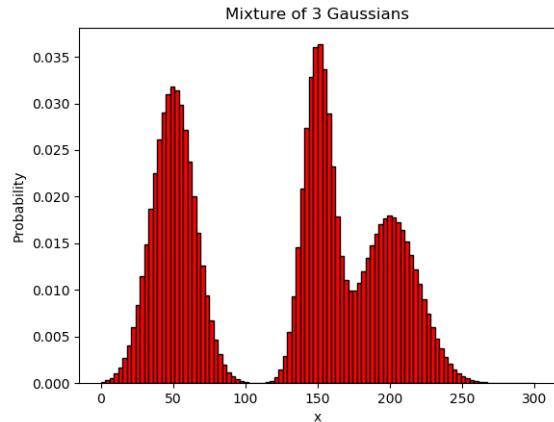


Figure 3: A mixture of 3 Gaussians with parameters $\mu = (50, 150, 200), \sigma = (15, 10, 20)$ and weights $(0.4, 0.3, 0.3)$ for bins of width 3 from 0 to 300

II.G.2 Gamma-Poisson Mixture Model

The Gamma-Poisson distribution has become an attractive distribution due to its ability to account for overdispersion from a Poisson distribution originally developed by *Bates and Neyman* [3] and has seen been applied to a regulating gene by *Amrhein et al.* [2] and mixtures to RNA-Seq and DNA methylation data by *Li et al.* [22]. The intuition behind the Gamma-Poisson model is to have a Poisson distribution $X|\lambda \sim \text{Poi}(\lambda)$ with it's rate parameter Gamma distributed $\lambda \sim \text{Gamma}(\alpha, \beta)$ with $X \in \mathbb{N}_0, \alpha > 0, \beta > 0$. Allowing us to change the variance and the mean of the model independently unlike with a Poisson distribution. The support is exactly the values that the distribution from our regulating gene model can take and most models in biology can take being \mathbb{N}_0 .

We can derive the PMF and CDF simply.

$$\begin{aligned}
f(x; \alpha, \beta) &= \int_0^\infty P(X = x | \lambda) P(\lambda | \alpha, \beta) d\lambda \\
&= \int_0^\infty \frac{\lambda^x e^{-\lambda}}{x!} \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} d\lambda \\
&= \frac{\beta^\alpha}{\Gamma(\alpha)\Gamma(x+1)} \int_0^\infty \lambda^{x+\alpha-1} e^{-\lambda(1+\beta)} d\lambda \\
&= \frac{\Gamma(\alpha+x)\beta^\alpha}{\Gamma(\alpha)\Gamma(x+1)(1+\beta)^{\alpha+x}} \\
F(x; \alpha, \beta) &= \int_0^\infty P(X \leq x | \lambda) P(\lambda | \alpha, \beta) d\lambda \\
&= \int_0^\infty \sum_{k=0}^{\lfloor x \rfloor} \frac{\lambda^k e^{-\lambda}}{k!} \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} d\lambda \\
&= \sum_{k=0}^{\lfloor x \rfloor} \frac{\beta^\alpha}{\Gamma(\alpha)\Gamma(k+1)} \int_0^\infty \lambda^{k+\alpha-1} e^{-\lambda(1+\beta)} d\lambda \\
&= \sum_{k=0}^{\lfloor x \rfloor} \frac{\Gamma(\alpha+k)\beta^\alpha}{\Gamma(\alpha)\Gamma(k+1)(1+\beta)^{\alpha+k}}
\end{aligned}$$

The integration comes from a substitution and applying the definition of the Gamma function, we use the identity $\Gamma(x) = (x-1)!$ for positive integers, the interchanging of the sum and the integral comes from Tonelli's theorem the requirements of which are satisfied by all distributions.

Definition II.2. Gamma-Poisson We define the Gamma-Poisson distribution with $\alpha \in (0, \infty)$, $\beta \in (0, \infty)$ and support $x \in \mathbb{N}_0$

$$f(x; \alpha, \beta) = \frac{\Gamma(\alpha+x)\beta^\alpha}{\Gamma(\alpha)\Gamma(x+1)(1+\beta)^{\alpha+x}}$$

To find the moments to show the variance isn't a function of the mean like Poisson we construct the moment generating function using that $X|\lambda \sim Poi(\lambda)$ and $\lambda \sim Gamma(\alpha, \beta)$ their known moment generating functions.

$$\begin{aligned}
M_{X|\lambda}(t) &= \exp(\lambda(e^t - 1)) \\
M_X(t) &= E(M_{X|\lambda}(t)) \\
\implies M_X(t) &= E(\exp(\lambda(e^t - 1))) \\
&= M_\lambda(e^t - 1) \\
&= \left(1 - \frac{e^t - 1}{\beta}\right)^{-\alpha}
\end{aligned}$$

We can use this to find the moments and show that we can control vary both the mean and variance with the parameters.

$$\begin{aligned}
E(X) &= \frac{\alpha}{\beta} \\
E(X^2) &= \frac{\alpha(\beta+1)+1}{\beta^2}
\end{aligned}$$

It is also important to note that this distribution is equivalent to the Negative Binomial distribution with a change of variables $\alpha = r$, $\beta = \frac{1-p}{p}$. Negative Binomial Mixture Models have seen use with some success.

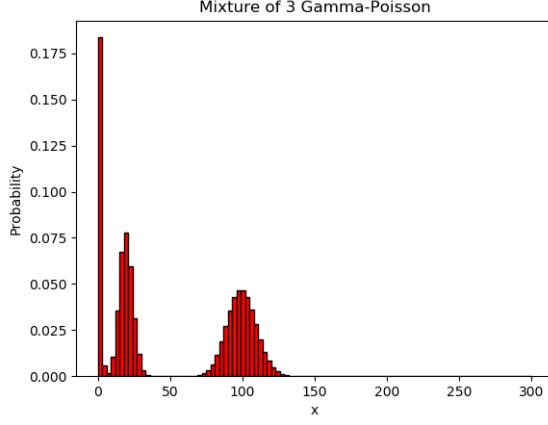


Figure 4: A mixture of 3 Gamma-Poissons with model parameters $\alpha = (5000, 300, 1000)$, $\beta = (50, 300, 50)$ and weights $(0.4, 0.3, 0.3)$ for bins of width 3 from 0 to 300

II.H Bivariate Gamma-Poisson Mixture Model

We can extend this to the bivariate case by using the product of two Gamma-Poisson distributions with their weight parameters weighted by p_1, p_2 as suggested in *Nelson* [24]. This leads to the probability distribution.

$$\begin{aligned} f(X_{12} = (i, j); \alpha, \beta, p) &= \int_0^\infty P(X_1 = i|p_1\lambda)P(X_2 = j|p_2\lambda)P(\lambda|\alpha, \beta)d\lambda \\ &= \frac{\Gamma(i + j + \alpha)\beta^\alpha p_1^i p_2^j}{\Gamma(\alpha)\Gamma(i + 1)\Gamma(j + 1)(1 + \beta)^{\alpha-(i+j)}} \end{aligned}$$

Definition II.3. Bivariate Gamma-Poisson We define the Bivariate Gamma-Poisson distribution with $\alpha \in (0, \infty)$, $\beta \in (0, \infty)$, $p_1 \in (0, 1]$, $p_2 \in (0, 1]$ with $p_1 + p_2 = 1$ and support $i, j \in \mathbb{N}_0$

$$f(i, j; \alpha, \beta, p) = \frac{\Gamma(i + j + \alpha)\beta^\alpha p_1^i p_2^j}{\Gamma(\alpha)\Gamma(i + 1)\Gamma(j + 1)(1 + \beta)^{\alpha-(i+j)}}$$

A suggested extension to this model from *Nelson* [24] that we use is the Dirichlet-Multinomial Gamma-Poisson Model. We model the weights p_1, p_2 with the Dirichlet distribution with parameters θ_1, θ_2 this leads to the the probability distribution. A graphical plot of this distribution is shown in figure 5.

$$\begin{aligned} f(X_{12} = (i, j); \alpha, \beta, \theta) &= \int_0^\infty \int_0^1 P(X_1 = i|p_1\lambda)P(X_2 = j|p_2\lambda)P(\lambda|\alpha, \beta)P(P|\theta)dPd\lambda \\ &= \frac{\Gamma(i + j + \alpha)\Gamma(\theta_1 + \theta_2)\Gamma(\theta_1 + i)\Gamma(\theta_2 + j)\Gamma(i + 1)\Gamma(j + 1)\beta^\alpha}{\Gamma(\alpha)\Gamma(\theta_1 + \theta_2 + i + j)\Gamma(\theta_1)\Gamma(\theta_2)(1 + \beta)^{\alpha-(i+j)}} \end{aligned}$$

Definition II.4. Dirichlet-Multinomial Gamma-Poisson We define the Dirichlet-Multinomial Gamma-Poisson distribution with $\alpha \in (0, \infty)$, $\beta \in (0, \infty)$, $\theta_1 \in (0, \infty)$, $\theta_2 \in (0, \infty)$ and support $i, j \in \mathbb{N}_0$

$$f(i, j; \alpha, \beta, \theta) = \frac{\Gamma(i + j + \alpha)\Gamma(\theta_1 + \theta_2)\Gamma(\theta_1 + i)\Gamma(\theta_2 + j)\Gamma(i + 1)\Gamma(j + 1)\beta^\alpha}{\Gamma(\alpha)\Gamma(\theta_1 + \theta_2 + i + j)\Gamma(\theta_1)\Gamma(\theta_2)(1 + \beta)^{\alpha-(i+j)}}$$

An important feature we would like is the ability to model the correlation between the two variates. We apply the law of total covariance to get.

$$\begin{aligned} COV(X_1, X_2) &= \mathbb{E}(COV(X_1, X_2|p_1\lambda, p_2\lambda)) + COV(\mathbb{E}(X_1|p_1\lambda), \mathbb{E}(X_2|p_2\lambda)) \\ &= \mathbb{E}(\mathbb{E}(X_1 X_2|p_1\lambda, p_2\lambda) - \mathbb{E}(X_1|p_1\lambda)\mathbb{E}(X_2|p_2\lambda)) + COV(p_1\lambda, p_2\lambda) \\ &= 0 + \mathbb{E}(\lambda^2)COV(p_1, p_2) \\ &= -\frac{\alpha(1 + \beta)}{\beta^2} \frac{\theta_1\theta_2}{(\theta_1 + \theta_2)^2(\theta_1 + \theta_2 + 1)} \end{aligned}$$

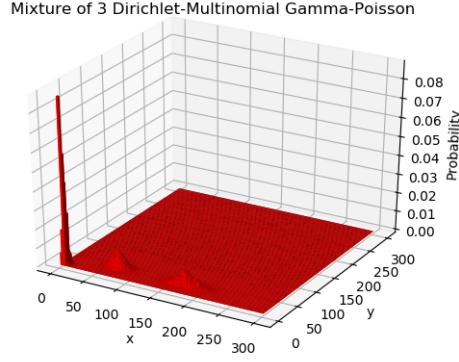


Figure 5: A mixture of 3 Dirichlet-Multinomial Gamma-Poissons with model parameters $\alpha = (1000, 5000, 60000)$, $\beta = (5, 50, 6000)$, $\theta_1 = (25171, 36770, 133977)$, $\theta_2 = (179284, 90339, 127516)$ and weights $(0.4, 0.3, 0.3)$ for bins of width 3 from 0 to 300 in both dimensions

II.I Poisson Log-Normal Mixture Model

A distribution that has seen use in modelling in biology with species abundance *Bulmer* [9], microbial counts in food *Gonzales-Barron and Butler* [17] and other settings is the Poisson Log-Normal distribution. Proposed by *AITCHISON and HO* [1] as a way to better model variation and correlation structures when applied to higher dimensions. We model the rate parameter λ of the Poisson distribution by a Log-Normal distribution. A graphical plot of a mixture of Poisson Log-Normals is shown in figure 6.

Definition II.5. Poisson Log-Normal We define the Poisson Log-Normal distribution with μ, σ and support $x \in \mathbb{N}_0$ where $P(\lambda|\mu, \sigma)$ is the PDF of the Log-Normal distribution as

$$f(x; \mu, \sigma) = \int_{\mathbb{R}^+} P(X = x|\lambda)P(\lambda|\mu, \sigma)d\lambda$$

Definition II.6. Log-Normal We define the Log-Normal distribution with μ, σ and support $x \in \mathbb{R}^+$

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\log(x) - \mu)^2}{2\sigma^2}\right)$$

Unfortunately this integral has no analytical solution so we rely on numerical methods to evaluate the PDF. An easy and effective method that works in the univariate case is quadrature based on the quartiles of the Log-Normal distribution. Sadly this doesn't work in higher dimensions due to the inverse CDF in the multivariate case having more than one solution. As the Log-Normal distribution is a transform of the Normal distribution we can define the inverse CDF as the inverse transform of the the inverse normal CDF.

$$\begin{aligned} \Phi_{LN}^{-1}(p) &= \exp(\Phi_N^{-1}(p)) \\ &= \exp(\mu + \sigma\sqrt{2}\operatorname{erf}^{-1}(2p - 1)) \end{aligned}$$

We then for n quadrature points take $n+1$ points p_i in $(0,1)$ we don't take $p=0,1$ as this can lead to inf or NaN values. We evaluate the LN distribution at these values and take the pairwise average to give λ_i . For higher dimensions as suggested in *AITCHISON and HO* [1] multivariate Hermitian integration can be used.

$$f(x; \mu, \sigma) \approx \sum_{i=1}^n \frac{1}{n} f_{Po}(x; \lambda_i)$$

We can find its n th moment easily by noticing if X is Log-Normal then $Y = \log(X)$ is normal

$$\mathbb{E}(X^n) = \mathbb{E}(\exp(nY)) = \int_{\mathbb{R}} \exp(ny) \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right) dy$$

This is the moment generating function of the normal distribution which we know to be

$$\mathbb{E}(X^n) = \exp\left(n\mu + \frac{1}{2}n^2\sigma^2\right)$$

The expectation and variance of this distribution is calculated easily as

$$\mathbb{E}(X) = \exp\left(\mu + \frac{1}{2}\sigma^2\right)$$

$$\begin{aligned} \text{var}(X) &= \exp(2\mu + 2\sigma^2) + \left(\exp\left(\mu + \frac{1}{2}\sigma^2\right)\right)^2 \\ &= \exp(2\mu + \sigma^2) (\exp(\sigma^2) + 1) \end{aligned}$$

Allowing us to change the variance and mean independently.

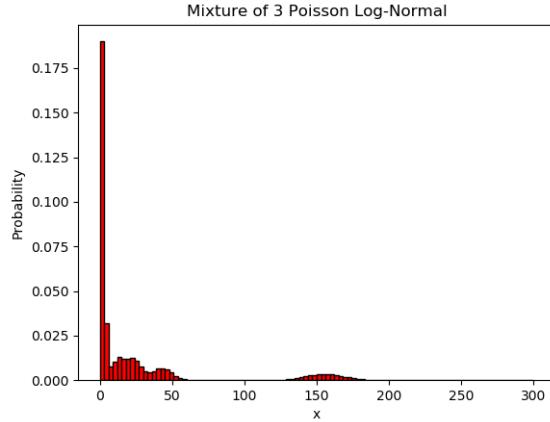
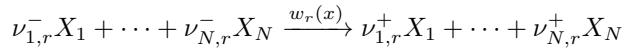


Figure 6: A mixture of 3 Poisson Log-Normal distributions with model parameters $\mu = (-1, 4, 2), \sigma = (3, 3, 1)$ and weights $(0.4, 0.3, 0.3)$ for bins of width 3 from 0 to 300

II.J Reaction Networks

Definition II.7. Reaction Network A reaction network consisting of a set of species X_1, \dots, X_N with abundances $\mathbf{x} = (x_1, x_2, \dots, x_N)^T$ that interact via a set of R reactions of the form



where $r = 1, \dots, R$, and the stoichiometric coefficients $(\nu_{ir}^\pm)_{i=1, \dots, N}$ are non-negative integers and $w_r(x) \geq 0$ is the reaction rate.

This leads us to a system of ODEs.

Definition II.8. Deterministic rate equation The deterministically modelled reaction network is given by the system of ODEs

$$\frac{dx_i}{dt} = \sum_{r=1}^R (\nu_{ir}^+ - \nu_{ir}^-) w_r(x) \text{ for } i = 1, \dots, N$$

called the rate equations which are solved subject to the initial condition $x(0) = x_0$

An important definition we will use in the derivation of analytical solutions is

Definition II.9. Mass-action Kinetics The vector of reaction rates $\mathbf{w}(x)$ using mass action kinetics are given by

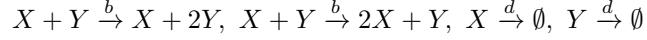
$$w_r(x) = k_r \prod_{i=1}^N x_i^{\nu_{ir}^-} \text{ for } r = 1, \dots, R$$

where k_r is the reaction rate constant

Not all systems are solvable analytically such as the regulating gene model and thus if we want a deterministic solution we would be required to use a numerical ODE solver.

II.J.1 Deterministic

We can think of a solution to our problem in a deterministic way without any stochasticity. An example where we can find an analytical deterministic solution is population explosion. Consider the two sexes (male and female) with equal rates of reproduction b and mortality d



Using Mass-action Kinetics this forms the system of ODEs

$$\frac{dx}{dt} = bxy - dx, \quad \frac{dy}{dt} = bxy - dy$$

If we chose the initial conditions $x(0) = y(0)$, clearly $x(t) = y(t) \forall t$ and there is no distinction between female and male populations. Thus

$$\frac{dx}{dt} = bx^2 - dx$$

which is a Bernoulli differential equation which has known exact solutions. In this case it is

$$x(t) = \frac{dx(0)}{(d - bx(0))e^{dt} + bx(0)}$$

However, this solution does not capture any of the stochastic nature of these systems which we would expect to exist. If we compared our analytical solution to real world data we would likely see heterogeneity that isn't accounted for by the deterministic setting.

II.J.2 Stochastic

We can seek to model the reactions in a more realistic stochastic formulation. As the number of events is necessarily discrete this is a counting process.

Definition II.10. Counting process $Y(t)$ is a counting process if $Y(t)$ is constant except for jumps of size ± 1 , whose paths are cadlag, i.e., they are right-continuous with left-limits $Y(t^-) = \lim_{s \rightarrow t} Y(s)$.

The simplest counting process is the Poisson Process and we can often interpret simple reaction networks as Poisson Processes.

Definition II.11. Poisson Process A counting process is a Poisson Process with rate λ if for $t < s$, it holds that

$$Pr(Y(s) - Y(t) = k) = \lambda^k \frac{(s-t)^k}{k!} e^{-\lambda(s-t)}$$

We also make a modification to the law of Mass-action for the stochastic case and give the Chemical Master Equation the proof of which comes from Dykin's formula.

Definition II.12. Mass-action Kinetics for stochastic reaction networks The vector of reaction rates $\mathbf{w}(x)$ using mass action kinetics are given by

$$w_r(x) = k_r \prod_{i=1}^N \frac{x_i!}{(x_i - \nu_{ir})!} \text{ for } r = 1, \dots, R$$

where k_r is the reaction rate constant

Definition II.13. Chemical Master Equation The time-dependent probability of the stochastic reaction network $p(z, t) = Pr(x(t) = z)$ satisfies the Chemical Master Equation:

$$\frac{dp(z, t)}{dt} = \sum_{r=1}^R [w_r(z - \nu_r)p(z - \nu_r) - w_r(z)p(z)]$$

The Chemical Master Equation is often very difficult to solve analytically if it is even possible. An example that exhibits a closed form solution is a Linear birth-death process under the detailed balance equation, thus giving us the stationary distribution.

$$\emptyset \xrightarrow{a} x, x \xrightarrow{b} \emptyset, w_+(x) = a, w_-(x) = bx$$

Where $w_+(x)$ and $w_-(x)$ are the birth and death rates respectively giving us the Chemical Master Equation:

$$\frac{dp(x,t)}{dt} = w_+(x-1)p(x-1) - w_+(x)p(x) + w_-(x+1)p(x+1) - w_-(x)p(x)$$

Assuming stationarity and using the detailed balance condition we get

$$w_+(x-1)p(x-1) = w_-(x)p(x)$$

Giving us

$$p(x) = p(0) \prod_{i=0}^x \frac{w_+(i-1)}{w_-(i)}$$

Which exists if

$$\frac{1}{p(0)} = \sum_{x=0}^{\infty} \prod_{i=0}^x \frac{w_+(i-1)}{w_-(i)} < \infty$$

Substituting in for the reaction rates we get

$$p(x) = \exp\left(-\frac{a}{b}\right) \frac{1}{x!} \left(\frac{a}{b}\right)^x$$

Which is Poisson with rate $\frac{a}{b}$. The solution if we looked at the Linear birth-death process as deterministic would be.

$$\begin{aligned} \frac{dx}{dt} &= a - bx \\ \implies x_t &= e^{-bt}x_0 + \frac{a}{b}(1 - e^{-bt}) \end{aligned}$$

Although, in general solutions like this aren't obtainable analytically such as in our case of a self-regulating gene and thus this leads us to the need for a simulation algorithm.

II.J.3 Gillespie Algorithm

A exact simulation algorithm exists in the Gillespie Allgorithm from *Gibson and Bruck* [16]. The algorithm exists in two equivalent formulations below we describe the algorithm for the Direct Method which calculates explicitly which reaction occurs next and when it occurs. We consider a system of r reactions each with an associated rate constant k_1, \dots, k_r these can be a true constant or have time variance. The system is in exactly one state at each time step. This leads that at a given time step the system has r possible transitions represented by each reaction from its state. Using a random number generator with the correct distribution we can choose the next state from the r possible transitions. We specify for the next reaction μ and its occurrence time τ a probability density $P(\mu, \tau)$. This can be shown to be

$$P(\mu, \tau)d\tau = \alpha_\mu \exp\left(-\tau \sum_j \alpha_j\right) d\tau$$

Where α_i is the propensity function of reaction i derived from the law of mass action. We can obtain the the probability of each reaction by integrating over 0 to ∞

$$Pr(\text{Reaction} = \mu) = \frac{\alpha_\mu}{\sum_j \alpha_j}$$

The distribution for the times can be found by summing over all μ

$$Pr(\tau)d\tau = \left(\sum_j \alpha_j\right) \exp\left(-\tau \sum_j \alpha_j\right) d\tau \sim Exp\left(\sum_j \alpha_j\right)$$

We formulate this into the following algorithm

Algorithm 1 Direct Method Gillespie

```
Initialise Set initial numbers of species, stop time  $t_{MAX}$  and  $t = 0$ 
while  $t < t_{MAX}$  do
    Calculate the propensity  $\alpha_i \forall i = 1, \dots, r$ 
    Choose the next reaction  $\mu$  according to the distribution  $Pr(\text{Reaction} = \mu)$ 
    Choose the next timestep  $\tau$  according to the distribution  $Exp\left(\sum_j \alpha_j\right)$ 
    Change the number of molecules to reflect reaction  $\mu$  occurring and  $t \leftarrow t + \tau$ 
end while
```

The equivalent and faster method of the Gillespie First Reaction Method using the previously defined for the Direct Method.

Algorithm 2 First Reaction Method Gillespie

```
Initialise Set initial numbers of species, stop time  $t_{MAX}$  and  $t = 0$ 
while  $t < t_{MAX}$  do
    Calculate the propensity  $\alpha_i \forall i = 1, \dots, r$ 
    for  $i = 1, \dots, r$  do
        Generate a putative time  $\tau_i$  according to  $Exp\left(\sum_j \alpha_j\right)$ 
    end for
    Choose  $\mu = \text{argmin}_{i \in \{1, \dots, r\}}(\tau_i)$  set  $\tau = \tau_\mu$ 
    Change the number of molecules to reflect reaction  $\mu$  occurring and  $t \leftarrow t + \tau$ 
end while
```

This allows us to generate samples from the distribution. Even for small reaction networks such as a regulating gene the generation of a useful number of independent samples is time consuming and for larger networks or for greater exploration of the reaction rate space it can be prohibitive. Faster non-exact algorithms exist such as tau-leaping shown in *Gibson and Bruck* [16] but these methods although faster still require a large amount of time to explore the reaction rate space of large reaction networks.

II.K Ensemble Models

We suggest a few possible ensembles as efforts to improve upon our models individual performance. An important requirement to build an ensemble with additional predictive power over the individual models is variation in the predictions of the submodels. This is rather obvious as if all our models predicted distributions that are arbitrarily close to each other there would be no need for an ensemble. As neural networks are stochastic in nature with random initial weights and using random batches to train on there exists variation between our trained models which is how an ensemble was constructed in *Wang et al.* [36] but in an effort to increase this variation we use bagging proposed by *Breiman* [8]. Bagging splits the dataset into a training data set and test data set as before but we train k submodels on a different random sample with replacement of the training dataset. This increases variation between the models by increasing and decreasing at random the occurrence of different patterns in the dataset. We then use a dense layer or multiple layers to take the predictions of each submodel and form a new prediction. Ensemble methods have frequently been shown to increase the predictive ability of an individual model provided there exists variation between the models as shown in *Wang et al.* [36] and *Breiman* [8].

III Alternative Emulation Methods

Alternative methods exist in the Bayesian setting with approximate Bayesian computation *Beaumont* [5]. These methods involve sampling from a parameter prior distribution then sampling a posterior distribution based on the prior and accepting or rejecting the simulated posterior based on summary statistics of the distribution or the full distribution. These methods require a large number of random samples which are often computationally expensive and the complexity increases greatly as dimensionality increases. Improvements upon these methods using Markov Chain Monte-Carlo *Beaumont* [4] and sequential Monte-Carlo *Sisson et al.* [32] techniques have been proposed but still require a large number of time consuming model simulations. Further suggested improvements include the use of regression or Gaussian Process emulation of these models where simulation is difficult *Wilkinson* [37].

The choice of parameters greatly affect the performance of such models. Also the choice of summary statistics can greatly affect performance and there is no guarantee we can find a finite dimensional number of summary statistics that encompass all the information of the distribution exist outside of the exponential family *Robert et al.* [28]. As such these methods can be too complex and time consuming for problems of interest as such an alternative method less effected by dimensionality and require generating random numbers would be useful.

IV Proposed Long Short Term Memory (LSTM) Model

We propose using a model based on *Wang et al.* [36] consisting first a series of 5 dense layers of 64, 128, 256, 128, 64 neurons to learn general model characteristics from the parameters. Then we split this off into a dense layer for each T time steps which feeds into one or two LSTM layers depending on if it is bidirectional. Taking in the p previous time steps predicted by the LSTM array. This allows the model to predict based off previous points in the curve. The outputs of the LSTM arrays are fed into a dense layer to give the T points. The model structure we use is shown in figure 7. We can use any activation we want for the dense layers at the start but for the final dense layer we use Softmax that gives an output that sums to 1. This model also allows for the use of both the Standard LSTM Cell and the Peephole LSTM Cell.

To form an ensemble of these models we train a given number of models n on random samples of a training dataset and then take their outputs to input into a dense layer with an input dimension of nT and output dimension of T . We use a Softmax activation function as our activation function as in the submodels.

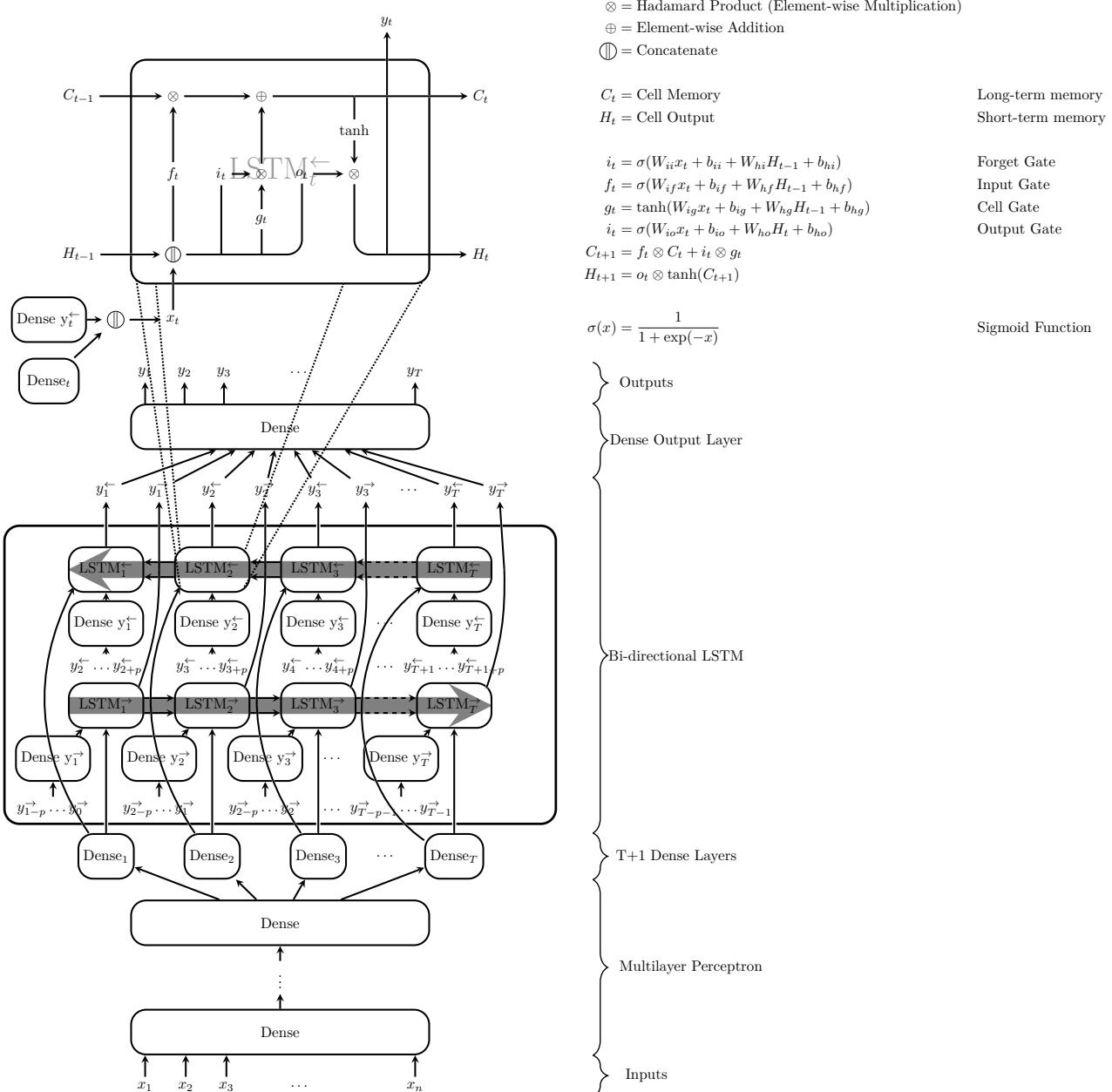


Figure 7: General structure for the Bidirectional LSTM Model

V Proposed Mixture Model

We utilise a feed foward neural network to predict the distribution parameters $\boldsymbol{\theta}$ and the distribution weights α given the input which in the case of the replicating gene model is the reaction rates. We cannot calculate a loss directly from the parameters so we construct our estimate for each bin using the CDF of our chosen distribution $\Phi(x; \boldsymbol{\theta}_i)$. We would most like to choose distributions that have the same or similar support to the distribution we are estimating and at least the ability to change the first and second moments without them being a function of each other that is the case for the proposed distributions.

Say our distribution takes p parameters ($\theta^{(1)} \dots \theta^{(p)} \in \mathbb{R}^n$) and there are r inputs for the neural network then a general diagram is shown in figure 8. The complexity of model is invariant to the number of bins which is useful if we want to train with lots of small bins. These models also extend to the multivariate case if we use a mixture of multivariate distributions. This would allow for the modelling of the joint distribution of multiple species if we wish with a minimal increase in the model's complexity. The dense layers before the output layer attempt to learn the general mapping from the inputs to the distribution parameters increasing its depth and size can increase the models ability to learn more complex mappings. We would therefore expect more complex reaction networks to require this to be deeper. In our experiments we use 5 layers of 64, 128, 256, 128, 64 neurons. The number of mixtures controls the maximum complexity of the output mixture thus increasing the number of mixtures may not increase the model's performance if the initial dense layers can't learn a good mapping or if the distribution is not appropriate. The activation of the dense layers can be any that works best on a given problem but the activation functions of the output layers must be chosen carefully to map to the values allowed for that parameter or a subset of those values. The activation for the layer predicting the weights α is Softmax as this enforces the summing to unity and values in the range [0, 1].

Using neural networks as a way to get the parameters of a Gaussian Mixture Model has seen some limited use but on rather different and more simplistic problems such as the mapping of an inversion of a function with noise to obtain the conditional probability of the input of the function given the output *Bishop* [7]. Taking a mean, mode or the mean of the most heavily weighted distribution in the mixture as the output depending in the application. It predicted the inputs of highly non-linear function given the outputs much better than a simple feed forward neural network.

To form ensembles of this model we use a series of dense layers that takes the outputs of all the submodels into a dense layer. This will be of dimension $n(p+1)d$ where d is the sum of the dimensionality of each model parameter. We can feed this through multiple dense layers before as in the individual models we have a dense layer for each parameter with appropriate activation functions.

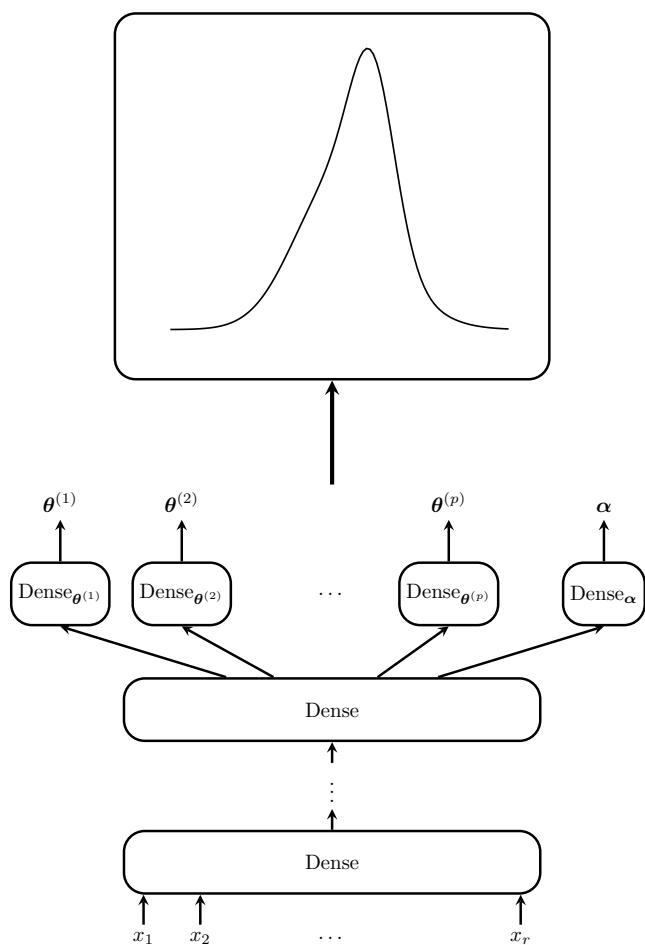


Figure 8: General structure of the Mixture Model

VI Loss Functions

To show the vanishing gradient problem caused by common loss functions and show alternatives which offer larger gradients for our problem we construct 3D surface plots to show their shape in figure 9. We use T points and a batch size of B when defining the loss functions.

VI.A Mean Squared Error

The mean squared error is commonly used in regression problems as it is twice differentiable and the gradient reduces as $x_t \rightarrow y_t$.

Definition VI.1. MSE The Mean Squared Error for $X, Y \in \mathbb{R}^{B \times T}$ of B batches and T points

$$MSE(X, Y) = \frac{1}{B} \frac{1}{T} \sum_{i=1}^B \sum_{t=1}^T (y_t^{(i)} - x_t^{(i)})^2$$

We can see visually that the gradient becomes very small for the values we expect in figure 9. Also remembering most values will be closer to 0 than 1 as we don't expect y_t to be uniformly distributed in its interval.

VI.B L1

One common alternative is the L1 Norm but a disadvantage of this approach is this function is not twice differentiable and thus can lead to the model bouncing around optimal solutions. The plot of this is shown in figure 9.

Definition VI.2. L1 The L1 Norm for $X, Y \in \mathbb{R}^{B \times T}$ of B batches and T points

$$L1 = \frac{1}{B} \frac{1}{T} \sum_{i=1}^B \sum_{t=1}^T |y_t^{(i)} - x_t^{(i)}|$$

VI.C Squared Hellinger Distance

We propose as an alternative measure of loss the Squared Hellinger Distance. This is a f-divergence measure *Eguchi* [13], it's twice differentiable and is a measure of distance due to it being symmetric. The attractiveness of this measure is it has similar features to the MSE with it's gradient becoming smaller as it approaches its minimum but its gradient is much larger than that of the MSE in the domain of our problems. We do however need to be careful that when the model predicts a 0 as when calculating the gradient during back propagation this leads to NaN due to a division by zero. This is easily resolved by not letting the prediction fall below a small $\epsilon > 0$ and also doing this for the true value y . This effectively means if they are both below ϵ we assume they are both equal and zero.

Definition VI.3. H^2 The Squared Hellinger Distance for $X, Y \in \mathbb{R}^{B \times T}$ of B batches and T points

$$H^2(X, Y) = \frac{2}{B} \sum_{i=1}^B \sum_{t=1}^T \left(\sqrt{y_t^{(i)}} - \sqrt{x_t^{(i)}} \right)^2$$

We can check this is a suitable loss by first checking the differential w.r.t x.

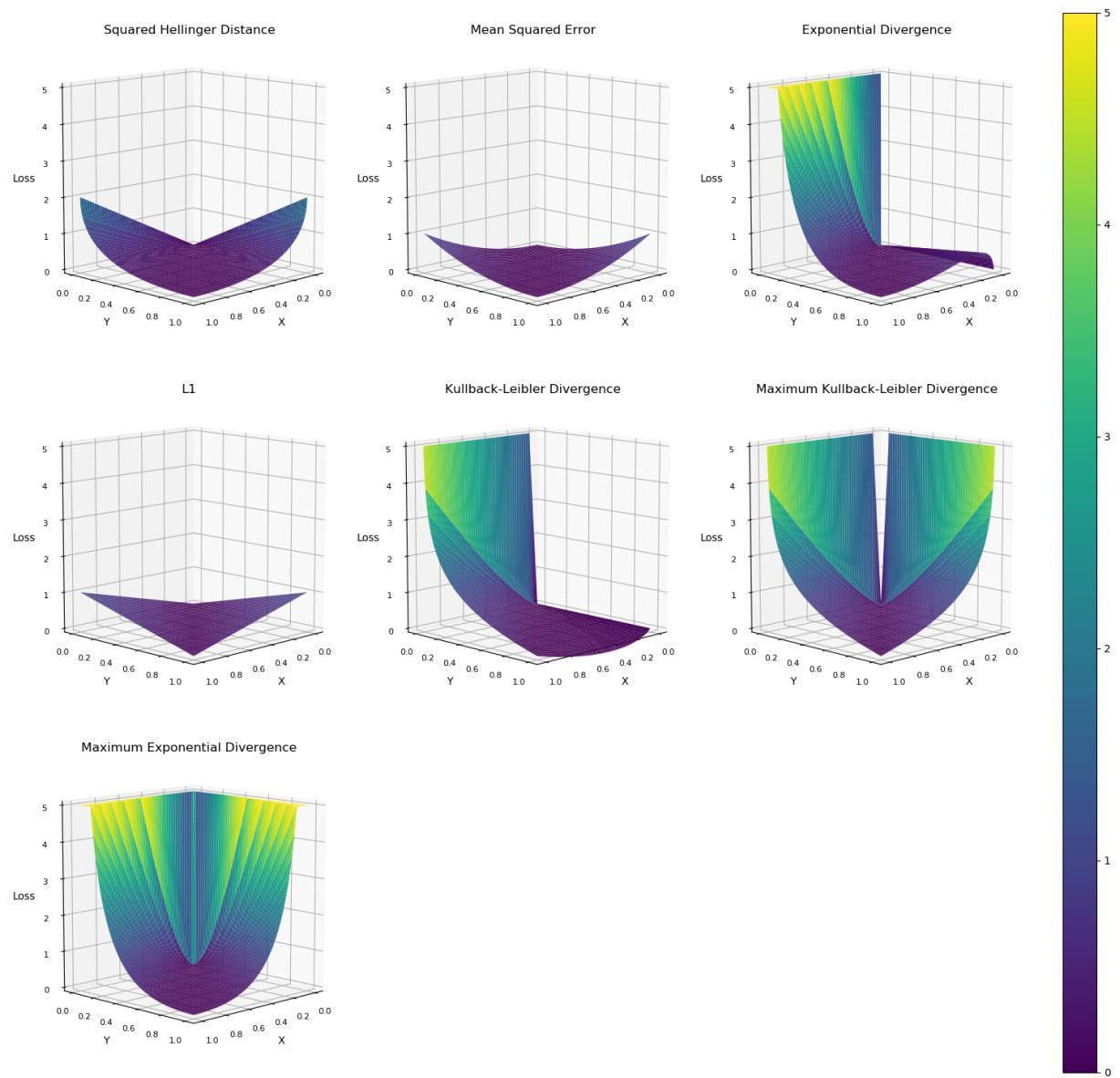


Figure 9: Proposed loss functions plotted for $x \in [0, 1], y \in [0, 1]$ plotted up to a maximum loss of 5. The colour map has the same scale for all the plots to show their differences in terms of magnitude. This also clearly shows the existance of other minimas in the divergence measures and their lack of symmetry.

$$\frac{\partial H^2}{\partial x_t} = \frac{1}{B} \sum_{i=1}^B \sum_{t=1}^T \left(1 - \frac{\sqrt{y_t^{(i)}}}{\sqrt{x_t^{(i)}}} \right)$$

$$\sum_{i=1}^B \sum_{t=1}^T y_t^{(i)} = \sum_{i=1}^B \sum_{t=1}^T x_t^{(i)} \implies \frac{\partial H^2}{\partial x_t} = 0$$

$$\sum_{i=1}^B \sum_{t=1}^T y_t^{(i)} > \sum_{i=1}^B \sum_{t=1}^T x_t^{(i)} \implies \frac{\partial H^2}{\partial x_t} < 0$$

$$\sum_{i=1}^B \sum_{t=1}^T y_t^{(i)} < \sum_{i=1}^B \sum_{t=1}^T x_t^{(i)} \implies \frac{\partial H^2}{\partial x_t} > 0$$

We also check that the loss function is convex by looking at the 2nd differential and showing it is greater than or equal to 0 for the entire domain this is also implied by it being an f-divergence measure *Eguchi* [13]. This shows the minimum is a global minimum for the entire domain.

$$\frac{\partial^2 H^2}{\partial x_t} = \frac{1}{B} \sum_{i=1}^B \sum_{t=0}^T \frac{\sqrt{y_t^{(i)}}}{2(x_t^{(i)})^{\frac{3}{2}}} \geq 0$$

We also show a plot of this loss function in figure 9.

VI.D Kullback-Leibler Divergence

We also propose the use of the KL Divergence. This is a twice differentiable function and as a f-divergence measure it's a useful metric for evaluating the performance of the model and is convex.

Definition VI.4. KLD The Kullback-Leibler Divergence for $X, Y \in \mathbb{R}^{B \times T}$ of B batches and T points

$$KLD(X, Y) = \frac{1}{B} \sum_{i=1}^B \sum_{t=1}^T (x_t^{(i)}) \log \left(\frac{x_t^{(i)}}{y_t^{(i)}} \right)$$

One problem with this loss function is it requires x and y to be from a proper probability distribution or it can give negative values leading to convergence to incorrect minima. To be able to use it as a loss function we must normalise the output of our model so it sums to 1. We also as we are taking logs we use clamp to enforce a minimum to prevent $\log(0)$. A plot of this function is shown in figure 9. It is important to note that this function is not symmetric so not a distance this can make its value harder to interpret. The existence of this other minimum makes it not useful as a loss function to give to our optimiser as we can get convergence to the wrong minimum. In the current implementation of Pytorch 1.5.0 there exists overflow errors causing negative values that are attributed to it requiring a log input which have been fixed in the nightly version and will be included in the next release.

VI.E Exponential Divergence

Another f-divergence measure which should offer larger gradients than the previous f-divergence measures is the Exponential Divergence (ED) *Eguchi* [13]. As it hasn't been used as a loss function before we calculate derivatives to check but these results are also implied through definition of an f-divergence measure.

Definition VI.5. ED The Exponential Divergence for $X, Y \in \mathbb{R}^{B \times T}$ of B batches and T points

$$ED(X, Y) = \frac{1}{B} \sum_{i=1}^B \sum_{t=1}^T (x_t^{(i)}) (\log(x_t^{(i)}) - \log(y_t^{(i)}))^2$$

$$\begin{aligned}
\frac{\partial ED}{\partial x_t} &= \frac{1}{B} \sum_{i=1}^B \sum_{t=1}^T \log \left(\frac{x_t^{(i)}}{y_t^{(i)}} \right) \\
&\quad \left(\log \left(\frac{x_t^{(i)}}{y_t^{(i)}} \right) + 2 \right) \\
\sum_{i=1}^B \sum_{t=1}^T y_t^{(i)} &= \sum_{i=1}^B \sum_{t=1}^T x_t^{(i)} \implies \frac{\partial ED}{\partial x_t} = 0 \\
\sum_{i=1}^B \sum_{t=1}^T y_t^{(i)} &> \sum_{i=1}^B \sum_{t=1}^T x_t^{(i)} \implies \frac{\partial ED}{\partial x_t} < 0 \\
\sum_{i=1}^B \sum_{t=1}^T y_t^{(i)} &< \sum_{i=1}^B \sum_{t=1}^T x_t^{(i)} \implies \frac{\partial ED}{\partial x_t} > 0
\end{aligned}$$

Looking at the 2nd derivative we get that the function is convex.

$$\begin{aligned}
\frac{\partial ED^2}{\partial x_t^2} &= \frac{1}{B} \sum_{i=1}^B \sum_{t=1}^T \frac{2 \left(\log \left(\frac{x_t^{(i)}}{y_t^{(i)}} \right) + 1 \right)}{x_t^{(i)}} \\
&\geq 0
\end{aligned}$$

Without the enforcing x and y from a proper probability distribution we get a non-convex function so could have our model converge to the wrong minima. We see this in the plot of the loss function in figure 9. It is important to note that this function is not symmetric like the KLD meaning $ED(X, Y) \neq ED(Y, X)$ for most X, Y because of this like KLD we would not want to give this to the optimiser to update the weights.

VI.F Converting A Divergence To A Distance Measure

We can however construct a symmetric distance measure from the Exponential Divergence and Kullback-Liebler Divergence to make its value more interpretable and also as a side effect eliminate the other minima. One approach is Intrinsic Discrepancy Loss applied to the KL divergence this offers reduced gradients over the normal KL divergence but is more robust than KLD and out performed the MSE *Bernardo et al.* [6]. For discrete distributions P and Q.

$$\delta(P, Q) = \min \left\{ \sum_{i=1}^n P(x_i) \log \left(\frac{P(x_i)}{Q(x_i)} \right), \sum_{i=1}^n Q(x_i) \log \left(\frac{Q(x_i)}{P(x_i)} \right) \right\}$$

This makes the function symmetric however it does allow for incorrect minima if P, Q are not true distributions which we may not always want to force. We suggest a modification of this function in the form of taking the maximum instead of the minimum to avoid this issue and extend this to include the exponential divergence as well as the KL and all f-divergence measures $D(\cdot, \cdot)$ in the Maximum Divergence Distance Loss $\rho(P, Q)$.

$$\rho(P, Q) = \max \{D(P, Q), D(Q, P)\}$$

The proof of symmetry comes from the symmetry of max

$$\rho(P, Q) = \max \{D(P, Q), D(Q, P)\} = \max \{D(Q, P), D(P, Q)\} = \rho(Q, P)$$

It is zero when $P = Q$ using that for all divergence measure by definition $P = Q \implies D(P, Q) = 0$.

$$\rho(P, Q) = \max \{D(P, P), D(Q, Q)\} = \max \{0, 0\} = 0$$

To satisfy the requirements of a loss function we note that the divergence measures are convex and the maximum of two convex functions is also convex so $\rho(X, Y)$ is convex. Proof that $f(x) = \max\{f_1(x), \dots, f_m(x)\}$ the maximum of convex functions f_1, \dots, f_m with $\text{dom}(f) = \text{dom}(f_1) \cap \dots \cap \text{dom}(f_m)$ is also convex. for any

$x, y \in \text{dom}(f), \lambda \in [0, 1]$ for some $j \in \{1, \dots, m\}$

$$\begin{aligned} f(\lambda x + (1 - \lambda)y) &= f_j(\lambda x + (1 - \lambda)y) \\ &\leq \lambda f_j(x) + (1 - \lambda)f_j(y) \\ &\leq \lambda\{f_1(x), \dots, f_m(x)\} + (1 - \lambda)\{f_1(y), \dots, f_m(y)\} \\ &= \lambda f(x) + (1 - \lambda)f(y) \end{aligned}$$

Plots of these two loss functions are in figure 9 and offer much steeper gradients than the other loss functions.

VI.G Maximum L1

We use as a loss function for evaluating performance the maximum L1 Norm across each batch.

Definition VI.6. L1 MAX The Maximum L1 Norm for $X, Y \in \mathbb{R}^{B \times T}$ of B batches and T points

$$L1M(X, Y) = \frac{1}{B} \sum_{i=1}^B \max_{t \in \{1, \dots, T\}} |y_t^{(i)} - x_t^{(i)}|$$

This allows us to see the maximum absolute difference between the 2 probability measures and thus giving us an interpretation of the largest deviation we would expect from the true distribution. We would expect this to help us capture when the model produces predictions with lower modality than that of the true distribution.

VII Experiments: Gaussian Mixture

As a simplistic and easy to generate example that we use to test hyper-parameters of the models. We use a mixture of three Gaussians. To be able to form bins we evaluate the CDF across the bin.

$$Y(x) = k_1\Phi(x, \mu_1, \sigma_1) + k_2\Phi(x, \mu_2, \sigma_2) + k_3\Phi(x, \mu_3, \sigma_3)$$

Where,

$$\Phi(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{x-\frac{\delta}{2}}^{x+\frac{\delta}{2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx$$

with $\sum_{i=1}^3 k_i = 1$ and δ representing the difference between each consecutive x . We sampled independently from uniform distribution $\mu \sim U(-3, 3)$, $\sigma \sim U(0.1, 5)$ and we sampled from the Dirichlet distribution their weighting $k \sim Dir(1, 1, 1)$. We used as our bins 100 bins of width 0.3 from -15 to 15. We then normalise the bins to sum to 1 and generate 10^5 random samples. As is common practice we normalise the inputs to the models μ, σ , to have 0 mean and variance 1.

We expect that the Mixture Model with a mixture of Gaussians will be able to perfectly predict the parameters as we use as the input the mean, variance and weighting of each Gaussian used to make the distribution the output. We use this example as a way to easily find the best hyperparameters for each of the models and check the can perform well on a simple experiment.

VII.A LSTM Model

We first explore the LSTM model and the effect of varying batch size on its performance with a unidirectional model, ELU activation for the initial dense layers, Sigmoid activation for the output of the LSTM array and Adam as optimiser with the Squared Hellinger Distance. The LSTM cell was the standard type with a encoder size of 512. We use a random 80/20 split between training and validation data.

We see in figure 10 plots for the KLD, L1 Norm and L1 Max and a table of their minimums of the validation dataset in figure 34. In general a reduction in L1 Loss as we reduce batch size. However, L1 Max appears to be unaffected suggesting possibly that the model struggles to identify multiple modes in all cases. The initial increase and negative KLD do suggest it is not training to an appropriate distribution. We see small batches leading to a slightly lower loss this is likely because the model over generalises in the case of larger batches. This

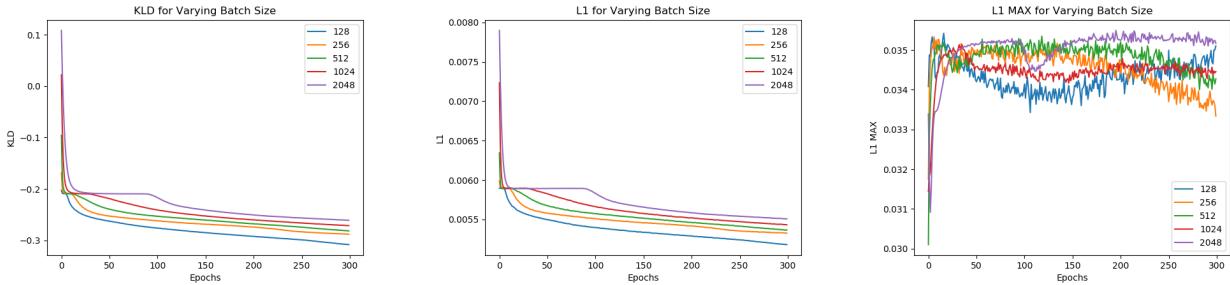


Figure 10: L1, L1 Max and KLD on the validation dataset for different batch sizes on the LSTM Model

is likely why in the source code provided in *Wang et al.* [36] they used a small batch size of 64. Looking in the table we see similar performance across the batches with a batch size of 512 achieving the lowest MSE. We therefore elect to use a batch size of 512 as the improvement for smaller batches seems small and the increase in train time is huge when we reduce the batch size.

We next explore the use of bidirectionality to increase the models predictive ability and perhaps identify modes of the distribution better in figure 11 we present their loss over epochs and in figure 35 their minimum loss across the validation data. We see that bidirectionality does not lead to a lower loss like we would expect although the difference is not significant and the loss over the training time is more consistent with the bidirectional model as the L1 Max and KLD becomes more negative as the unidirectional model trains.

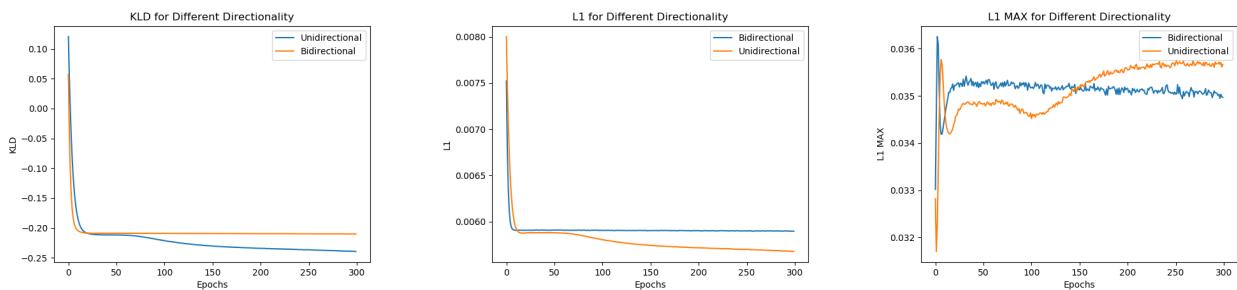


Figure 11: L1, L1 Max and KLD on the validation dataset for Unidirectional and Bidirectional LSTM models

We next test various combinations of loss function and optimiser to see which achieves the lowest loss over a given number of epochs with a unidirectional model. No results are shown for stochastic gradient descent as they all failed to train diverging to NaN in all cases. The validation loss using the Adam optimiser for various loss functions are shown in figure 12 and a table of their minimums in figure 36. They all converged to negative KLD values except in the case of the Maximum Exponential Divergence Distance Measure this is likely as this loss function is more strict in enforcing the requirement of the output being a proper distribution. Contrary to this performed worst in terms of the L1 and Maximum L1 loss. The MSE performs best in terms of the L1 Max but the L1 followed closely by the Squared Hellinger Distance perform best in terms of the L1 loss. The values of the L1 Max are high enough to suggest that the model is failing to capture some of the modes of the distribution at about 3%. We opt to use the Squared Hellinger Distance from here on but could have argued to use the L1 loss but the L1 loss appeared to be stationary where as H^2 is decreasing. We would suspect this is due to L1 not being twice differentiable so is bouncing around its optimal weights.

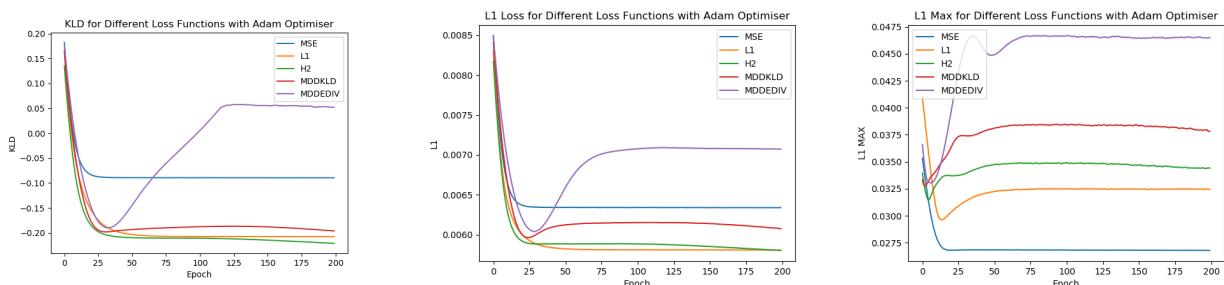


Figure 12: L1, L1 Max and KLD on the validation dataset for different loss functions with the Adam optimiser

Next we test the various choices of activation function for the initial dense layers and the activation function of the output of the LSTM array on a unidirectional LSTM model. The loss across the validation dataset is shown in figure 13 and the minimums it achieved in figure 39. We see a negative KLD again, similar L1 Max for all combinations but in terms of the L1 Loss the ELU activation function for the dense layers and Softplus for the output gives a better L1 loss with the others performing similarly. We elect to use the ELU activation and Softplus activation functions.

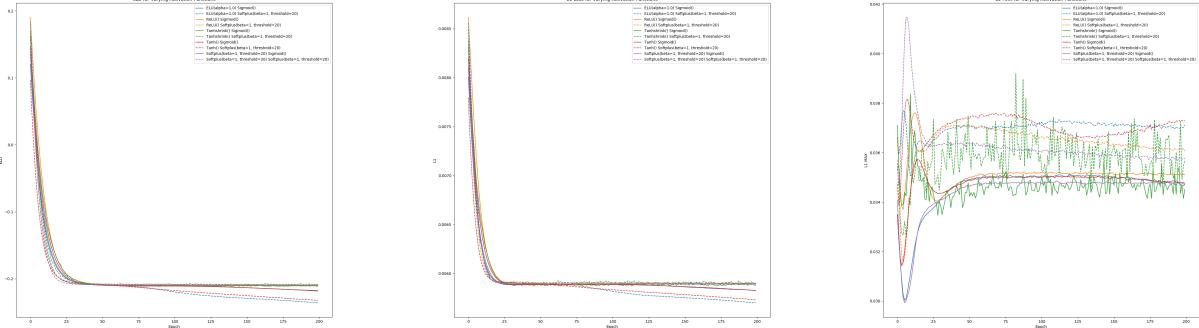


Figure 13: L1, L1 Max and KLD on the validation dataset for different activation functions. The colour and the first name in each lines label represents its initial dense layers activation function and the line type and the second name represents the activation for the output of the LSTM array

We move on to test varying LSTM types either standard or with peephole connection and their sizes. The validation loss is show in figure 14 with a table of minimums in figure 38. The 1024 encoder size standard LSTM would diverge which is why its results are not present. We see similar performance for both types of cell and the tested LSTM encoder sizes with the Peephole LSTM performing slightly better in terms of L1 Loss but worse in terms of L1 Max. The peephole cells also had a more negative KLD indicating convergence to an improper distribution. Although the difference is not very large so we opt to continue using the standard LSTM cell with an encoder size of 512.

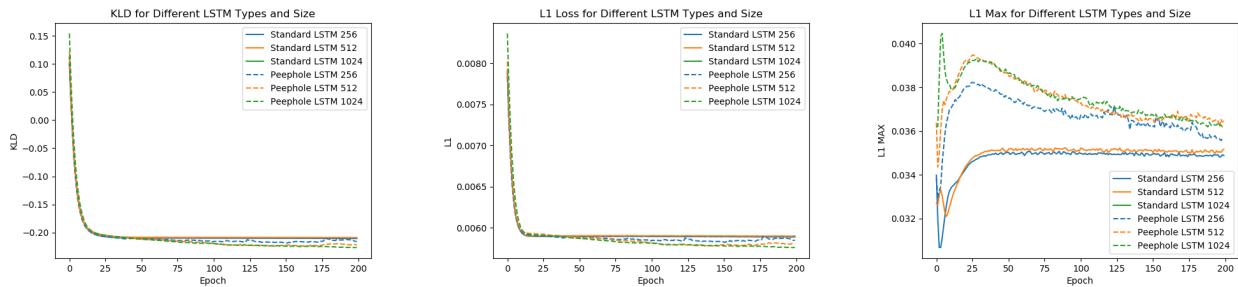


Figure 14: L1, L1 Max and KLD on the validation dataset for different LSTM Types. The colour and the number in each lines label represents its size and the line type and the name represents the type of LSTM array

To see the best performance we can expect we train for 1000 epochs our best performing LSTM model with the Squared Hellinger Distance and Bidirectionality the performance on the validation dataset is shown in figure 15 and the minimas therein achieved in figure 39. We see the model quickly reaches a set of weights such that the KLD remains negative and the L1 remain consistent. We see however the L1 Max decreases slowly as we train showing that if we trained for a longer time we could see some improvement in the predictions. The difference is not massive but the lowest value is still lower than any previously trained model.

To try and improve our performance we next train a ensemble of 5 Bidirectional LSTM models each for 1000 epochs and form an ensemble by training a dense layer to take the outputs of these 5 models and form a new prediction. The performance on the validation dataset is shown in figure 16 and it minimum loss in figure 40. We see a small increase in performance over the individual models although not massive in terms of the L1 Max but not the other 2 metrics. Which suggests that the models have similar predictions for most inputs so we don't gain much extra information from 5 separate models.

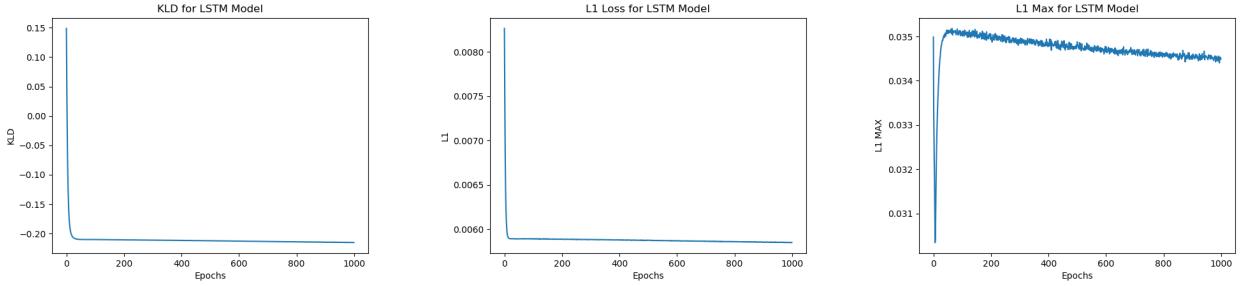


Figure 15: L1, L1 Max and KLD on the validation dataset for the LSTM Model for 1000 epochs

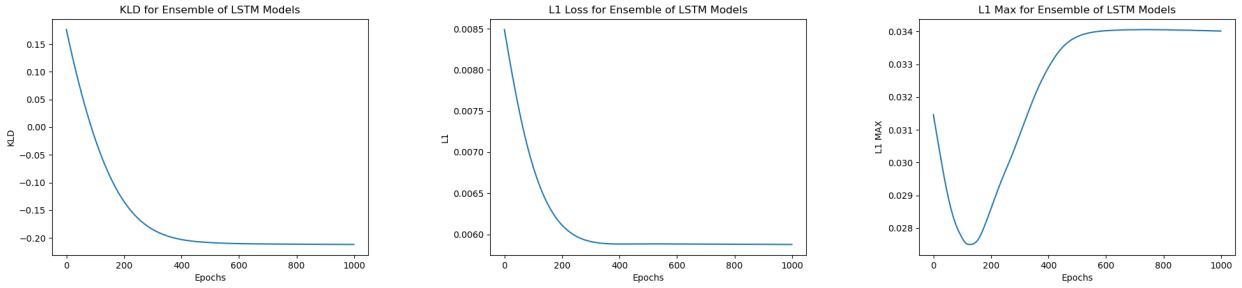


Figure 16: L1, L1 Max and KLD on the validation dataset for an ensemble of 5 LSTM Models

VII.B Mixture Model

We only use the Gaussian Mixture for our distribution for this question as it is the only applicable distribution. We use the Adam optimiser with H^2 as our loss function. As our activation functions we use the identity function for μ as we want to map to any value in \mathbb{R} and for σ we use the exponential function as we need a positive real number but we impose a minimum of 10^{-12} to prevent division by zero errors. Similarly to the LSTM Model we use ELU as the activation function for the initial dense layers. We first explore the optimal number of distributions in our mixture by performing Kfold cross validation with $K = 5$ over 500 epochs for 1, 2, 3, 4, 5, 6, 8, 16, 32, 64 Gaussians in our mixture. We expect that we should find that the optimal number is 3 or close to 3.

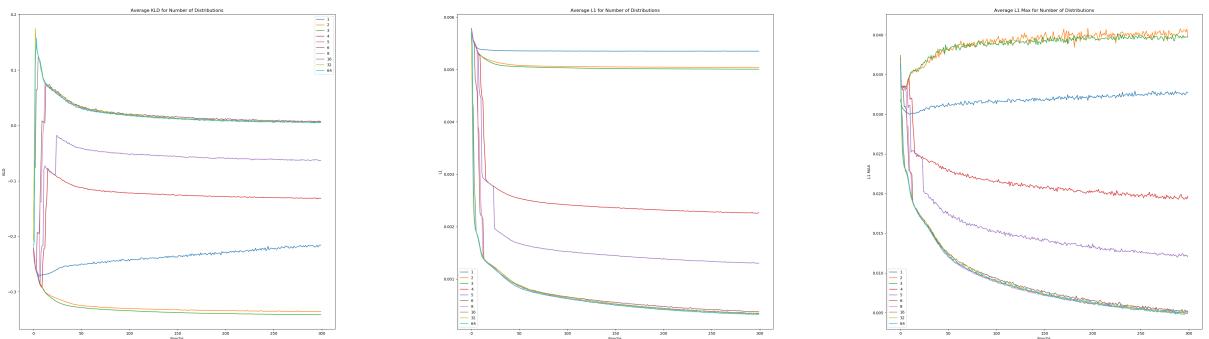


Figure 17: L1, L1 Max and KLD on the validation datasets averaged across 5 folds for different numbers of distributions in the mixture

Plots of the loss for the number of distributions are shown in figure 17 and the minimum loss in figure 41. We see that the optimal value is 6 higher than we might first expect however, this is likely down to the re-scaling of the data as the estimates of the mixture model aren't normalised. This is shown by the negative KLD for the lower numbers of distributions which is negative when the predicted values do not sum to 1. This requirement is not enforced across the bins with this model. We could try to normalise the output of the mixture across the bins however this is likely to induce errors in backpropagation caused by NaN gradients. A potential solution is to use a SoftMax activation for the values given by the mixture but this would change the distribution we predict to no longer be Gaussian. The largest drop occurs at 6 distributions but we see for 16 and 32 the divergence is a

considerable amount lower. The large jumps in minimum loss going from 1 to 2 to 3 to 4 distributions indicate the model is learning to predict the additional modes. From here on we use 6 Gaussians.

We next want to explore the effect of batch size on the models performance we do this simply by training models for different batch sizes. The validation loss for these are shown in figure 18 and their minimas in figure 42. Ideally we would also use kfold cross validation here but for small batch sizes it can take prohibitively long to train the model. We see the model is largely invariant to the size of the batch with all converging but lower losses are observed with smaller batches. Either a batch of 1024 or 512 samples appears to be optimal so we select 1024 from here on.

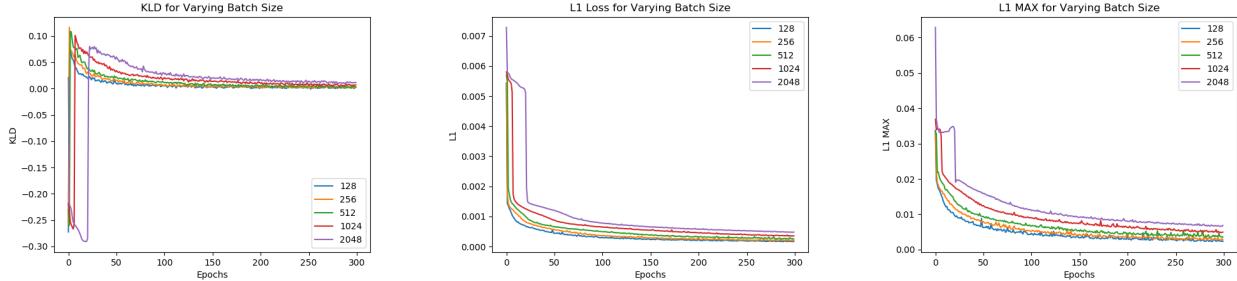


Figure 18: L1, L1 Max and KLD on the validation datasets for different batch sizes with the Mixture Model

We now explore the effects of two optimisers Stochastic Gradient Descent and Adam with different loss functions. We look at the results on the validation dataset and are shown in 19 and their minimas in figure 43. Similar to the LSTM model in that the Adam optimiser significantly outperforms SGD with all loss functions but SGD does not diverge to NaN like it did with the LSTM Model. The Maximum Exponential Distance performs the best with all loss functions. Although looking at how it reaches this we see a massive jump at 200 epochs would suggest perhaps this loss is unstable and might not always work best. It could be best used when we have a pretrained model we want to train further to squeeze out more performance. The squared Hellinger Distance has the closest performance closely followed be L2 and is more consistent so from here on we use it with the Adam optimiser for this model.

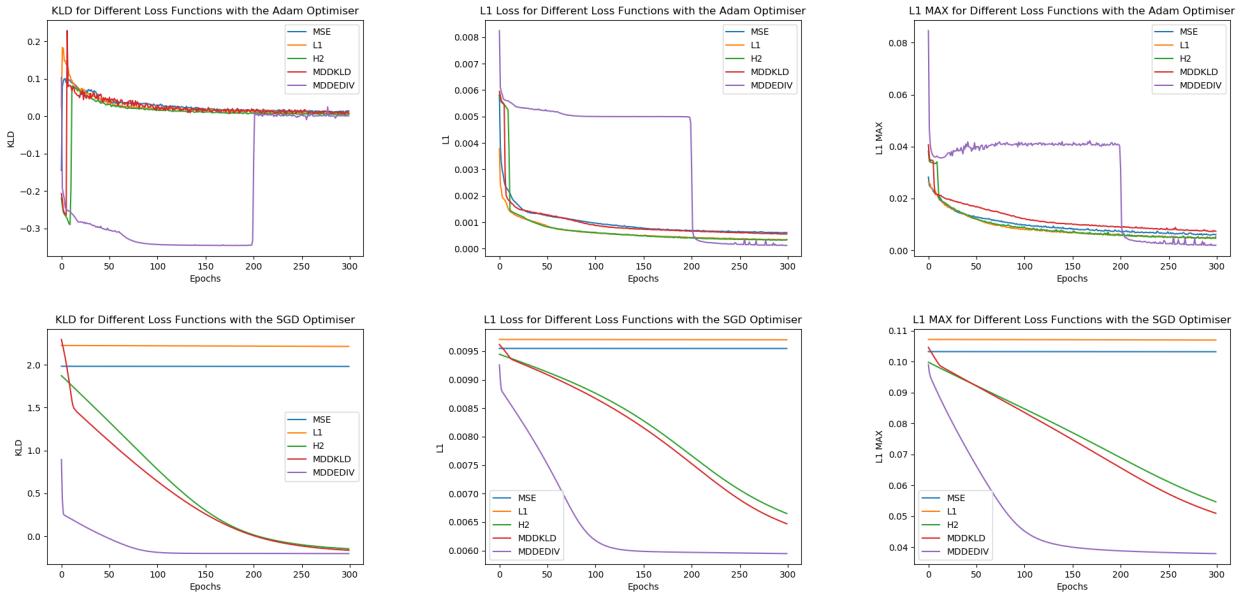


Figure 19: L1, L1 Max and KLD on the validation datasets for different optimisers. Adam was used to train the first row and SGD the second row with each line representing a different loss function as indicated in the legend.

We next want to explore different activation function for the dense layers and the output for σ . The activation function for the weights is fixed as we require it to sum to unity and for μ we use the identity function as we know it can take any value in \mathbb{R} . The results on the validation dataset are in figure 20 the line colour represents dense activation and the line style the activation of the σ output layer. The minimum loss achieved across the validation data set is shown in figure 44. The Clamped function represents $ReLU(x) + \epsilon$ where $\epsilon = 10^{-12}$. We see a general trend of the best performing activation for the σ layer being the exponential function and 2 combinations failed to converge Tanhshrink with Exponential and Softplus. Overall the best performing combination in terms of

the KLD, L1 Loss and L1 Max is the hyperbolic tangent function as the activation for the dense layers and the Exponential function on the layer the outputs σ . This is why we elect to use this combination from here on.

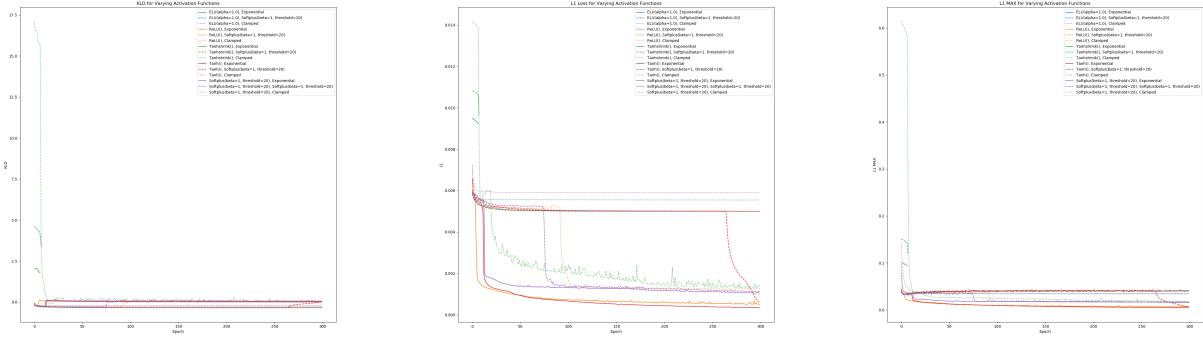


Figure 20: L1 and KLD on the validation datasets for different batch sizes

We now train for a large number of epochs with first the Squared Hellinger Distance as loss function then with the Maximum Exponential Divergence to explore the performance after the jump compared to the next best model. The performance on the validation data set for H^2 is in figure 45 and its minimum values in figure 45 and for the Maximum Exponential Divergence in figure 22 and its minimum values in figure 22. We see the problem we observed before with the Maximum Exponential Divergence with a large jump in loss now this time at 2000 epochs showing despite offering a low loss it is inconsistent perhaps due to the steepness of the loss function causing bouncing around the optimal value. The model trained with the Squared Hellinger Distance achieves very low loss but the Maximum Exponential Divergence still outperforms it with half the L1 and L1 Max. We do choose to continue using H^2 as the large jumps seem to happen randomly and this is not an ideal characteristic.

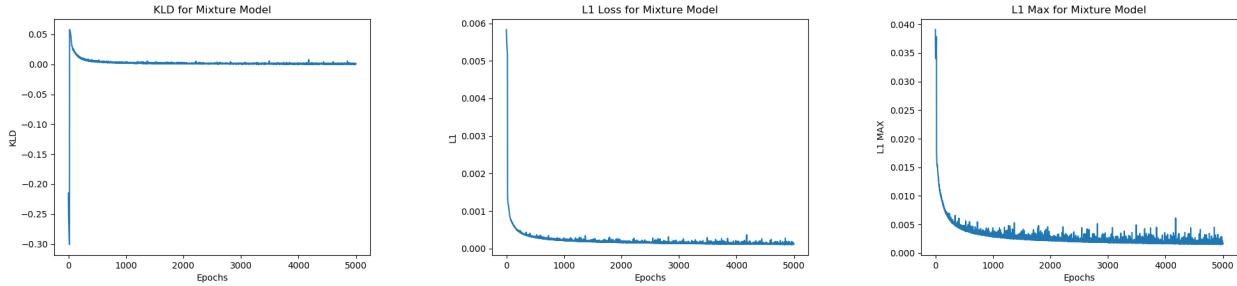


Figure 21: L1, L1 Max and KLD on the validation dataset trained with the Squared Hellinger Distance for 5000 epochs

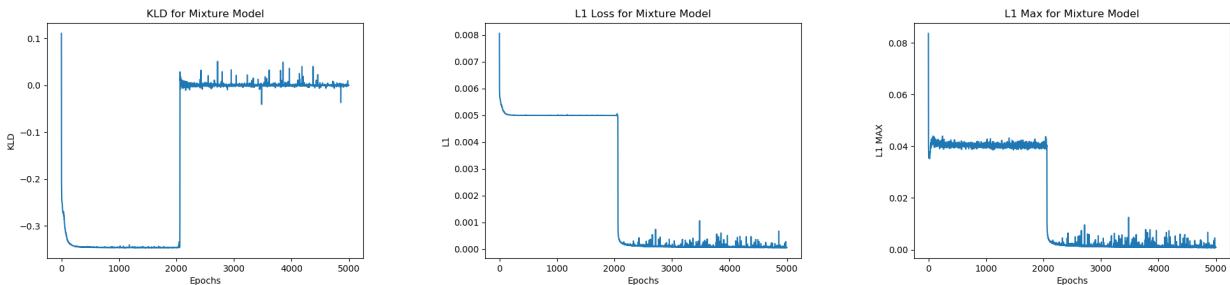


Figure 22: L1, L1 Max and KLD on the validation dataset trained with the Maximum Exponential Divergence for 5000 epochs

Finally we utilise an ensemble method we train 5 Mixture Models on a random sample of the training dataset for 1000 epochs. We then use the predictions of the models that achieved the lowest validation loss on the data to set to produce a Gaussian Mixture from the Mixtures predicted by the 5 Models. The performance on the common validation dataset is shown in figure 23 and the minimum they achieved in figure 47. We achieve slightly worse performance to the model we trained for a large number of epochs but better than models trained in previous tests. We suspect this is due to the submodels requiring longer to train but the difference in loss between the model and the performance of other previous models is small suggesting not much is gained from using an ensemble as the predictions are already very accurate.

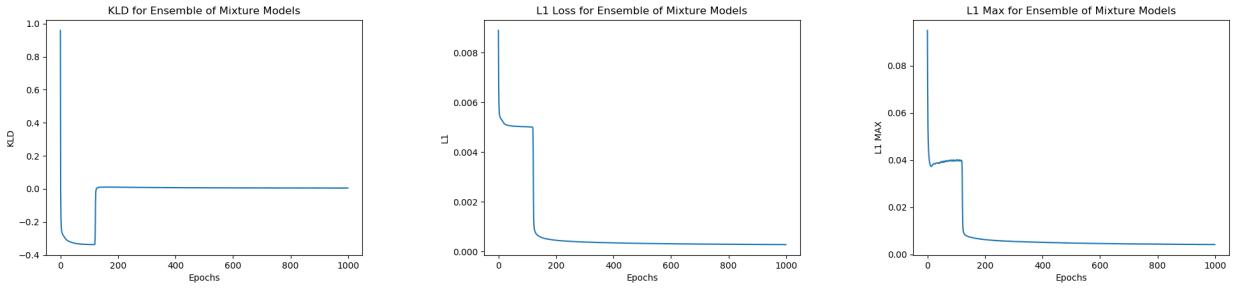


Figure 23: L1, L1 Max and KLD on the validation datasets for Ensemble of 5 Gaussian Mixture Models which were each trained for 1000 epochs on a random sample of the training dataset

VII.C Comparison

The Mixture Model achieved significantly lower loss than that of the LSTM array and due to it not having the LSTM arrays, its complexity is not scaled by the number of bins T and as it is able to use a larger batch size it trains much faster than the LSTM Model. The mixture with this data has very little it needs to learn from its inputs as it is predicting a Gaussian Mixture from the parameters of a Gaussian mixture so its performance on other real world problem could be worse and very dependent on the matching of an appropriate distribution. The LSTM Model does not have this problem and does not need this adjustment but for large numbers of bins the training quickly becomes very long and requires a significant amount of memory to be able to perform back propagation on the 1 or 2 LSTM arrays.

VIII Experiments: Protein Distribution

As previously mentioned a real world example is the production of proteins from a self-regulating gene which is shown in figure 24. We fix the average number of proteins to allow us to capture the distribution with a fixed number of bins of a fixed size with the average number of proteins given by

$$P = \frac{v_1 v_0 k_0}{d_1 d_0 (k_0 + k_1)}$$

We also fix $d_0 = 1$ as we can scale the other rate parameters to have the same effect as if this was a free parameter. We sample k_0, k_1 from the log uniform distribution $LOG_{10}U(0.01, 1)$, d_1, v_0 from the uniform distribution $U(0.01, 1)$ and finally we choose v_1 such that the average number of proteins is 100. Using Gillespie's Algorithm we generate the timeseries of proteins for given reaction rates then take independent samples to get the distribution of proteins. To get independent samples for the rate parameters r we sample every $\max(\max(1/r), 1)$. The binning procedure is we take a number T of equal size bins on the interval $[0, 300]$ then normalise so they sum to 1. The standard binning functions given in numpy and scipy don't allow for fast vectorised binning so we make some improved binning functions that allow for vectorisation. Even in this simplistic model the binning and generation of the timeseries of proteins is rather long at 45.638 seconds on average per set of parameters. Even using parallelisation across nodes the generation and binning for the 10^5 samples requires a lot of time and access to significant computing resources motivating the need for a faster generation method.

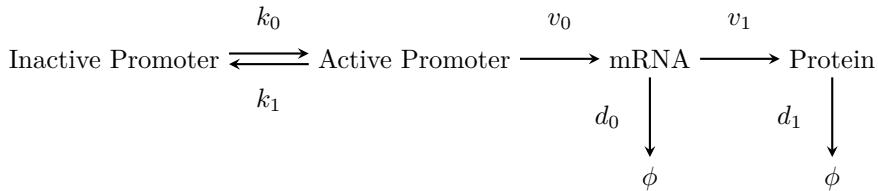


Figure 24: Gene switching model

VIII.A LSTM Model

We first explore the performance of the Bidirectional LSTM model trained with Adam and H^2 on this problem for 1000 epochs. Its performance on the validation dataset is shown in figure 25 and its minimum performance

in figure 48. The KLD here is not negative as it was on the Gaussian data but is high and L1 loss is low and slowly decreasing but the L1 Max is rather large suggesting the model is not predicting significant peaks in the distribution of proteins.

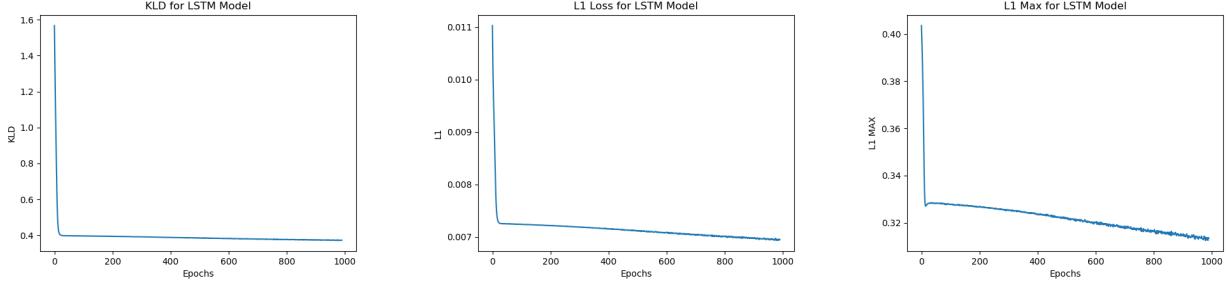


Figure 25: L1,L1 MAX and KLD on the validation dataset for the LSTM model using the distribution of Proteins in a self-regulating gene model

We next see if an ensemble of 5 LSTM models trained for 500 epochs as used before can increase the performance with this data. The performance on the validation dataset is shown in figure 26 and its minimum loss in figure 49. The ensemble has much lower L1 and L1 Max but these values are still much higher than we would like.

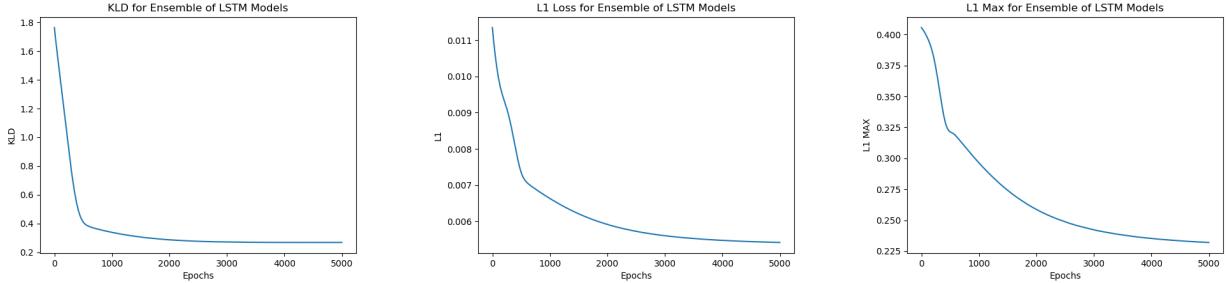


Figure 26: L1, L1 Max and KLD on the validation dataset for an ensemble of 5 LSTM Models

VIII.B Mixture Model

We first explore the possible distributions Gaussian, Gamma-Poisson, Poisson Log-Normal and the number of them first. We used the Adam optimiser and H^2 for activation functions the Gamma-Poisson distribution we used the clamped exponential for both as $\alpha, \beta > 0$ and for the Poisson Log-Normal we used the same activations as with the Gaussian distribution. The Log-Normal distribution we do not display results for as it produced NaNs. This was due to precision errors caused by evaluating the Poisson Distribution at the quadrature points. I tried implementing methods to prevent the NaNs being output such as assigning them as zero but this worked better but still resulted in NaNs. Perhaps a more strict activation function that prevents it evaluating models with very large μ or σ at the risk of hindering the range of distributions we can match. We could use double precision in the evaluation of the CDF but this will result in a large increase in computation time, or casting to a new set of bins that will lead to smaller appropriate parameters. Here this could be done with 100 bins in $[0, 100]$ instead of $[0, 300]$. Lastly perhaps another quadrature method that avoids using Pytorch's Poisson distribution. The results for 5 folds on the validation datasets for a mixture of Gaussians are show in figure 27 with its minimums in figure 27 and for a mixture of Gamma-Poissons in 28 with its minimums in figure 28.

We see that the Gamma-Poisson Model out performs the Gaussian Model on the regulating gene model by a significant margin for the KLD, L1 and L1 Max more so with the KLD and L1 Max. This is what we expected with the support being more similar to the distribution we are trying to emulate. Both mixtures perform significantly better when the number of distributions is greater than 1. The Gaussian Model performed best with 6 distributions with more offering only slightly greater performance. The Gamma-Poisson Model has decreasing loss as we increase the number of distributions however beyond 3/4 distributions the reduction in loss is small. We therefore select for this problem the Gamma-Poisson distribution with 4 distributions in the mixture.

We then train a Mixture Model of 4 Gamma-Poisson distribution for 3000 epochs with Adam and H^2 to see its performance after a large number of epochs. The performance of this model on validation data is shown in figure

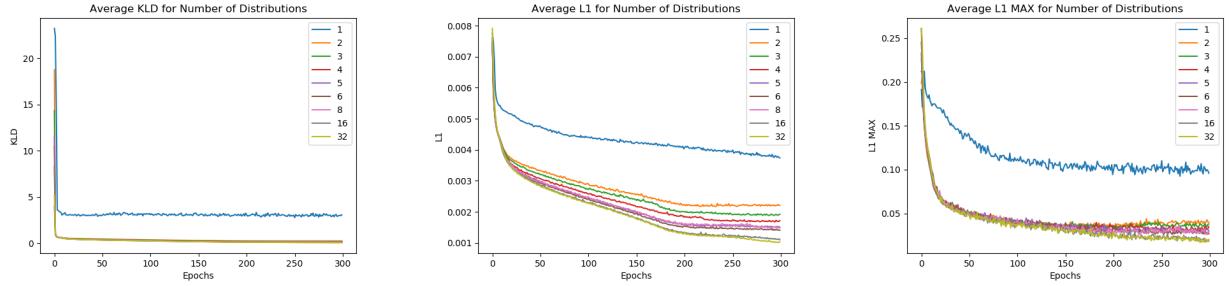


Figure 27: L1,L1 MAX and KLD on the validation datasets averaged across 5 folds for 300 epochs on different numbers of Gaussian distributions on the distribution of Proteins in a self-regulating gene model

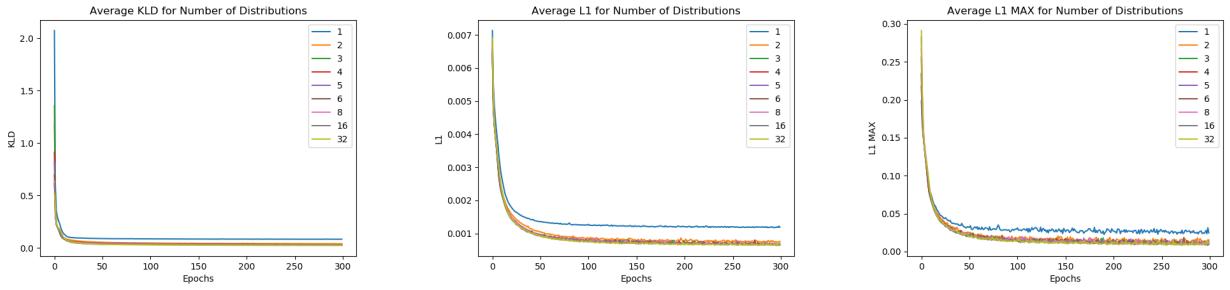


Figure 28: L1,L1 MAX and KLD on the validation datasets averaged across 5 folds for 300 epochs on different numbers of Gamma-Poisson distributions on the distribution of Proteins in a self-regulating gene model

29 and its minimum loss in figure 52. We see its performance after relatively few epochs is rather good but not as good as it was on the Gaussian data as we would expect and we see a reduction in loss coming from the results on the 5 folds.

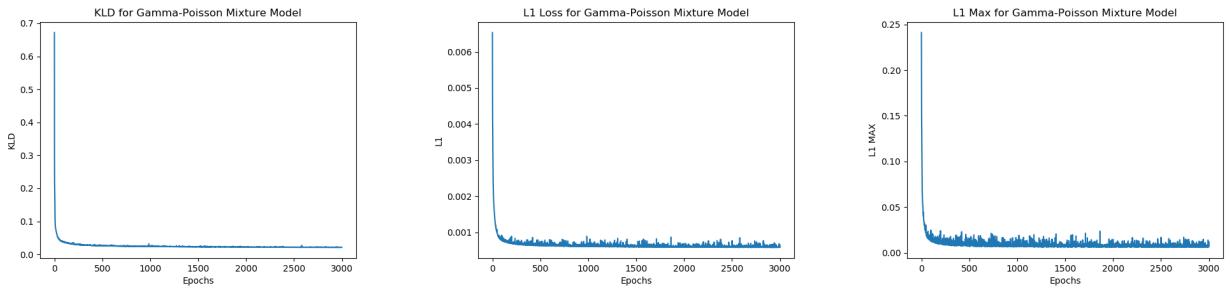


Figure 29: L1,L1 MAX and KLD on the validation dataset for a Mixture Model of 4 Gamma-Poisson distributions on the distribution of Proteins in a self-regulating gene model

As before we create an ensemble of 5 of the previous models trained for 1000 epochs the results on the common validation dataset the results of which are shown in figure 30. We achieve a slight decrease in performance over the individual model.

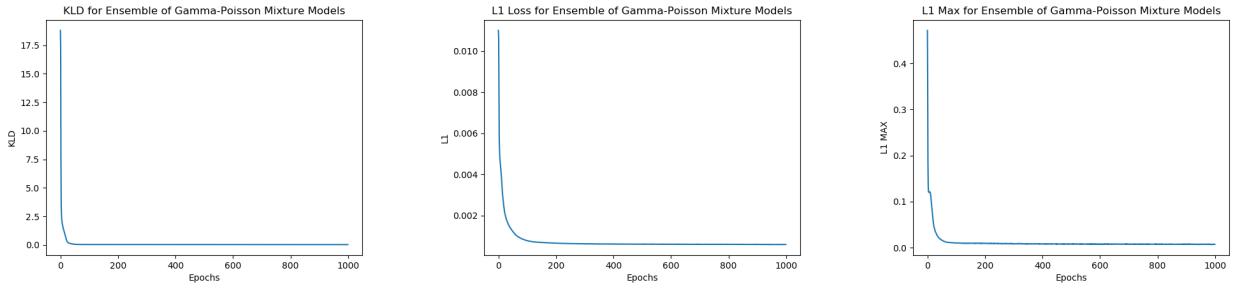


Figure 30: L1,L1 MAX and KLD on the validation datasets for an Ensemble of 5 Mixture Models consisting of 4 Gamma-Poisson distributions on the distribution of Proteins in a self-regulating gene model

VIII.C Comparison

The Mixture Model with the Gamma-Poisson distribution significantly outperforms the LSTM model in all metrics and trains significantly faster with less memory usage. We do however with the mixture lose some generality as we need to find an appropriate distribution but the Gaussian mixture trained across 5 folds still had lower validation loss than the LSTM Model.

IX Experiments: mRNA and Protein Joint Distribution

To show the increased versatility of the Mixture Model we use it to predict the bivariate joint distribution of proteins and mRNA. The creation of the data is almost the same as the univariate case of the protein distribution except we also store the number of mRNA at each sampled point and fix v_0 such that the average number of mRNAs is 50. The formula for the average number of mRNA is shown below.

$$M = \frac{v_0 k_0}{d_0(k_0 + k_1)}$$

We only show this with the Mixture Model as to use the LSTM model we would need a model with a prohibitively large number of LSTM arrays for the bivariate case.

IX.A Mixture Model

We first evaluate the performance across 5 folds as we increase the number of Dirichlet-Multinomial Gamma-Poisson distributions in our mixture. Using hyperbolic tangent activation for the dense layers and the Exponential as the activation for the parameters. The results on the validation dataset averaged over the kfolds are shown in figure 31 and the minimums in figure 54. The training of this model takes much longer than before due to the extra dimension which is why we only checked using 1, 2, 4 distributions. The KLD with 1 distribution is high and the L1 and L1 Max are also much higher than that of 2 or 4 distributions suggesting it fails to capture additional peaks. We see similar performance for 2 and 4 distributions with a slight increase when moving to 4 from 2. We opt to use 4 distributions from here on.

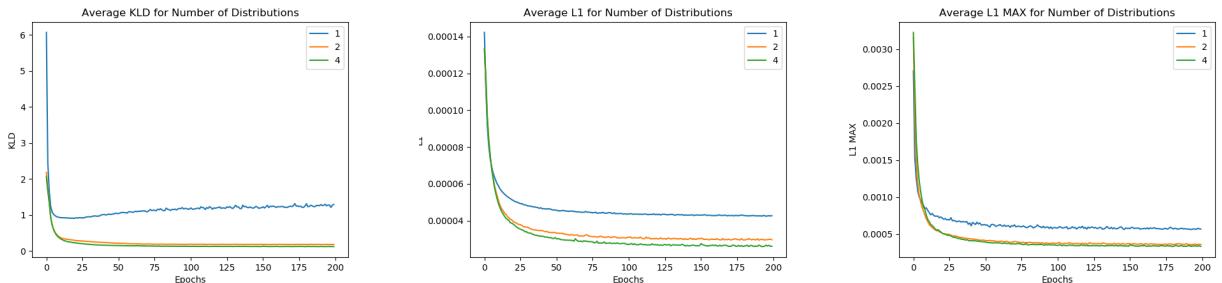


Figure 31: L1,L1 MAX and KLD on the validation dataset averaged across 5 folds for different numbers of Dirichlet-Multinomial Gamma-Poisson distributions

To examine the performance further we train the Mixture Model with 4 distributions for 800 epochs. The performance on the validation dataset is shown in figure 32 and the minimum loss achieved in figure 55. We see that the models performance doesn't improve much after the first 200 epochs but we see the loss is lower than it was averaged over the 5 folds.

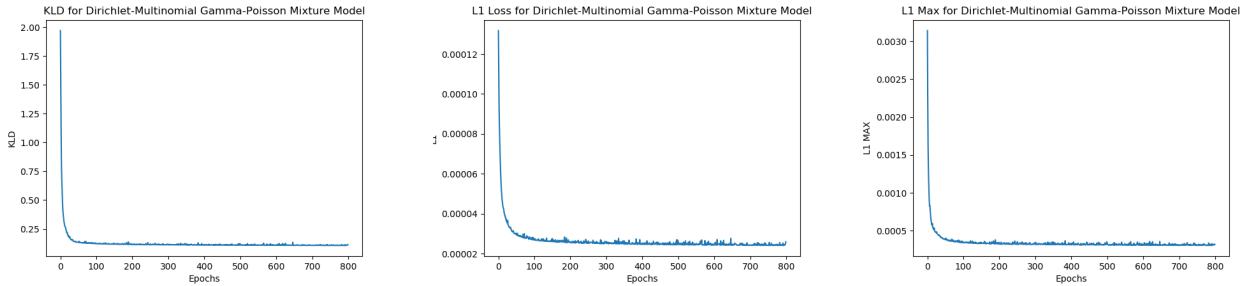


Figure 32: L1,L1 MAX and KLD on the validation dataset for a Mixture Model using 4 Dirichlet-Multinomial Gamma-Poisson distributions trained for 800 epochs

As previously we construct an ensemble using 5 of these models trained for 500 epochs. The loss on the validation dataset are in figure 33 and the minimum loss in figure 56. We see minimum loss very close to that of the individual model and thus we haven't gained a significant increase in emulation accuracy from using an ensemble.

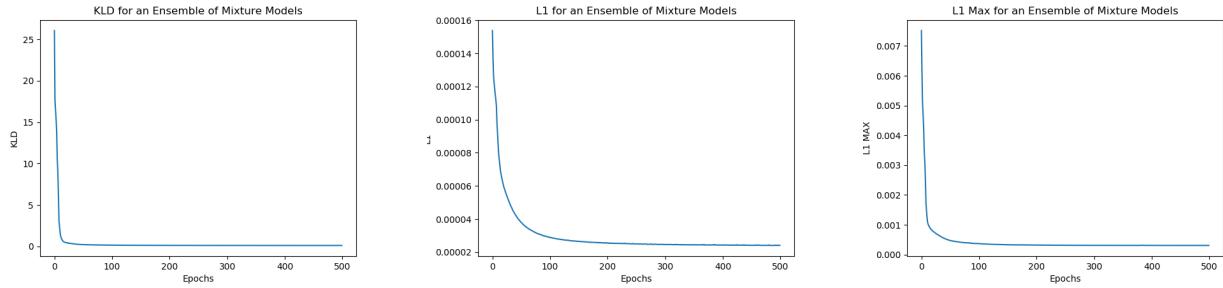


Figure 33: L1,L1 MAX and KLD on the validation dataset for a ensemble of 5 Mixture Models using 4 Dirichlet-Multinomial Gamma-Poisson distributions trained for 500 epochs

X Concluding Remarks

Using neural networks to produce emulating mixture distributions provides a fast and effective method that can be scaled easily to higher dimension problems. An exploration of its performance on smaller datasets and on more complex networks would provide an insight into whether it can replace Approximate Bayesian Computation (ABC) in larger stochastic reaction networks. We suggest in such problems using a deeper neural network and exploring other distributions which could further increase its emulating ability with much more complex problems such as the growth of cancerous tumours. In higher dimensions the Poisson Log-Normal distribution, provided an effective gradient preserving numerical approximation technique is used, could provide emulation of n-variate distributions. Stochastic reaction networks that exhibit a wider range of complex behaviours may increase the performance of ensembles of Mixture Models where they can emulate rare complex behaviours. Their lack of improvement in predictive ability is likely due to the individual models performance being good enough that they very rarely produce incorrect predictions.

The main advantage the LSTM based model has is in its increased generality as we aren't required to select a distribution to fit the problem. Although its performance here was significantly worse than that of the mixture model, both taking longer to train and having a significantly worse emulating ability. Improvements could be found increasing the depth of the network or in the application of clockwork RNN's *Koutník et al.* [20] that can learn more complex relationships in series but increase the number of hyperparameters that need to be tuned.

In the training of our models it was clear the the MSE offered poor performance and functions with larger gradients for values close to zero but the L1 Norm and the Squared Hellinger Distance offer an attractive alternative. The Squared Hellinger Distance also is twice differentiable which can be important when we want to achieve very accurate models. Our Exponential Divergence Distance Measure despite offering the lowest validation loss is too

unstable to train models except maybe in the case of further training of an already well-trained model. Exploring the performance of other f-divergence measures like those stated in *Eguchi* [13] modified to become distance measures could provide other loss functions of interest.

References

- [1] J. AITCHISON and C. H. HO. The multivariate Poisson-log normal distribution. *Biometrika*, 76(4):643–653, 12 1989. ISSN 0006-3444. doi: 10.1093/biomet/76.4.643. URL <https://doi.org/10.1093/biomet/76.4.643>.
- [2] L. Amrhein, K. Harsha, and C. Fuchs. A mechanistic model for the negative binomial distribution of single-cell mRNA counts. *bioRxiv*, 2019. doi: 10.1101/657619. URL <https://www.biorxiv.org/content/early/2019/06/05/657619>.
- [3] G. E. Bates and J. Neyman. Contributions to the theory of accident proneness. 1. an optimistic model of the correlation between light and severe accidents. 1952.
- [4] M. A. Beaumont. Estimation of population growth or decline in genetically monitored populations. *Genetics*, 164(3):1139–1160, 2003. ISSN 0016-6731. URL <https://www.genetics.org/content/164/3/1139>.
- [5] M. A. Beaumont. Approximate bayesian computation in evolution and ecology. *Annual Review of Ecology, Evolution, and Systematics*, 41(1):379–406, 2010. doi: 10.1146/annurev-ecolsys-102209-144621. URL <https://doi.org/10.1146/annurev-ecolsys-102209-144621>.
- [6] J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith, M. W. (eds), J. M. Bernardo, and M. A. Juárez. Intrinsic estimation. In *Bayesian Statistics 7*, pages 456–476. Oxford University Press, 2003.
- [7] C. Bishop. Mixture density networks. Technical Report NCRG/94/004, January 1994. URL <https://www.microsoft.com/en-us/research/publication/mixture-density-networks/>.
- [8] L. Breiman. Bagging predictors. Technical report, Department of Statistics, University of California Berkeley, September 1994. URL <https://www.stat.berkeley.edu/~breiman/bagging.pdf>.
- [9] M. G. Bulmer. On fitting the poisson lognormal distribution to species-abundance data. *Biometrics*, 30(1): 101–110, 1974. ISSN 0006341X, 15410420. URL <http://www.jstor.org/stable/2529621>.
- [10] D. Clevert, T. Unterthiner, and S. Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). In Y. Bengio and Y. LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL <https://arxiv.org/abs/1511.07289>.
- [11] K. Cranmer, J. Brehmer, and G. Louppe. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 2020. ISSN 0027-8424. doi: 10.1073/pnas.1912789117. URL <https://www.pnas.org/content/early/2020/05/28/1912789117>.
- [12] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12(null):2121–2159, July 2011. ISSN 1532-4435.
- [13] S. Eguchi. A differential geometric approach to statistical inference on the basis of contrast functionals. *Hiroshima Mathematical Journal*, 15(2):341–391, 1985. doi: 10.32917/hmj/1206130775. URL <http://projecteuclid.org/euclid.hmj/1206130775>.
- [14] F. A. Gers and J. Schmidhuber. Recurrent nets that time and count. In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, volume 3, pages 189–194 vol.3, 2000.
- [15] F. A. Gers, J. Schmidhuber, and F. Cummins. Learning to forget: Continual prediction with lstm. *Neural computation*, 12(10):2451–2471, 2000. doi: 10.1162/089976600300015015. URL <https://doi.org/10.1162/089976600300015015>.
- [16] M. A. Gibson and J. Bruck. Efficient exact stochastic simulation of chemical systems with many species and many channels. *The Journal of Physical Chemistry A*, 104(9):1876–1889, Mar 9, 2000. doi: 10.1021/jp993732q. URL <http://dx.doi.org/10.1021/jp993732q>.

- [17] U. Gonzales-Barron and F. Butler. A comparison between the discrete poisson-gamma and poisson-lognormal distributions to characterise microbial counts in foods. *Food Control*, 22(8):1279 – 1286, 2011. ISSN 0956-7135. doi: <https://doi.org/10.1016/j.foodcont.2011.01.029>. URL <http://www.sciencedirect.com/science/article/pii/S0956713511000491>.
- [18] J. Guo, S. Lu, H. Cai, W. Zhang, Y. Yu, and J. Wang. Long text generation via adversarial training with leaked information, 2018. URL <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16360/16061>.
- [19] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Y. Bengio and Y. LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- [20] J. Koutník, K. Greff, F. J. Gomez, and J. Schmidhuber. A clockwork RNN. *CoRR*, abs/1402.3511, 2014. URL <http://arxiv.org/abs/1402.3511>.
- [21] Y. LeCun, L. Bottou, G. B. Orr, and K. R. Müller. *Efficient BackProp*, pages 9–50. Springer Berlin Heidelberg, Berlin, Heidelberg, 1998. ISBN 978-3-540-49430-0. doi: 10.1007/3-540-49430-8_2. URL https://doi.org/10.1007/3-540-49430-8_2.
- [22] Q. Li, A. Cassese, M. Guindani, and M. Vannucci. Bayesian negative binomial mixture regression models for the analysis of sequence count and methylation data. *Biometrics*, 75(1):183–192, 2019. doi: 10.1111/biom.12962. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/biom.12962>.
- [23] F. Murtagh. Multilayer perceptrons for classification and regression. *Neurocomputing*, 2(5):183 – 197, 1991. ISSN 0925-2312. doi: [https://doi.org/10.1016/0925-2312\(91\)90023-5](https://doi.org/10.1016/0925-2312(91)90023-5). URL <http://www.sciencedirect.com/science/article/pii/0925231291900235>.
- [24] J. F. Nelson. Multivariate gamma-poisson models. *Journal of the American Statistical Association*, 80(392): 828–834, Dec 1, 1985. doi: 10.1080/01621459.1985.10478190. URL <http://www.tandfonline.com/doi/abs/10.1080/01621459.1985.10478190>.
- [25] H. D. Nguyen and G. McLachlan. On approximations via convolution-defined mixture models. *Communications in Statistics - Theory and Methods*, 48(16):3945–3955, Aug 18, 2019. doi: 10.1080/03610926.2018.1487069. URL <http://www.tandfonline.com/doi/abs/10.1080/03610926.2018.1487069>.
- [26] T. T. Nguyen, H. D. Nguyen, F. Chamroukhi, and G. J. McLachlan. Approximation by finite mixtures of continuous density functions that vanish at infinity. *Cogent Mathematics Statistics*, 7(1), Jan 1, 2020. doi: 10.1080/25742558.2020.1750861. URL <http://www.tandfonline.com/doi/abs/10.1080/25742558.2020.1750861>.
- [27] S. J. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. *CoRR*, abs/1904.09237, 2019. URL <http://arxiv.org/abs/1904.09237>.
- [28] C. P. Robert, J.-M. Cornuet, J.-M. Marin, and N. S. Pillai. Lack of confidence in approximate bayesian computation model choice. *Proceedings of the National Academy of Sciences*, 108(37):15112–15117, 2011. ISSN 0027-8424. doi: 10.1073/pnas.1102900108. URL <https://www.pnas.org/content/108/37/15112>.
- [29] M. Schuster and K. K. Paliwal. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11):2673–2681, 1997.
- [30] V. Shahrezaei and P. S. Swain. The stochastic nature of biochemical networks. *Current Opinion in Biotechnology*, 19(4):369–374, 2008. doi: 10.1016/j.copbio.2008.06.011. URL <https://www.clinicalkey.es/playcontent/1-s2.0-S0958166908000815>.
- [31] S. Siami-Namini and A. S. Namin. Forecasting economics and financial time series: ARIMA vs. LSTM. *CoRR*, abs/1803.06386, 2018. URL <http://arxiv.org/abs/1803.06386>.
- [32] S. A. Sisson, G. W. Peters, M. Briers, and Y. Fan. A note on target distribution ambiguity of likelihood-free samplers, 2010. URL <http://arxiv.org/abs/1005.5201>.
- [33] T. Székely and K. Burrage. Stochastic simulation in systems biology. *Computational and Structural Biotechnology Journal*, 12(20):14 – 25, 2014. ISSN 2001-0370. doi: <https://doi.org/10.1016/j.csbj.2014.10.003>. URL <http://www.sciencedirect.com/science/article/pii/S2001037014000403>.
- [34] R. T. Trevor Hastie and J. Friedman. *The Elements of Statistical Learning*. Springer, New York, second edition, 2009. ISBN 9780387848570. doi: 10.1007/978-0-387-84858-7.

- [35] O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, M. Yeo, A. Makhzani, H. Küttler, J. Agapiou, J. Schrittwieser, J. Quan, S. Gaffney, S. Petersen, K. Simonyan, T. Schaul, H. van Hasselt, D. Silver, T. P. Lillicrap, K. Calderone, P. Keet, A. Brunasso, D. Lawrence, A. Ekermo, J. Repp, and R. Tsing. Starcraft II: A new challenge for reinforcement learning. *CoRR*, abs/1708.04782, 2017. URL <http://arxiv.org/abs/1708.04782>.
- [36] S. Wang, K. Fan, N. Luo, Y. Cao, F. Wu, C. Zhang, K. A. Heller, and L. You. Massive computational acceleration by using neural networks to emulate mechanism-based biological models. *Nature Communications*, 10(1):1–9, Dec 2019. doi: 10.1038/s41467-019-12342-y. URL <https://search.proquest.com/docview/2297119101>.
- [37] R. D. Wilkinson. Accelerating abc methods using gaussian processes, 2014. URL <http://arxiv.org/abs/1401.1436>.

XI Data Tables

XI.A Gaussian Mixture

XI.A.1 LSTM Model

Batch Size	H2	MSE	DIV	L1 MAX	L1
128	0.227146	0.000109	-0.308390	0.033424	0.005180
256	0.255286	0.000091	-0.288083	0.031748	0.005326
512	0.264446	0.000080	-0.281693	0.030103	0.005364
1024	0.278667	0.000082	-0.271466	0.031440	0.005432
2048	0.293843	0.000083	-0.260971	0.030916	0.005507

Figure 34: Minimum Validation Loss for different batch sizes

Bidirectional	H2	MSE	DIV	L1 MAX	L1
False	0.324661	0.000085	-0.239323	0.031700	0.005676
True	0.367802	0.000087	-0.210061	0.033016	0.005895

Figure 35: Minimum Validation Loss for bidirectionality

Optimiser	Loss	H2	MSE	DIV	L1 MAX	L1
Adam	H2	0.351781	0.000082	-0.221339	0.031448	0.005801
Adam	L1	0.379517	0.000078	-0.207952	0.029633	0.005806
Adam	MDDEDIV	0.392923	0.000098	-0.190069	0.033060	0.006040
Adam	MDDKLD	0.371288	0.000092	-0.198279	0.032741	0.005962
Adam	MSE	0.612590	0.000074	-0.089754	0.026787	0.006337
SGD	H2	NaN	NaN	NaN	NaN	NaN
SGD	L1	NaN	NaN	NaN	NaN	NaN
SGD	MDDEDIV	NaN	NaN	NaN	NaN	NaN
SGD	MDDKLD	NaN	NaN	NaN	NaN	NaN
SGD	MSE	NaN	NaN	NaN	NaN	NaN

Figure 36: Minimum Validation Loss for different loss functions

Dense Activation	Output Activation	H2	MSE	DIV	L1 MAX	L1
ELU(alpha=1.0)	Sigmoid()	0.356720	0.000080	-0.217725	0.030024	0.005828
ELU(alpha=1.0)	Softplus(beta=1, threshold=20)	0.328169	0.000090	-0.236009	0.035070	0.005701
ReLU()	Sigmoid()	0.368639	0.000089	-0.209631	0.033706	0.005890
ReLU()	Softplus(beta=1, threshold=20)	0.369432	0.000085	-0.208994	0.031530	0.005902
Softplus(beta=1, threshold=20)	Sigmoid()	0.367826	0.000080	-0.210125	0.029928	0.005890
Softplus(beta=1, threshold=20)	Softplus(beta=1, threshold=20)	0.367821	0.000100	-0.210016	0.034692	0.005889
Tanh()	Sigmoid()	0.355825	0.000085	-0.218434	0.031420	0.005826
Tanh()	Softplus(beta=1, threshold=20)	0.333291	0.000090	-0.232506	0.033869	0.005730
Tanhshrink()	Sigmoid()	0.366828	0.000084	-0.211328	0.033465	0.005871
Tanhshrink()	Softplus(beta=1, threshold=20)	0.370049	0.000085	-0.209338	0.032591	0.005882

Figure 37: Minimum Validation Loss for different activation functions

LSTM Type	Encoder Size	H2	MSE	DIV	L1 MAX	L1
Standard	1024	NaN	NaN	NaN	NaN	NaN
Standard	256	0.366975	0.000083	-0.210589	0.030676	0.005889
Standard	512	0.369820	0.000083	-0.208983	0.032110	0.005899
Peephole	1024	0.342686	0.000094	-0.226813	0.036187	0.005760
Peephole	256	0.355613	0.000085	-0.218492	0.032919	0.005825
Peephole	512	0.346774	0.000091	-0.224123	0.034359	0.005781

Figure 38: Minimum Validation Loss for different LSTM types and sizes

H2	MSE	DIV	L1 MAX	L1
0.36021	0.000081	-0.215393	0.030346	0.005847

Figure 39: Minimum Validation Loss for the best performing LSTM model

H2	MSE	DIV	L1 MAX	L1
0.365106	0.000076	-0.211963	0.0275	0.005877

Figure 40: Minimum Validation Loss for an ensemble of 5 LSTM models

XI.A.2 Mixture Model

Number of Distributions	H2	MSE	DIV	L1 MAX	L1
1	0.242886	0.000108	-0.271140	0.029988	0.005349
2	0.185605	0.000112	-0.336140	0.033075	0.005039
3	0.178443	0.000112	-0.341469	0.033409	0.005007
4	0.074123	0.000048	-0.292148	0.019272	0.002257
5	0.039166	0.000025	-0.285674	0.011980	0.001300
6	0.004963	0.000002	-0.260716	0.005065	0.000372
8	0.004098	0.000002	-0.274334	0.004872	0.000341
16	0.004160	0.000002	-0.207477	0.005007	0.000343
32	0.003968	0.000002	-0.209130	0.004879	0.000332
64	0.003714	0.000002	-0.217570	0.004746	0.000318

Figure 41: Minimum Validation Loss for different numbers of Gaussians

Batch Size	H2	MSE	DIV	L1 MAX	L1
128	0.001002	5.431698e-07	-0.272449	0.002245	0.000163
256	0.001445	7.430181e-07	-0.257867	0.002724	0.000194
512	0.002397	1.255268e-06	-0.259909	0.003521	0.000254
1024	0.004255	2.106039e-06	-0.267012	0.004788	0.000351
2048	0.007447	3.747046e-06	-0.290907	0.006572	0.000475

Figure 42: Minimum Validation Loss for different batch sizes

Optimiser	Loss	H2	MSE	DIV	L1 MAX	L1
Adam	H2	0.003647	1.931740e-06	-0.289956	0.004535	0.000318
Adam	L1	0.005569	2.449547e-06	0.004191	0.004605	0.000330
Adam	MDDEDIV	0.000454	3.273433e-07	-0.345629	0.001878	0.000121
Adam	MDDKLD	0.010005	4.498752e-06	-0.264407	0.007208	0.000550
Adam	MSE	0.014379	3.517463e-06	-0.144875	0.005877	0.000600
SGD	H2	0.437010	2.242867e-04	-0.145695	0.054646	0.006649
SGD	L1	0.965049	6.209610e-04	2.212706	0.106964	0.009700
SGD	MDDEDIV	0.355855	1.507516e-04	-0.202062	0.037870	0.005945
SGD	MDDKLD	0.413791	2.062625e-04	-0.162953	0.050932	0.006470
SGD	MSE	0.931416	5.906887e-04	1.979915	0.103181	0.009548

Figure 43: Minimum Validation Loss for different loss functions

Dense Activation	Parameter Activation	H2	MSE	DIV	L1 MAX	L1
ELU(alpha=1.0)	Clamped	0.292729	0.000117	-0.258812	0.034004	0.005541
ELU(alpha=1.0)	Exponential	0.175812	0.000112	-0.343054	0.033598	0.004997
ELU(alpha=1.0)	Softplus(beta=1, threshold=20)	0.176226	0.000110	-0.342970	0.032654	0.005000
ReLU()	Clamped	0.038657	0.000020	-0.311674	0.016242	0.001219
ReLU()	Exponential	0.008960	0.000004	-0.235422	0.006167	0.000506
ReLU()	Softplus(beta=1, threshold=20)	0.178130	0.000113	-0.341622	0.033533	0.005000
Softplus(beta=1, threshold=20)	Clamped	0.368433	0.000121	-0.209523	0.034099	0.005894
Softplus(beta=1, threshold=20)	Exponential	0.030627	0.000019	-0.262770	0.015570	0.001045
Softplus(beta=1, threshold=20)	Softplus(beta=1, threshold=20)	0.033786	0.000020	-0.300614	0.016318	0.001114
Tanh()	Clamped	0.178223	0.000111	-0.341584	0.035044	0.005002
Tanh()	Exponential	0.004261	0.000002	-0.280361	0.004700	0.000343
Tanh()	Softplus(beta=1, threshold=20)	0.014528	0.000004	-0.341571	0.006930	0.000579
Tanhshrink()	Clamped	0.049479	0.000023	-0.110702	0.017988	0.001294
Tanhshrink()	Exponential	NaN	NaN	NaN	NaN	NaN
Tanhshrink()	Softplus(beta=1, threshold=20)	NaN	NaN	NaN	NaN	NaN

Figure 44: Minimum Validation Loss for different activation functions

H2	MSE	DIV	L1 MAX	L1
0.000538	2.961174e-07	-0.300871	0.001511	0.000108

Figure 45: Minimum Validation Loss for Mixture Model trained with the Squared Hellinger Distance

H2	MSE	DIV	L1 MAX	L1
0.000095	6.200201e-08	-0.346132	0.000766	0.000054

Figure 46: Minimum Validation Loss for Mixture Model trained with the Maximum Exponential Divergence Distance

H2	MSE	DIV	L1 MAX	L1
0.003457	0.000002	-0.337599	0.004106	0.000273

Figure 47: Minimum Validation Loss for ensemble of 5 Mixture Models

XI.B Protein Distribution

XI.B.1 LSTM Model

H2	MSE	DIV	L1 MAX	L1
0.469361	0.001254	0.370608	0.312515	0.00693

Figure 48: Minimum Validation Loss for LSTM Model

H2	MSE	DIV	L1 MAX	L1
0.314929	0.000853	0.268331	0.231965	0.005414

Figure 49: Minimum Validation Loss for an ensemble of 5 LSTM Models

XI.B.2 Mixture Model

Number of Distributions	H2	MSE	DIV	L1 MAX	L1
1	0.327024	0.000267	2.763582	0.092562	0.003736
2	0.077675	0.000036	0.247811	0.033208	0.002169
3	0.052161	0.000029	0.191391	0.032986	0.001885
4	0.036786	0.000021	0.154179	0.030207	0.001673
5	0.026860	0.000016	0.124969	0.026716	0.001505
6	0.023504	0.000015	0.111989	0.026096	0.001407
8	0.025782	0.000016	0.123313	0.025397	0.001475
16	0.014816	0.000008	0.078032	0.017465	0.001108
32	0.012609	0.000007	0.070254	0.016728	0.001007

Figure 50: Minimum Validation Loss for different numbers of Gaussians

Number of Distributions	H2	MSE	DIV	L1 MAX	L1
1	0.018553	0.000011	0.081763	0.022607	0.001179
2	0.008533	0.000004	0.038923	0.011393	0.000739
3	0.007913	0.000003	0.030430	0.009762	0.000678
4	0.007860	0.000003	0.029604	0.009391	0.000669
5	0.007865	0.000003	0.029268	0.009045	0.000666
6	0.007748	0.000003	0.027307	0.008933	0.000654
8	0.007708	0.000003	0.026384	0.009249	0.000652
16	0.007555	0.000003	0.024800	0.008063	0.000640
32	0.007338	0.000003	0.022266	0.008377	0.000641

Figure 51: Minimum Validation Loss for different numbers of Gamma-Poissons

H2	MSE	DIV	L1 MAX	L1
0.007295	0.000002	0.020879	0.005897	0.000579

Figure 52: Minimum Validation Loss for Gamma-Poisson Mixture Model

H2	MSE	DIV	L1 MAX	L1
0.007323	0.000002	0.023029	0.006814	0.000591

Figure 53: Minimum Validation Loss for ensemble of 5 Gamma-Poisson Mixture Models

XI.C mRNA and Protein Joint Distribution

Number of Distributions	H2	MSE	DIV	L1 MAX	L1
1	0.316246	3.691391e-07	0.901115	0.000561	0.000042
2	0.139137	5.578090e-08	0.177821	0.000355	0.000029
4	0.100825	4.854157e-08	0.118900	0.000331	0.000026

Figure 54: Minimum Validation Loss for different numbers of Dirichlet-Multinomial Gamma-Poisson distributions

H2	MSE	DIV	L1 MAX	L1
0.089727	3.815524e-08	0.106211	0.000306	0.000024

Figure 55: Minimum Validation Loss for Mixture Model of 4 Dirichlet-Multinomial Gamma-Poisson distributions

H2	MSE	DIV	L1 MAX	L1
0.091038	3.841848e-08	0.107542	0.000307	0.000024

Figure 56: Minimum Validation Loss for a ensemble of 5 Mixture Models of 4 Dirichlet-Multinomial Gamma-Poisson distributions

XII Source Code

Source code for the simulation and all the models is provided below.

XII.A Regulating Gene Model Simulation

XII.A.1 Gillespie simulation joint Protein distribution

```

1 import gillespy2 as gp
2 from gillespy2.solvers.cpp import SSACSolver
3 import numpy as np
4 import os
5 from sklearn.model_selection import ParameterGrid
6 import mpi4py.MPI as MPI
7 import h5py
8 import time
9
10
11 class geneModel(gp.Model):
12     """
13     Gene Model gillespie simulator
14     """
15     def __init__(self):
16         gp.Model.__init__(self,name='geneModel')# inherit
17
18         self.paramGrid = None
19         self.n=0
20
21     def setValues(self,n):
22         # set model parameters
23         parameters = self.paramGrid[n]
24         self.delete_all_parameters()
25         self.delete_all_reactions()
26         self.delete_all_species()
27
28         # set reaction rates
29         k0 = gp.Parameter(name='k0', expression=parameters[0])
30         k1 = gp.Parameter(name='k1', expression=parameters[1])
31         v0 = gp.Parameter(name='v0', expression=parameters[2])
32         v1 = gp.Parameter(name='v1', expression=parameters[-1])#v1Value
33         d0 = gp.Parameter(name='d0', expression=1.0)#1.0

```

```

34     d1 = gp.Parameter(name='d1', expression=parameters[3])
35     self.add_parameter([k0,k1,v0,v1,d0,d1])
36
37     # define species
38     ip = gp.Species(name='ip', initial_value=0)
39     ap = gp.Species(name='ap', initial_value=1)
40     m = gp.Species(name='m', initial_value=int(parameters[-2]))
41     p = gp.Species(name='p', initial_value=100)
42     self.add_species([ip,ap,m,p])
43
44     # define reactions
45     r_k0 = gp.Reaction(name="on", rate=k0, reactants={ip:1}, products={ap:1})
46     r_k1 = gp.Reaction(name="off", rate=k1, reactants={ap:1}, products={ip:1})
47     r_v0 = gp.Reaction(name="m_creation", rate=v0, reactants={ap:1}, products={m:1,ap:1})
48     r_d0 = gp.Reaction(name="m_destruction", rate=d0, reactants={m:1}, products={})
49     r_v1 = gp.Reaction(name="p_creation", rate=v1, reactants={m:1}, products={p:1,m:1})
50     r_d1 = gp.Reaction(name="p_destruction", rate=d1, reactants={p:1}, products={})
51     self.add_reaction([r_k0,r_k1,r_v0,r_d0,r_v1,r_d1])
52
53     # get independent sample dists
54     self.sample = int(np.maximum(np.ceil((1/parameters[np.arange(len(parameters))-4])).max(),1))
55     # get end time
56     self.t = self.sample*10**5
57     self.timespan(np.linspace(0,self.t,10**5+1))# set times to return results for
58
59     def generateParams(self,n):
60         # set parameters
61         self.n = n
62         k01 = np.power(10,np.random.uniform(-2,1,(self.n,2)))
63         div0 = np.random.uniform(0.01,1,(self.n,2))
64         self.paramGrid = np.ones((self.n,6))
65         self.paramGrid[:,0:2] = k01
66         self.paramGrid[:,2:4] = div0
67         self.paramGrid[:, -2] = self.paramGrid[:,2]* self.paramGrid[:,0]/(( self.paramGrid[:,0]+ self.paramGrid[:,1]))
68         self.paramGrid[:, -1]= 100* self.paramGrid[:,3]/self.paramGrid[:, -2]
69
70
71
72     def MPE_DECOMP1D(n,numprocs,myid):
73         """
74             Python function to get start end indicies for mpi proc
75
76         """
77         nlocal = n // numprocs
78         s = myid * nlocal
79         deficit = n % numprocs
80         s = s + min(myid,deficit)
81         if myid < deficit:
82             nlocal = nlocal +1
83         e = s + nlocal
84         if e > n or myid==numprocs-1:
85             e = n
86         return s,e
87
88
89
90     if __name__=="__main__":
91         comm = MPI.COMM_WORLD
92         rank = comm.Get_rank()
93         nprocs = comm.Get_size()
94         model = geneModel()
95         n = 10**5 # number of samples
96         if rank==0:# generate parameters on root
97             model.generateParams(n)
98             temp =model.paramGrid
99         else:
100             temp=None
101
102
103
104         start,end = MPE_DECOMP1D(n,nprocs,rank) # get indicies
105         temp = comm.bcast(temp,root=0) # bcast parameters to all processes
106         if rank != 0:
107             model.paramGrid = temp[start:end] # set paramGrids of model
108         del temp
109
110         sample = np.empty((end-start,10**5),dtype="f") # create array to store results
111         if rank ==0:# start timer
112             tic = time.perf_counter()
113             for i in range(end-start):
114
115                 model.setValues(i)# set rates
116                 sample[i]= model.run(number_of_trajectories=1,solver=SSACSolver())["p"] [1:] # run gillespie and store protein series
117                 # gather indicies on root
118                 ends = np.array(comm.gather(int(end),root=0))
119                 starts = np.array(comm.gather(int(start),root=0))
120                 comm.Barrier()
121                 if rank == 0:
122                     toc = time.perf_counter()# stop timer
123                     temp = np.empty((n,10**5),dtype="f")
124                     temp[start:end] = sample

```

```

125     with h5py.File("modelData.hdf5","w") as f:# create file to store results
126         f.create_dataset("results",data=temp,dtype="f")
127         f.create_dataset("parameters",data=model.paramGrid)
128         f.create_dataset("avg_time",data=np.array([(toc-tic)*nprocs/n]))
129     del temp
130     for i in range(1,nprocs):# send samples from each process individually due to Gather limit on elements
131         if rank == i:
132             comm.Send([sample,MPI.FLOAT],dest=0)
133         elif rank == 0:
134             temp = np.empty((ends[i]-starts[i],10**5),dtype="f")
135             comm.Recv(temp,source=i)
136             with h5py.File("modelData.hdf5","a") as f:# save sample from proc
137                 f["results"][starts[i]:ends[i]] = temp
138
139
140
141
142

```

XII.A.2 Gillespie simulation joint mRNA and Protein distribution

```

1  from GeneModelSimulation import *
2  class biGeneModel(geneModel):
3      """
4          Joint distribution of mrna and proteins
5      """
6      def __init__(self):
7          super().__init__()# inherit
8      def generateParams(self,n):
9          # generate parameters
10         self.n = n
11         k01 = np.power(10,np.random.uniform(-2,1,(self.n,2)))
12         d1 = np.random.uniform(0.01,1,(self.n))
13         self.paramGrid = np.ones((self.n,6))
14         self.paramGrid[:,0:2] = k01
15         self.paramGrid[:,2] = 50*(( self.paramGrid[:,0]+ self.paramGrid[:,1]))/self.paramGrid[:,0]
16         self.paramGrid[:,3] = d1
17         self.paramGrid[:,2] = 50
18         self.paramGrid[:,3] = 100*self.paramGrid[:,3]/self.paramGrid[:,2]
19
20
21     if __name__=="__main__":
22         comm = MPI.COMM_WORLD
23         rank = comm.Get_rank()
24         nprocs = comm.Get_size()
25         model = biGeneModel()
26         n = 10**5 # number of samples
27         if rank==0:
28             model.generateParams(n)# generate params on root
29             temp =model.paramGrid
30         else:
31             temp=None
32
33
34
35     start,end = MPE_DECOMP1D(n,nprocs,rank) # get indicies
36     temp = comm.bcast(temp,root=0) #bcast params
37     if rank != 0:
38         model.paramGrid = temp[start:end] # set param grid
39     del temp
40     # create array to store proteins and mrna
41     sampleP = np.empty((end-start,10**5),dtype="f")
42     sampleM = np.empty((end-start,10**5),dtype="f")
43     if rank ==0:# time
44         tic = time.perf_counter()
45     for i in range(end-start):
46
47         model.setValues(i)# set rates
48         res = model.run(number_of_trajectories=1,solver=SSACSolver())# run gillespie
49         # store samples
50         sampleP[i]=res["p"][:1]
51         sampleM[i]=res["m"][:1]
52     # get indicies on root
53     ends = np.array(comm.gather(int(end),root=0))
54     starts = np.array(comm.gather(int(start),root=0))
55     comm.Barrier()
56     if rank == 0:
57         toc = time.perf_counter()#stop timer
58         temp1 = np.empty((n,10**5),dtype="f")
59         temp1[start:end] = sampleP
60         temp2 = np.empty((n,10**5),dtype="f")
61         temp2[start:end] = sampleM
62     # save data
63     with h5py.File("modelDataBi.hdf5","w") as f:
64         f.create_dataset("resultsP",data=temp1,dtype="f")
65         f.create_dataset("resultsM",data=temp2,dtype="f")
66         f.create_dataset("parameters",data=model.paramGrid)
67         f.create_dataset("avg_time",data=np.array([(toc-tic)*nprocs/n]))

```

```

68     del temp1
69     del temp2
70 for i in range(1,nprocs):# we send from each proc seperately rather than using Gather as Gather has a limit on the number of
    ↵ elements
71     if rank == i:# send proteins from i to root
72         comm.Send([sampleP,MPI.FLOAT],dest=0)
73     elif rank == 0:
74         temp = np.empty((ends[i]-starts[i],10**5),dtype="f")
75         comm.Recv(temp,source=i)
76         with h5py.File("modelDataBi.hdf5","a") as f:# save into file
77             f["resultsP"][starts[i]:ends[i]] = temp
78     if rank == i:# send mrna from i to root
79         comm.Send([sampleM,MPI.FLOAT],dest=0)
80     elif rank == 0:
81         temp = np.empty((ends[i]-starts[i],10**5),dtype="f")
82         comm.Recv(temp,source=i)
83         with h5py.File("modelDataBi.hdf5","a") as f:# save into file
84             f["resultsM"][starts[i]:ends[i]] = temp

```

XII.A.3 Vectorised Binning

```

1 import numpy as np
2 import h5py
3 import time
4 from tqdm import tqdm
5 from sklearn.preprocessing import scale
6 import mpi4py.MPI as MPI
7 from GeneModelSimulation import MPE_DECOMP1D
8 from numba import jit
9
10 @jit(nopython=True)
11 def vbins(x,b,r):
12
13     bins = np.linspace(r[0],r[1],b+1)# get bins
14     dig = np.searchsorted(bins, x, side='right')# get bins
15     res = np.empty((b,x.shape[0]),dtype=np.int32)
16     for i in range(b):# count numbers in each bin
17         res[i] = np.sum(dig==i+1, axis=-1)
18     res = res.T
19     s = res.sum(-1)
20     # stop division by zero and normalise
21     s[s==0] = 1
22     return res/s.reshape(x.shape[0],1)
23
24
25
26 if __name__=="__main__":
27     comm = MPI.COMM_WORLD
28     rank = comm.Get_rank()
29     nprocs = comm.Get_size()
30     t = 100 # number of time points
31     comm.Barrier()
32     if rank == 0:# start timer
33         tic = time.perf_counter()
34     with h5py.File("modelData.hdf5","r") as f:
35         tx = f["results"]
36         s = tx.shape
37         n= s[0]
38         start,end = MPE_DECOMP1D(n,nprocs,rank)# get index
39         it = np.append(np.arange(0,end-start,1536,dtype=int),[int(end-start)])# only load 1536 rows at once due to mem constraints
40         x = np.empty((end-start,t),dtype="f")
41         for i in range(len(it)-1):
42
43             x[it[i]:it[i+1]] = vbins(tx[(it[i]+start):(it[i+1]+start)],t,(0,300))# bin
44
45             if rank ==0:# in root store reaction rates and time
46                 y = f["parameters"][:,np.arange(6)!=4]
47                 times = f["avg_time"][:]
48             # gather send counts
49             sendcounts = np.array(comm.gather((end-start)*t, root=0))
50             disp = np.array(comm.gather((start)*t, root=0))
51             if rank == 0:
52                 r = np.empty((n,t),dtype="f")
53
54             else:
55                 r = None
56             # use gatherv to put results in r
57             comm.Gatherv(x,[r,sendcounts,disp,MPI.FLOAT],root=0)
58             comm.Barrier()
59             if rank ==0:
60
61                 toc =time.perf_counter() # stop timer
62                 # get scaling parameters
63                 scalesM = np.mean(y, axis=0)
64                 scalesS = np.std(y, axis=0)
65
66                 with h5py.File(f"modelDataBin_{t}.hdf5","w") as f: # save binned data
67                     f.create_dataset("results", data=r)

```

```

68     f.create_dataset("scalesM",data=scalesM)
69     f.create_dataset("scalesS",data=scalesS)
70
71     f.create_dataset("parameters",data=scale(y))
72     f.create_dataset("avg_time",data=times+(toc-tic)*nprocs/(s[0]))

```

XII.A.4 Bivariate Vectorised Binning

```

1  from binning import *
2
3  @jit(nopython=True)
4  def vbinsbi(x1,x2,b1,b2,r1,r2):
5      # bivariate version of binner
6      # get bin edges
7      bins1 = np.linspace(r1[0],r1[1],b1+1)
8      bins2 = np.linspace(r2[0],r2[1],b2+1)
9      # get bin of each
10     dig1 = np.searchsorted(bins1, x1, side='right')
11     dig2 = np.searchsorted(bins2, x2, side='right')
12     res = np.empty((x1.shape[0],b1,b2),dtype=np.int32)
13     for i in range(b1):# for all bins for x1
14         for j in range(b2):# for all bins for x2
15             res[:,i,j] = np.sum((dig1==i+1) & (dig2==j+1),axis=-1)
16     s = res.sum(-1).sum(-1)
17     # normalise
18     s[s==0] = 1
19     return res/s.reshape(x1.shape[0],1,1)
20
21
22 if __name__=="__main__":
23     comm = MPI.COMM_WORLD
24     rank = comm.Get_rank()
25     nprocs = comm.Get_size()
26     t = 100 # set number of bins
27     comm.Barrier()
28     if rank ==0:
29         tic = time.perf_counter()
30         with h5py.File("modelDataBi.hdf5","r") as f:
31             xP = f["resultsP"]
32             xM = f["resultsM"]
33             s = xP.shape
34             n= s[0]
35             start,end = MPE_DECOMPID(n,nprocs,rank) # get indicies
36             it = np.append(np.arange(0,end-start,768,dtype=int),[int(end-start)]) # only load 768 rows at a time
37             x = np.empty((end-start,t,t),dtype="f")
38             for i in range(len(it)-1): # bin data
39
40                 x[it[i]:it[i+1]] = vbinsbi(xM[(it[i]+start):(it[i+1]+start)],xP[(it[i]+start):(it[i+1]+start)],t,t,(0,300),(0,300))
41
42             if rank ==0:# in root store reaction rates and time
43                 y = f["parameters"][:,np.arange(6)!=4]
44                 times = f["avg_time"][:]
45             # send send counts and disp to root
46             sendcounts = np.array(comm.gather((end-start)*t*t, root=0))
47             disp = np.array(comm.gather((start)*t*t, root=0))
48             if rank == 0:
49                 r = np.empty((n,t,t),dtype="f")
50
51             else:
52                 r = None
53             # gather results in root
54             comm.Gatherv(x,[r,sendcounts,disp,MPI.FLOAT],root=0)
55             comm.Barrier()
56             if rank ==0:
57                 toc =time.perf_counter() # stop timer
58                 # get scaling parameters
59                 scalesM = np.mean(y,axis=0)
60                 scalesS = np.std(y,axis=0)
61
62             with h5py.File(f"modelDataBiBin_{t}.hdf5","w") as f:# save binned results
63                 f.create_dataset("results",data=r)
64                 f.create_dataset("scalesM",data=scalesM)
65                 f.create_dataset("scalesS",data=scalesS)
66
67                 f.create_dataset("parameters",data=scale(y))
68                 f.create_dataset("avg_time",data=times+(toc-tic)*nprocs/(s[0]))

```

XII.B LSTM Model

XII.B.1 LSTM Cells

```

1  import torch
2  import torch.nn as nn
3

```

```

4
5  class LSTMCell(nn.Module):
6      """
7          Standard LSTM cell
8      """
9      def __init__(self, input_size : int, hidden_size : int):
10         super().__init__()
11         # set initial parameters
12         self.input_size = input_size
13         self.hidden_size = hidden_size
14         self.weight_ih = nn.Parameter(torch.randn(input_size,4 * hidden_size))
15         self.weight_hh = nn.Parameter(torch.randn(hidden_size,4 * hidden_size))
16         self.bias = nn.Parameter(torch.randn(4 * hidden_size))
17     @staticmethod
18     def __str__():
19         return "LSTM"
20     def forward(self,x,hx,cx):
21         # forward method of LSTM
22         # do gate matrix operations
23         gates = (torch.mm(x, self.weight_ih) + self.bias +
24                   torch.mm(hx, self.weight_hh))
25         ingate, forgetgate, cellgate, outgate = gates.chunk(4, 1)
26
27         # apply activation functions
28         ingate = torch.sigmoid(ingate)
29         forgetgate = torch.sigmoid(forgetgate)
30         cellgate = torch.tanh(cellgate)
31         outgate = torch.sigmoid(outgate)
32
33         # update cell and hidden state
34         cy = (forgetgate * cx) + (ingate * cellgate)
35         hy = outgate * torch.tanh(cy)
36
37         return hy, cy
38
39 class PeepholeLSTMCell(nn.Module):
40     """
41         Peephole LSTM
42     """
43     def __init__(self, input_size, hidden_size):
44         super().__init__()
45         # initialise parameters
46         self.input_size = input_size
47         self.hidden_size = hidden_size
48         self.weight_ih = nn.Parameter(torch.randn(input_size,3 * hidden_size))
49         self.weight_hh = nn.Parameter(torch.randn(hidden_size,3 * hidden_size))
50         self.bias = nn.Parameter(torch.randn(3*hidden_size))
51         self.weight_c =nn.Parameter(torch.randn(input_size,hidden_size))
52         self.bias_c = nn.Parameter(torch.randn(hidden_size))
53     @staticmethod
54     def __str__():
55         return "PeepholeLSTM"
56     def forward(self, x,hx,cx):
57         # forward of cell takes hidden state despite it not being used so it can be interchanged with standard LSTM
58         # do matrix operations
59         gates = (torch.mm(x, self.weight_ih) + self.bias +
60                   torch.mm(cx, self.weight_hh) )
61         ingate, forgetgate, outgate = gates.chunk(3, 1)
62
63         # activations
64         ingate = torch.sigmoid(ingate)
65         forgetgate = torch.sigmoid(forgetgate)
66
67         outgate = torch.sigmoid(outgate)
68
69         # get new cell and hidden state
70         cy = (forgetgate * cx) + (ingate * torch.sigmoid(torch.mm(x, self.weight_c)+self.bias_c))
71         hy = outgate*cy
72
73         return hy,cy
74 class GRUCell(nn.Module):
75     """
76         Grated Recurrent Unit Cell
77     """
78     def __init__(self, input_size, hidden_size):
79         super().__init__()
80         self.input_size = input_size
81         self.hidden_size = hidden_size
82         self.weight_ih = nn.Parameter(torch.randn(input_size,2 * hidden_size))
83         self.weight_hh = nn.Parameter(torch.randn(hidden_size,2 * hidden_size))
84         self.bias =nn.Parameter(torch.randn(2*hidden_size))
85         self.weight_ic =nn.Parameter(torch.randn(input_size,hidden_size))
86         self.weight_hc =nn.Parameter(torch.randn(hidden_size,hidden_size))
87         self.bias_c =nn.Parameter(torch.randn(hidden_size))
88     @staticmethod
89     def __str__():
90         return "GRU"
91     def forward(self, x,hx,cx):
92         gates = (torch.mm(x, self.weight_ih) + self.bias +
93                   torch.mm(hx, self.weight_hh) )
94         updategate, resetgate = gates.chunk(2, 1)

```

```

95     updategate = torch.sigmoid(updategate)
96     resetgate = torch.sigmoid(resetgate)
97
98     hy = (updategate * hx) +
99         (1.0 - updategate) * torch.tanh(torch.mm(x, self.weight_ic) + torch.mm(resetgate * hx, self.weight_hc) + self.bias_c)
100
101    return hy, cx

```

XII.B.2 LSTM Model

```

1 import numpy as np
2 import os
3 from tqdm import tqdm
4 import torch.nn as nn
5 import torch
6 from torch.utils.data import DataLoader
7 from torch.utils.tensorboard import SummaryWriter
8 import torch.nn.functional as F
9 from scipy.stats import norm
10 import pickle
11 from sklearn.model_selection import train_test_split
12 from sklearn.model_selection import ParameterGrid
13 from sklearn.preprocessing import scale
14 import matplotlib.pyplot as plt
15 from multiprocessing import Pool
16 from functools import partial
17 import h5py
18 import time
19 #from copy import deepcopy
20 import pickle
21 import gc
22 from torch.distributions.normal import Normal
23 from cells import LSTMCell, PeepholeLSTMCell, GRUCell
24
25
26
27 def gaussianMix(x, points, dt):
28     g1 = norm(loc=x[3], scale=x[0])
29     g2 = norm(loc=x[4], scale=x[1])
30     g3 = norm(loc=x[5], scale=x[2])
31     return np.diff(x[6] * (g1.cdf(points)) + x[7] * (g2.cdf(points)) + x[8] * (g3.cdf(points))) / 2
32
33 def h2_dist(x, y, ep=1e-12):
34     return 2 * ((y.clamp(ep).sqrt() - x.clamp(ep).sqrt()).pow(2).sum(axis=tuple(np.linspace(1, x.dim() - 1, x.dim() - 1, dtype=np.int))).mean())
35
36 def e_div(x, y, ep=1e-12):
37     return (x * (((x.clamp(ep) / (y.clamp(ep))).log()).pow(2)).sum(axis=tuple(np.linspace(1, x.dim() - 1, x.dim() - 1, dtype=np.int))).mean())
38
39 def e_div2(x, y, ep=1e-12):
40     return (y * (((y.clamp(ep) / (x.clamp(ep))).log()).pow(2)).sum(axis=tuple(np.linspace(1, x.dim() - 1, x.dim() - 1, dtype=np.int))).mean())
41
42 def gaussianLoss(x, y, sigma=0.5):
43     dist = Normal(x, sigma)
44     l = dist.log_prob(y)
45     ind = l.ne(1)
46     l[ind] = 0.0
47     return -l.sum(axis=-1).mean()
48
49 def polynomial(x, points):
50     return np.diff(np.polynomial.polynomial.polyval(points, x))
51 kld = nn.KLDivLoss(reduction="batchmean")
52 def MDDKLD(x, y):
53     d1 = kld((x.clamp(1e-12)).log(), y)
54     d2 = kld((y.clamp(1e-12)).log(), x)
55     return torch.max(d1, d2)
56 def MDDEDIV(x, y):
57     d1 = e_div(x, y)
58     d2 = e_div2(x, y)
59     return torch.max(d1, d2)
60 def LOGMSE(x, y, ep=1e-12):
61     return (y.clamp(ep) / x.clamp(ep)).log().pow(2).mean()
62
63
64
65 def dataLoader(T=100, data="Gene"):
66     """
67         loads data from files into python and if it doesn't exist in the case
68         of the gaussian data it will generate it.
69     """
70     if data == "Gene":
71         print(f"Using modelDataBin_{T}.hdf5")
72         #load gene data
73         with h5py.File(f"modelDataBin_{T}.hdf5", "r") as f:
74             xall = f["parameters"][:]
75             yall = f["results"][:]
76             scalesM = f["scalesM"][:]

```

```

77     scalesS = f["scalesS"][:]
78 elif data=="BiGene":
79     print(f"Using modelDataBiBin_{T}.hdf5")
80     #load bigene data
81     with h5py.File(f"modelDataBiBin_{T}.hdf5","r") as f:
82         xall = f["parameters"][:]
83         yall = f["results"][:]
84         scalesM = f["scalesM"][:]
85         scalesS = f["scalesS"][:]
86
87
88 elif data == "Gaussian":
89
90     path = f"gaussian_{T}.hdf5"
91     if os.path.exists(path):
92         #load gaussian data
93         print(f"Loading gaussian mixture from {path}")
94         with h5py.File(path,"r") as f:
95             xall=f["parameters"][:]
96             scalesM =f["scalesM"][:]
97             scalesS =f["scalesS"][:]
98             yall = f["results"][:]
99     else:
100
101     print("Generating mixture of gaussians...")
102     n = 10**5 # samples
103     mu = np.random.uniform(-3,3,size=(n,3))
104     var = np.random.uniform(0.1,5,size=(n,3))
105     k = np.zeros((n,4))
106     k[:,3] = 1
107     k[:,1:3] = np.sort(np.random.uniform(0,1,size=(n,2)),axis=-1)
108     k = np.diff(k)
109     x = np.hstack((var,mu,k))
110
111     points = np.linspace(-15,15,T+1)
112
113     dt = 30/T
114
115     # openmp parallelised calls function to get the bin probability
116     with Pool(os.cpu_count()-1) as p:
117         y = p.map(partial(gaussianMix,points=points,dt=dt),x)
118     y= np.array(y)
119     # get scaling parameters
120     scalesM = np.mean(x,axis=0)
121     scalesS = np.std(x,axis=0)
122     xall = scale(x) # scale input
123     yall = y
124     # save gaussian data
125     with h5py.File(path,"w") as f:
126         f.create_dataset("results",data=yall)
127         f.create_dataset("scalesM",data=scalesM)
128         f.create_dataset("scalesS",data=scalesS)
129         f.create_dataset("parameters",data=xall)
130
131     return xall,yall,scalesM,scalesS
132
133 class model(nn.Module):
134     """
135     LSTM Model class
136     """
137     def __init__(self,input_dim=9,
138                  T=100,prev=32,lstm_size=512,device_index=0,
139                  batch_size=1024,e_learning_rate=1e-5,bidir=True,
140                  output_act=nn.Sigmoid(),output_act_dense=nn.Softmax(dim=1),dense_act=nn.ELU(),celltype=LSTMCell):
141
142         super().__init__()
143         # set model constants
144         self.input_dim = input_dim
145         self.T = T
146         self.prev = prev
147         self.batch_size = batch_size
148         self.e_learning_rate = e_learning_rate
149         self.enc_size = lstm_size
150         self.output_act_type = output_act
151         self.output_dense_act_type = output_act_dense
152         self.dense_act_type = dense_act
153         self.bidir = bidir
154         self.celltype=celltype
155         self.device_index = device_index
156
157         self._build_net()
158     def _build_net(self):
159         # build layers
160         self.use_cuda = torch.cuda.is_available()
161         self.device = torch.device(f"cuda:{self.device_index}" if self.use_cuda else "cpu")
162         self.index_tensor = torch.tensor([[t]*self.batch_size for t in
163                                         range(self.T+1)]).to(self.device,non_blocking=self.use_cuda)
164         self.input_1_1 = nn.Linear(self.input_dim,32).to(self.device,non_blocking=self.use_cuda)
165         self.input_act = self.dense_act_type.to(self.device,non_blocking=self.use_cuda)
166         self.input_1_2 = nn.Linear(32,64).to(self.device,non_blocking=self.use_cuda)
167         self.input_1_3 = nn.Linear(64,128).to(self.device,non_blocking=self.use_cuda)
168         self.input_1_4 = nn.Linear(128,256).to(self.device,non_blocking=self.use_cuda)

```

```

167     self.input_1_5 = nn.Linear(256,128).to(self.device,non_blocking=self.use_cuda)
168     self.input_1_x = nn.ModuleList([nn.Linear(128,64) for _ in range(self.T)]).to(self.device,non_blocking=self.use_cuda)
169
170
171     self.output_layer=nn.Linear((self.T)*(self.bidir+1),self.T).to(self.device,non_blocking=self.use_cuda)
172     self.output_act = self.output_act_type.to(self.device,non_blocking=self.use_cuda)
173     self.loss = torch.zeros(1,dtype=torch.float32).to(self.device,non_blocking=self.use_cuda)
174     # build lstm array
175     if self.bidir:
176
177         ↪ self.lstmArray("array1",self.input_dim,self.T,self.prev,self.enc_size,self.batch_size,self.dense_act_type,self.output_act_type
178         ↪ )
179
180         ↪ self.lstmArray("array2",self.input_dim,self.T,self.prev,self.enc_size,self.batch_size,self.dense_act_type,self.output_act_type
181         ↪ )
182     else:
183
184         ↪ self.lstmArray("array1",self.input_dim,self.T,self.prev,self.enc_size,self.batch_size,self.dense_act_type,self.output_act_type
185         ↪ )
186
187     def lstmArray(self,array,input_dim=4,
188                 T=100,prev=32,lstm_size=512,
189                 batch_size=128,dense_act=nn.ReLU(),output_act=nn.Sigmoid(),celltype=LSTMCell):
190
191         ys = [nn.Linear(self.enc_size,1) for _ in range(self.T)]
192         setattr(self,array+"_ys",nn.ModuleList(ys).to(self.device,non_blocking=self.use_cuda))
193         ys_sigmoid = [output_act for _ in range(self.T)]
194         setattr(self,array+"_ys_sigmoid",nn.ModuleList(ys_sigmoid).to(self.device,non_blocking=self.use_cuda))
195         y_prev_dense = [nn.Linear(self.prev,128) for _ in range(self.T)]
196         setattr(self,array+"_y_prev_dense",nn.ModuleList(y_prev_dense).to(self.device,non_blocking=self.use_cuda))
197         y_prev_dense_act = [dense_act for _ in range(self.T)]
198         setattr(self,array+"_y_prev_dense_act",nn.ModuleList(y_prev_dense_act).to(self.device,non_blocking=self.use_cuda))
199         #setattr(self,array+"_dropouts",nn.ModuleList(nn.Dropout(0.1) for _ in
200         #↪ range(self.T+1))).to(self.device,non_blocking=self.use_cuda)
201
202         ↪ setattr(self,array+"_lstm_enc",torch.jit.script(self.celltype(192,self.enc_size)).to(self.device,non_blocking=self.use_cuda))
203
204     def _embedding_input(self,x):
205         # forward function for initial dense layers
206         x = self.input_1_1(x)
207         x = self.input_act(x)
208         x = self.input_1_2(x)
209         x = self.input_act(x)
210         x = self.input_1_3(x)
211         x = self.input_act(x)
212         x = self.input_1_4(x)
213         x = self.input_act(x)
214         x = self.input_1_5(x)
215         x = self.input_act(x)
216
217         return x
218
219     def forward(self,x):
220
221         # put through initial dense layers
222         xe = self._embedding_input(x)
223         # Feed through LSTM array
224         ys = self._forwardLoop1(xe)
225         if self.bidir: # if bidirectional feed through other lstm array and concatenate
226             ys = torch.cat((self._forwardLoop2(xe),ys),1)
227         # feed through final dense layer and softmax
228         x=ys
229
230         x = self.output_layer(x)
231         x = self.output_dense_act_type(x)
232
233         return x
234
235     def _forwardLoop1(self,xe):
236         # forward LSTM array
237         enc_state_0,enc_state_1 = self.initial_hidden() # initialise hidden layers
238         # initialise list to store previous outputs
239         y_prev = [torch.zeros([self.batch_size,1],dtype=torch.float32).to(self.device,non_blocking=self.use_cuda)] * self.prev
240         # initialise result tensor
241         ys = torch.zeros([self.batch_size,self.T],dtype=torch.float32).to(self.device,non_blocking=self.use_cuda)
242
243         for t in range(0,self.T):# go through time
244             # feed previous results into dense layer
245             x = self.array1_y_prev_dense[t](torch.cat(tuple(y_prev),dim=1))
246             x = self.array1_y_prev_dense_act[t](x)
247             # feed output initial of dense layers into dense layer
248             xt = self.input_1_x[t](xe)
249             xt = self.input_act(xt)
250             x = torch.cat((xt,x),1)# concatenate
251             # feed concatenated results into lstm array
252             enc_state_0,enc_state_1 = self.array1_lstm_enc(x,enc_state_0,enc_state_1)
253             # put output into dense layer
254             x = self.array1_ys[t](enc_state_0)
255             x = self.array1_ys_sigmoid[t](x)
256             ys.scatter_(1,self.index_tensor[t],x) # put results into results tensor
257             # update previous tensor
258             y_prev = y_prev[1:] +[x]
259
260

```

```

250     return ys
251
252     def _forwardLoop2(self,xe):
253         enc_state_0,enc_state_1 = self.initial_hidden() # initialise hidden states of lstm
254         # initialise list to store previous results
255         y_prev = [torch.zeros([self.batch_size,1],dtype=torch.float32).to(self.device,non_blocking=self.use_cuda)] * self.prev
256         # initialise tensor to store output
257         ys = torch.zeros([self.batch_size,self.T],dtype=torch.float32).to(self.device,non_blocking=self.use_cuda)
258
259         for t in range(0,self.T): # go through time index in for loop reverses order
260             # feed previous states into dense layer
261             x = self.array2_y_prev_dense[t](torch.cat(tuple(y_prev),dim=1))
262             x = self.array2_y_prev_dense_act[t](x)
263             # feed output of initial dense layers into dense layer
264             xt = self.input_1_x[-t-1](xe)
265             xt = self.input_act(xt)
266             x = torch.cat((xt,x),1)#concatenate
267             # feed into lstm array
268             enc_state_0,enc_state_1 = self.array2_lstm_enc(x,enc_state_0,enc_state_1)
269             # feed output into dense layer and apply softmax
270             x = self.array2_ys[t](enc_state_0)
271             x = self.array2_ys_sigmoid[t](x)
272             # cast to results tensor
273             ys.scatter_(1,self.index_tensor[t],x)
274             # update previous results list
275             y_prev = y_prev[1:] +[x]
276
277         return ys
278
279     @staticmethod
280     def get_yprev23(y_prev,t,ys):
281         return torch.cat((y_prev[:,1:],ys[:,t-1].unsqueeze(1)),1)
282
283     def initial_hidden(self):
284         # create initial states for lstm arrays
285         enc_state_0 = torch.zeros([self.batch_size,self.enc_size],dtype=torch.float32).to(self.device,non_blocking=self.use_cuda)
286         enc_state_1 = torch.zeros([self.batch_size,self.enc_size],dtype=torch.float32).to(self.device,non_blocking=self.use_cuda)
287         return (enc_state_0,enc_state_1)
288
289     @Torch.jit.export
290     def train_model(self,epochs=100,data="Gene",
291                     optim="Adam",loss_fn="MSE",custom_data=None,
292                     model_name=None):
293         # create directory to store models
294         self._directory_maker()
295         # set constants
296         self.epochs = epochs
297         self.loss_fn = loss_fn
298         self.optim = optim
299
300         if custom_data is None:
301             # load data
302             xall,yall=self.load_data(data=data)
303             # create train and test split
304             x_train, x_test,y_train,y_test = train_test_split(xall,yall,train_size=0.8)
305         else:
306             # if custom data is given use that
307             x_train, x_test,y_train,y_test =
308                 custom_data["x_train"],custom_data["x_test"],custom_data["y_train"],custom_data["y_test"]
309             self.scalesM,self.scalesS = custom_data["scalesM"],custom_data["scalesS"]
310             print(f"Loading model onto {self.device}")
311             print(f"Loading tensorboard writer")
312             # create model name
313             if model_name is None:
314                 oact = str(self.output_act_type).split(",1")[0]
315                 dact = str(self.dense_act_type).split(",1")[0]
316                 odact = str(self.output_dense_act_type).split(",1")[0]
317                 celltype = self.celltype.__str__()
318                 t = time.strftime("%d-%m-%H-%M-%S",time.localtime())
319                 self.model_name = f"cell={celltype}_b={self.batch_size}_loss={self.loss_fn}_oact={oact}_dact={dact}_odact={odact}"
320             else:
321                 self.model_name = model_name
322             # start tensorboard log
323             writer = SummaryWriter(comment=self.model_name)
324             # create optimiser
325             self._init_optim()
326             print("Loading dataset")
327             # create data loaders for train and test data
328             xtrainTensor=torch.from_numpy(x_train).type(torch.float32)
329             ytrainTensor=torch.from_numpy(y_train).type(torch.float32)
330             dataset = torch.utils.data.TensorDataset(xtrainTensor,ytrainTensor)
331             loader = DataLoader(dataset, batch_size=self.batch_size,drop_last=True,shuffle=True,
332                                 pin_memory=self.use_cuda,num_workers=2)
333             xtestTensor=torch.from_numpy(x_test).type(torch.float32)
334             ytestTensor=torch.from_numpy(y_test).type(torch.float32)
335             valDataset = torch.utils.data.TensorDataset(xtestTensor,ytestTensor)
336             valLoader = DataLoader(valDataset, batch_size=self.batch_size,drop_last=True,shuffle=False,
337                                 pin_memory=self.use_cuda,num_workers=2)
338             # create matrix to store loss
339             self._init_loss()

```

```

340     del valDataset,dataset,xtestTensor,ytestTensor,x_train,x_test,y_train,y_test
341     # start train loop
342     print("Train Loop Starting...")
343     for i in range(self.epochs):
344         #nvtx.range_push("Epoch " + str(i+1))
345         self.train() # put model in train mode
346         print(f"STEP: {i+1}/{self.epochs}")
347         #nvtx.range_push("Training")
348         # train
349         with tqdm(total=train_batches,ascii=True,desc="Training") as pbar:# loading bar
350             for idx, (local_batch, local_labels) in enumerate(loader):
351                 #nvtx.range_push("Batch "+str(idx))
352
353                 #nvtx.range_push("Copying to device")
354                 # get batch and send to device
355                 y=local_labels.to(self.device,non_blocking=self.use_cuda)
356                 x_batch=local_batch.to(self.device,non_blocking=self.use_cuda)
357                 #nvtx.range_pop()
358                 #nvtx.range_push("Forward pass")
359                 # forward pass
360                 x=self(x_batch)
361                 #nvtx.range_pop()
362
363                 #nvtx.range_push("Backward pass")
364                 # update loss and backpropigate
365                 self._update_loss_train(x,y,i)
366
367
368                 #nvtx.range_pop()
369                 #nvtx.range_pop()
370
371
372             pbar.update(1)
373
374             #nvtx.range_pop()
375             # average
376             self.train_loss[i] = self.train_loss[i]/train_batches
377             if np.isnan(self.train_loss[i,0]):# if NaN stop training
378                 print("GOT NaN")
379                 break
380             # update tensorbard and print
381             for l,s in enumerate(self.loss_names):
382                 name=f'Loss {s}/train'
383                 writer.add_scalar(name, float(self.train_loss[i,l]), i)
384                 print(f'{name}: {self.train_loss[i,l]}')
385
386
387
388             #nvtx.range_push("Evaluation")
389             # get val loss
390             with torch.no_grad():
391                 self.eval() # set to eval
392                 with tqdm(total=test_batches,ascii=True,desc="Validating") as pbar:# loading bar
393                     for local_batch, local_labels in valLoader:
394                         # get batch
395                         y = local_labels.to(self.device,non_blocking=self.use_cuda)
396                         x_batch=local_batch.to(self.device,non_blocking=self.use_cuda)
397                         # forward pass
398                         pred = self(x_batch)
399                         # update val loss
400                         self._update_loss_val(pred,y,i)
401                         pbar.update(1)
402
403             #nvtx.range_pop()
404             # average
405             self.val_loss[i] = self.val_loss[i]/test_batches
406             # update tensorboard and print results
407             for l,s in enumerate(self.loss_names):
408                 name=f'Loss {s}/test'
409                 writer.add_scalar(name, float(self.val_loss[i,l]), i)
410                 print(f'{name}: {self.val_loss[i,l]}')
411             #nvtx.range_pop()
412             # save model if val loss is lower than previous best
413             if self.val_loss[i,0] < old_test_loss:
414                 #nvtx.mark("Saving best model")
415                 torch.save({ "state_dict":self.state_dict(),
416                             "scalesM":self.scalesM,
417                             "scalesS":self.scalesS,
418                             "epoch":i+1,
419                             "optim":self.optimiser,
420                             "val_loss":self.val_loss,
421                             "dense_act":str(self.dense_act_type),
422                             "output_act":str(self.output_act_type),
423                             "batch_size": self.batch_size,
424                             "enc_size":self.enc_size,
425                             "bidir":self.bidir,
426                             "T": self.T,
427                             "output_dense_act":self.output_dense_act_type,
428                             "prev": self.prev,
429                             "input_dim": self.input_dim,
430                             "loss_fn":self.loss_fn}, f"./best/bestModel_{self.model_name}.pth")

```

```

431         old_test_loss = self.val_loss[i,0]
432     # save final model
433     torch.save({"state_dict":self.state_dict(),
434                 "scalesM":self.scalesM,
435                 "scalesS":self.scalesS,
436                 "epoch":self.epochs,
437                 "optim":self.optimiser,
438                 "val_loss":self.val_loss,
439                 "dense_act":str(self.dense_act_type),
440                 "output_act":str(self.output_act_type),
441                 "batch_size": self.batch_size,
442                 "enc_size":self.enc_size,
443                 "bidir":self.bidir,
444                 "T": self.T,
445                 "output_dense_act":self.output_dense_act_type,
446                 "prev": self.prev,
447                 "input_dim": self.input_dim,
448                 "loss_fn":self.loss_fn}, f"./final/finalModel_{self.model_name}.pth")
449     return self.train_loss,self.val_loss
450 def _directory_maker(self):
451     # create directories for models
452     try:
453         print("Making directory for best models")
454         os.mkdir("./best")
455     except FileExistsError:
456         print("Directory for best models exists")
457     try:
458         print("Making directory for final models")
459         os.mkdir("./final")
460     except FileExistsError:
461         print("Directory for final models exists")
462
463 def _init_optim(self):
464     # create optimiser
465     if self.optim == "Adam":
466         self.optimiser = torch.optim.Adam(self.parameters(),
467                                         lr=self.e_learning_rate,
468                                         betas=(0.5,0.999))#,weight_decay=0.1
469     elif self.optim == "SGD":
470         self.optimiser = torch.optim.SGD(self.parameters(),
471                                         lr=self.e_learning_rate)
472                                         #momentum=0.1),weight_decay=0.1
473
474 def _init_loss(self):
475     # initialise loss functions
476     self.mse_fn = nn.MSELoss()
477     self._kld = nn.KLDivLoss(reduction="batchmean")
478     self._div_fn = lambda x,y: self._kld(x.clamp(1e-20).log(),y)
479     self.l1LossFN_t = nn.L1Loss(reduction="none")
480     self.l1LossFN = nn.L1Loss()
481     self.h2 = h2_dist
482     self.ediv = e_div
483     self.ediv2 = e_div2
484     self.MDDKLD = MDDKLD
485     self.MDDEDIV = MDDEDIV
486     self.LOGMSE = LOGMSE
487     #self.edivc = e_div_c
488     #self.a_mse =lambda x,y : adj_mse(x,y)+ dist_loss(x,y)
489
490     self.l1_max_fn = lambda x,y : self.l1LossFN_t(x,y).max(dim=-1)[0].mean()
491     allLoss = np.array(["MSE","DIV","L1 MAX","L1","H2","EDIV","EDIV2","MDDKLD","MDDEDIV","LOGMSE"])
492     lossFns =
493     ↪ np.array([self.mse_fn,self.div_fn,self.l1_max_fn,self.l1LossFN,self.h2,self.ediv,self.ediv2,self.MDDKLD,self.MDDEDIV,self.LOGMSE])
494     ind1 = np.where(allLoss==self.loss_fn)
495     ind2 = np.where(allLoss!=self.loss_fn)
496     s = allLoss.size
497     self.loss_names = np.concatenate((allLoss[ind1],allLoss[ind2]))
498     self.loss_fns = np.concatenate((lossFns[ind1],lossFns[ind2]))
499     # create array to store loss
500     self.train_loss = np.zeros((self.epochs,s))
501     self.val_loss = np.zeros((self.epochs,s))
502
503 def _update_loss_train(self,x,y,i):
504     # get loss
505     loss = self.loss_fns[0](x,y)
506     self.optimiser.zero_grad()
507     # backpropigate
508     loss.backward()
509     self.optimiser.step()
510     # update loss
511     self.train_loss[i,0] += float(loss.detach().cpu())
512     self.train_loss[i,1:] += [float(s(x.detach(),y).cpu()) for s in self.loss_fns[1:]]
513
514 def _update_loss_val(self,x,y,i):
515     # update validation loss
516     self.val_loss[i] += [float(s(x.detach(),y).cpu()) for s in self.loss_fns]
517
518 def load_data(self,data="Gene"):
519     #load data using data loader function
520     xall,yall,self.scalesM,self.scalesS = dataLoader(T=self.T,data=data)

```

```

521     return xall,yall
522
523     @torch.jit.export
524     def load_model(self,PATH):
525         # load model
526         model = torch.load(PATH,map_location=self.device)
527         self.scalesM = model["scalesM"]
528         self.scalesS = model["scalesS"]
529         self.output_act_type = eval("nn."+model["output_act"])
530         self.dense_act_type = eval("nn."+model["dense_act"])
531         self.batch_size = model["batch_size"]
532         self.bidir = model["bidir"]
533         self.enc_size = model["enc_size"]
534         self.T = model["T"]
535         self.prev = model["prev"]
536         self.input_dim=model["input_dim"]
537         #self._build_net()
538         self.load_state_dict(model["state_dict"])
539
540
541
542
543
544     @torch.jit.export
545     def predict(self,x):
546         # predict outputs
547         self.eval()#eval mode
548         #scale
549         x = (np.atleast_2d(x) - self.scalesM)/self.scalesS
550         p = int(np.ceil(x.shape[0] /self.batch_size))
551         t = x.shape[0]
552         r = self.batch_size - (t % self.batch_size)
553         x = np.pad(x,((0,r),(0,0)))
554         xs = np.array_split(x,p)
555         res = [None] * p
556         for i in range(p):# feed x in batch size chunks
557             res[i] =
558                 ↪ self(torch.from_numpy(xs[i]).type(torch.float32).to(self.device,non_blocking=self.use_cuda)).detach().cpu().numpy()
559
560         res[-1] = res[-1][:-r]
561         return np.vstack(tuple(res))
562
563 if __name__ == "__main__":
564     # learning_rate=1e-4
565     use_cuda = torch.cuda.is_available()
566     device = torch.device("cuda" if use_cuda else "cpu")
567     cell_types = [LSTMCell,PeepholeLSTMCell,GRUCell]
568     lstm_size=[512,1024]
569     dense_act = [nn.ELU(),nn.ReLU(),nn.Tanhshrink(),nn.Tanh(),nn.Softplus()]
570     output_act = [nn.Sigmoid(),nn.Softplus()]
571     output_act_dense =[nn.Sigmoid(),nn.Softplus(),nn.Softmax()]
572     loss = ["MSE","L1","H2","MDDKLD","MDDEDIV","LOGMSE"]
573     optim = ["Adam","SGD"]
574     prev = [1,4,8,16,32,64]
575     bidir = [True,False]
576     g = ParameterGrid({"optim":optim,"loss":loss})
577
578     val_loss = [[None,i] for i in g]
579     train_loss = [[None,i] for i in g]
580
581     for v,i in enumerate(g):
582         try:
583             print(i)
584             m = model(bidir=False,e_learning_rate=1e-4,batch_size=128).to(device)
585             train_loss[v][0],val_loss[v][0] = m.train_model(data="Gaussian",optim=i["optim"],loss_fn=i["loss"],epochs=100)
586
587         except Exception as inst:
588             print("ERROR!!!!")
589             print(type(inst))
590             print(inst)
591         with open('train_loss.pickle', 'wb') as f:
592             pickle.dump(train_loss, f, protocol=pickle.HIGHEST_PROTOCOL)
593         with open('val_loss.pickle', 'wb') as f:
594             pickle.dump(val_loss, f, protocol=pickle.HIGHEST_PROTOCOL)
595
596
597

```

XII.B.3 Ensemble LSTM Model

```

1  from pytorchModel import *
2
3  class ensembleLSTM(model):
4      """
5          ensemble of lstm models
6      """
7  def __init__(self,T,n,batch_size=1024,output_activation=nn.Softmax(dim=1),learning_rate=1e-5):

```

```

8      super(model,self).__init__()# inherit
9      # set consts
10     self.use_cuda = torch.cuda.is_available()
11     self.device = torch.device("cuda" if self.use_cuda else "cpu")
12     self.n = n
13     self.T = T
14     self.batch_size=batch_size
15     self.e_learning_rate=learning_rate
16     # create layer
17     self.output_activation=output_activation
18     self.linear = nn.Linear(self.n*(self.T),self.T).to(self.device,non_blocking=self.use_cuda)
19
20     def forward(self,x):
21         # forward method
22         return self.output_activation(self.linear(x))
23     @torch.jit.export
24     def train_model(self,loss_fn,optim,epochs,data):
25         self._directory_maker() # create directory to store models
26         # set const
27         self.epochs = epochs
28         self.loss_fn = loss_fn
29         self.optim = optim
30         # create date
31         self.create_data(data)
32         # set model name
33         self.model_name = f"ensemble_n={self.n}"
34         # create tensorboard writer
35         writer = SummaryWriter(comment=self.model_name)
36         # create optimiser
37         self._init_optim()
38         print("Loading dataset")
39         # create data loader
40         xtrainTensor=torch.from_numpy(self.x_train).type(torch.float32)
41         ytrainTensor=torch.from_numpy(self.y_train).type(torch.float32)
42         dataset = torch.utils.data.TensorDataset(xtrainTensor,ytrainTensor)
43         loader = DataLoader(dataset, batch_size=self.batch_size,drop_last=True,shuffle=True,
44                             pin_memory=self.use_cuda,num_workers=2)
45         xtestTensor=torch.from_numpy(self.x_test).type(torch.float32)
46         ytestTensor=torch.from_numpy(self.y_test).type(torch.float32)
47         valDataset = torch.utils.data.TensorDataset(xtestTensor,ytestTensor)
48         valLoader = DataLoader(valDataset, batch_size=self.batch_size,drop_last=True,shuffle=False,
49                             pin_memory=self.use_cuda,num_workers=2)
50         # create loss array
51         self._init_loss()
52         old_test_loss = np.inf
53         train_batches = self.x_train.shape[0] // self.batch_size
54         test_batches = self.x_test.shape[0] // self.batch_size
55         print("Train Loop Starting...")
56         # start train loop
57         for i in range(self.epochs):
58             #nutx.range_push("Epoch " + str(i+1))
59             self.train()
60             print(f"STEP: {i+1}/{self.epochs}")
61             #nutx.range_push("Training")
62             with tqdm(total=train_batches,ascii=True,desc="Training") as pbar:#loading bar
63                 for idx, (local_batch, local_labels) in enumerate(loader):
64                     #nutx.range_push("Batch "+str(idx))
65
66                     #nutx.range_push("Copying to device")
67                     # get batches
68                     y=local_labels.to(self.device,non_blocking=self.use_cuda)
69                     x_batch=local_batch.to(self.device,non_blocking=self.use_cuda)
70                     #nutx.range_pop()
71                     #nutx.range_push("Forward pass")
72                     # forward pass
73                     x=self(x_batch)
74                     #nutx.range_pop()
75
76                     #nutx.range_push("Backward pass")
77                     # update loss and back prop
78                     self._update_loss_train(x,y,i)
79
80                     #nutx.range_pop()
81                     #nutx.range_pop()
82
83
84
85                     pbar.update(1)
86
87                     #nutx.range_pop()
88                     # average
89                     self.train_loss[i] = self.train_loss[i]/train_batches
90                     if np.isnan(self.train_loss[i,0]):# break if nan
91                         print("GOT NaN")
92                         break
93
94                     # write loss to tensorboard log and print loss
95                     for l,s in enumerate(self.loss_names):
96                         name=f'Loss {s}/train'
97                         writer.add_scalar(name, float(self.train_loss[i,1]), i)
98                         print(f"{name}: {self.train_loss[i,1]}")

```

```

99
100
101
102     #nvtx.range_push("Evaluation")
103     # get val loss
104     with torch.no_grad():
105         self.eval() # set to eval mode
106         with tqdm(total=test_batches,ascii=True,desc="Validating") as pbar:# loading bar
107             for local_batch, local_labels in valLoader:
108                 #get batch
109                 y = local_labels.to(self.device,non_blocking=self.use_cuda)
110                 x_batch=local_batch.to(self.device,non_blocking=self.use_cuda)
111                 # forward pass
112                 pred = self(x_batch)
113                 # update val loss
114                 self._update_loss_val(pred,y,i)
115                 pbar.update(1)
116
117     #nvtx.range_pop()
118     # average
119     self.val_loss[i] = self.val_loss[i]/test_batches
120     # write tensorboard log print loss
121     for l,s in enumerate(self.loss_names):
122         name=f'Loss {s}/test'
123         writer.add_scalar(name, float(self.val_loss[i,l]), i)
124         print(f'{name}: {self.val_loss[i,l]}')
125     #nvtx.range_pop()
126     # if model has the lowest val loss save
127     if self.val_loss[i,0] < old_test_loss:
128         torch.save({"state_dict":self.state_dict(),
129                     "n":self.n,
130                     "epoch":i+1,
131                     "optim":self.optimiser,
132                     "val_loss":self.val_loss,
133                     "batch_size": self.batch_size,
134                     "T": self.T,
135                     "loss_fn":self.loss_fn}, f"./best/bestModel_{self.model_name}.pth")
136     # save final model
137     torch.save({"state_dict":self.state_dict(),
138                     "n":self.n,
139                     "epoch":i+1,
140                     "optim":self.optimiser,
141                     "val_loss":self.val_loss,
142                     "batch_size": self.batch_size,
143                     "T": self.T,
144                     "loss_fn":self.loss_fn}, f"./final/finalModel_{self.model_name}.pth")
145
146     return self.train_loss,self.val_loss
147
148     def _load_model(self,i):
149         # load sub models and set to eval
150         m = model(bidir=True,batch_size=512, input_dim=self.inp, T=100, dense_act=nn.ELU(),output_act=nn.Softplus(),
151                    ↪ output_act_dense=nn.Softmax()).to(self.device)
152         m.load_model(f"bestModel_ensemble_model_lstm_{self.data}_{i+1}.pth")
153         m.eval()
154
155         return m
156
157     def _submodel_predict(self,x):
158         # load each model and get its prediction in series
159         # done in series to save memory
160         out = [None] * self.n
161         for i in range(self.n):
162             out[i] = self._load_model(i).predict(x)
163         return np.hstack(tuple(out))
164
165     def create_data(self,data):
166         # set const
167         self.data=data
168         if data == "gau":
169             self.inp = 9
170         else:
171             self.inp = 5
172
173         # load data
174         with h5py.File(f"ensemble_data_{data}.hdf5","r") as f:
175             x_train = f["x_train"][:]
176             x_test = f["x_test"][:]
177             y_train = f["y_train"][:]
178             y_test = f["y_test"][:]
179             self.x_train,self.x_test,self.y_train,self.y_test
180             ↪ =self._submodel_predict(x_train),self._submodel_predict(x_test),y_train,y_test
181
182     if __name__=="__main__":
183         use_cuda = torch.cuda.is_available()
184         device = torch.device("cuda" if use_cuda else "cpu")
185
186         m = ensemble(n=10,T=100).to(device)
187         m(torch.randn((1024,10*101)).to(device))
188         train_loss,val_loss=m.train_model("EDIV2","adam",500)
189         with open('train_loss.pickle', 'wb') as f:
190             pickle.dump(train_loss, f, protocol=pickle.HIGHEST_PROTOCOL)
191         with open('val_loss.pickle', 'wb') as f:
192             pickle.dump(val_loss, f, protocol=pickle.HIGHEST_PROTOCOL)

```

XII.C Mixture Model

XII.C.1 Mixture Model

```

1  from torch.distributions.normal import Normal
2  from torch.distributions.log_normal import LogNormal
3  from torch.distributions.poisson import Poisson
4  from torch.distributions.exp_family import ExponentialFamily
5  from torch.distributions import constraints
6  from torch.distributions.utils import broadcast_all
7  from torch.distributions.categorical import Categorical
8  from numbers import Number
9  from pytorchModel import *
10 import itertools
11
12
13 # define these wrappers so functions have a readable name when converted to a string
14
15 class Exponential():
16     def __init__(self):
17         return None
18     def __call__(self,x):
19         return x.exp().clamp(1e-12)
20     def __str__(self):
21         return "Exponential"
22
23 class Clamped():
24     def __init__(self,ep=1e-12):
25         self.ep = ep
26         return None
27     def __call__(self,x):
28         return x.clamp(self.ep)
29     def __str__(self):
30         return "Clamped"
31
32
33
34 class gammaPoisson(ExponentialFamily):
35     """
36     create gamma-poisson distribution class
37     """
38
39
40     arg_constraints = {'alpha': constraints.positive,
41                        'beta': constraints.positive}
42     support = constraints.nonnegative_integer
43     def __init__(self,alpha,beta,validate_args=None):
44         # set parameters
45         self.alpha,self.beta = broadcast_all(alpha,beta)
46         if isinstance(alpha, Number):
47             batch_shape = torch.Size()
48             self.device = torch.device("cpu")
49             self.batch_size = torch.tensor([1])
50         else:
51             batch_shape = self.alpha.size()
52             self.device = self.alpha.device
53             self.batch_size = batch_shape[:-1]
54
55         super(gammaPoisson, self).__init__(batch_shape, validate_args=validate_args)
56     def cdf(self,x):
57         if self._validate_args:
58             self._validate_sample(x)
59         x, = broadcast_all(x)
60         n = x.floor().long() # cast to integers
61         # create tensor of all numbers up to max
62         r = torch.linspace(0,n.max(),n.max()+1).to(self.device)
63         res = torch.zeros(*self.batch_size,n.max()+2).to(self.device)
64         t1 = (r+self.alpha).lgamma()
65         t2 = -self.alpha.lgamma()-(r+1.0).lgamma()
66         t3 = self.alpha*self.beta.log() - (r+self.alpha)*(1.0+self.beta).log()
67         # evaluate cdf and take cumsum
68         res[...,-1] = (t1+t2+t3).exp().cumsum(dim=-1).clamp(0.0,1.0)
69         # return relevant values
70         ind = n +1
71         ind[n<0] = 0
72         return res[...,-ind]
73     @property
74     def mean(self):
75         return self.alpha/self.beta
76
77     @property
78     def variance(self):
79         return (self.alpha * (1.0 + self.beta)+1.0)/self.beta.pow(2)
80
81 class bivariateGammaPoisson(ExponentialFamily):

```

```

82 """
83 bivariate multinomial-dirichlet gamma-poisson distribution class
84 """
85 def __init__(self, alpha, beta, theta0, theta1, validate_args=None):
86     # initialise parameters of distribution
87     self.alpha, self.beta, self.theta0, self.theta1 = broadcast_all(alpha, beta, theta0, theta1)
88     self.alpha, self.beta, self.theta0, self.theta1 = self.alpha.reshape(self.alpha.shape[0], 1, 1),
89     ↪ self.beta.reshape(self.alpha.shape[0], 1, 1), self.theta0.reshape(self.alpha.shape[0], 1, 1),
90     ↪ self.theta1.reshape(self.alpha.shape[0], 1, 1)
91     self.thetaT = self.theta0 + self.theta1
92     if isinstance(alpha, Number):
93         batch_shape = torch.Size()
94         self.device = torch.device("cpu")
95         self.batch_size = torch.tensor([1])
96     else:
97         batch_shape = self.alpha.size()
98         self.device = self.alpha.device
99         self.batch_size = batch_shape[:2]
100    super(bivariateGammaPoisson, self).__init__(batch_shape, validate_args=validate_args)
101   def cdf(self, x, y):
102       # returns cdf
103       if self._validate_args:
104           self._validate_sample(x)
105           self._validate_sample(y)
106           x, y = broadcast_all(x, y)
107           # cast to integers
108           nx = x.floor().long()
109           ny = y.floor().long()
110
111       # create grid of all points up to max
112       rx = torch.linspace(0, nx.max(), nx.max() + 1).to(self.device)
113       ry = torch.linspace(0, ny.max(), ny.max() + 1).to(self.device)
114       rx, ry = torch.meshgrid([rx, ry])
115
116       # calculate log pdf
117       rx = rx.repeat(self.batch_size[0], 1).reshape(self.batch_size[0], *rx.shape)
118       ry = ry.repeat(self.batch_size[0], 1).reshape(self.batch_size[0], *ry.shape)
119       t1 = self.alpha * (self.beta.log() - (self.beta + 1).log()) + (self.alpha + rx + ry).lgamma()
120       t2 = -self.alpha.lgamma() - (rx + ry + 1).lgamma() - (rx + ry) * (self.beta + 1).log()
121       t3 = (rx + ry + 1).lgamma() - (rx + 1).lgamma() - (ry + 1).lgamma() + self.thetaT.lgamma() +
122           ↪ (rx + self.theta0).lgamma() + (ry + self.theta1).lgamma()
123       t4 = -(self.thetaT + rx + ry).lgamma() - self.theta0.lgamma() - self.theta1.lgamma()
124       res = torch.zeros(self.batch_size[0], nx.max() + 2, ny.max() + 2).to(self.device)
125       # using cumsum get cdf
126       res[:, 1:, 1:] = (t1 + t2 + t3 + t4).exp().cumsum(-1).cumsum(-2)
127       # get relevant indicies
128       ind1 = nx + 1
129       ind1[nx < 0] = 0
130       ind2 = ny + 1
131       ind2[ny < 0] = 0
132       return res[:, :, ind1, ind2].clamp(0.0, 1.0)
133   def _cdf(self, x, y):
134       # get cdf across bins
135       if self._validate_args:
136           self._validate_sample(x)
137           self._validate_sample(y)
138           # cast to integers
139           x, y = broadcast_all(x, y)
140           nx = x.floor().long()
141           ny = y.floor().long()
142           # get bin sizes
143           sx = (x[1:] - x[:-1]).long()
144           sy = (y[1:] - y[:-1]).long()
145
146           # get indicies to sum over
147           idx = torch.zeros(sx.sum() + 1, dtype=torch.long).to(self.device)
148           idx[torch.cumsum(sx, dim=0)] = torch.tensor(1).to(self.device)
149           idx = torch.cumsum(idx, dim=0)[-1]
150
151           indy = torch.zeros(sy.sum() + 1, dtype=torch.long).to(self.device)
152           indy[torch.cumsum(sy, dim=0)] = torch.tensor(1).to(self.device)
153           indy = torch.cumsum(indy, dim=0)[-1]
154           # create grid of points up to max
155           rx = torch.linspace(0, nx.max(), nx.max() + 1).to(self.device)
156           ry = torch.linspace(0, ny.max(), ny.max() + 1).to(self.device)
157           rx, ry = torch.meshgrid([rx, ry])
158
159           # calculate log pdf
160           rx = rx.repeat(self.batch_size[0], 1).reshape(self.batch_size[0], *rx.shape)
161           ry = ry.repeat(self.batch_size[0], 1).reshape(self.batch_size[0], *ry.shape)
162           t1 = self.alpha * (self.beta.log() - (self.beta + 1).log()) + (self.alpha + rx + ry).lgamma()
163           t2 = -self.alpha.lgamma() - (rx + ry + 1).lgamma() - (rx + ry) * (self.beta + 1).log()
164           t3 = (rx + ry + 1).lgamma() - (rx + 1).lgamma() - (ry + 1).lgamma() + self.thetaT.lgamma() +
165               ↪ (rx + self.theta0).lgamma() + (ry + self.theta1).lgamma()
166           t4 = -(self.thetaT + rx + ry).lgamma() - self.theta0.lgamma() - self.theta1.lgamma()
167           temp = torch.zeros(self.batch_size[0], len(x) - 1, ny.max() + 1).to(self.device)
168           # sum pdf over x for bins
169           temp.index_add_(1, idx, (t1 + t2 + t3 + t4).exp())
170           res = torch.zeros(self.batch_size[0], len(x) - 1, len(y) - 1).to(self.device)
171           # sum pdf over y for bins
172           res.index_add_(2, indy, temp).cumsum(-1).cumsum(-2)

```

```

169         return res.clamp(0.0,1.0)
170     class poissonLogNormal(ExponentialFamily):
171         """
172             poisson log-normal distribution class
173             """
174         def __init__(self,mu,sigma,quad_points=8,validate_args=None):
175             # set parameters
176             self.mu,self.sigma = broadcast_all(mu,sigma)
177             self.quad_points = quad_points
178             if isinstance(mu, Number):
179                 batch_shape = torch.Size()
180                 self.device = torch.device("cpu")
181                 self.batch_size = torch.tensor([1])
182             else:
183                 batch_shape = self.mu.size()
184                 self.device = self.mu.device
185                 self.batch_size = batch_shape[0]
186             # get quadrature points
187             grid = self.quadrature()
188             # create poissons
189             self.poisson = Poisson(grid.exp())
190             # calculate logits
191             self.logits = -np.log(self.quad_points)
192
193
194
195         super(poissonLogNormal, self).__init__(batch_shape, validate_args=validate_args)
196     def quadrature(self):
197         # calculate quad points
198         edges = torch.linspace(0.,1., steps=self.quad_points+3).to(self.device)[1:-1] # dont use 0,1 as this may give inf
199         logNormal = LogNormal(self.mu,self.sigma)
200         # use icdf to get quad points
201         quantiles = logNormal.icdf(edges)
202         # take midpoints
203         grid = (quantiles[:,1:]+quantiles[:,:-1])/2
204         return grid
205     def log_prob(self,x):
206         # return log prob
207         if self._validate_args:
208             self._validate_sample(x)
209             x, = broadcast_all(x)
210             return torch.logsumexp(self.logits + self.poisson.log_prob(x[:,None,None]),dim=-1)
211     def cdf(self,x):
212         if self._validate_args:
213             self._validate_sample(x)
214             x, = broadcast_all(x)
215             # cast to integers
216             n = x.floor().long()
217             # create vector of all points up to maximum
218             r = torch.linspace(0,n.max(),n.max()+1).to(self.device)
219             # calculate cdf from log prob
220             res = self.log_prob(r).exp().t().cumsum(-1)
221             res[res.ne(res)]=0.0 # replace nans with 0
222             return res[:,n].clamp(0.0,1.0)
223
224     class bivariatePoissonLogNormal(ExponentialFamily):
225         """
226             bivariate poisson log-normal distribution class
227             (not working)
228             """
229         def __init__(self,mu,sigma,rho,quad_points=10,validate_args=None):
230             self.mu,self.sigma = broadcast_all(mu,sigma)
231             self.rho, = broadcast_all(rho)
232             self.quad_points= quad_points
233             if isinstance(mu, Number):
234                 batch_shape = torch.Size()
235                 self.device = torch.device("cpu")
236                 self.batch_size = torch.tensor([1])
237             else:
238                 batch_shape = self.mu.size()
239                 self.device = self.mu.device
240                 self.batch_size = batch_shape[0]
241             c = self.sigma.sqrt().prod(dim=-1)[:,None]*self.rho
242
243             self.covmatrix = torch.cat((self.sigma[:,0,None], c, c, self.sigma[:,1,None]),dim=-1).view(self.batch_size,2,2)
244             self.sigmanM = self.covmatrix.cholesky()
245             self.quadrature()
246
247             self.poissonx = Poisson(self.lam1)
248             self.poissony = Poisson(self.lam2)
249
250
251         super(bivariatePoissonLogNormal, self).__init__(batch_shape, validate_args=validate_args)
252     def quadrature(self):
253         p,w = np.polynomial.hermite.hermgauss(self.quad_points)
254         #p,w = torch.from_numpy(p).type(torch.float32).to(self.device),torch.from_numpy(w).type(torch.float32).to(self.device)
255         xn = torch.tensor([i for i in itertools.product(p, repeat=2)]).type(torch.float32).to(self.device)
256         self.wn =torch.tensor([i for i in itertools.product(w, repeat=2)]).type(torch.float32).to(self.device).prod(-1).log()
257         r = (torch.matmul(self.sigmanM,xn.t())).transpose(1,2)*np.sqrt(2) + self.mu.unsqueeze(1).repeat(1,100,1)).exp()
258         lr = torch.matmul(self.sigmanM,(1./r).transpose(1,2)).transpose(1,2)
259         self.lam1,self.lam2 = r[:,0], r[:,1]

```

```

260     self.llam1,self.llam2 = lr[...,:], lr[...,1]
261     self.const = self.sigmaM.det()/(np.pi)
262
263     return
264   def prob(self,x,y):
265     if self._validate_args:
266       self._validate_sample(x)
267       self._validate_sample(y)
268     x,y = broadcast_all(x,y)
269
270     return (self.const*(self.wn+(self.llam1*self.llam2).log() +
271     ↪ self.poissonx.log_prob(x.unsqueeze(-1).unsqueeze(-1))+self.poissony.log_prob(y.unsqueeze(-1).unsqueeze(-1)) -
272     ↪ (self.lam1*self.lam2).log()).exp().sum(-1)).reshape(self.batch_size,*x.shape)
273   def cdf(self,x,y):
274     if self._validate_args:
275       self._validate_sample(x)
276     x,y = broadcast_all(x,y)
277     nx = x.floor().long()
278     ny = y.floor().long()
279     rx = torch.linspace(0,nx.max(),nx.max()+1).to(self.device)
280     ry = torch.linspace(0,ny.max(),ny.max()+1).to(self.device)
281     rx,ry=torch.meshgrid([rx,ry])
282     rx,ry=rx.flatten(),ry.flatten()
283
284     res = torch.zeros((self.batch_size,nx.max()+2,ny.max()+2)).to(self.device)
285     res[:,1:,1:] = self.prob(rx,ry).view(self.batch_size,nx.max()+1,ny.max()+1).cumsum(-1).cumsum(-2)
286     return res
287 class MixD(model):
288   """
289   Mixture Model
290   """
291   def __init__(self,T,num_mix,input_dim,bins,
292    dense_act=nn.ELU(),mix_dist="Gaussian",
293    p_act=(Exponential(),Exponential(),Exponential(),Exponential()),
294    learning_rate=1e-4,batch_size=1024):
295     super(model,self).__init__() # inherit
296     # set constants
297     self.use_cuda = torch.cuda.is_available()
298     self.device = torch.device(f"cuda" if self.use_cuda else "cpu")
299     self.batch_size = batch_size
300     self.T = T
301     self.input_dim = input_dim
302     self.p0_act = p_act[0]
303     self.p1_act = p_act[1]
304     self.mix_dist = mix_dist
305     self.bivariate = False
306     self.num_mix = num_mix
307     # create layers
308     self.p0 = nn.Linear(64,self.num_mix).to(self.device,non_blocking=self.use_cuda)
309     self.p1 = nn.Linear(64,self.num_mix).to(self.device,non_blocking=self.use_cuda)
310
311     self.e_learning_rate = learning_rate
312     self.embedding1 = nn.Linear(self.input_dim,64).to(self.device,non_blocking=self.use_cuda)
313     self.embedding1Act = dense_act
314     self.embedding2 = nn.Linear(64,128).to(self.device,non_blocking=self.use_cuda)
315     self.embedding2Act = dense_act
316     self.embedding3 = nn.Linear(128,256).to(self.device,non_blocking=self.use_cuda)
317     self.embedding3Act = dense_act
318     self.embedding4 = nn.Linear(256,128).to(self.device,non_blocking=self.use_cuda)
319     self.embedding4Act = dense_act
320     self.embedding5 = nn.Linear(128,64).to(self.device,non_blocking=self.use_cuda)
321     self.embedding5Act = dense_act
322
323
324
325     self.ai = nn.Linear(64,self.num_mix)
326     self.aiAct = nn.Softmax(dim=1)
327     self.mix = torch.arange(0,self.num_mix,dtype=torch.int).to(self.device,non_blocking=self.use_cuda)
328     self.mix_dist = mix_dist
329     # set which function to use for getting predictions and set bins
330     if mix_dist == "Gaussian":
331       self.distPredict = self.normalPredict
332       self.bins = bins.to(self.device,non_blocking=self.use_cuda)
333     elif mix_dist == "GammaPoisson":
334       self.distPredict = self.gammaPoissonPredict
335       self.bins = (bins-1).to(self.device,non_blocking=self.use_cuda)
336     elif mix_dist == "PoissonLogNormal":
337       self.distPredict = self.poissonLogNormalPredict
338       self.bins = (bins-1).to(self.device,non_blocking=self.use_cuda)
339     elif mix_dist == "bivariateGammaPoisson":
340       self.distPredict = self.biGammaPoissonPredict
341
342       self.binsx = (bins[0]-1).to(self.device,non_blocking=self.use_cuda)
343       self.binsy = (bins[1]-1).to(self.device,non_blocking=self.use_cuda)
344       # create additional layers required for this distribution
345       self.p2 = nn.Linear(64,self.num_mix).to(self.device,non_blocking=self.use_cuda)
346       self.p3 = nn.Linear(64,self.num_mix).to(self.device,non_blocking=self.use_cuda)
347       self.p2_act = p_act[2]
348       self.p3_act = p_act[3]

```

```

349         self.bivariate = True
350     else:
351         raise ValueError("Invalid mixture distribution")
352
353     def forward(self, x):
354         # forward function
355         # initial dense layers
356         x = self.embedding1Act(self.embedding1(x))
357         x = self.embedding2Act(self.embedding2(x))
358         x = self.embedding3Act(self.embedding3(x))
359         x = self.embedding4Act(self.embedding4(x))
360         x = self.embedding5Act(self.embedding5(x))
361         # get predictions of parameters
362         if self.bivariate:
363             p0 = self.p0_act(self.p0(x))
364             p1 = self.p1_act(self.p1(x))
365             p2 = self.p2_act(self.p2(x))
366             p3 = self.p3_act(self.p3(x))
367             p = (p0,p1,p2,p3)
368         else:
369             p0 = self.p0_act(self.p0(x))
370             p1 = self.p1_act(self.p1(x))
371             p = (p0,p1)
372         # get weights
373         pr = self.aiAct(self.ai(x))
374
375         return p,pr
376     def normalPredict(self,x):
377
378         (m,s),p = x
379         y = 0.0
380         for i in self.mix: # for each dist in mixture evaluate cdf at bins
381             dist = Normal(m[:,i,None],s[:,i,None])
382             t = dist.cdf(self.bins)
383             y += (t[:,1:]-t[:,:-1]).clamp(0)*p[:,i,None]
384         return y
385     def poissonLogNormalPredict(self,x):
386
387         (m,s),p = x
388         y = 0.0
389         for i in self.mix:# for each dist in mixture evaluate cdf at bins
390             dist = poissonLogNormal(m[:,i,None],s[:,i,None])
391             t = dist.cdf(self.bins)
392             y += (t[:,1:]-t[:,:-1]).clamp(0)*p[:,i,None]
393         return y
394     def gammaPoissonPredict(self,x):
395
396         (alpha,beta),p = x
397         y = 0.0
398         for i in self.mix:# for each dist in mixture evaluate cdf at bins
399             dist = gammaPoisson(alpha[:,i,None],beta[:,i,None])
400             t = dist.cdf(self.bins)
401             y += (t[:,1:]-t[:,:-1]).clamp(0)*p[:,i,None]
402         return y
403     def biGammaPoissonPredict(self,x):
404
405         (alpha,beta,theta0,theta1),p = x
406         y = 0.0
407         for i in self.mix:# for each dist in mixture evaluate cdf at bins
408             dist = bivariateGammaPoisson(alpha[:,i,None],beta[:,i,None],theta0[:,i,None],theta1[:,i,None])
409             y += dist._cdf(self.binsx,self.binsy).clamp(0)*p[:,i,None,None]
410         return y
411
412     @torch.jit.export
413     def train_model(self,epochs=100,optim="Adam",data="Gene",loss_fn="H2",custom_data=None,model_name=None):
414         self._directory_maker() # create directory
415         # set consts
416         self.loss_fn = loss_fn
417         self.epochs = epochs
418         self.optim = optim
419
420         if custom_data is None:
421             # load data and split into train and test data
422             xall,yall=self.load_data(data=data)
423             x_train, x_test,y_train,y_test = train_test_split(xall,yall,train_size=0.8)
424         else:
425             # if custom data is given use that
426             x_train, x_test,y_train,y_test =
427                 custom_data["x_train"],custom_data["x_test"],custom_data["y_train"],custom_data["y_test"]
428             self.scalesM,self.scalesS = custom_data["scalesM"],custom_data["scalesS"]
429         # set name of model
430         if model_name == None:
431             self.model_name = f"mix_{self.mix_dist}_{loss_fn}_{optim}_{self.num_mix}"
432         else:
433             self.model_name=model_name
434         # create tensorboard writer
435         writer = SummaryWriter(comment=self.model_name)
436         # create optimiser
437         self._init_optim()
438         print("Loading dataset")
439         # cast from numpy to tensor

```

```

439 xtrainTensor=torch.from_numpy(x_train).type(torch.float32)
440 ytrainTensor=torch.from_numpy(y_train).type(torch.float32)
441
442 if self.mix_dist=="bivariateGammaPoisson":
443     # use no additional workers for loading in the bivariate case
444     # to stop the pytorch data loader loading too many batches and causing
445     # a memory overflow
446     num_workers = 0
447 else:
448     # for the univariate distribution use extra loaders as we dont
449     # get a memory overflow
450     num_workers = 2
451
452
453     # create datat loaders
454 dataset = torch.utils.data.TensorDataset(xtrainTensor,ytrainTensor)
455 loader = DataLoader(dataset, batch_size=self.batch_size,drop_last=True,shuffle=True,
456                     pin_memory=self.use_cuda,num_workers=num_workers)#self.use_cuda
457 xtestTensor=torch.from_numpy(x_test).type(torch.float32)
458 ytestTensor=torch.from_numpy(y_test).type(torch.float32)
459 valDataset = torch.utils.data.TensorDataset(xtestTensor,ytestTensor)
460 valLoader = DataLoader(valDataset, batch_size=self.batch_size,drop_last=True,shuffle=False,
461                     pin_memory=self.use_cuda,num_workers=num_workers)#self.use_cuda
462 old_test_loss = np.inf
463 train_batches = x_train.shape[0] // self.batch_size
464 test_batches = x_test.shape[0] // self.batch_size
465
466 # initialise loss arrays
467 self._init_loss()
468 del valDataset,dataset,xtestTensor,ytestTensor,x_train,x_test,y_train,y_test
469
470 print("Train Loop Starting...")
471 for i in range(self.epochs):
472     #nvtx.range_push("Epoch " + str(i+1))
473     self.train() # put in train mode
474     print(f"STEP: {i+1}/{self.epochs}")
475     with tqdm(total=train_batches,ascii=True,desc="Training") as pbar:#loading bar
476         for idx, (local_batch, local_labels) in enumerate(loader):
477             #nvtx.range_push("Batch "+str(idx))
478
479             #nvtx.range_push("Copying to device")
480             # get batches
481             y=local_labels.to(self.device,non_blocking=self.use_cuda)
482             x_batch=local_batch.to(self.device,non_blocking=self.use_cuda)
483             #nvtx.range_pop()
484             #nvtx.range_push("Forward pass")
485             # forward pass
486             x=self(x_batch)
487             #nvtx.range_pop()
488
489             #nvtx.range_push("Backward pass")
490             # get output of mixture for bins and update loss and back prop
491             self._update_loss_train(self.distPredict(x),y,i)
492
493             #nvtx.range_pop()
494             #nvtx.range_pop()
495
496
497             pbar.update(1)
498
499
500             #nvtx.range_pop()
501             # average
502             self.train_loss[i] = self.train_loss[i]/train_batches
503             if np.isnan(self.train_loss[i,0]):# break if we get nan
504                 print("GOT NaN")
505                 break
506             # update tensorboard log and print loss
507             for l,s in enumerate(self.loss_names):
508                 name=f'Loss {s}/train'
509                 writer.add_scalar(name, float(self.train_loss[i,1]), i)
510                 print(f"{name}: {self.train_loss[i,1]}")
511             # get val loss
512             with torch.no_grad(): # dont keep gradients
513                 self.eval()# set to eval
514                 with tqdm(total=test_batches,ascii=True,desc="Validating") as pbar:
515                     for local_batch, local_labels in valLoader:
516                         # get batch
517                         y = local_labels.to(self.device,non_blocking=self.use_cuda)
518                         x_batch=local_batch.to(self.device,non_blocking=self.use_cuda)
519                         # forward pass
520                         pred = self(x_batch)
521                         # update val loss
522                         self._update_loss_val(self.distPredict(pred),y,i)
523                         pbar.update(1)
524             #nvtx.range_pop()
525             # average
526             self.val_loss[i] = self.val_loss[i]/test_batches
527             # update tensorboard print val loss
528             for l,s in enumerate(self.loss_names):
529                 name=f'Loss {s}/test'

```

```

530     writer.add_scalar(name, float(self.val_loss[i,1]), i)
531     print(f"{{name}}: {self.val_loss[i,1]}")
532     #nvtx.range_pop()
533     # if achieves lowest val loss save
534     if self.val_loss[i,0] < old_test_loss:
535         #nvtx.mark("Saving best model")
536         torch.save({"state_dict":self.state_dict(),
537                     "scalesM":self.scalesM,
538                     "scalesS":self.scalesS,
539                     "epoch":i+1,
540                     "optim":self.optimiser,
541                     "val_loss":self.val_loss,
542                     "T": self.T,
543                     "input_dim": self.input_dim}, f"./best/mdBestModel_{self.model_name}.pth")
544     old_test_loss = self.val_loss[i,0]
545     # save final model
546     torch.save({"state_dict":self.state_dict(),
547                     "scalesM":self.scalesM,
548                     "scalesS":self.scalesS,
549                     "epoch":self.epochs,
550                     "optim":self.optimiser,
551                     "val_loss":self.val_loss,
552                     "T": self.T,
553                     "input_dim": self.input_dim}, f"./final/mdFinalModel_{self.model_name}.pth")
554     return self.train_loss, self.val_loss
555     @torch.jit.export
556     def predict(self,x,raw=True):
557         # get predictions
558         self.eval() # set to eval
559         # scale x
560         x = (np.atleast_2d(x) - self.scalesM)/self.scalesS
561         p = int(np.ceil(x.shape[0] /self.batch_size))
562         t = x.shape[0]
563         r = self.batch_size - (t % self.batch_size)
564         x = np.pad(x,((0,r),(0,0)))
565         xs = np.vsplit(x,p)
566         res = [None] * p
567         if raw: # return parameters
568             for i in range(p):# feed x in batch size chunks
569                 # forward
570                 a,p = self(torch.from_numpy(xs[i]).type(torch.float32).to(self.device,non_blocking=self.use_cuda))
571                 # detach and send to numpy
572                 p = p.detach().cpu().numpy()
573                 a = torch.cat(a,dim=-1).detach().cpu().numpy()
574                 res[i] = np.hstack((p,a))
575
576         else:# return bin probs
577             for i in range(p):# feed x in batch size chunks
578                 # forward
579                 x = self(torch.from_numpy(xs[i]).type(torch.float32).to(self.device,non_blocking=self.use_cuda))
580                 # get dist predictions
581                 x = self.distPredict(x)
582
583                 res[i] = x.detach().cpu().numpy()
584             res[-1] = res[-1][:-r]
585
586             return np.vstack(tuple(res))
587             @torch.jit.export
588             def load_model(self,PATH):
589                 # load model
590                 model = torch.load(PATH,map_location=self.device)
591                 self.scalesM = model["scalesM"]
592                 self.scalesS = model["scalesS"]
593
594                 self.load_state_dict(model["state_dict"])
595
596
597
598
599
600
601     if __name__=="__main__":
602         use_cuda = torch.cuda.is_available()
603         device = torch.device("cuda" if use_cuda else "cpu")
604         num_mix = [2*i for i in range(7)]
605         loss = ["MSE","L1","H2","MDDKLD","MDDEDIV"]
606         optim = ["Adam", "SGD"]
607         gau = torch.linspace(-15,15,100+1)
608         real = torch.linspace(0,300,100+1)
609         g = ParameterGrid({"loss":loss,"optim":optim})
610
611         val_loss = [[None,i] for i in g]
612         train_loss = [[None,i] for i in g]
613         for v,i in enumerate(g):
614             try:
615                 print(i)
616                 m = MixD(T=100,num_mix=6,input_dim=9,bins=gau,mix_dist="Gaussian").to(device)
617
618                 train_loss[v][0],val_loss[v][0] = m.train_model(epochs=200,data="Gaussian",loss_fn=i["loss"],optim=i["optim"])
619
620             except Exception as inst:

```

```

621     print("ERROR!!!!")
622     print(type(inst))
623     print(inst)
624     with open('train_loss.pickle', 'wb') as f:
625         pickle.dump(train_loss, f, protocol=pickle.HIGHEST_PROTOCOL)
626     with open('val_loss.pickle', 'wb') as f:
627         pickle.dump(val_loss, f, protocol=pickle.HIGHEST_PROTOCOL)
628
629

```

XII.C.2 Ensemble Mixture Model

```

1  from MIXD import *
2
3  class ensembleMIX(MixD):
4      """
5          model to make ensemble of n mixture models
6      """
7      def __init__(self,n,T,bins,num_mix,
8          mix_dist="Gaussian",batch_size=1024,
9          p_act=(Exponential(),Exponential(),Exponential(),Exponential()),
10         dense_act=nn.ELU(),learning_rate=1e-5):
11         super(model,self).__init__().__inherit
12         # set consts
13         self.use_cuda = torch.cuda.is_available()
14         self.device = torch.device("cuda" if self.use_cuda else "cpu")
15         self.n = n
16         self.T = T
17         self.bini = bins
18         self.batch_size=batch_size
19         self.e_learning_rate=learning_rate
20         self.dense_act=dense_act
21         self.num_mix = num_mix
22         const = 3 if mix_dist != "bivariateGammaPoisson" else 5
23         self.input_dim = const*n*self.num_mix
24         # create layers
25         self.embedding1 = nn.Linear(self.input_dim, self.input_dim*2).to(self.device,non_blocking=self.use_cuda)
26         self.embedding1Act = dense_act
27         self.embedding2 = nn.Linear(self.input_dim*2, self.input_dim*4).to(self.device,non_blocking=self.use_cuda)
28         self.embedding2Act = dense_act
29         self.embedding3 = nn.Linear(self.input_dim*4, self.input_dim*8).to(self.device,non_blocking=self.use_cuda)
30         self.embedding3Act = dense_act
31         self.embedding4 = nn.Linear(self.input_dim*8, self.input_dim*4).to(self.device,non_blocking=self.use_cuda)
32         self.embedding4Act = dense_act
33         self.embedding5 = nn.Linear(self.input_dim*4, self.input_dim).to(self.device,non_blocking=self.use_cuda)
34         self.embedding5Act = dense_act
35         self.p0_act = p_act[0]
36         self.p1_act = p_act[1]
37         self.bivariate = False
38
39
40         self.p0 = nn.Linear(self.input_dim,self.num_mix).to(self.device,non_blocking=self.use_cuda)
41         self.p1 = nn.Linear(self.input_dim,self.num_mix).to(self.device,non_blocking=self.use_cuda)
42
43         self.ai = nn.Linear(self.input_dim,self.num_mix)
44         self.aiAct = nn.Softmax(dim=1)
45         self.mix = torch.arange(0,self.num_mix,dtype=torch.int).to(self.device,non_blocking=self.use_cuda)
46         self.mix_dist = mix_dist
47         # set mixture distribution and bins
48         if mix_dist == "Gaussian":
49             self.distPredict = self.normalPredict
50             self.bins = bins.to(self.device,non_blocking=self.use_cuda)
51         elif mix_dist == "GammaPoisson":
52             self.distPredict = self.gammaPoissonPredict
53             self.bins = (bins-1).to(self.device,non_blocking=self.use_cuda)
54         elif mix_dist == "PoissonLogNormal":
55             self.distPredict = self.poissonLogNormalPredict
56             self.bins = (bins-1).to(self.device,non_blocking=self.use_cuda)
57         elif mix_dist == "bivariateGammaPoisson":
58             self.distPredict = self.biGammaPoissonPredict
59             self.binsx = (bins[0]-1).to(self.device,non_blocking=self.use_cuda)
60             self.binsy = (bins[1]-1).to(self.device,non_blocking=self.use_cuda)
61             # create additional layers required by this distribution
62             self.p2 = nn.Linear(self.input_dim,self.num_mix).to(self.device,non_blocking=self.use_cuda)
63             self.p3 = nn.Linear(self.input_dim,self.num_mix).to(self.device,non_blocking=self.use_cuda)
64             self.p2_act = p_act[2]
65             self.p3_act = p_act[3]
66             self.bivariate = True
67         else:
68             raise ValueError("Invalid mixture distribution")
69
70
71     def forward(self,x):
72         # forward method
73         # initial dense layers
74         x = self.embedding1Act(self.embedding1(x))
75         x = self.embedding2Act(self.embedding2(x))
76         x = self.embedding3Act(self.embedding3(x))

```

```

77     x = self.embedding4Act(self.embedding4(x))
78     x = self.embedding5Act(self.embedding5(x))
79     # predict parameters
80     if self.bivariate:
81         p0 = self.p0_act(self.p0(x))
82         p1 = self.p1_act(self.p1(x))
83         p2 = self.p2_act(self.p2(x))
84         p3 = self.p3_act(self.p3(x))
85         p = (p0,p1,p2,p3)
86     else:
87         p0 = self.p0_act(self.p0(x))
88         p1 = self.p1_act(self.p1(x))
89         p = (p0,p1)
90     # predict weights
91     pr = self.aiAct(self.ai(x))
92
93     return p,pr
94 @torch.jit.export
95 def train_model(self,epochs=100,optim="Adam",data="gau",loss_fn="H2"):
96     # create directory to store models
97     self._directory_maker()
98     # set consts
99     self.loss_fn = loss_fn
100    self.epochs = epochs
101    self.optim = optim
102    # get data
103    self.create_data(data)
104    self.model_name = f"mix_{self.mix_dist}_{loss_fn}_{optim}_{self.num_mix}"
105    # set tensorboard writer
106    writer = SummaryWriter(comment=self.model_name)
107    # initialise optimiser
108    self._init_optim()
109    print("Loading dataset")
110
111    # set number of workers for data loaders
112    if self.mix_dist=="bivariateGammaPoisson":
113        num_workers = 0
114    else:
115        num_workers = 2
116
117    # create data loaders
118    xtrainTensor=torch.from_numpy(self.x_train).type(torch.float32)
119    ytrainTensor=torch.from_numpy(self.y_train).type(torch.float32)
120    dataset = torch.utils.data.TensorDataset(xtrainTensor,ytrainTensor)
121    loader = DataLoader(dataset, batch_size=self.batch_size,drop_last=True,shuffle=True,
122                         pin_memory=self.use_cuda,num_workers=num_workers)
123    xtestTensor=torch.from_numpy(self.x_test).type(torch.float32)
124    ytestTensor=torch.from_numpy(self.y_test).type(torch.float32)
125    valDataset = torch.utils.data.TensorDataset(xtestTensor,ytestTensor)
126    valLoader = DataLoader(valDataset, batch_size=self.batch_size,drop_last=True,shuffle=False,
127                          pin_memory=self.use_cuda,num_workers=num_workers)
128    # create array to store loss
129    self._init_loss()
130    old_test_loss = np.inf
131    train_batches = self.x_train.shape[0] // self.batch_size
132    test_batches = self.x_test.shape[0] // self.batch_size
133    del valDataset,dataset,xtestTensor,ytestTensor,xtrainTensor,ytrainTensor
134    # start train loop
135    print("Train Loop Starting...")
136    for i in range(self.epochs):
137        #nvtz.range_push("Epoch " + str(i+1))
138        self.train()# set train mode
139        print(f"STEP: {i+1}/{self.epochs}")
140        with tqdm(total=train_batches,ascii=True,desc="Training") as pbar:# loading bar
141            for idx, (local_batch, local_labels) in enumerate(loader):
142                #nvtz.range_push("Batch "+str(idx))
143
144                #nvtz.range_push("Copying to device")
145                # load batch
146                y=local_labels.to(self.device,non_blocking=self.use_cuda)
147                x_batch=local_batch.to(self.device,non_blocking=self.use_cuda)
148                #nvtz.range_pop()
149                #nvtz.range_push("Forward pass")
150                # forward pass
151                x=self(x_batch)
152                #nvtz.range_pop()
153
154                #nvtz.range_push("Backward pass")
155                #evaluate mixture and update loss and back prop
156                self._update_loss_train(self.distPredict(x),y,i)
157
158
159                #nvtz.range_pop()
160                #nvtz.range_pop()
161
162
163
164                pbar.update(1)
165
166
167            #nvtz.range_pop()

```

```

168     # average
169     self.train_loss[i] = self.train_loss[i]/train_batches
170     if np.isnan(self.train_loss[i,0]):# if nan break
171         print("GOT NaN")
172         break
173     # write tensorboard log and print loss
174     for l,s in enumerate(self.loss_names):
175         name=f'Loss {s}/train'
176         writer.add_scalar(name, float(self.train_loss[i,1]), i)
177         print(f'{name}: {self.train_loss[i,1]}')
178     # get val loss
179     with torch.no_grad():
180         self.eval()# set eval mode
181         with tqdm(total=test_batches,ascii=True,desc="Validating") as pbar:# loading bar
182             for local_batch, local_labels in valLoader:
183                 # load batch
184                 y = local_labels.to(self.device,non_blocking=self.use_cuda)
185                 x_batch=local_batch.to(self.device,non_blocking=self.use_cuda)
186                 # forward pass
187                 pred = self(x_batch)
188                 # evaluate mixture and get loss
189                 self._update_loss_val(self.distPredict(pred),y,i)
190             pbar.update(1)
191     #nvtx.range_pop()
192     # average
193     self.val_loss[i] = self.val_loss[i]/test_batches
194     # update tensorbard and print loo
195     for l,s in enumerate(self.loss_names):
196         name=f'Loss {s}/test'
197         writer.add_scalar(name, float(self.val_loss[i,1]), i)
198         print(f'{name}: {self.val_loss[i,1]}')
199     #nvtx.range_pop()
200     # if it has the lowest val loss print
201     if self.val_loss[i,0] < old_test_loss:
202         #nvtx.mark("Saving best model")
203         torch.save({"state_dict":self.state_dict(),
204                     "scalesM":self.scalesM,
205                     "scalesS":self.scalesS,
206                     "epoch":i+1,
207                     "optim":self.optimiser,
208                     "val_loss":self.val_loss,
209                     "T": self.T,
210                     "input_dim": self.input_dim}, f"./best/mdEnsBestModel_{self.model_name}.pth")
211     old_test_loss = self.val_loss[i,0]
212     # save final model
213     torch.save({"state_dict":self.state_dict(),
214                     "scalesM":self.scalesM,
215                     "scalesS":self.scalesS,
216                     "epoch":self.epochs,
217                     "optim":self.optimiser,
218                     "val_loss":self.val_loss,
219                     "T": self.T,
220                     "input_dim": self.input_dim}, f"./final/mdRnsFinalModel_{self.model_name}.pth")
221     return self.train_loss,self.val_loss
222 def _load_model(self,i):
223     # load sub model and put in eval mode
224     m =
225     ↪ MixD(T=self.T,input_dim=self.inp,num_mix=self.num_mix,bins=self.bini,mix_dist=self.mix_dist,dense_act=nn.Tanh()).to(self.device)
226     m.load_model(f"mdBestModel_ensemble_model_{self.data}_{i+1}.pth")
227     m.eval()
228
229     return m
230
231 def _submodel_predict(self,x,n):
232     # get prediction of n submodels in series
233     out = [None] * n
234     for i in range(n):
235         out[i] = self._load_model(i).predict(x)
236     return np.hstack(tuple(out))
237
238 def create_data(self,data):
239     # set const
240     self.data = data
241     if data == "gau":
242         self.inp = 9
243     else:
244         self.inp = 5
245     # load data
246     with h5py.File(f"ensemble_data_{data}.hdf5","r") as f:
247         x_train = f["x_train"][:]
248         x_test = f["x_test"][:]
249         y_train = f["y_train"][:]
250         y_test = f["y_test"][:]
251
252         self.x_train,self.x_test,self.y_train,self.y_test
253         ↪ =self._submodel_predict(x_train,self.n),self._submodel_predict(x_test,self.n),y_train,y_test
254         xall = np.vstack((self.x_train,self.x_test))
255         self.scalesM = np.mean(xall,0)
256         self.scalesS = np.std(xall,0)
257         self.x_train,self.x_test = (self.x_train-self.scalesM)/self.scalesS,(self.x_test-self.scalesM)/self.scalesS

```

This is my own work except where otherwise stated