

**ISTANBUL TECHNICAL UNIVERSITY**

**FACULTY OF SCIENCE AND LETTERS**

**Advanced Physics Project Report**

**İTÜ**



**Machine Learning and Nonlinear Schrödinger Equation**

**Hüseyin Talha Şenyaşa**

**Department : Physics Engineering**

**Student ID : 090120132**

**Advisor : Assoc. Prof. A. Levent Subaşı**

**FALL 2017**

## Summary

We train an artificial neural network to estimate the ground state energy of a one-dimensional Bose-Einstein condensate in harmonic trapping potential. Such a system can be described by the solution of a non-linear Schrödinger equation also called a Gross-Pitaevskii equation. We also use the method for the inverse problem of predicting the non-linearity parameter using the ground state density profile for a given harmonic trapping potential.

# Contents

<b>1</b>	<b>Introduction and Motivation</b>	<b>1</b>
<b>2</b>	<b>Gross Pitaevskii Equation</b>	<b>1</b>
2.1	Numerical Solution and XMDS Framework . . . . .	4
<b>3</b>	<b>Problem Statement and Dataset Generation</b>	<b>5</b>
3.1	Dataset and Dataset Generation . . . . .	9
<b>4</b>	<b>Machine Learning for GPE</b>	<b>10</b>
4.1	Architecture . . . . .	10
4.2	Hyperparameters . . . . .	11
4.3	Training Results . . . . .	13
4.3.1	Non-interacting System . . . . .	13
4.3.2	Interacting Systems . . . . .	16
4.3.3	Inverse Problem Prediction of $g$ . . . . .	28
<b>5</b>	<b>Conclusion</b>	<b>30</b>
<b>A</b>	<b>APPENDIX A</b>	<b>32</b>
A.1	Comment on Python Codes . . . . .	32

# 1 Introduction and Motivation

Machine learning is a growing research area, not only its theoretical base but also its share in applications in different areas. One of the reason why machine learning is used increasingly today is that artificial neural networks used in machine learning can approximate any continuous function within desired accuracy. This means that if we take  $\epsilon > 0$  as desired accuracy,  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  output of the network and  $\mathbf{g}(\mathbf{x})$  real value of the function, then it is guaranteed that there exists a network that satisfies the relation  $|\mathbf{g}(\mathbf{x}) - \mathbf{f}(\mathbf{x})| < \epsilon$ . From this point of view, a process or calculation that can be thought of as a function can be represented by a corresponding neural network that mimics this function in desired accuracy [1]. With this in mind, neural networks have the potential to learn general functions and be exploited for their advantages. Approximate value of the quantity can be determined for different scenarios.

Many different kind of applications of machine learning have already been implemented in physics<sup>1</sup>. For example, In [3, 4] machine learning is applied to quantum many body problems. A machine learning method called Unsupervised Learning to detect patterns in big datasets is used for discovering phase transitions [5]. There are also developed techniques in machine learning inspired from physics such as quantum tensor networks [6]. Relation between physics and machine learning also caused foundation of a new branch called Quantum Machine Learning which aims to implement a quantum software to make machine learning faster than its classical version [7].

In [8], machine learning approach is applied to a 2D Schrödinger Equation with random potentials. The authors built a convolutional deep neural network, and trained it to predict ground state energy of the system under four different two dimensional confining potentials including some random potentials. Their study showed that machine learning is a promising alternative in electronic structure calculations of quantum systems. In our study inspired from the article mentioned above, we apply machine learning method to the nonlinear Schrödinger equation to predict ground state energy of a Bose-Einstein condensate at absolute zero temperature in one dimensional harmonic trapping potentials.

# 2 Gross Pitaevskii Equation

A Bose-Einstein Condensate (BEC) is described by Gross Pitaevskii Equation (GPE) also known as non-linear Schrodinger Equation (NLSE). General form of

---

<sup>1</sup>For a detailed list, see [2]

GPE in three dimension is given as,

$$i\hbar \frac{\partial \Psi}{\partial t} = \frac{-\hbar^2}{2m} \nabla^2 \Psi + V(\mathbf{r}, t) \Psi + g|\Psi|^2 \Psi \quad (1)$$

where  $\hbar$  is Planck constant,  $\mathbf{r}$  is position vector,  $t$  is time,  $\Psi(\mathbf{r}, t)$  is the wave function,  $m$  is mass,  $\nabla^2$  is the Laplacian operator,  $V$  is the potential,  $g$  is a constant interaction parameter and it is defined as

$$g = \frac{4\pi\hbar^2 a_s}{m} \quad (2)$$

where,  $a_s$  is the s wave scattering length. Non-linearity of GPE comes from the interaction term. If there is no interaction GPE reduces to the Schrödinger Equation (SE) and becomes a linear equation. By definition,  $g$  can be positive or negative. If  $g > 0$ , it represents repulsive interaction, and if  $g < 0$ , it represents attractive interactions[9].

In Eq. (1), the first term on the right side is the kinetic term, the second term represents the potential energy, and the last term is the interaction term of the Hamiltonian. The expectation value of the total energy can be written as;

$$\langle E \rangle = \int \Psi^* \hat{H} \Psi d^3 \mathbf{r} \quad (3)$$

$$E = \int \left( \frac{\hbar^2}{2m} |\nabla \Psi|^2 + V|\Psi|^2 + g|\Psi|^4 \right) d^3 \mathbf{r} \quad (4)$$

If potential is not function of time, then, total energy of the system is conserved.

We assume a time-independent harmonic trapping potential. In such a case, the potential is only function of position  $V(\mathbf{r})$  and it can be written as

$$V(x, y, z) = \frac{1}{2} m (\omega_x^2 x^2 + \omega_y^2 y^2 + \omega_z^2 z^2) \quad (5)$$

where  $\omega$  is the angular frequency. Then method of separation of variables can be applied and it can be obtained that there is a solution in stationary form. In this case the factorized wave function is,

$$\Psi(\mathbf{r}, t) = \psi(\mathbf{r}) e^{-i\mu t/\hbar} \quad (6)$$

Here,  $\mu$  is the chemical potential. Time independent GPE becomes,

$$\mu\psi = \frac{-\hbar^2}{2m} \nabla^2 \psi + V(\mathbf{r})\psi + g|\psi|^2\psi \quad (7)$$

BECs can also be studied in one dimension by making angular frequencies asymmetric such that  $\omega_x, \omega_y \gg \omega_z$  and keeping their energy order much greater

than condensate's energy  $\hbar(\omega_x\omega_y)^{1/2} \gg \mu$ . This enables us to confine the dynamics of the system in one dimension and it can be described with a corresponding one dimensional GPE equation [9]. To obtain one dimensional GPE one can rewrite wave function,

$$\psi(x, y, z, t) = \psi_z(z, t)\psi_x(x)\psi_y(y) \quad (8)$$

$\psi_x$  and  $\psi_y$  are the ground state solution of the transverse potential;

$$\begin{aligned}\psi_x(x) &= \frac{1}{(\pi l_x^2)^{1/4}} e^{-x^2/2l_x^2} \\ \psi_y(y) &= \frac{1}{(\pi l_y^2)^{1/4}} e^{-y^2/2l_y^2}\end{aligned} \quad (9)$$

where  $l_x = \sqrt{\hbar/m\omega_x}$  and  $l_y = \sqrt{\hbar/m\omega_y}$  are the harmonic oscillator lengths that related to the density distribution of the condensate. Normalization of the wave function is chosen such that  $\int \psi_x(x)dx = \int \psi_y(y)dy = 1$ , and  $\int \psi_z(z, t) = N$ . If we plug in the factorized wave function to Eq. (7), then, the equation becomes;

$$\mu'\psi_z = \frac{-\hbar^2}{2m} \frac{d^2\psi_z}{dz^2} + \frac{1}{2}m\omega_z^2 z^2 \psi_z + g'|\psi_z|^2\psi_z \quad (10)$$

where  $\mu'$  and  $g'$  one dimensional effective chemical potential and interaction strength respectively

$$\mu' = \mu - \frac{\hbar}{2}(\omega_x - \omega_y), \quad g' = \frac{g}{2\pi l_x l_y} \quad (11)$$

From now on, we refer to Eq. (10) as GPE.

It is convenient to work with dimensionless quantities, to make equation dimensionless, one can define  $z = a_0 \tilde{z}$ . In this case Eq. (10) becomes,

$$\mu\psi = \frac{-\hbar^2}{2ma_0^2} \frac{d^2\psi}{d\tilde{z}^2} + \frac{1}{2}m\omega^2 a_0^2 \tilde{z}^2 \psi + g|\psi|^2\psi \quad (12)$$

(Here we dropped primes). After that, setting  $\hbar\omega = \hbar^2/ma_0^2$ , we find  $a_0^2 = \hbar\omega/m$ . Plugging  $a_0^2$  in to the Eq. (12), we get

$$\mu\psi = -\frac{1}{2}m\omega^2 \frac{d^2\psi}{d\tilde{z}^2} + \frac{\hbar^2}{2m} \tilde{z}^2 \psi + g|\psi|^2\psi \quad (13)$$

If we divide Eq. (13) by  $\hbar\omega$  and define  $\tilde{\mu} = \mu/\hbar\omega$ , we obtain

$$\tilde{\mu}\psi = -\frac{1}{2\sqrt{a_0}} \frac{d^2\psi}{d\tilde{z}^2} + \frac{1}{2\sqrt{a_0}} \tilde{z}^2 \psi + \frac{g|\psi|^2\psi}{\hbar\omega} \quad (14)$$

From the normalization condition,

$$\int |\psi|^2 dz = \int |\psi|^2 a_0 d\tilde{z} = N \quad (15)$$

we can define  $\tilde{\psi} = (\sqrt{a_0}/\sqrt{N})\psi$ . The third term on the right in Eq. (14) becomes,

$$\frac{g|\psi|^2}{\hbar\omega} = \frac{ga_0}{N\hbar\omega} |\tilde{\psi}|^2 \quad (16)$$

Finally, defining

$$\frac{ga_0}{N\hbar\omega_z} = \tilde{g} \quad (17)$$

the dimensionless GPE can be written as,

$$\tilde{\mu}\tilde{\psi} = -\frac{1}{2} \frac{d^2\tilde{\psi}}{d\tilde{z}^2} + \frac{1}{2}\tilde{z}^2\tilde{\psi} + \tilde{g}|\tilde{\psi}|^2\tilde{\psi} \quad (18)$$

There is no known analytic solution of GPE for harmonic trapping potential except the case in which interaction parameter is zero. In this case, GPE reduces to the SE and solution is well known with the ground state energy;

$$\tilde{\mu} = \frac{1}{2} \quad (19)$$

we used this case to check our numerical solutions.

## 2.1 Numerical Solution and XMDS Framework

In the previous section we stated that there is no general analytic solution of GPE. It is known for only few cases such as for a uniform condensate  $V = 0$  constant. Most of the time GPE is solved with numerical techniques or approximations. In this study, we used a framework called XMDS [10], implemented specifically to solve differential equation systems with well optimized numerical methods. In this framework the partial differential equations system can be described by a markup language called XML. When equation system is declared properly, XMDS produces a source code written in C++ that solves the equation with pre-specified numerical method. In our program, XMDS takes angular frequency, shift of equilibrium point, value of interaction parameter externally. It generates trapping potential and other quantities internally with supplied expressions.

Cross-check of the framework has done in two ways. First, since there exists an analytical solution for  $g = 0$ , we compared the results of the XMDS with analytical ones, and we also used another program here [11], and compared few results for situations where analytical solution does not exist.

### 3 Problem Statement and Dataset Generation

GPE can be solved with stated methods above. These methods must be applied every time when there is a change in parameters of the equation such as applying different trapping potential or changing interaction parameter. For example, when interaction parameter is changed, if there is no analytic solution, numerical method must be reapplied. With machine learning these steps can be encoded in network and information about the system can be obtained instantly with only one time cost which is training process of the network. In[8], authors apply deep learning method to the Schrödinger Equation to obtain ground state energy of the system under different potentials.

We try to built an artificial neural network to predict the ground state energy of a BEC for a given harmonic trapping potential and interaction parameter. We also try the inverse problem which is the prediction of interaction parameter for a given potential and density.

An artificial neural network is made of layers and these layers contain simple units called neurons. These neurons were inspired by the neurons in the human brain. A neuron takes one or more input and generates an output. To do this, it uses its internal variables which are called weights and biases. In general, the number of weights in a neuron is equal to the size of the input and there is one bias value. There is no restriction on the range of weights and biases except computational, they can even take on complex values [12]. Input output relation of a neuron can be described in the following way; let  $\mathbf{x} = [x_1, x_2, x_3, \dots, x_n]$  be the input of the neuron and  $f$  be the function that describes behavior of the neuron.

$$f(\mathbf{x}, \boldsymbol{\omega}, b) = \sum_{i=1}^n \omega_i x_i + b \quad (20)$$

Here,  $\omega_i$  are the weights and  $b$  is the bias. A layer of a network involves neurons that implement with this function. In our example,  $x_n$  inputs represent the potential expression respect to the position or represents density for the inverse problem. The range of the function  $f$  is infinite and it can be used directly. However, there are various ways to interpret result of the function  $f$ . The most primitive version of interpretation is sending the result of  $f$  as an argument to unit step function to generate 0 (False) or 1 (True). In this case,

$$a(f) = u_1(f(\mathbf{x}, \boldsymbol{\omega}, b)) \quad (21)$$

where  $a$  is called as the activation function of the neuron and it is a unit step function in this case. However, the usage of the step function introduces discontinuity and does not allow to use differential methods to calculate behavior of the

neurons when a small change in weights or biases is applied. Therefore, a more common activation function is sigmoid function denoted by  $\sigma$  which is continuous version of the step function and it is defined as,

$$\sigma(f) = \frac{1}{1 + e^{-f}} \quad (22)$$

If a neuron defined in this way such that its output given by Eq. (22), then it is called a sigmoid neuron. In the former case which is the step function, the neuron is called as perceptron.

If we denote  $j^{th}$  neuron in the  $l^{th}$  layer with  $f_j^l$ , and the weight from  $k^{th}$  neuron in the  $(l-1)^{th}$  layer to the  $j^{th}$  neuron in the layer  $l^{th}$  with  $\omega_{jk}^l$ , finally the bias of the  $j^{th}$  neuron in the  $l^{th}$  layer with  $b_j^l$ , then, the output of the  $j^{th}$  neuron in the  $l^{th}$  layer will be;

$$f_j^l(\sigma(f_k^{l-1}), \boldsymbol{\omega}, b) = \sum_{k=1}^n \omega_{jk}^l \sigma(f_k^{l-1}) + b_j^l \quad (23)$$

and the activation of the output can be written as;

$$\sigma(f_j^l(\sigma(f_k^{l-1}), \boldsymbol{\omega}, b)) \quad (24)$$

Here we use sigmoid function as activation function but it can be any suitable activation function such as hyperbolic tangent  $\tanh(f)$  or Rectified Linear Unit (ReLU).

Connection complexity of the neurons may vary. Output of each neuron in one layer can be input of every neuron in next layer. Such networks are called Fully Connected Networks (FCN) or multilayer perceptrons (MLP)[1]. If the result of one layer is directly sent to the next layer without any circulation or feedback, then the network is called as Feed Forward Network [1].

We are going to change these weights and biases in such a way that the difference between output generated by network and real value of the corresponding quantity will be minimized. To do that, we are going to use a proxy relation between output and real value which is called as Cost Function [1]. There are different kind of cost functions such as quadratic cost function, cross entropy cost function etc. We are going to use quadratic cost function also called as mean squared error to show the general mechanism and to introduce few notions that directly effects the behavior of the network.

The quadratic cost function is defined as;

$$C(\boldsymbol{\omega}, b) = \frac{1}{2n} \sum_x \|\mathbf{y} - a(f_L(\mathbf{x}))\|^2 \quad (25)$$

where  $n$ ,  $\omega$ ,  $b$  are number of examples in the training set, weights and biases, respectively. The subscript  $L$  indicates the last layer which is the output, so  $f_L(\mathbf{x})$  output produced by network and  $\mathbf{y}$  is real value of the quantity. To minimize  $C$ , we can take derivative of the function and can find extremum values, but this method is extremely costly because the total number of neurons in the network is enormous [1]. This problem can be overcame by an iterative algorithm called Gradient Descent. Since  $C$  is also a scalar field the algorithm tries to determine a direction which points to the decrement and moves in that direction with small steps. In each iteration, algorithm repeats itself to reach minimum value.

A small displacement in arbitrary direction which corresponds to a small change in weight or bias or both can be written as;

$$\Delta C = \frac{\partial C}{\partial \omega} \Delta \omega + \frac{\partial C}{\partial b} \Delta b \quad (26)$$

This expression gives information about what happens when small displacement is made. In order to determine direction of decrement we can rewrite this expression in terms of the gradient since it points to direction of maximum rate of increase.

We can define gradient operator as;

$$\nabla C = \left( \frac{\partial C}{\partial \omega}, \frac{\partial C}{\partial b} \right)^T \quad (27)$$

In this case, Eq. (26) can be written as;

$$\Delta C = \nabla C \cdot \Delta \mathbf{l} \quad (28)$$

where  $\mathbf{l} = (\omega, b)^T$ . We want to minimize the cost function, thus; in each iteration value of the  $C$  must be smaller than the previous one which means that  $\Delta C$  must be smaller than zero. Therefore, the direction of displacement must be in the opposite direction of the gradient. Rather than calculating the opposite direction for each iteration we can define  $\Delta \mathbf{l}$  in such a way so that it always points to opposite direction of the gradient.

$$\Delta \mathbf{l} = -\eta \nabla C \quad (29)$$

In this case, Eq. (28) becomes;

$$\Delta C = -\eta \|\nabla C\|^2 \quad (30)$$

where  $\eta$  is positive definite number and it is called the learning rate. Since it is guaranteed that  $\|\nabla C\|^2 \geq 0$ , therefore;  $\Delta C \leq 0$ . In this situation, learning rate

is the step size in each iteration. From Eq. (30) it is intuitive that  $\eta$  can not be a large number because in such a case, the algorithm can miss the minimum point. This situation also brings up the subject that gradient descent does not give guarantee to find the minimum point.

Gradient Descent algorithm has an statistical version called Stochastic Gradient Descent to increase the speed of training process. It is assumed that gradient of  $m$  randomly selected examples in the dataset is nearly equal to gradient of the whole dataset and it can be expressed as;

$$\frac{1}{m} \sum_{k=1}^m \nabla C_{X_j} \approx \frac{1}{n} \sum_{x=1}^n \nabla C_x = \nabla C \quad (31)$$

where  $\nabla C_x$  is the gradient of a single training input and  $m$  is also called as mini-batch size.

It is worth noting that approaching a minimum value is proportional to the number of iterations which equals to number of examples in the dataset, therefore; it is a corollary that there can be two situations<sup>2</sup>. First, the algorithm cannot reach the minimum point due to lack of training examples, and secondly, the algorithm reaches a global or local minimum point and starts to oscillate around it[13]. There are mechanisms to prevent such cases and other optimization problems like direction sensitivity and they are called as adaptive learning rate [8]. In our network we used an algorithm called Adam [14] provided by the framework we used.

The final step in the gradient descent is updating the weights and biases. We use Eq. (30) to perform this task. If we combine Eq. (30) with Eq. (31) and write it in open form, the expression to update the weight and bias can be written as;

$$\begin{aligned} \omega' &= \omega - \frac{\eta}{m} \sum_{j=1}^m \frac{\partial C_{X_j}}{\partial \omega} \\ b' &= b - \frac{\eta}{m} \sum_{j=1}^m \frac{\partial C_{X_j}}{\partial b} \end{aligned} \quad (32)$$

Numerical calculation of the partial derivatives in the equation is one of the costliest computational operation in the training process. To reduce computational cost, there is a well known algorithm called backpropagation [15]. We use Pytorch Framework in our study and it provides similar mechanism called automatic differentiation.

In this section, we reviewd the general concept in artificial neural networks

---

<sup>2</sup>The third case is that algorithm may diverge but we ignore this situation here.

and introduced notions we are going to use in our example.

### 3.1 Dataset and Dataset Generation

We solved GPE for four different interaction parameters, therefore; there are four different corresponding datasets and each of them contains 10000 elements. We used 8500 of them as training examples and 1500 of them for testing if otherwise is specified. An element of a dataset involves an array containing potential values as a function of position and an another array containing interaction, kinetic, potential and total energy values respectively.

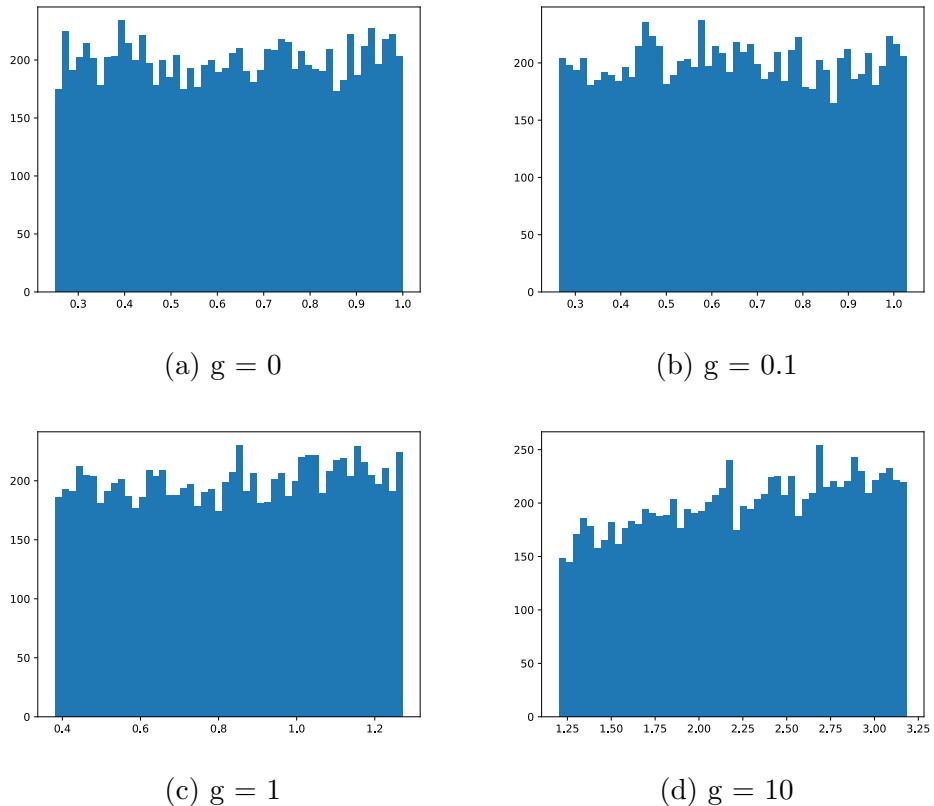


Figure 1: Total energy distributions of the generated solution for different interaction parameter values.

The array containing potential energy is generated according to harmonic trap expression in XMDS, thus; we only supplied angular frequency and shift of equilibrium to the XMDS randomly within predetermined limits. The dimensionless angular frequency can take values between 0.5 and 2. Shift of the equilibrium point is determined due to boundary conditions since density must be zero at the boundaries. In the numerical solution of the GPE, domain of the  $z$  dimension is between -10 to 10. Taking the maximum magnitude of the shift as  $\mp 5$  is enough to ensure that density function goes to zero at infinity. In both angular frequency

and shift, we used uniformly generated random numbers. Total energy distribution is given in Fig 1. Another type distribution is not required in our problem since the potential type is fixed but there are examples such as here [8], which different distribution is used.

## 4 Machine Learning for GPE

We used Pytorch Framework [16] to build our neural networks. It allows the client codes to work on both CPU and GPU via its internal python object called Tensor. If any CUDA supported GPU is available then Pytorch can use GPU without any change in the code. Code have three main parts, first one is dataloaders; it reads train and test data for specified interaction parameter from corresponding file and generates tensor dataset object. In this part, manipulation on dataset can be applied such as normalization, shuffling etc. Second part is implementation of the network. Architecture of the network is represented with a python class inherited from Pytorch's base class for networks called Module. This class also includes a forward method which is responsible for how data will be sent to the next layer. The last part is training and testing. In the training part, the network is iterated with training dataset and loss (result of cost function) is calculated at the end of each iteration. After that, weights and bias are updated accordingly.

Our aim is to show that the machine learning methods can be used to bypass GPE equation and physical features of the BEC can be determined directly. To do that, we implemented two different types of neural network; fully connected network (FCN) and convolutional neural network (CNN). CNN is more complex than FCN and it is developed to handle data which have grid-like topology [15]. Since our dataset involves translation, we expect that CNN can handle such a situation better than FCN.

### 4.1 Architecture

FCN involves 128 input neurons as input layer, next layer is the first hidden layer with 30 neurons, the second is another hidden layer same as the first, the next one is last hidden layer with 10 neurons and the last layer is the output layer. Totally there are 5 layers in our FCN and it will be denoted as FCN[128, 30, 30, 10, 4]. The most common activation function in feedforward networks rectified linear unit (ReLU) [8], is used for each forward feed except output. No operation is applied to the output neuron. Learning rate of the FCN is fixed and it is 0.001. Cost function is mean squared error (MSE) and optimization is

done with Adam.

CNN has two convolution layers, two maxpool layers and three fully connected layer and the last layer of the fully connected part is the output layer. Maxpooling is applied to output of the first and second convolution layers. ReLu is also applied each forward except output neuron same as in the FCN. Fully connected part of the CNN is FCN[310, 100, 20, 1]. Learning rate, cost function and optimizer of the CNN is same as FCN which are 0.001, MSE and Adam respectively.

## 4.2 Hyperparameters

We started with small dataset which has 800 elements for training 200 for test to see response of the architecture and to obtain a range for learning rate. The size of this dataset is small but it can give an idea about the range of learning rate since stochastic gradient descent is a statistical approach. We first set learning rate,  $\eta$  to 3 and start to decrease it, until 0.01 prediction of the network was extremely poor. For lower values, predictions were more encouraging. In this section, we show the effect of learning rate for two different architectures and interaction parameters.

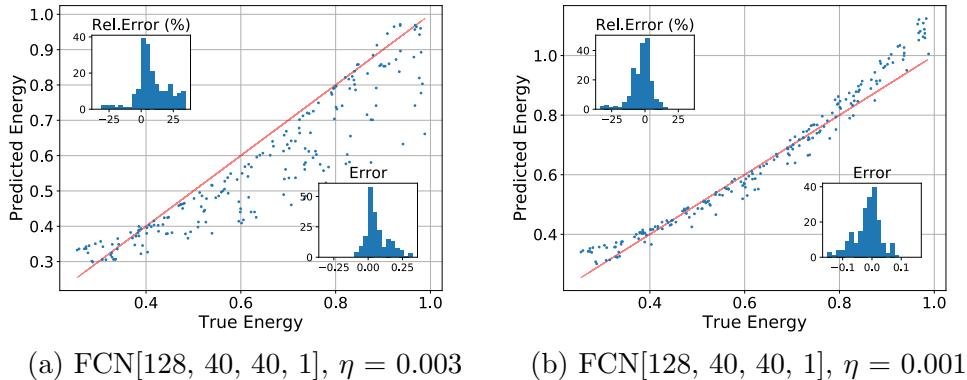


Figure 2: True Energy of Ground State vs Predicted Ground State Energy. Their hyperparameters are identical except the learning rate. Total number of epoch is 20 and batch size is 10. In this example interaction parameter  $g$ , is zero. Even the precision of the learning rate 0.001 seems higher, it is trivial to expect that change in interaction parameter will effect the result.

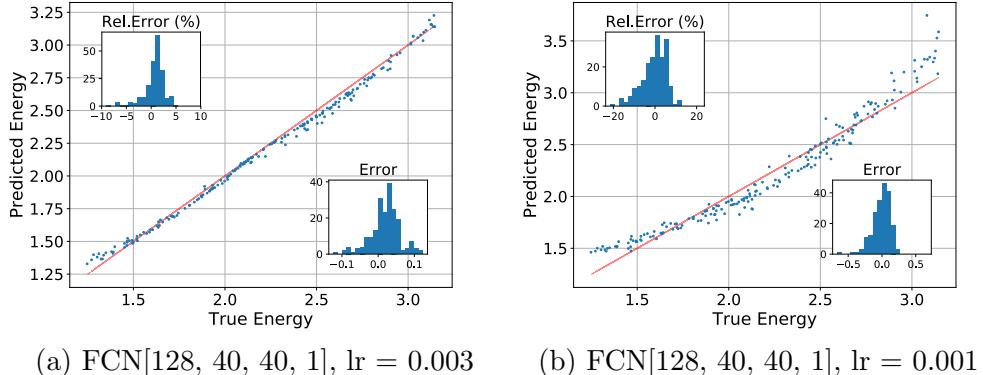


Figure 3: Here the same network with different interaction parameter which is 10.

It is clear that precision of the network with learning rate 0.003 is higher than 0.001. Of course, dramatic increase in the training dataset length may affect them both to converge same precision but our intention here to show that small change in learning rate causes different results.

After determination of learning rate we added extra one hidden layer to the network and tried some different number of neurons combination in the layers.

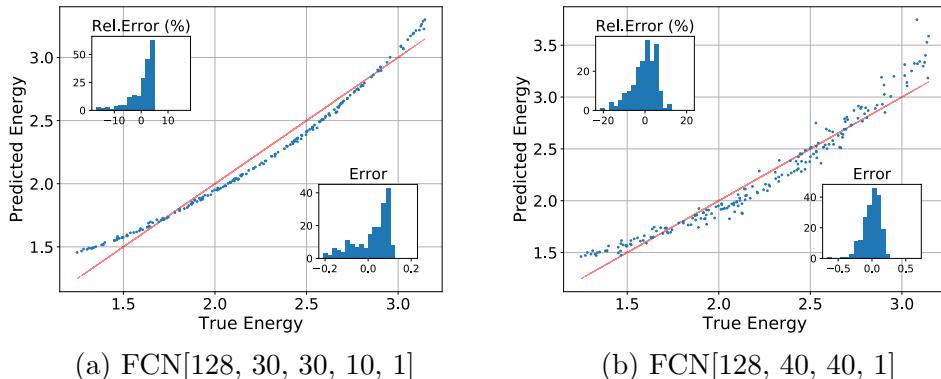


Figure 4: Interaction paramater  $g$ , is 10. Batch size is 10 and total number of epoch is 20 for this example. Precision of the network of 5 layer (a) is higher than 4 layer (b). Bias in (a) can be eliminated by increasing training dataset length.

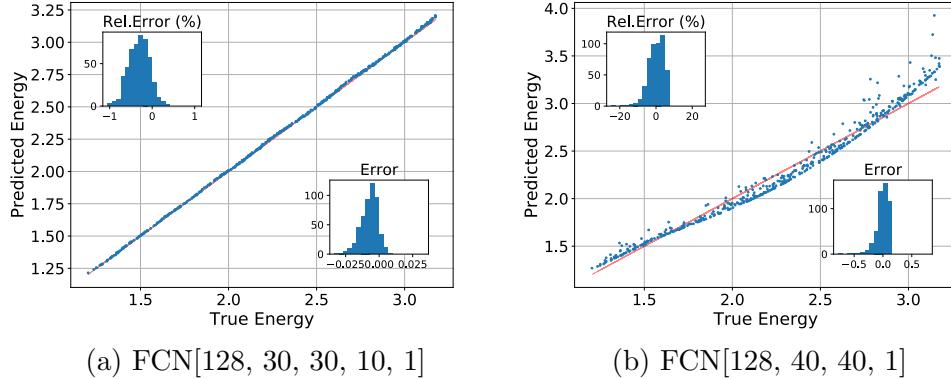


Figure 5: Total training dataset length is increased to 3500 and test dataset length is 500. Batch size, Learning rate are same as example given in Fig. 4 but the epoch is 30. The bias is eliminated. Precision and accuracy of the network of 5 layers has improved.

In conclusion, we chose learning rate  $\eta$  as 0.001 and started the training process. The other two hyperparameters epoch and mini batch size are determined in the training process. After number of different tries, we set the total number of epochs to 60, and mini batch size to 30.

### 4.3 Training Results

We give the result of total energy predictions per 20 epochs for both FCN and CNN also to show how predictions change. Results of FCN[128, 30, 30, 10, 4] used in prediction of interaction, potential, kinetic and total energy separately are given directly. For the inverse problem, three different scenarios are presented. All trainings are done with identical networks and only the training and test datasets are changed.

### 4.3.1 Non-interacting System

In non-interacting systems, GPE reduces to SE which does not involve non-linear term.

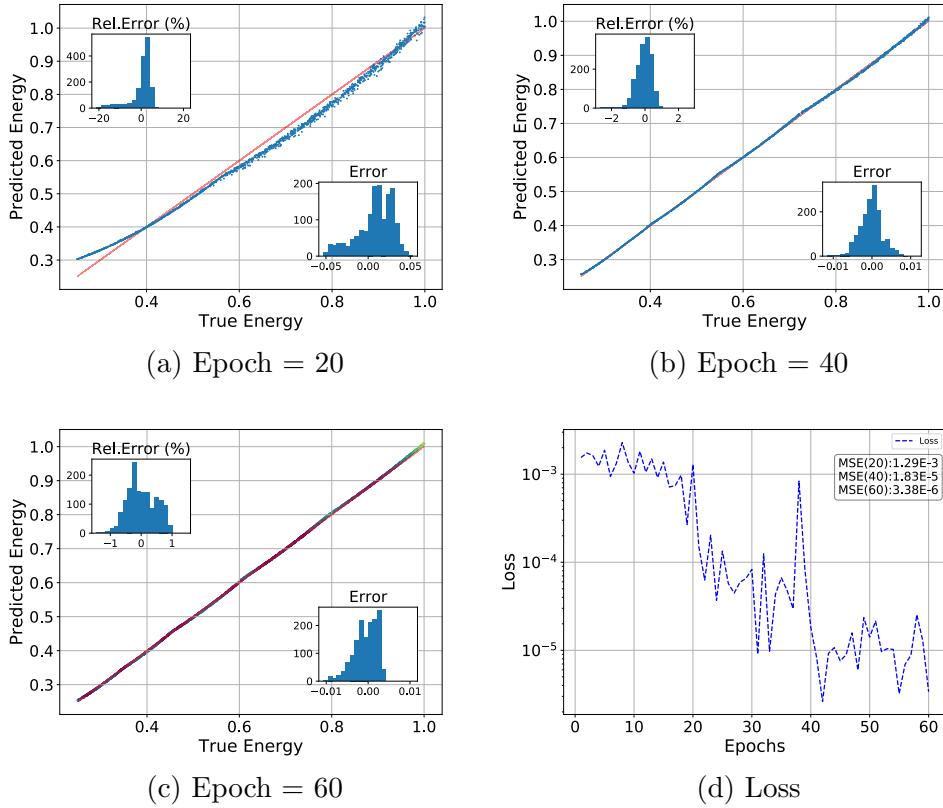


Figure 6: Result of FCN[128, 30, 30, 10, 1] for non-interacting system.

In figures,  $x$  axis represents true dimensionless total energy values, and  $y$  is predicted energy by trained network. Inset histogram at the left upper corner is relative error given as percentage and the inset histogram at the right corner is difference error.

Since the system has no interaction parameter, problem is relatively easy when it is compared to systems that involve interaction. In Fig. 6(a) it is obvious that the network requires more training examples and it reaches a satisfactory level in Fig. 6(b). Supplying another 20 more epochs reduces both relative error and MSE as shown in Fig. 6(d)

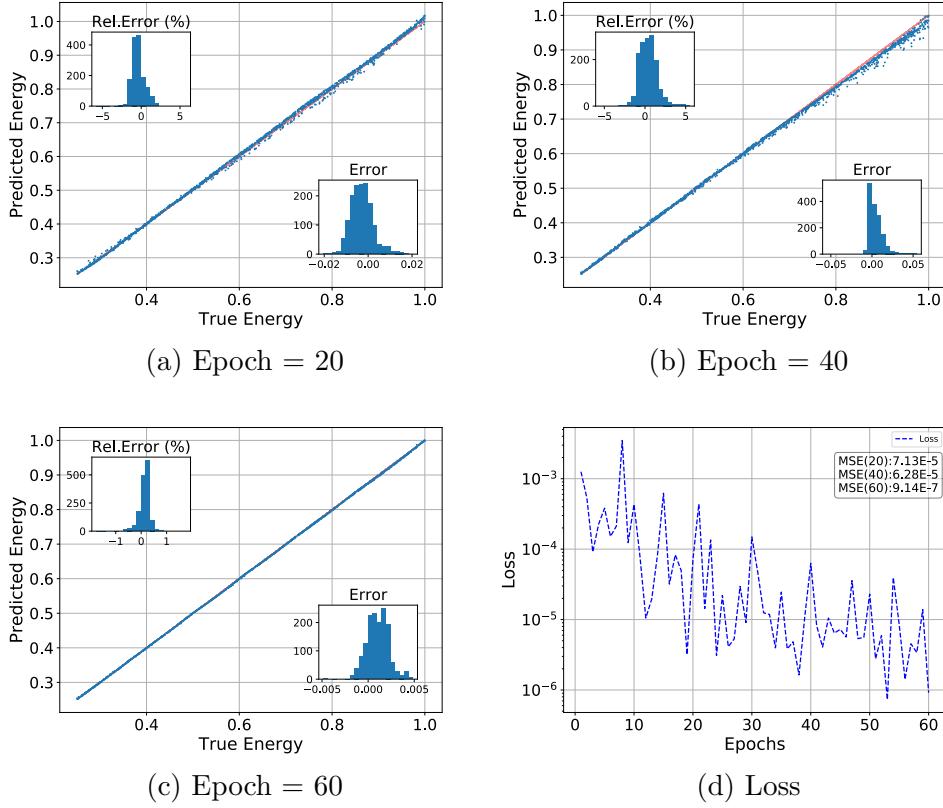


Figure 7: CNN results for  $g = 0$

The bias in Fig. 6(a) does not appear in Fig. 7(a). One reason why there is less or no bias in CNN predictions could be because CNN's have the ability to handle translations in dataset. In our dataset, translation corresponds to shift of equilibrium point.

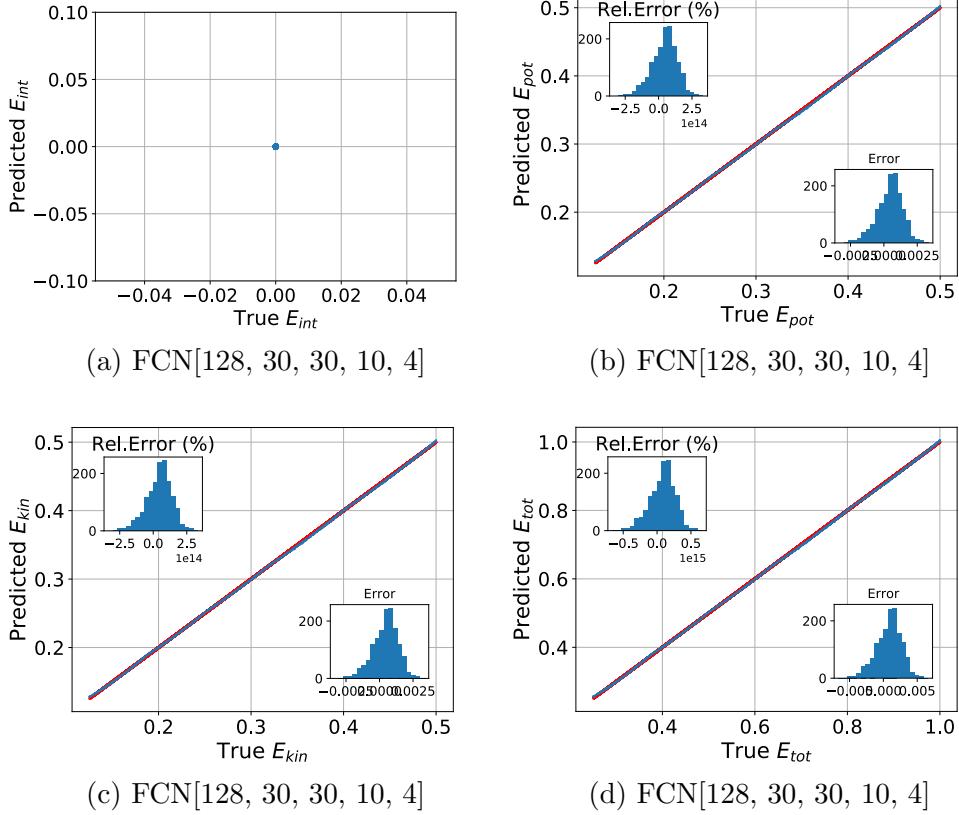


Figure 8: FCN[128, 30, 30, 10, 4] results for  $g = 0$ . Here (a) is interaction energy, (b) potential energy, (c) kinetic energy and (d) total energy predictions.

Only number of neurons in the output layer is increased to 4, and still network predicts total energy of the system within same sensitivity as former FCN network.

#### 4.3.2 Interacting Systems

In interacting systems, values of interaction parameter  $g$  are 0.1, 1, 10, 20 and the results are given in this order. In this stage, the values of the interaction parameter  $g$  can be thought directly as a coefficient of a nonlinear partial differential equation rather than its physical meaning. In  $g = 1$  scenario, contribution of interaction parameter to Eq. (18) will be more vivid. (We do not take into account the attractive interaction case which is represented by negative  $g$ . In such a situation, BEC collapses [9].)

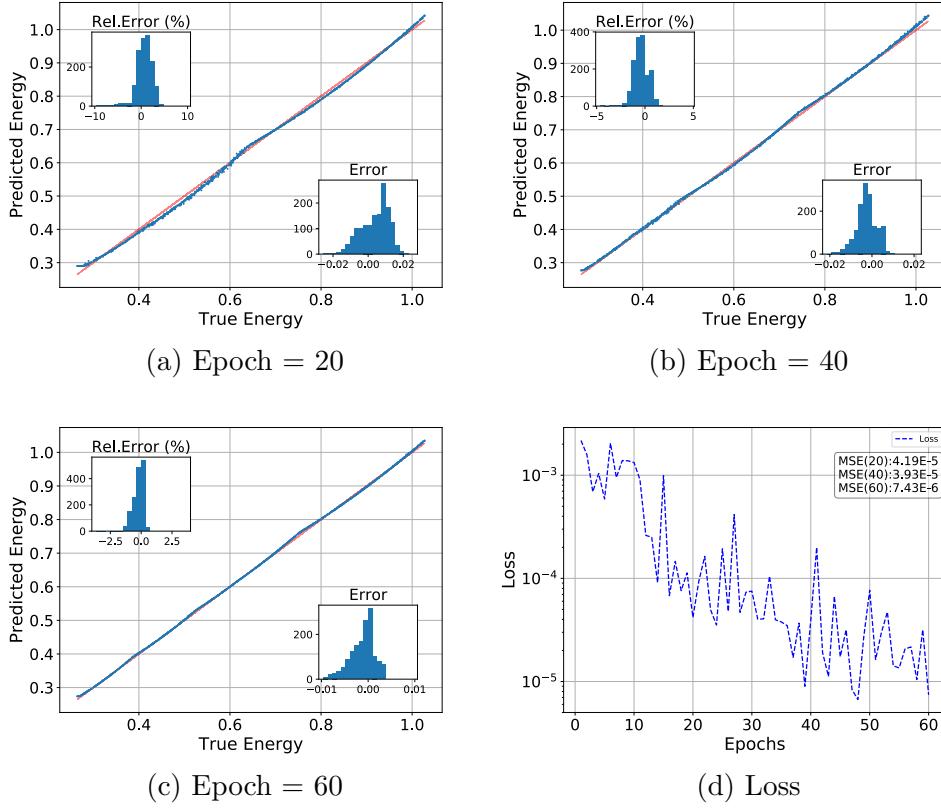


Figure 9: FCN[128, 30, 30, 10, 1] results for  $g = 0.1$ .

In first 20 epoch Fig. 9(a), there are relatively higher deviations in lower energies. In Fig. 9(b), deviations become smaller but still visible. In Fig. 9(c), prediction line is smoother but still deviations remains in lower energy levels.

Although the accuracy of the network is high, loss figure Fig. 9(d) still has fluctuation.

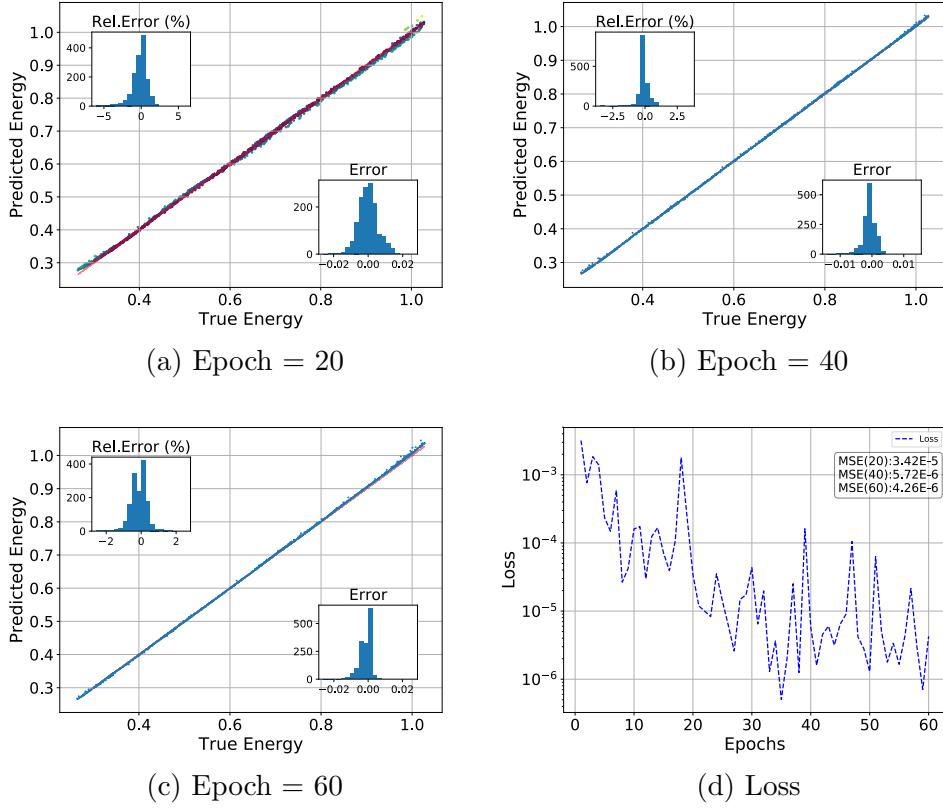


Figure 10: CNN results for  $g = 0.1$

In Fig. 10(a) deviation has nearly same characteristic with FCN in lower energies. After 60 epochs Fig. 10(c), there is a visible shift in higher energies but this does not occur in Fig. 10(b). Local peaks can be seen on Fig. 10(d) which increases error. Different techniques may be applied to recover the state at Fig. 10(b). For instance, before final test, network can check itself and compare different losses at different epochs and can select the minimum case.

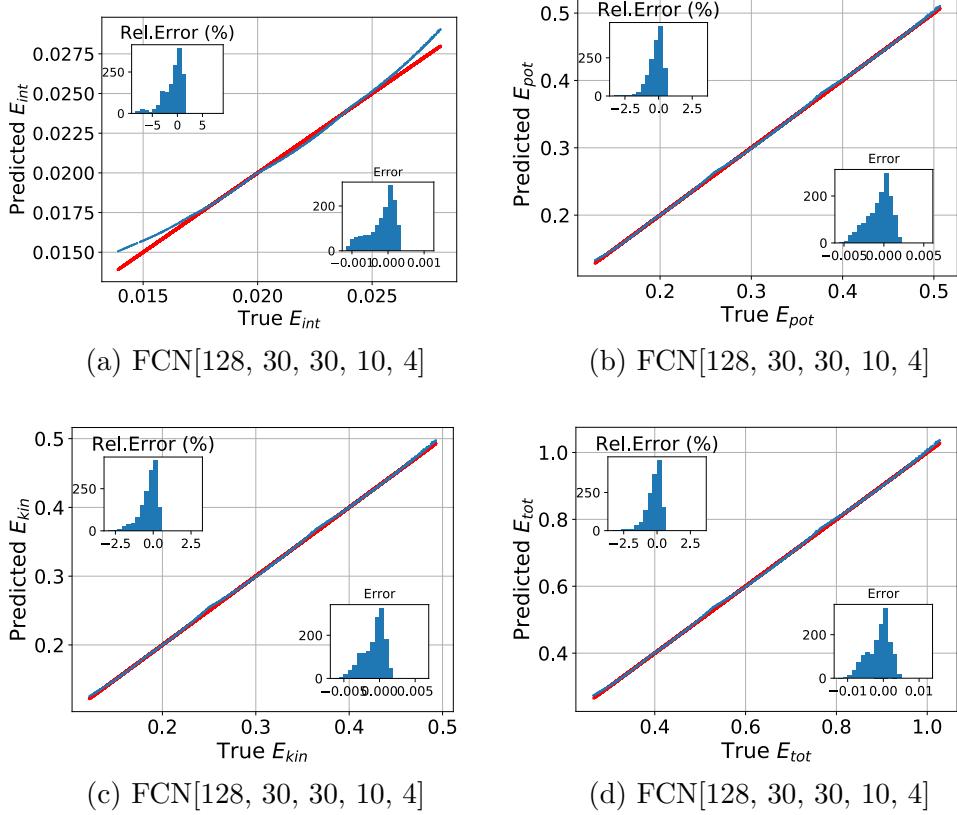


Figure 11: Separate Energy Predictions for  $g = 0.1$

Output layer of the network involves four layer and since our network is fully connected, some of the weights at previous layers are shared indirectly. The order of interaction energy is much smaller than the others. This is the reason of bias in Fig 11(a). Because, the order of incoming inputs are same as other energy levels and network has not enough example to reduce weights and bias values in order to increase prediction accuracy. Applying a proper normalization to the variables may reduce the error. In this time, incoming inputs will be in the same order.

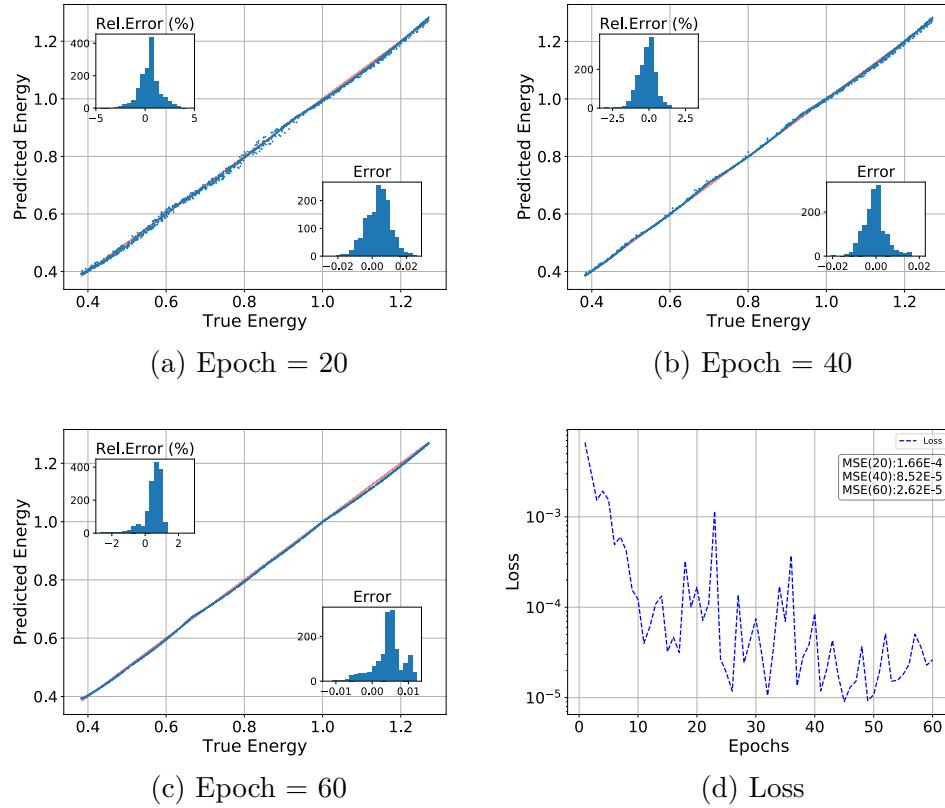


Figure 12: FCN[128, 30, 30, 10, 1] results for  $g = 1$

Relative error gets smaller after each 20 epoch and network goes into saturation after 40 epoch.

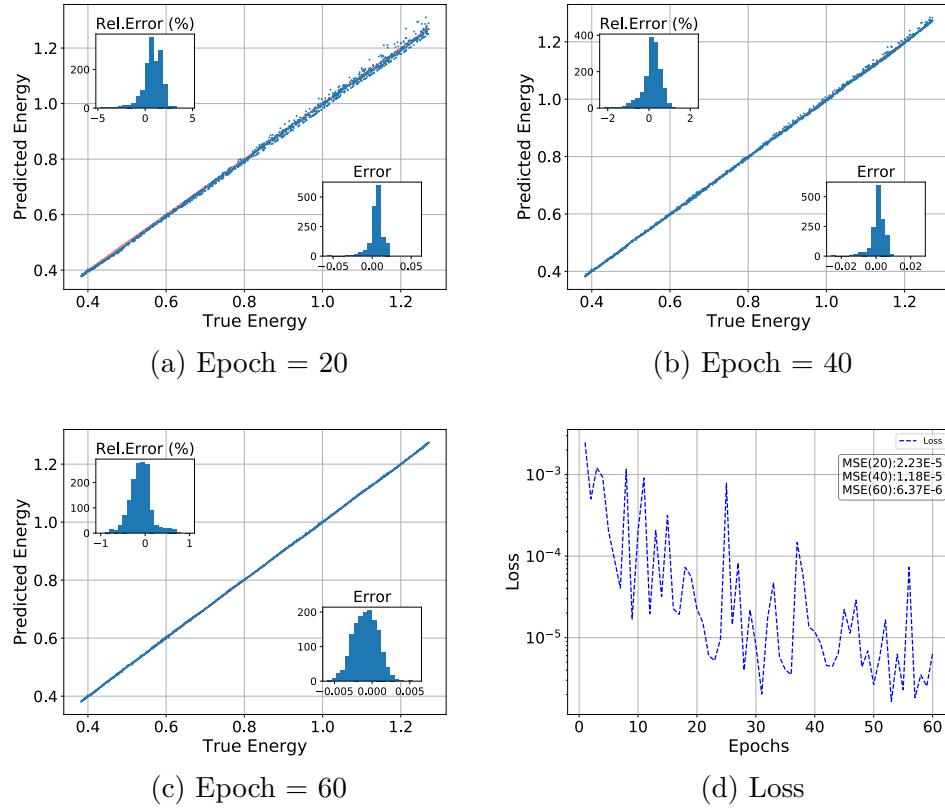


Figure 13: CNN results for  $g = 1$

Even the loss of the CNN decreases slower than FCN Fig. 12(d) and involves local peaks, its relative error at epoch 40 and 60 is smaller.

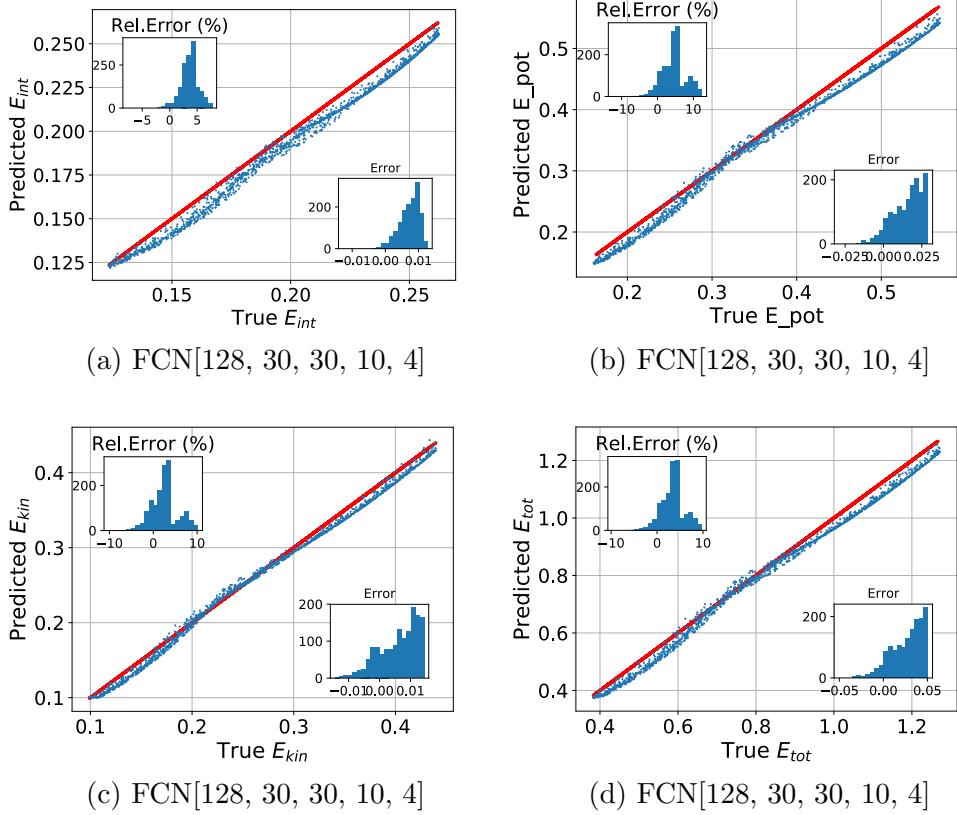


Figure 14: FCN[128, 30, 30, 10, 4], Separate energy predictions for  $g = 1$ .

This is the worst case in our training results. In first sight, it can be thought that a proper normalization may reduce the error and increase the accuracy but the difference between energy values are not as higher as in the case  $g = 0.1$ , so the problem is not due to normalization. Here, contribution of interaction parameter to Eq. (18) is in same order with other terms. So, it can be said that the degree of freedom of the network is not enough to handle such a situation.

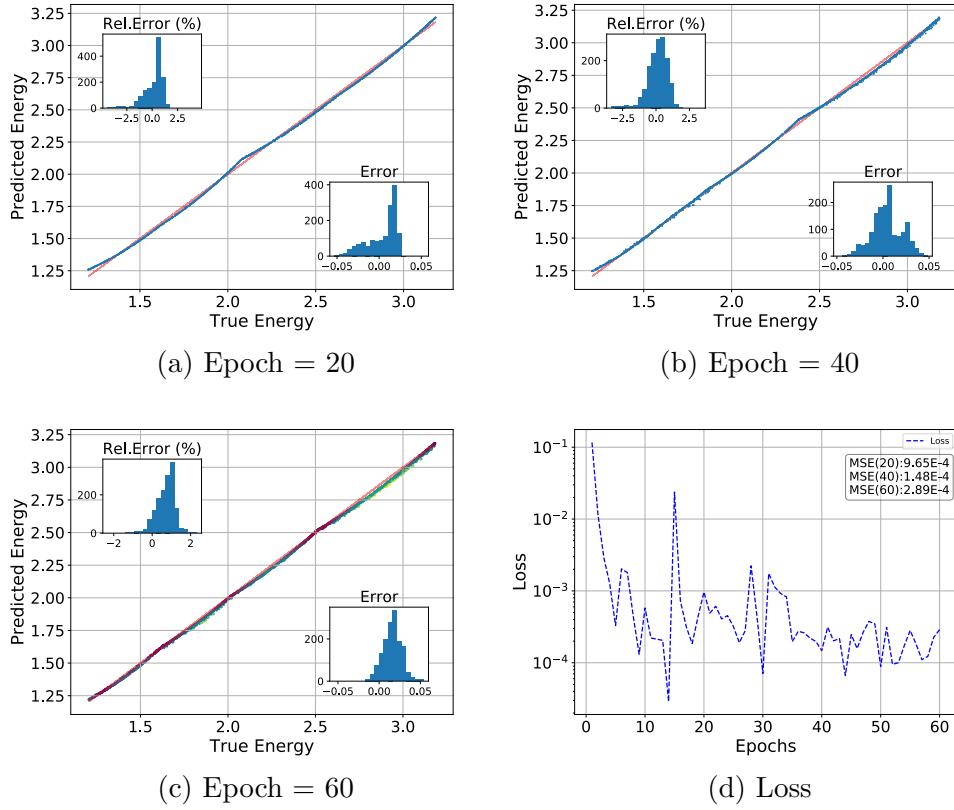


Figure 15: FCN[128, 30, 30, 10, 1] predictions for  $g = 10$

In Fig. 15(a) same deviation occurs in lower and higher energies as in previous results. Loss of the network decreases rapidly but order of it greater than the case of which  $g = 1$

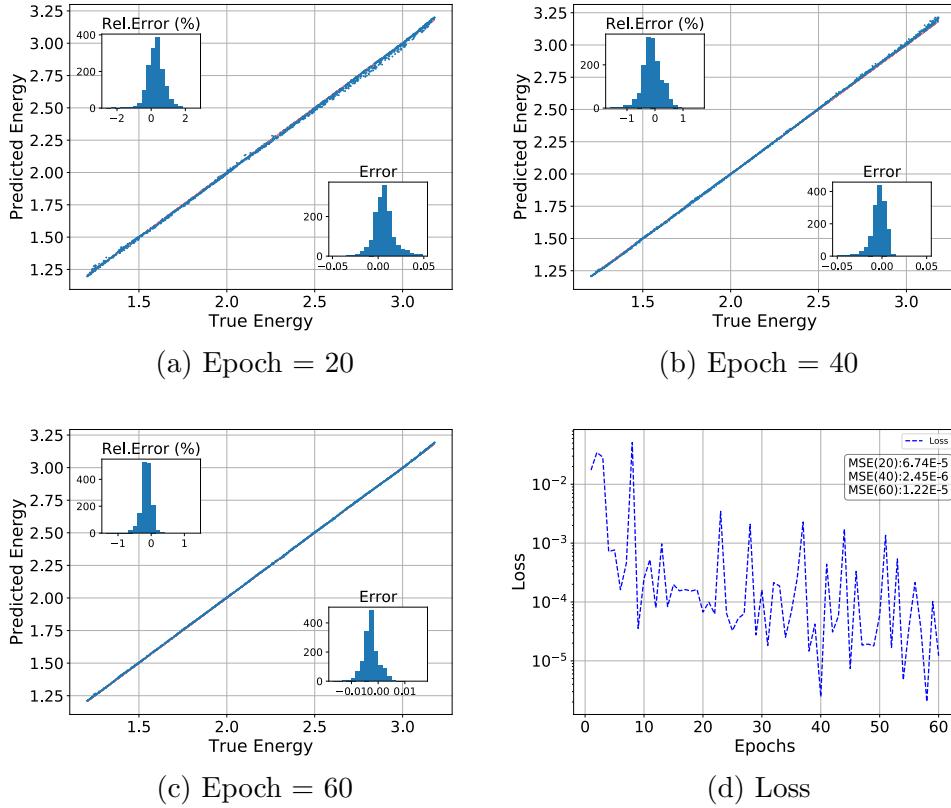


Figure 16: CNN results for  $g = 10$

Accuracy and precision of the CNN are greater than FCN's, both relative error and MSE are smaller than FCN predictions. Also unlike FCN, there is no bias in first 20 epochs. Another important subject is that last 20 epochs increases MSE but reduces relative error slightly.

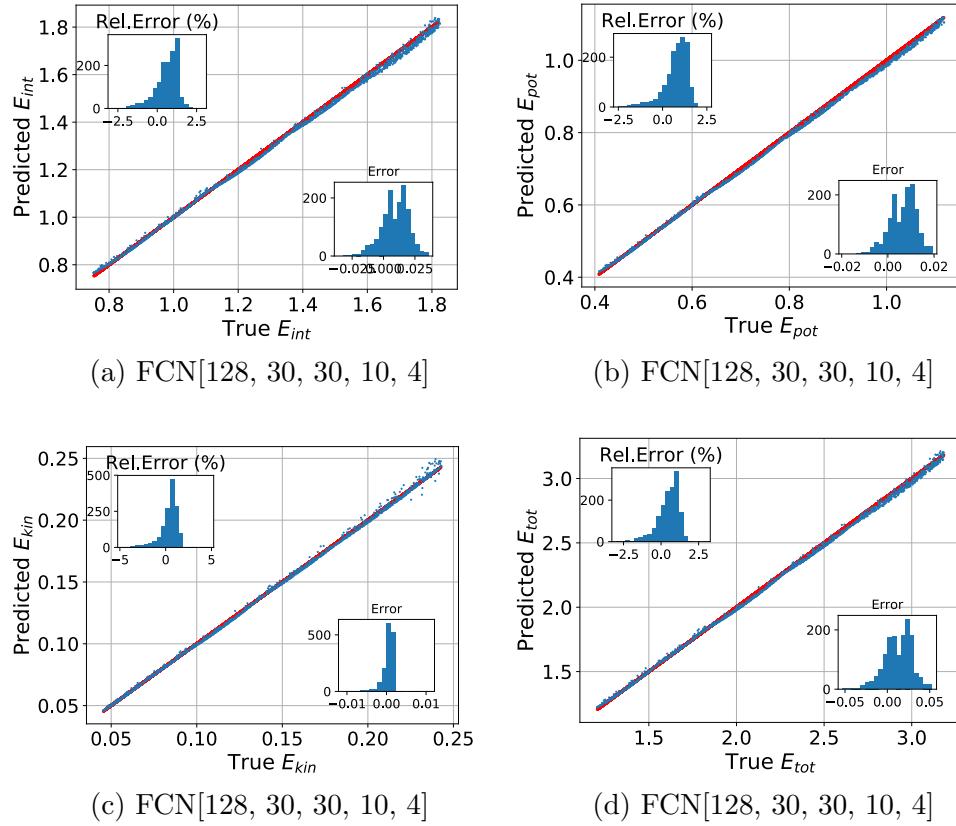


Figure 17: Separate prediction  $g=10$

Result of separate energy prediction for  $g = 10$  is not poor as  $g = 1$ . This also shows that dominance of a term in the equation directly affects the prediction.

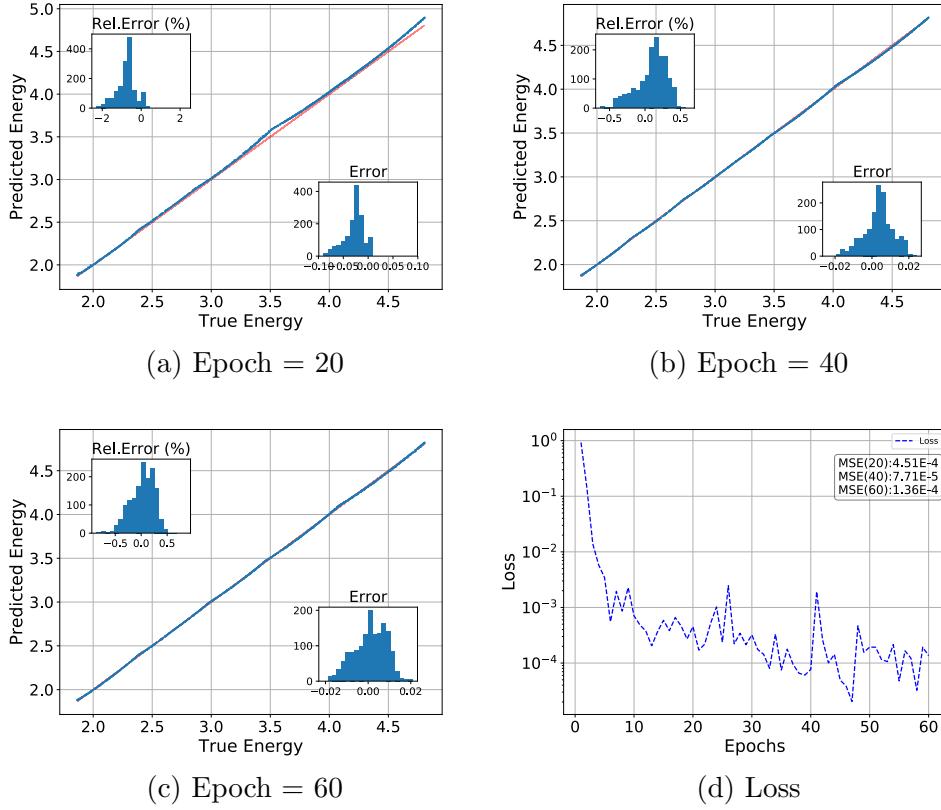


Figure 18: FCN[128, 30, 30, 10, 1] predictions for  $g = 20$

In Fig. 18(d), the network goes into saturation before 10 epochs. This shows that how the arrangement of the dataset and initialization of the network affects the predictions. A shuffling on both dataset and network initialization may produce different loss graph.

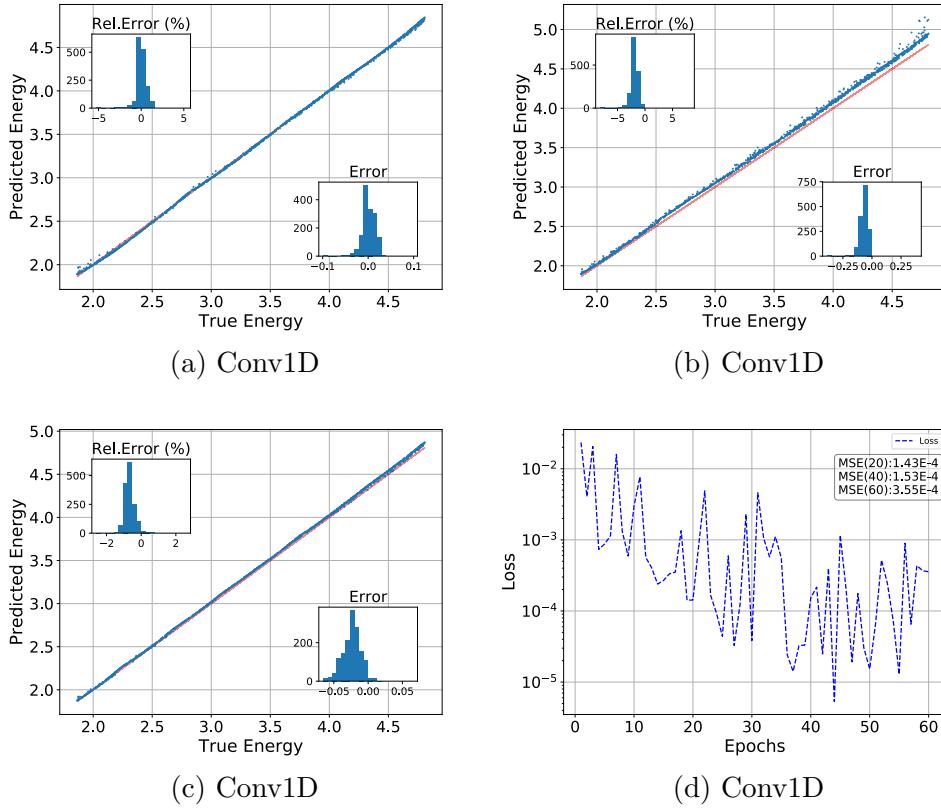


Figure 19: CNN predictions for  $g = 20$

Unbiased predictions in first 20 epochs becomes biased after 40 epochs. It is because of energy distribution in the dataset. The number of examples at higher energies is greater than number of lowers.

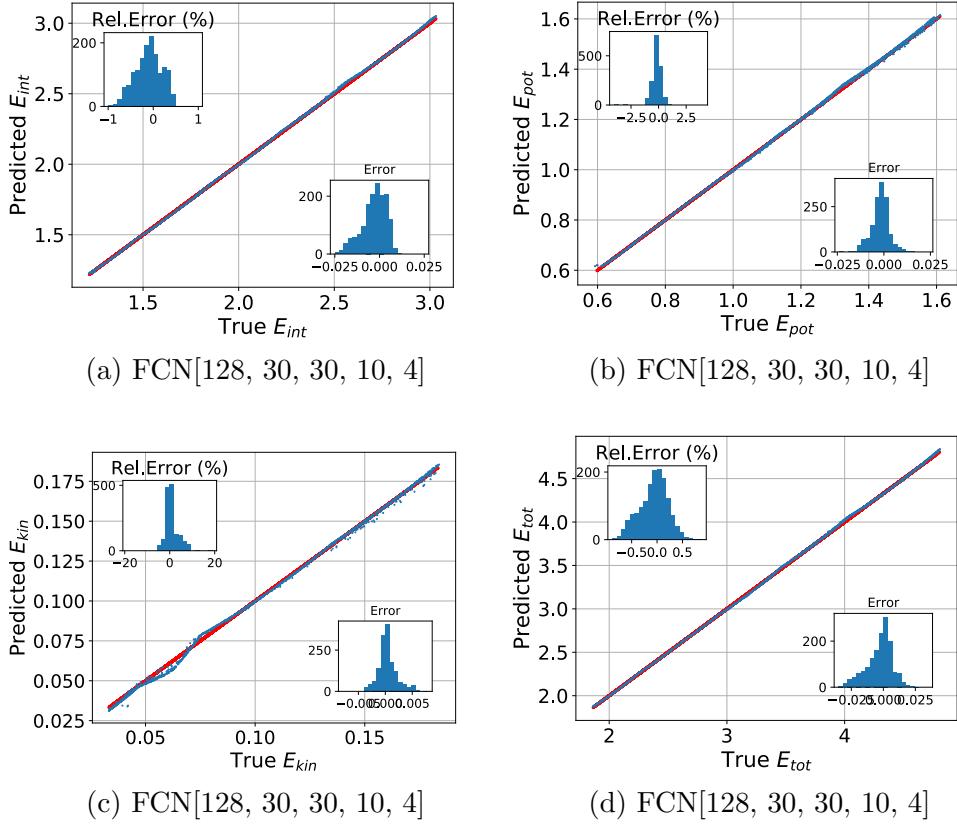


Figure 20: Separate predictions for  $g=20$

The problem in Fig. 20(c) is same as in the case  $g = 0.1$  which is shown in Fig. 11(a). The order of kinetic energy is lower than other energies.

#### 4.3.3 Inverse Problem Prediction of $g$

In this section, we try to predict the interaction parameter for the density. We accept that  $g$  is positive and the value of it can vary between 0 to 10. We used modified version of CNN to predict interaction parameter. Modification only involves a change in input layer of the network. We increased the number of channels to 2 and one of the channels represents potential the other is density. Dataloaders also updated to handle new input representation.

There are three different cases. The first one is that both potential and shift are fixed, only the interaction parameter varies. In the second one we used random angular frequencies. In final case, we both randomized angular frequency and shift.

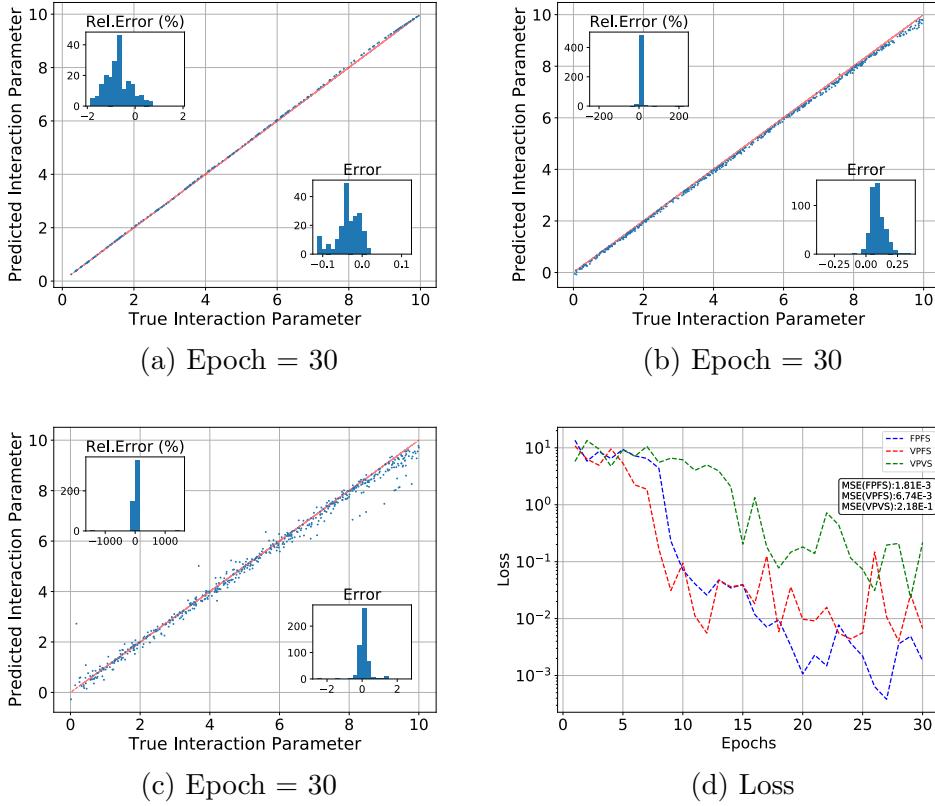


Figure 21: In (a) angular frequency is fixed and there is no shift (FPFS). In (b) angular frequency is randomized and no shift (VPFS). In (c) both angular frequency and shift are randomized (VPVS).

When there is no randomization in parameters the network handles the situation easily but when angular frequency and shift starts to vary, the network produces poorer results. It is an expected situation and it does not show that the network is not capable of to predict interaction parameter. As shown in Fig. 21(c), precision of the network is not good but the accuracy is. Expansion in the training dataset and increase in total number of epochs will fix the problem.

## 5 Conclusion

In this study, we implemented an artificial neural network to show a proof of concept that machine learning methods can also be applied to non-linear Schrödinger equation. Our results show that when the harmonic trapping potential is given, ground state energy of the Bose-Einstein can be obtained without solving the corresponding Gross-Pitaevskii equation. We also managed to predict interaction, kinetic, potential and total energy of a system separately. Error in our results differs case to case. The worst result occurred in separate energy predictions when the interaction parameter is 1 and the relative error came out about %10. Because of the contribution of the interaction parameter to the Gross-Pitaevskii equation, all terms become same order and the degree of freedom of the network is not enough to make predictions. However, in most cases except separate energy predictions relative error is lower than %2. In inverse problem, we built another neural network to predict interaction parameter when potential and density are given. The prediction accuracy decreases when angular frequency of the trapping potential and shift of equilibrium point is not fixed but mean of the relative error for it is 1.15.

Our work only involves harmonic trapping potential and limited interaction parameters. Each trained network only accept inputs for a fixed interaction parameter. Our future work will focus on more generic situations such that the network will be able to predict ground state energy of a condensate when both potential and interaction parameter is random. It is also possible to work with two dimensional trapping potential. Our aim is first successful implementation of random interaction parameter and one dimensional random potentials. After that, we will try two dimensional case.

## References

- [1] M. A. Nielsen, “Neural networks and deep learning,” 2015.
- [2] PhysicsML, “Papers.”
- [3] G. Carleo and M. Troyer, “Solving the quantum many-body problem with artificial neural networks,” *Science*, vol. 355, no. 6325, pp. 602–606, 2017.
- [4] Z. Cai, “Approximating quantum many-body wave-functions using artificial neural networks,” *arXiv preprint arXiv:1704.05148*, 2017.
- [5] L. Wang, “Discovering phase transitions with unsupervised learning,” *Physical Review B*, vol. 94, no. 19, p. 195105, 2016.
- [6] E. M. Stoudenmire and D. J. Schwab, “Supervised learning with quantum-inspired tensor networks,” *arXiv preprint arXiv:1605.05775*, 2016.
- [7] J. Biamonte, P. Wittek, N. Pancotti, P. Rebentrost, N. Wiebe, and S. Lloyd, “Quantum machine learning, 2016,” *arXiv preprint arXiv:1611.09347*.
- [8] K. Mills, M. Spanner, and I. Tamblyn, “Deep learning and the schrödinger equation,” *Physical Review A*, vol. 96, no. 4, p. 042113, 2017.
- [9] C. F. Barenghi and N. G. Parker, *A primer on quantum fluids*. No. arXiv: 1605.09580, Springer, 2016.
- [10] G. R. Dennis, J. J. Hope, and M. T. Johnsson, “Xmds2: Fast, scalable simulation of coupled stochastic partial differential equations,” *Computer Physics Communications*, vol. 184, no. 1, pp. 201–208, 2013.
- [11] P. Muruganandam and S. K. Adhikari, “Fortran programs for the time-dependent gross–pitaevskii equation in a fully anisotropic trap,” *Computer Physics Communications*, vol. 180, no. 10, pp. 1888–1912, 2009.
- [12] H.-G. Zimmermann, A. Minin, and V. Kushnerbaeva, “Comparison of the complex valued and real valued neural networks trained with gradient descent and random search algorithms.,” in *ESANN*, 2011.
- [13] M. D. Zeiler, “Adadelta: an adaptive learning rate method,” *arXiv preprint arXiv:1212.5701*, 2012.
- [14] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [15] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.

- [16] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in pytorch,” 2017.

## A APPENDIX A

### A.1 Comment on Python Codes

There are two neural network codes and three utility codes. Code of FCN is feedforwardnetwork.py, CNN is cnnnetwork1d.py. analyzer.py is kind of a network tracer. All information about network internals and iteration information is saved and traced by class inside this module. readdata.py reads the data generated by XMDS and produces proper representations of the data. sample-trainloader.py takes the data from readdata.py and produces Pytorch’s tensors. run\_training.py automatizes the training process.

Codes are not optimized and consistency of the naming convention is poor. But volume of the code is small and it is self explanatory. Code can be obtained from