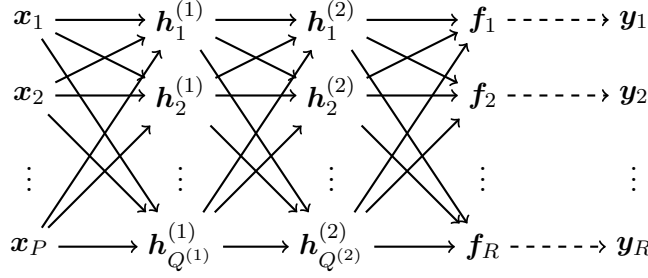


# DEEP LEARNING TUTORIAL

MAZIAR RAISSI

**1. Neural Networks [2].** Consider the following deep neural network with two hidden layers.



Here,  $\mathbf{x}_p \in \mathbb{R}^{N \times 1}$  denotes dimension  $p = 1, \dots, P$  of the input data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_P] \in \mathbb{R}^{N \times P}$ . The first hidden layer is given by

$$\mathbf{h}_q^{(1)} = h\left(\sum_{p=1}^P a_{pq}^{(1)} \mathbf{x}_p + b_q^{(1)}\right), \quad q = 1, \dots, Q^{(1)}, \quad (1.1)$$

where  $\mathbf{h}_q^{(1)} \in \mathbb{R}^{N \times 1}$  denotes dimension  $q = 1, \dots, Q^{(1)}$  of the matrix  $\mathbf{H}^{(1)} = [\mathbf{h}_1^{(1)}, \dots, \mathbf{h}_{Q^{(1)}}^{(1)}] \in \mathbb{R}^{N \times Q^{(1)}}$ . In matrix-vector notations we obtain  $\mathbf{H}^{(1)} = h(\mathbf{X}A^{(1)} + b^{(1)})$  with  $A^{(1)} = [a_{pq}^{(1)}]_{p=1, \dots, P, q=1, \dots, Q^{(1)}}$  being a  $P \times Q^{(1)}$  matrix of multipliers and  $b^{(1)} = [b_1^{(1)}, \dots, b_{Q^{(1)}}^{(1)}] \in \mathbb{R}^{1 \times Q^{(1)}}$  denoting the bias vector. Here,  $h(x)$  is the activation function given explicitly by  $h(x) = \tanh(x)$ . Similarly, the second hidden layer  $\mathbf{H}^{(2)} \in \mathbb{R}^{N \times Q^{(2)}}$  is given by  $\mathbf{H}^{(2)} = h(\mathbf{H}^{(1)}A^{(2)} + b^{(2)})$  where  $A^{(2)}$  is a  $Q^{(1)} \times Q^{(2)}$  matrix of multipliers and  $b^{(2)} \in \mathbb{R}^{1 \times Q^{(2)}}$  denotes the bias vector. The output of the neural network  $\mathbf{F} \in \mathbb{R}^{N \times R}$  is given by  $\mathbf{F} = \mathbf{H}^{(2)}A + b$  with  $A \in \mathbb{R}^{Q^{(2)} \times R}$  and  $b \in \mathbb{R}^{1 \times R}$ . Moreover, we assume a Gaussian noise model

$$\mathbf{y}_r = \mathbf{f}_r + \epsilon_r, \quad r = 1, \dots, R, \quad (1.2)$$

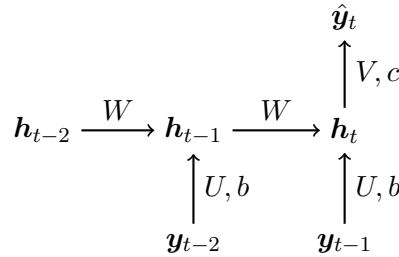
with mutually independent  $\epsilon_r \sim \mathcal{N}(\mathbf{0}, \sigma_r^2 \mathbf{I})$ . Letting  $\mathbf{y} := [\mathbf{y}_1; \dots; \mathbf{y}_R] \in \mathbb{R}^{NR \times 1}$ , we obtain the following likelihood

$$\mathbf{y} \sim \mathcal{N}(\mathbf{y}|\mathbf{f}, \Sigma \otimes \mathbf{I}), \quad (1.3)$$

where  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_S^2)$  and  $\mathbf{f} := [\mathbf{f}_1; \dots; \mathbf{f}_R] \in \mathbb{R}^{NR \times 1}$ . One can train the parameters of the neural network by minimizing the resulting negative log likelihood.

**1.1. Illustrative Example.** Figure 1.1 depicts a neural network fit to a synthetic dataset generated by random perturbations of a simple one dimensional function.

**2. Recurrent Neural Networks [2].** Let us consider a time series dataset of the form  $\{\mathbf{y}_t : t = 1, \dots, T\}$ . We can employ the following recurrent neural network



to model the next value  $\hat{\mathbf{y}}_t$  of the variable of interest as a function of its own lagged values  $\mathbf{y}_{t-1}$  and  $\mathbf{y}_{t-2}$ ; i.e.,  $\hat{\mathbf{y}}_t = f(\mathbf{y}_{t-1}, \mathbf{y}_{t-2})$ . Here,  $\hat{\mathbf{y}}_t = \mathbf{h}_t V + c$ ,  $\mathbf{h}_t = \tanh(\mathbf{h}_{t-1} W + \mathbf{y}_{t-1} U + b)$ ,  $\mathbf{h}_{t-1} = \tanh(\mathbf{h}_{t-2} W + \mathbf{y}_{t-2} U + b)$ , and  $\mathbf{h}_{t-2} = \mathbf{0}$ . The parameters  $U, V, W, b$ , and  $c$  of the recurrent neural network can be trained by minimizing the mean squared error

$$\mathcal{MSE} := \frac{1}{T-2} \sum_{t=3}^T |\mathbf{y}_t - \hat{\mathbf{y}}_t|^2.$$

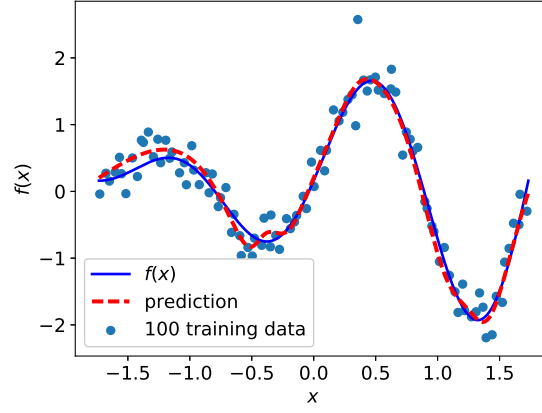


FIG. 1.1. Neural network fitting a simple one dimensional function.

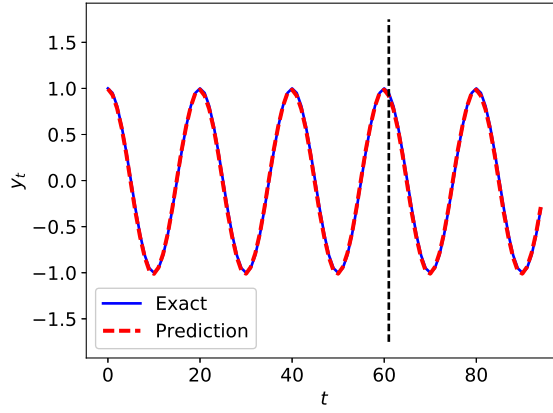


FIG. 2.1. Recurrent neural network predicting the dynamics of a simple sine wave.

**2.1. Illustrative Example.** Figure 2.1 depicts a recurrent neural network (with 5 lags) learning and predicting the dynamics of a simple sine wave.

**3. Long Short Term Memory (LSTM) Networks [2].** A long short term memory (LSTM) network replaces the units  $\mathbf{h}_t = \tanh(\mathbf{h}_{t-1}W + \mathbf{y}_{t-1}U + b)$  of a recurrent neural network with

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{s}_t),$$

where

$$\mathbf{o}_t = \sigma(\mathbf{h}_{t-1}W_o + \mathbf{y}_{t-1}U_o + b_o),$$

is the output gate and

$$\mathbf{s}_t = \mathbf{f}_t \odot \mathbf{s}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{s}}_t,$$

is the cell state. Here,

$$\tilde{\mathbf{s}}_t = \tanh(\mathbf{h}_{t-1}W_s + \mathbf{y}_{t-1}U_s + b_s).$$

Moreover,

$$\mathbf{i}_t = \sigma(\mathbf{h}_{t-1}W_i + \mathbf{y}_{t-1}U_i + b_i)$$

is the external input gate while

$$\mathbf{f}_t = \sigma(\mathbf{h}_{t-1}W_f + \mathbf{y}_{t-1}U_f + b_f),$$

is the forget gate.

**3.1. Illustrative Example.** Figure 3.1 depicts a long short term memory network (with 10 lags) learning and predicting the dynamics of a simple sine wave.

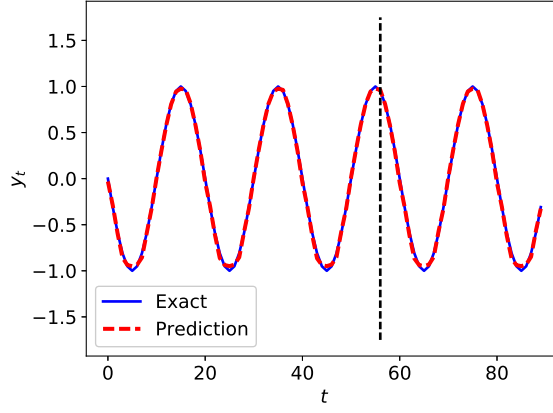


FIG. 3.1. Long short term memory network predicting the dynamics of a simple sine wave.

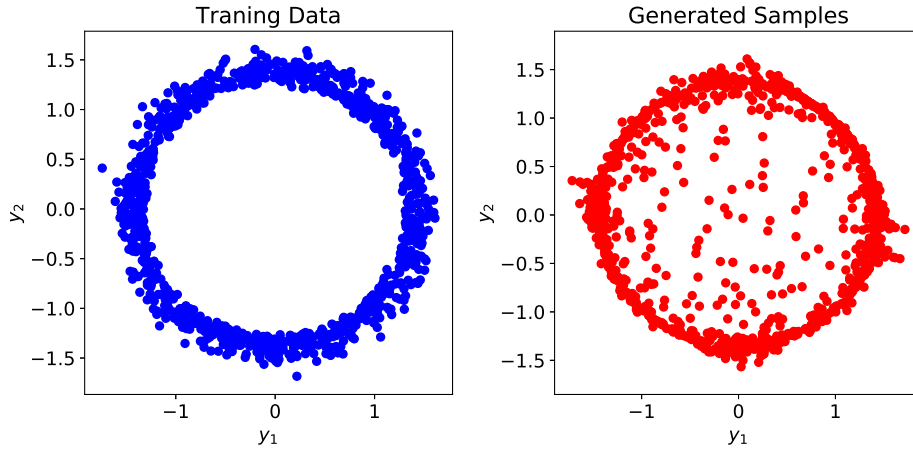


FIG. 3.2. Training data and samples generated by a variational auto-encoder.

#### 4. Variational Auto-encoders [1].

Let us start by the prior assumption that

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, I),$$

where  $\mathbf{z}$  is a latent variable. Moreover, let us assume

$$p(\mathbf{y}|\mathbf{z}) = \mathcal{N}(\mathbf{y}|\mu_2(\mathbf{z}), \Sigma_2(\mathbf{z})),$$

where  $\mu_2(\mathbf{z})$  and  $\Sigma_2(\mathbf{z})$  are modeled as deep neural networks. Here,  $\Sigma_2(\mathbf{z})$  is constrained to be a diagonal matrix. We are interested in minimizing the negative log likelihood  $-\log p(\mathbf{y})$ , where  $p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ . However,  $-\log p(\mathbf{y})$  is not analytically tractable. To deal with this issue, one could employ a variational distribution  $q(\mathbf{z}|\mathbf{y})$  and compute the following Kullback-Leibler divergence; i.e.,

$$\mathbb{KL}[q(\mathbf{z}|\mathbf{y}) || p(\mathbf{z}|\mathbf{y})] = \int \log \frac{q(\mathbf{z}|\mathbf{y})}{p(\mathbf{z}|\mathbf{y})} q(\mathbf{z}|\mathbf{y}) d\mathbf{z} = \int [\log q(\mathbf{z}|\mathbf{y}) - \log p(\mathbf{z}|\mathbf{y})] q(\mathbf{z}|\mathbf{y}) d\mathbf{z}.$$

Using the Bayes rule  $p(\mathbf{z}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{y})}$ , one obtains

$$\mathbb{KL}[q(\mathbf{z}|\mathbf{y}) || p(\mathbf{z}|\mathbf{y})] = \int [\log q(\mathbf{z}|\mathbf{y}) - \log p(\mathbf{y}|\mathbf{z}) - \log p(\mathbf{z}) + \log p(\mathbf{y})] q(\mathbf{z}|\mathbf{y}) d\mathbf{z}.$$

Therefore,

$$\mathbb{KL}[q(\mathbf{z}|\mathbf{y}) || p(\mathbf{z}|\mathbf{y})] = \log p(\mathbf{y}) + \mathbb{KL}[q(\mathbf{z}|\mathbf{y}) || p(\mathbf{z})] - \int \log p(\mathbf{y}|\mathbf{z}) q(\mathbf{z}|\mathbf{y}) d\mathbf{z}.$$

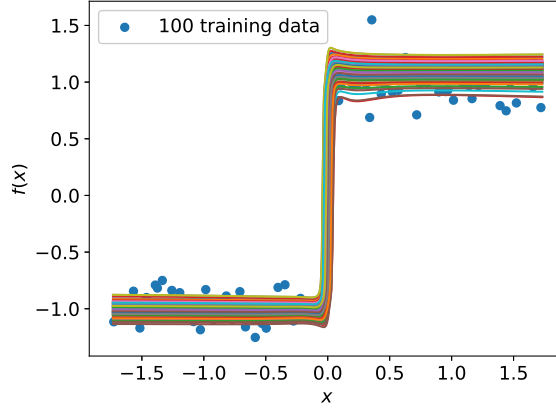


FIG. 5.1. Training data and samples generated by a conditional variational auto-encoder.

Rearranging the terms yields

$$-\log p(\mathbf{y}) + \mathbb{KL}[q(\mathbf{z}|\mathbf{y}) \parallel p(\mathbf{z}|\mathbf{y})] = - \int \log p(\mathbf{y}|\mathbf{z})q(\mathbf{z}|\mathbf{y})d\mathbf{z} + \mathbb{KL}[q(\mathbf{z}|\mathbf{y}) \parallel p(\mathbf{z})].$$

A variational auto-encoder proceeds by minimizing the terms on the right hand side of the above equation. Moreover, let us assume that

$$q(\mathbf{z}|\mathbf{y}) = \mathcal{N}(\mathbf{z}|\mu_1(\mathbf{y}), \Sigma_1(\mathbf{y})),$$

where  $\mu_1(\mathbf{y})$  and  $\Sigma_1(\mathbf{y})$  are modeled as deep neural networks. Here,  $\Sigma_1(\mathbf{y})$  is constrained to be a diagonal matrix. One can use

$$\mu_1(\mathbf{y}) + \epsilon \Sigma_1(\mathbf{y})^{1/2}, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, I).$$

to generate samples from  $q(\mathbf{z}|\mathbf{y})$ .

**4.1. Illustrative Example.** Figure 3.2 depicts the training data and the samples generated by a variational auto-encoder.

**5. Conditional Variational Auto-encoders.** Conditional variational auto-encoders, rather than making the assumption that  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, I)$ , start by assuming that

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu_0(\mathbf{x}), \Sigma_0(\mathbf{x})),$$

where  $\mu_0(\mathbf{x})$  and  $\Sigma_0(\mathbf{x})$  are modeled as deep neural networks. Here,  $\Sigma_0(\mathbf{x})$  is constrained to be a diagonal matrix.

**5.1. Illustrative Example.** Figure 5.1 depicts the training data and the samples generated by a conditional variational auto-encoder.

## REFERENCES

- [1] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- [2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.