

(Draft : March 18, 2018)

Technische Universität Berlin

Institut für Softwaretechnik und Theoretische Informatik
Maschinelles Lernen / Intelligente Datenanalyse

Fakultät IV – Elektrotechnik und Informatik
Franklinstrasse 28-29
10587 Berlin
<http://www.ml.tu-berlin.de>



Master Thesis

Designing Recurrent Neural Networks for Explainability

Pattarawat Chormai

Matriculation Number: 387441
31.03.2018

Supervisor

Prof. Dr. Klaus-Robert Müller

Advisor

Dr. Grégoire Montavon

Acknowledgement

First of all, I would like to thank Prof. Dr. Klaus-Robert Müller and Dr. Grégoire Montavon for this research opportunity and invaluable guidance throughout the course of conducting the thesis as well as facilitating me at TU Berlin, Machine Learning group with a great research environment . I would also like to thank Prof. Dr. Klaus Obermayer and staffs of Neural Information Processing group for organizing Machine Intelligence I & II and Neural Information Project. These courses provide me necessary knowledge to conduct the thesis.

Secondly, I would like to thank to my family for always providing me financial and mental support.

Thanks EIT colleagues and . In particular, Someone for proofreading..

I would like to thank EIT Master School as well as TU/e and TUB staffs that involve in this study program. Your support and advice are always helpful.

Lastly, I would like to acknowledge Amazon AWS for generous educational credits and Auction Club S.a.r.L, Luxembourg for lending me a powerful laptop to use in my study. None of experiments would have been possible without these computational resource supports.

Hereby I declare that I wrote this thesis myself with the help of no more than the mentioned literature and auxiliary means.

Berlin, 15.03.2018

.....

Pattarawat Chormai

Abstract

Standard (non-LSTM) recurrent neural networks have been challenging to train, but special optimization techniques such as heavy momentum makes this possible. However, the potentially strong entangling of features that results from this difficult optimization problem can cause deep Taylor or LRP-type to perform rather poorly due to their lack of global scope. LSTM networks are an alternative, but their gating function make them hard to explain by deep Taylor LRP in a fully principled manner. Ideally, the RNN should be expressible as a deep ReLU network, but also be reasonably disentangled to let deep Taylor LRP perform reasonably. The goal of this thesis will be to enrich the structure of the RNN with more layers to better isolate the recurrent mechanism from the representational part of the model.

Zusammenfassung

Da die meisten Leuten an der TU deutsch als Muttersprache haben, empfiehlt es sich, das Abstract zusätzlich auch in deutsch zu schreiben. Man kann es auch nur auf deutsch schreiben und anschließend einem Englisch-Muttersprachler zur Übersetzung geben.

Contents

Notation	xiii
List of Figures	xv
List of Tables	xvii
1 Introduction	1
2 Related Work	3
3 Background	5
3.1 Neural Networks	6
3.1.1 Loss functions	7
3.1.2 Learning Algorithm : Gradient Descent and Backpropagation	8
3.1.3 Convolutional Neural Networks	11
3.1.4 Recurrent Neural Networks	11
3.2 Explainability of Neural Networks	15
3.2.1 Sensitivity Analysis	16
3.2.2 Guided Backpropagation	16
3.2.3 Simple Taylor Decomposition	17
3.2.4 Layer-wise Relevance Propagation	18
3.2.5 Deep Taylor Decomposition	19
4 Experiments	25
4.1 Dataset	25
4.1.1 MNIST	25
4.1.2 Fashion-MNIST	25
4.2 RNN Cell Architectures	26
4.2.1 Shallow Cell	27
4.2.2 Deep Cell Architecture	27
4.3 Setup	28
4.4 Experiment 1 : Sequence Classification	31
4.4.1 Problem Formulation	31
4.4.2 Result	31
4.4.3 Summary	35
4.5 Experiment 2 : Majority Sample Classification	35
4.5.1 Problem Formulation	35

4.5.2	Result	36
4.5.3	Summary	37
4.6	Experiment 3 : Improve Relevance Distribution	38
4.6.1	Proposal 1 : Gating units	38
4.6.2	Proposal 2 : Stationary Dropout	38
4.6.3	Proposal 3 : Convolutional layers with literal connections	38
4.6.4	Result	38
4.6.5	Summary	38
5	Conclusion	41
5.1	Summary	41
5.2	Challenges	41
5.3	Future work	41
	References	43

Notation

θ	Parameters of a neural network
$\mathbf{x}, \mathbf{x}^{(\alpha)}$	A vector representing an input sample
σ	An activation function
$\{a_j\}_L$	A vector of activations of neurons in layer L
a_j	Activation of neuron j
b_k	Bias of neuron k
R_j	Relevance score of neuron j
$R_{j \leftarrow k}$	Relevance score distributed from neuron k to neuron j
w_{jk}	Weight between neuron j to neuron k
x_i	Feature i of input sample x

List of Figures

3.1	xxx	6
3.2	xxx	6
3.6	RNN Structure	12
3.7	LSTM Structure	14
3.8	Comparison between Global and Local Analysis	16
3.9	The LOF caption	19
3.10	An illustration of R_k functional view and root point candidates	22
3.11	R_k functional view and root points from z^+ -rule	23
3.12	R_k functional view and root points from z^β -rule with $-1 < a_j < 1$	24
4.1	MNIST Dataset	25
4.2	Fashion-MNIST Dataset TODO : Figure use same size as MNIST figure	26
4.3	Shallow and Deep Cell Architecture	27
4.4	DeepV2 and ConvDeep Cell Architecture	28
4.5	Number of neurons in each layer for each cell architecture	29
4.6	General setting of RNN classifiers in this thesis	31
4.7	Relevance heatmaps from Shallow and Deep Cell trained on MNIST with different sequence lengths	32
4.8	Relevance heatmaps from Shallow and Deep Cell trained on FashionMNIST with different sequence lengths	33
4.9	DTD relevance heatmaps of MNIST <i>Class 1</i> and FashionMNIST <i>Class Trouser</i> samples from Shallow-7 and Deep-7.	34
4.10	Pixel intensity and DTD relevance distribution from Shallow-7 and Deep-7 averaged over MNIST <i>Class 1</i> and FashionMNIST <i>Class Trouser</i> test population.	34
4.11	Majority Sample Classification Problem	35
4.12	Relevance heatmaps for MSC problem	37
4.13	Percentage of relevance in data region	38
4.14	Variational Dropout[Gal and Ghahramani, 2016]	38
4.15	ConvDeep with Literal connections	39

List of Tables

3.1	LRP rules from Deep Taylor Decomposition	24
4.1	Hyperparameter Summary	29
4.2	Total variables in each architecture and sequence length	30
4.3	Classification Accuracy Criteria	30
4.4	Dimensions of \mathbf{x}_t for each dataset and sequence length	32
4.5	Model Accuracy	32
4.6	Number of trainable variables and model accuracy for Majority Sample Classification experiment.	36

1 Introduction

2 Related Work

3 Background

3.1 Neural Networks

Neural networks (NNs) are a type of machine learning algorithms that inspired by how human brain works. In particular, NNs have units called neurons connecting together similar to the way of neurons our brain do. These connections allow NNs to hierarchically build representations that are necessary to perform an objective task. Figure 3.1 illustrates a simple coordination reaction task by our brain.

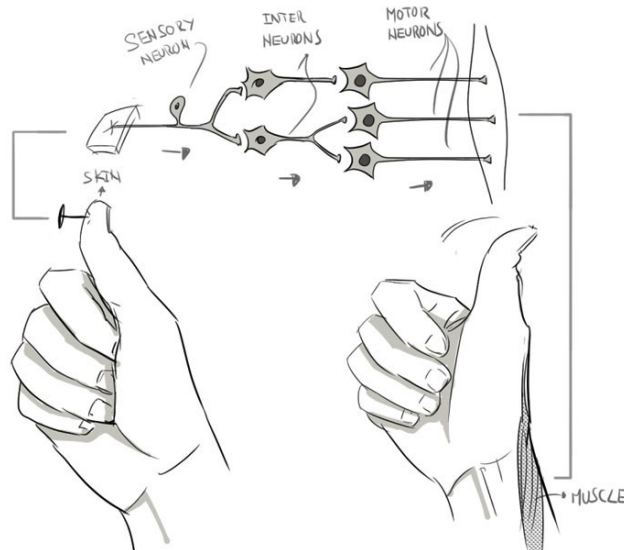


Figure 3.1: An illustration of how neurons in our brain cooperate together to sense the pain and move the thumb. ¹

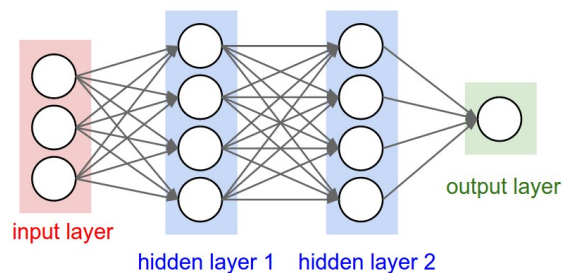


Figure 3.2

Figure 3.2 illustrates a general structure of NNs. It has input layer, output layer and hidden layers, which are analogously similar to sensory, motor and inter neurons in

¹source: Eugenio N. Leon, <https://becominghuman.ai/making-a-simple-neural-network-2ea1de81ec20>

Figure 3.1. The goal is to build connections between these neurons such that the NN is able to perform an objective task with high accuracy. This process is called *training*. In this example, the objective task to classify what is the given digit.

Consider a given set of p training samples $\mathcal{D} = \{\mathbf{x}^{(\alpha)}, y^{(\alpha)}\}_{\alpha=1}^p$, there are 3 primary components to train a NN, namely

1. **Network architecture** : it defines how neurons connect and communicate to each other. More precisely, these are corresponding to connection weights and biases $\boldsymbol{\theta}$ and neuron's activation functions illustrated in Figure 3.2. Mathematically, it is a function f with parameters $\boldsymbol{\theta}$ that nonlinearly transforms an input $\mathbf{x}^{(\alpha)} \in \mathbb{R}^d$ to some values.
2. **Loss function** L : it is a measurement corresponding to the objective task that quantifies whether output $f(\mathbf{x}^{(\alpha)})$ from the NN matches the true output $y^{(\alpha)}$.
3. **Learning algorithm** : it is responsible to find the network's parameters $\hat{\boldsymbol{\theta}}$ to optimize the loss function $L(\text{Empirical Error})$ as a proxy for optimizing loss or error for unforeseen $\mathbf{x} \notin \mathcal{D}$ (*Generalization Error*).

Hence, the goal of this empirical training process can be summarized as follows :

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_{\alpha=1}^p L(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)}) \quad (3.1)$$

3.1.1 Loss functions

Choosing loss function is depend on the objective of problem that the NN is being trianed to to solve. For classification problems, such as digit classification, which the goal is to categorize \mathbf{x} into K categories C , $f : \mathbf{x} \in \mathbb{R}^d \mapsto C \in \{C_k\}$, *Cross-Entropy*(CE) is the loss function for this purpose.

$$l_{\text{CE}} = - \sum_i y_k \log \hat{y}_k,$$

where $y_i \in [0, 1]$ and $\hat{y}_i \in [0, 1]$ are true and predicted probability that \mathbf{x} belongs to C_k respectively. Consider a K -class classification problem, $\mathbf{z} = f(\mathbf{x}) \in \mathbb{R}^K$, \hat{y}_k is computed via *softmax* activation function :

$$\hat{y}_k = \frac{e^{z_k}}{\sum_{k=1}^K e^{z_k}}$$

For regression problems, $f : \mathbf{x} \in \mathbb{R}^d \mapsto \mathbb{R}$, such as forecasting price, requires *Mean Squared Error*(MSE) loss function.

$$l_{\text{MSE}} = (f(\mathbf{x}) - y)^2$$

This is a brief introduction to loss functions widely used in machine learning. More loss functions do exist and are beyond scope of the thesis to cover.

3.1.2 Learning Algorithm : Gradient Descent and Backpropagation

Consider a function $L(\theta) = \theta^2$ on Figure 3.3 as a simplified version of loss function averaged over all $\mathbf{x}^{(\alpha)} \in \mathcal{D}$ from a NN with a parameter θ . In this case, $\hat{\theta}$ can trivially be computed by solving

$$\frac{dL(\theta)}{d\theta} \stackrel{!}{=} 0 \quad (3.2)$$

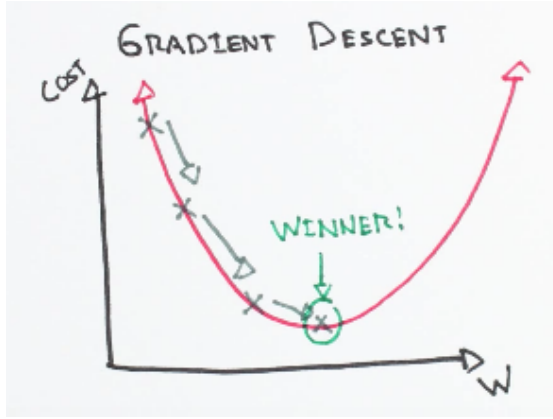


Figure 3.3: A Toy example

However, a NN usually contains thousands of parameters, hence $L(\theta)$ becomes a high dimensional function and solving (3.2) is not a trivial task. Nonetheless, Figure 3.3 shows an insight how we can reach the minimum location. There, we can see that if we move θ in the opposite direction of gradient $-\frac{dL(\theta)}{d\theta}$ with a proper step size λ (*learning rate*), we will eventually reach the minimum. Therefore, instead of directly solve (3.2), we rather slightly adjust θ in the same direction of the gradient for some iterations, for example until the loss function converges. This is called *Gradient Descent*.

$$\forall \theta_i \in \theta : \theta_i \leftarrow \theta_i - \frac{\lambda}{p} \sum_{\alpha=1}^p \frac{\partial l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})}{\partial \theta_i}$$

Let's further consider a simple NN shown in Figure 3.4 with $\theta = \{\forall i, j, k, l : w_{i \rightarrow j}^{(1)}, w_{j \rightarrow k}^{(2)}, w_{k \rightarrow l}^{(3)}\}$ and biases are omitted. $f(\mathbf{x}, \theta)$ is calculated as follows

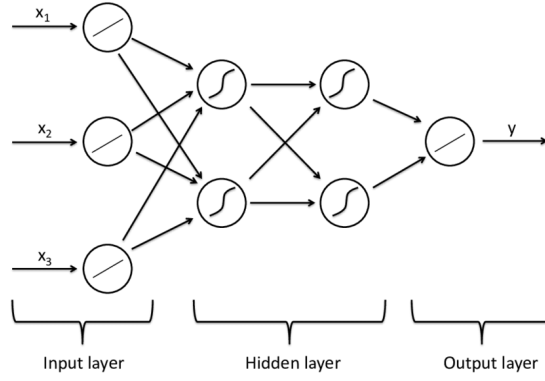


Figure 3.4: A Toy example

$$\begin{aligned}
 h_j^{(1)} &= \sum_i w_{i \rightarrow j}^{(1)} x_i & a_j^{(1)} &= \sigma(h_j^{(1)}) \\
 h_k^{(2)} &= \sum_j w_{j \rightarrow k}^{(2)} a_j^{(1)} & a_k^{(2)} &= \sigma(h_k^{(2)}) \\
 h_l^{(3)} &= \sum_k w_{k \rightarrow l}^{(3)} a_k^{(2)} & a_l^{(3)} &= \sigma(h_l^{(3)}) \\
 f(\mathbf{x}) &= [a_1^{(3)}, \dots, a_L^{(3)}]^T & L &= \frac{1}{p} \sum_{\alpha=1}^p l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})
 \end{aligned}$$

Consider a sample $(\mathbf{x}^{(\alpha)}, y^{(\alpha)})$. The gradients can be computed by recursively applying

chain rule to l from the last to the first layer, hence the name *Backpropagation*.

$$\frac{\partial l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})}{\partial w_{k \rightarrow l}^{(3)}} = \frac{\partial l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})}{\partial a_l^{(3)}} \frac{\partial a_l^{(3)}}{\partial w_{k \rightarrow l}^{(3)}} \quad (3.3)$$

$$= \underbrace{\frac{\partial l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})}{\partial a_l^{(3)}} \sigma'(h_l^{(3)}) a_k^{(2)}}_{\delta_l^{(3)}} \quad (3.4)$$

$$\frac{\partial l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})}{\partial w_{j \rightarrow k}^{(2)}} = \sum_{l'=1}^L \frac{\partial l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})}{\partial a_{l'}^{(3)}} \frac{\partial a_{l'}^{(3)}}{\partial w_{j \rightarrow k}^{(2)}} \quad (3.5)$$

$$= \sum_{l'=1}^L \frac{\partial l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})}{\partial a_{l'}^{(3)}} \sigma'(h_{l'}^{(3)}) \frac{\partial h_{l'}^{(3)}}{\partial w_{j \rightarrow k}^{(2)}} \quad (3.6)$$

$$= \sum_{l'=1}^L \delta_{l'}^{(3)} w_{k \rightarrow l'}^{(3)} \frac{\partial a_k^{(2)}}{\partial w_{j \rightarrow k}^{(2)}} \quad (3.7)$$

$$= a_j^{(1)} \sigma'(h_k^{(2)}) \underbrace{\sum_{l'=1}^L \delta_{l'}^{(3)} w_{k \rightarrow l'}^{(3)}}_{\delta_k^{(2)}} \quad (3.8)$$

$$\frac{\partial l(f(\mathbf{x}^{(\alpha)}), y^{(\alpha)})}{\partial w_{i \rightarrow j}^{(1)}} = x_i \sigma'(h_j^{(1)}) \sum_{k'=1}^K \delta_{k'}^{(2)} w_{j \rightarrow k'}^{(2)} \quad (3.9)$$

As shown in the derivation, *Backpropagation* allows us to efficiently compute the gradients by reusing calculated quantities from computation of later layer, for example $\delta_l^{(3)}, \delta_k^{(2)}$. Moreover, $\delta_l^{(3)}, \delta_k^{(2)}$ can interpreted as error of the neuron.

In practice, because the training set usually contains several thousand samples, the gradient update in (3.3) would require significant amount of computation to update one step, not to mention that it could also result in small gradient update step leading to slow convergence to desire objective performance. Therefore, the training data D is usually divided into batches \tilde{D}_i with equal size and perform the gradient update for every \tilde{D}_i . For example, the size of \tilde{D}_i is usually chosen between 32 and 512 samples. This refers to *Mini-Batch Gradient Descent*.

Lastly, because noise in the data and potentially highly non-smooth of the loss function, learning rate λ is a great influential to the training process. More precisely, it should not be too small or too large. This requires some effort and experience in order to get the value right. Some work have proposed alternative update rules aiming to make the training process more stable. For example, Adaptive Moment Estimation(Adam)[Kingma and Ba, 2014] uses adaptive learning rate and incorporates accumulated direction and speed of the previous gradients (momentum) into the update, hence more consistent gradient and fast convergence. Other similar proposals are RMSProp [Tieleman and Hinton, 2012] and Adadelta [Zeiler, 2012].

3.1.3 Convolutional Neural Networks

Convolutional Neural Networks(CNNs) refer to neural networks that employ convolutional operators to process information in early layers instead of fully-connected layers(weighted sum). Typically, a convolutional operator is followed by a pooling operator. Using this convolutional and pooling operators allows the NN to extract hierarchical features that are spatially invariant [Zeiler and Fergus, 2013], hence having higher predictive capability of the NN comparing to traditional fully-connected layers with the same number of parameters.

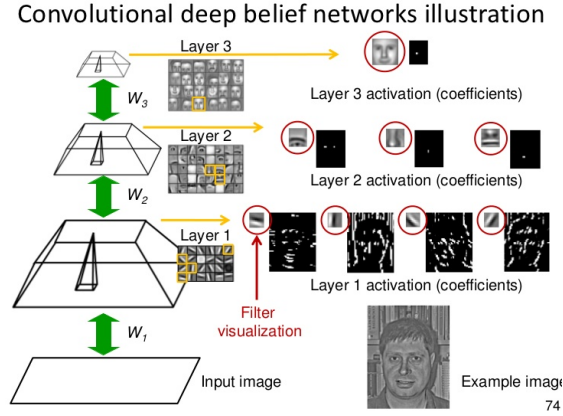


Figure 3.5: An intuition of what features CNN learns

Figure 3.5 illustrates an intuition behind hierarchical structures that CNN’s neurons learn to detect. More precisely, in this example, neurons in the first learn to detect low level features, such as edges, and neurons in middle layer then use knowledge to detect higher level features, for example nose, mouth or eyes, and so on.

Since [LeCun et al., 2001] successfully applied CNNs to a document classification problem, CNNs have become the first choice of architectures in many domains. Particularly, in computer vision, CNNs are the core component of state-of-the-art results in various contests. Such successful results are : AlexNet[Krizhevsky et al., 2012] that archive the remarkable results on ImageNet Large-Scale Visual Recognition Challenge 2012(ILSVRC 2012) followed by the achievement of VGG[Simonyan and Zisserman, 2014] and GoogleLenet [Szegedy et al., 2014] architecture in ILSVRC 2014 and ResNet[He et al., 2015] that won ILSVRC 2015.

3.1.4 Recurrent Neural Networks

Recurrent Neural Networks(RNNs) are neural networks whose computed outputs are repeatedly incorporated into the next computation. Figure 3.6 illustrates this idea of recurrent computation by unfolding RNN into steps. Let’s consider \mathbf{x} a sequence of x_1, \dots, x_t . At step t , RNN takes r_{t-1} and x_t to compute r_t and \hat{y}_t . This recurrent

connections can be interpreted as accumulating information from the past, hence RNNs are capable of processing sequential data, possibly coming with different length. Natural Language Processing(NLP) and Machine Translation(MT) are some of the fields that RNNs are widely applied to.

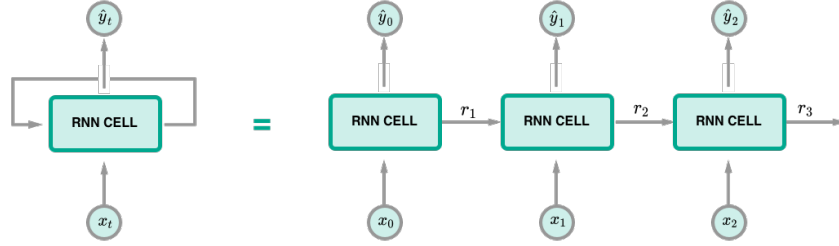


Figure 3.6: RNN Structure

Backpropagation Through Time

As the number of computation steps in RNNs is depend on the length of samples, which can be different in principle, one needs to organize data in such a way that samples in the same batch have the same steps of computations before training a RNN. In this case, training RNNs can be viewed as training a feedforward neural network where we apply backpropagation, except that variables are shared across steps.

Consider again the RNN in Figure 3.6 with $\mathbf{x} = \{x_1, \dots, x_T\}$ and $r_0 = 0$. Assume that only \hat{y}_T determines the value of the loss function and the computations are defined as follows

$$h_1 = w_{rx}x_1 + w_{rr}r_0 \quad r_1 = \sigma(h_1) \quad (3.10)$$

$$h_2 = w_{rx}x_2 + w_{rr}r_1 \quad r_2 = \sigma(h_2) \quad (3.11)$$

$$\vdots \quad \vdots \quad (3.12)$$

$$h_{T-1} = w_{rx}x_{T-1} + w_{rr}r_{T-2} \quad r_{T-1} = \sigma(h_{T-1}) \quad (3.13)$$

$$\hat{y} = \sigma(w_{yx}x_T + w_{yr}r_{T-1}) \quad (3.14)$$

The gradients can be computed by

$$\frac{\partial l}{\partial w_{yx}} = \sigma'(w_{yx}x_T + w_{yr}r_{T-1})x_T \quad (3.15)$$

$$\frac{\partial l}{\partial w_{yr}} = \sigma'(w_{yx}x_T + w_{yr}r_{T-1})r_{T-1} \quad (3.16)$$

$$\frac{\partial l}{\partial w_{rx}} = w_{yr}\sigma'(w_{yx}x_T + w_{yr}r_{T-1})\frac{\partial r_{T-1}}{\partial w_{rx}} \quad (3.17)$$

$$= w_{yr}\sigma'(w_{yx}x_T + w_{yr}r_{T-1})\left[\sigma'(h_{T-1})\left(x_{T-1} + w_{rr}\frac{\partial r_{T-2}}{\partial w_{rx}}\right)\right] \quad (3.18)$$

$$\frac{\partial l}{\partial w_{rr}} = w_{yr}\sigma'(w_{yx}x_T + w_{yr}r_{T-1})\frac{\partial r_{T-1}}{\partial w_{rr}} \quad (3.19)$$

$$= w_{yr}\sigma'(w_{yx}x_T + w_{yr}r_{T-1})\left[\sigma'(h_{T-1})\left(\frac{\partial r_{T-2}}{\partial w_{rr}}\right)\right] \quad (3.20)$$

$$(3.21)$$

However, as we unfold the computations, we can see that there are 2 possibilities that might happen to the gradients of the shared parameters w_{rx} and w_{rr} , namely

- Exploding Gradient : this scenario happens if the gradient is derived from shared weights, for example w_{rr} in (3.18), whose absolute value is greater than one. The recursive multiplication will result in a large value of the gradient leading to unreliable training. [Pascanu et al., 2012] have proposed Gradient Clipping to alleviate the problem.
- Varnishing Gradient : in contrast, when the values are smaller than one, the gradient will be very small causing slow learning. More precisely, RNNs would require enormous of time to learn long term dependencies. The next section discusses techniques to mitigate this problem

Long Short-Term Memory and Gated RNNs

Varnishing Gradient is a major problem that causes RNNs to learn long term memories with slow progress. This is due to how the computation of recurrent connections are constructed. In particular, as in (3.10), standard RNNs compute those connections with weighted sum at every step t leading to recursive multiplication terms in the gradient's computation.

Alternatively, [Hochreiter and Schmidhuber, 1997] have proposed *Long Short-Term Memory* (LSTM) network that employs a gating mechanism and additive updates in the calculation of the recurrent connections. This mechanism decreases number of damping factors involved in the gradients' computation, hence it allows the network to learn long memories better.

More precisely, as shown in Figure 3.7, LSTM utilizes 3 gates, namely input i_g , forget f_g and o_g output gate, to control the information flow through the LSTM cell. In

particular, i_g and f_g decides how to accumulate information in the cell state C_t from previous cell state C_{t-1} , previous output h_{t-1} and current input x_t , while o_g determines to the information in C_t leaks to outside h_t . Formally,

$$i_g = \sigma(w_{ix}x_t + w_{ih}h_{t-1}) \quad f_g = \sigma(w_{fx}x_t + w_{fh}h_{t-1}) \quad (3.22)$$

$$o_g = \sigma(w_{ox}x_t + w_{oh}h_{t-1}) \quad \tilde{C}_t = \tanh(w_{cx}x_t + w_{ch}h_{t-1}) \quad (3.23)$$

$$C_t = f_g \odot C_{t-1} + i_g \odot \tilde{C}_t \quad h_t = o_g \odot \tanh(C_t) \quad (3.24)$$

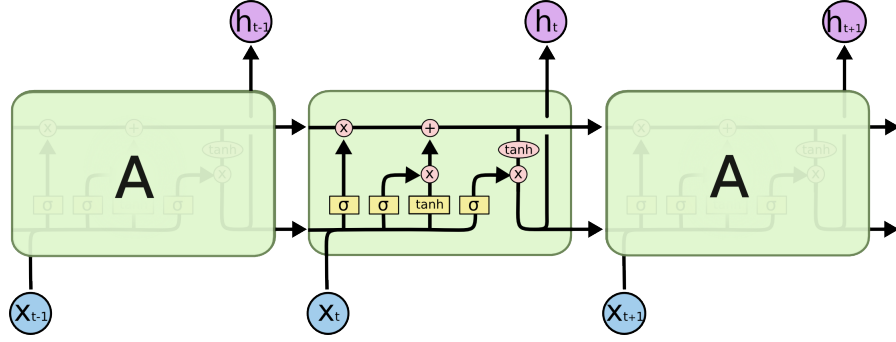


Figure 3.7: LSTM Structure

Since the work published, LSTM has successfully contributed to many state-of-the-art results, such as MT [Greff et al., 2017] have shown that the forget and output gate are the crucial parts of the network. [Cho et al., 2014] have proposed *Gated Recurrent Unit*(GRU) that employs only 2 gates, however [Jozefowicz et al., 2015] have conducted several benchmarking tasks and found no significant difference in performance between LSTM and GRU.

3.2 Explainability of Neural Networks

Neural networks have become one of major machine learning algorithms used in many applications, for example computer vision and medicine. Despite those achievements, they are still considered as a blackbox process that is difficult to interpret its results or find evidences that the networks use to make such accurate decisions for further analysis.

Moreover, it is always important to verify whether the trained network properly utilize data or what information it uses to make decisions. Literatures usually call this process as “Explainability”. [Bach et al., 2016] show that there are situations that the networks exploit artifacts in the data to make decisions. This discovery emphasizes the importance of having explainable neural networks, not to mention the fact that we are applying more and more neural networks to domains that human’s life involved, such as medicine or self-driving car.

There are 2 approaches towards explaining neural network, namely *Global* and *Local* analysis. Given a classification problem of \mathcal{C} classes classification problem and a trained network, global method aims to find an input \mathbf{x}^* that is the most representative to a class $c \in \mathcal{C}$. “Activation Maximization[Erhan et al., 2010]” is such method.

$$\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} \mathbb{P}(c|\mathbf{x}, \theta)$$

On the other hand, local analysis focuses on finding relevant information in \mathbf{x} that causes the network predicting class c_i . For example, consider an image classification problems, we can interpret a pixel $x_i \in \mathbf{x}$ as a feature, this local analysis tries to find relevance score of each pixel in respect to the classification decision. This process usually results in heatmap intensity, often called relevance heatmap $R(\mathbf{x})$.

The difference between the 2 approaches can be analogously described by formulating questions as follows : Consider \mathbf{x} is an image in category “car”.

- Global analysis : “what does the usual car look like?”
- Local analysis : “which area in the image make it look like car?” wheels, windows?

In the following, I will leave content of global analysis aside and discuss only approaches in local analysis further. In particular, I will start with properties that literatures usually use to analysis goodness of a method. These properties can imply the quality of relevance heatmap. Then, I will discuss 3 methods, namely Sensitivity Analysis, Simple Taylor Decomposition and Layer-wise Relevance Propagation.

Consider $f(\mathbf{x})$ is an output from a neural network classifier that is corresponding to the class prediction, for example the value at the final layer before applying softmax function.

Definition 3.2.1. Conservation Property

$$\forall \mathbf{x} : f(\mathbf{x}) = \sum_i R_i$$

Sum of relevance score of each pixel x_i should equal to the total relevance that the network outputs.



Figure 3.8: Comparison between Global and Local Analysis

Definition 3.2.2. Positivity Property

Definition 3.2.3. Consistency

3.2.1 Sensitivity Analysis

Sensitivity analysis[Simonyan et al., 2013] is a local analysis that derives relevance R_i of pixel x_i , from the partial derivative of $f(\mathbf{x})$ respect to x_i . In particular, literature usually formulates the calculation as

$$R_i = \left(\frac{\partial f(\mathbf{x})}{\partial x_i} \right)^2$$

Hence,

$$\sum_i R_i = \|\nabla f(\mathbf{x})\|^2$$

Although this technique can be easily implemented via automatic differentiation provided in modern deep learning frameworks, such as TensorFlow[Abadi et al., 2016], the derivation of $\sum_i R_i$ above implies that sensitivity analysis instead seeks to explain R_i from the aspect of variation magnitudes, not the actual relevance quantity.

3.2.2 Guided Backpropagation

Guided backpropagation is a extended version of sensitivity analysis where gradients are propagated in a controlled manner. It is designed specifically for neural network using piecewise linear activations. In particular, Springenberg et al.[Springenberg et al., 2014] reformulate the definition of ReLU function as:

$$\sigma(x) = x \mathbb{1}[x > 0],$$

where $\mathbb{1}[\cdot]$ is an indicator function. With the new formulation, [Springenberg et al., 2014] proposes a new derivative of a ReLU neuron j as:

$$\frac{\partial_* f(\mathbf{x})}{\partial a_j} = \mathbb{1}\left[a_j > 0\right] \mathbb{1}\left[\frac{\partial f(\mathbf{x})}{\partial a_j} > 0\right] \frac{\partial f(\mathbf{x})}{\partial a_j}$$

The 2 indicator functions control whether original gradients are propagated back, hence the name ‘‘Guided Backpropagation’’. Hence, the relevance score for x_i is:

$$R_i = \left(\frac{\partial_* f(\mathbf{x})}{\partial x_i} \right)^2$$

With this result, one can see that x_i is relevant to the problem if activations a_j that it supplies are active and positively contribute to $f(\mathbf{x})$.

3.2.3 Simple Taylor Decomposition

This method decomposes $f(\mathbf{x})$ into terms of relevance scores R_i via Taylor Decomposition. Formally,

$$f(\mathbf{x}) = f(\tilde{\mathbf{x}}) + \sum_i \underbrace{\frac{\partial f}{\partial x_i} \Big|_{x_i=\tilde{x}_i}}_{R_i} (x_i - \tilde{x}_i) + \zeta,$$

where ζ is the second and higher order terms of Taylor series and $\tilde{\mathbf{x}}$ is a root point where $f(\tilde{\mathbf{x}}) = 0$. To find such $\tilde{\mathbf{x}}$, one need to optimize :

$$\min_{\xi \in \mathcal{X}} \|\xi - \mathbf{x}\|^2 \quad \text{such that } f(\xi) = 0,$$

where \mathcal{X} represents the input distribution. However, this optimization is time consuming and ξ might potentially be close to or diverge from \mathbf{x} leading to non informative R_i .

Nonetheless, Montavon et al.[Montavon et al., 2017b] demonstrate that neural networks whose activations $\sigma(x)$ are piecewise linear functions with $\sigma(tx) = t\sigma(x), \forall t \geq 0$ property, for example a deep Rectified Linear Unit (ReLU) network without biases, $\tilde{\mathbf{x}}$ can be found in approximately the same flat region as \mathbf{x} , $\tilde{\mathbf{x}} = \lim_{\epsilon \rightarrow 0} \epsilon \mathbf{x}$, yielding

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \Big|_{\mathbf{x}=\tilde{\mathbf{x}}} = \frac{\partial f(\mathbf{x})}{\partial x_i} \Big|_{\mathbf{x}=\mathbf{x}}$$

Hence, the decomposition can be simplified to :

$$f(\mathbf{x}) = \sum_i \underbrace{\frac{\partial f(\mathbf{x})}{\partial x_i} \Big|_{\mathbf{x}=\mathbf{x}}}_{R_i} x_i$$

Moreover, the result suggests the relationship between sensitivity analysis and Taylor decomposition. Specifically, x_i has high relevance score if x_i activates and its variation positively affects $f(x)$ and vice versa.

3.2.4 Layer-wise Relevance Propagation

The methods mentioned so far derive R_i directly from $f(\mathbf{x})$ and do not use important knowledge about the network itself, such as architecture or activation values. Alternatively, Bach et al. [Binder et al., 0906] propose Layer-wise Relevance Propagation (LRP) framework that leverages this known information to distribute relevance scores to x_i . In particular, LRP propagates relevance scores backward from layer to layer, similar to the back-propagation algorithm of gradient descent.

Consider the neural network illustrated in Figure 3.9. R_j and R_k are relevance score of neurons j, k in successive layers. [Binder et al., 0906] formulates the general form of relevance propagation as :

$$R_j = \sum_k \delta_{j \leftarrow k} R_k, \quad (3.25)$$

where $\delta_{j \leftarrow k}$ defines a proportion that R_k contributes to R_j . Consider further that activity a_k of neuron k is computed by

$$a_k = \sigma \left(\sum_j w_{jk} a_j + b_k \right),$$

where w_{jk}, b_k are the corresponding weight and bias between neuron j and k , and σ is a monotonic increasing activation function. [Binder et al., 0906] suggests

$$\delta_{j \leftarrow k} = \alpha \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} - \beta \frac{a_j w_{jk}^-}{\sum_j a_j w_{jk}^-}, \quad (3.26)$$

where w_{jk}^+, w_{jk}^- are $\max(0, w_{jk}), \min(0, w_{jk})$, and α, β are parameters with $\alpha - \beta = 1$ condition. Together with Equation 3.25, [Binder et al., 0906] calls $\alpha\beta$ -rule.

$$R_j = \sum_k \left(\alpha \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} - \beta \frac{a_j w_{jk}^-}{\sum_j a_j w_{jk}^-} \right) R_k \quad (3.27)$$

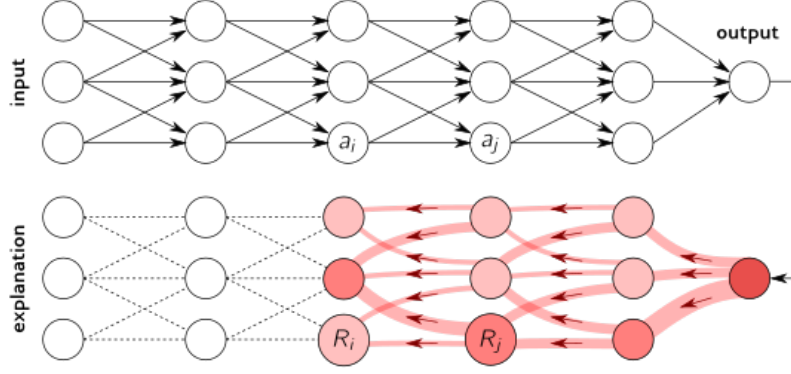
This $\alpha\beta$ -rule ensures that conservation property is satisfied as $f(\mathbf{x})$ is distributed through the network.

$$\sum_i R_i = \sum_j R_j = \sum_k R_k = f(\mathbf{x})$$

Moreover, if we rewrites $\alpha\beta$ -rule as

$$R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} \hat{R}_k + \frac{a_j w_{jk}^-}{\sum_j a_j w_{jk}^-} \check{R}_k,$$

where $\hat{R}_k = \alpha R_k$ and $\check{R}_k = -\beta R_k$. Then, we can intuitively interpret this propagation as

Figure 3.9: LRP Framework²

“Relevance \hat{R}_k ” should be redistributed to the lower-layer neurons $\{a_j\}_j$ in proportion to their excitatory effect on a_k . “Counter-relevance” \check{R}_k should be redistributed to the lower-layer neurons $\{a_j\}_j$ in proportion to their inhibitory effect on a_j - Section 5.1 [Montavon et al., 2017b]

However, it seems that there is no clear relationship between values of α, β and the structure of the heatmap. [Montavon et al., 2017b, Binder et al., 0906] demonstrate that the values are depend on the architecture of the network. In particular, [Montavon et al., 2017b] observes that $\alpha = 1, \beta = 0$ works well for deep architectures, such as GoogleNet [Szegedy et al., 2014], while $\alpha = 2, \beta = 1$ for shallower architectures, such as BVLC CaffeNet [Jia et al., 2014].

Algorithm 1: LRP Algorithm

```

 $f(\mathbf{x}), \{\{a\}_{l_1}, \{a\}_{l_2}, \dots, \{a\}_{l_n}\} = \text{forward\_pass}(\mathbf{x}, \boldsymbol{\theta});$ 
 $R_k = f(\mathbf{x});$ 
for  $\text{layer} \in \text{reverse}(\{l_1, l_2, \dots, l_n\})$  do
     $\text{prev\_layer} \leftarrow \text{layer} - 1$  ;
    for  $j \in \text{neurons}(\text{prev\_layer}), k \in \text{neurons}(\text{layer})$  do
         $R_j \leftarrow \alpha\beta\text{-rule}(R_k, \{a\}_j, \{w\}_{j,k});$ 
    end
end

```

3.2.5 Deep Taylor Decomposition

Deep Taylor Decomposition (DTD) is a explanation technique that decomposes R_k as a sum of R_j from previous layer using Simple Taylor Decomposition. Montavon et al. [Montavon et al., 2017a] proposes the method to explain decisions of neural networks

²Source: <http://heatmapping.org>

with piece-wise linear activations. Similar to LRP, DTD decomposes R_k and propagates the quantity backward to R_j . In particular, R_k is decomposed as follows :

$$R_k = R_k \Big|_{\tilde{a}_j} + \sum_j \frac{\partial R_k}{\partial a_j} \Big|_{a_j=\tilde{a}_j} (a_j - \tilde{a}_j) + \zeta_k \quad (3.28)$$

Assume further that there exists a root point \tilde{a}_j such that $R_k = 0$, and the second and higher terms $\zeta_k = 0$. Then, Equation 3.28 is simplified to

$$R_k = \sum_j \underbrace{\frac{\partial R_k}{\partial a_j} \Big|_{a_j=\tilde{a}_j}}_{R_{j \leftarrow k}} (a_j - \tilde{a}_j) \quad (3.29)$$

Moreover, as the relevance propagated back, [Montavon et al., 2017a] shows that R_j is an aggregation of $R_{j \leftarrow k}$ from neuron k in the next layer that neuron j contributes to. Formally, this is

$$R_j = \sum_k R_{j \leftarrow k} \quad (3.30)$$

Combining Equation 3.29 and 3.30 yields :

$$\begin{aligned} R_j &= \sum_k R_{j \leftarrow k} \\ \sum_j R_j &= \sum_j \sum_k R_{j \leftarrow k} \\ \sum_j R_j &= \sum_k \sum_j R_{j \leftarrow k} \\ \sum_j R_j &= \sum_k R_k \end{aligned} \quad (3.31)$$

Demonstrated by [Montavon et al., 2017a], Equation 3.31 holds for all j, k and all subsequent layers. Hence, this results in conservation property which guarantee that no relevance loss during the propagations.

$$\sum_i R_i = \cdots = \sum_j R_j = \sum_k R_k = \cdots = f(\mathbf{x}) \quad (3.32)$$

To find a root point \tilde{a}_j , consider a neural network whose R_k is computed by :

$$R_k = \max \left(0, \sum_j a_j w_{jk} + b_k \right), \quad (3.33)$$

where $b_k \leq 0$.

One can see that there are 2 cases to be analyzed, namely $R_k = 0$ and $R_k \geq 0$. For $R_k = 0$, \mathbf{a}_j is already the root point. For the latter, one can find such point by performing line search in a direction \mathbf{v}_j and magnitude t .

$$\tilde{\mathbf{a}}_j = \mathbf{a}_j - t\mathbf{v}_j, \quad (3.34)$$

The root point is then the intersection point between Equation 3.34 and 3.33. Hence,

$$0 = \sum_j (\mathbf{a}_j - t\mathbf{v}_j)w_{jk} + b_k \quad (3.35)$$

$$t \sum_j v_j w_{jk} = R_k \quad (3.36)$$

$$t = \frac{R_k}{\sum_j v_j w_{jk}} \quad (3.37)$$

$$(3.38)$$

Therefore, R_j can be computed by :

$$R_j = \sum_k \left. \frac{\partial R_k}{\partial \mathbf{a}_j} \right|_{\mathbf{a}_j - \tilde{\mathbf{a}}_j} (\mathbf{a}_j - \tilde{\mathbf{a}}_j) \quad (3.39)$$

$$= \sum_k w_{jk} t v_j \quad (3.40)$$

$$= \sum_k \frac{v_j w_{jk}}{\sum_j v_j w_{jk}} R_k \quad (3.41)$$

Noticing here is that $\tilde{\mathbf{a}}_j$ is not necessary the closest point to the line $R_k = 0$ in Euclidean distance, because $\tilde{\mathbf{a}}_j$ might not be in the same domain as \mathbf{a}_j , hence v_j needs to be chosen according to the domain of \mathbf{a}_j . Consider an example on Figure 3.10, if $\mathbf{a}_j \in \mathbb{R}^+$, then $\tilde{\mathbf{a}}_j$ must be also in \mathbb{R}^+ . In the following, I will summarize how \mathbf{a}_j can be computed for each domain of \mathbf{a}_j .

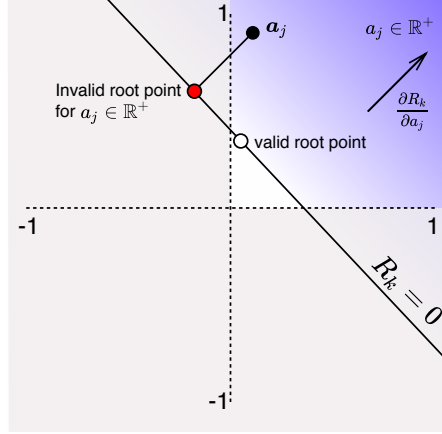
Case $\mathbf{a}_j \in \mathbb{R}$: w^2 -rule

Trivially, the search direction \mathbf{v}_j is just the direction of gradient $\frac{\partial R_k^{(l+1)}}{\partial \mathbf{a}_j^{(l)}}$:

$$\mathbf{v}_j = \mathbf{w}_{jk}$$

Hence,

$$R_j = \sum_k \frac{w_{jk}^2}{\sum_j w_{jk}^2} R_k$$

Figure 3.10: An illustration of R_k functional view and root point candidates**Case $a_j \geq 0$: z^+ -rule**

The root point is on the line segment $(\mathbf{a}_j \mathbb{1}\{w_{jk} < 0\}, \mathbf{a}_j)$. In particular, as shown on Figure 3.11, R_k has a root at $\mathbf{a}_j \mathbb{1}\{w_{jk} < 0\}$, because of:

$$R_k = \max \left(\sum_j a_j w_{jk} + b_k, 0 \right) \quad (3.42)$$

$$= \max \left(\sum_j a_j \mathbb{1}\{w_{jk} < 0\} w_{jk} + b_k, 0 \right) \quad (3.43)$$

$$= \max \left(\sum_j a_j w_{jk}^- + b_k, 0 \right) \quad (3.44)$$

$$= 0 \quad (3.45)$$

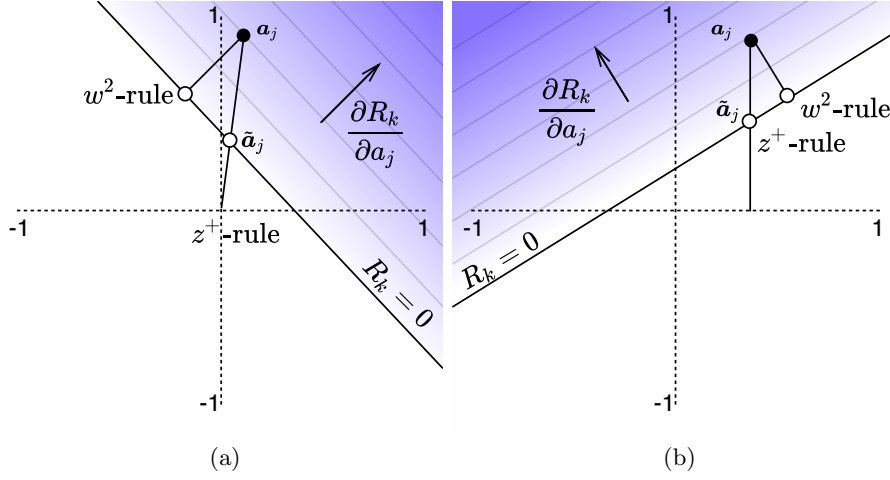
The last step uses the fact that $a_j \in R^+$ and $b_k \leq 0$ from the assumption. Hence, the search direction is:

$$\begin{aligned} v_j &= a_j - a_j \mathbb{1}\{w_{jk} < 0\} \\ &= a_j \mathbb{1}\{w_{jk} \geq 0\} \end{aligned}$$

Therefore,

$$\begin{aligned} R_j &= \sum_k \frac{w_{jk} a_j \mathbb{1}\{w_{jk} \geq 0\}}{\sum_j w_{jk} a_j \mathbb{1}\{w_{jk} \geq 0\}} R_k \\ &= \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k \end{aligned}$$

Moreover, one can also see that z^+ -rule is equivalent to LRP's $\alpha\beta$ -rule when $\alpha = 1, \beta = 0$.

Figure 3.11: R_k functional view and root points from z^+ -rule

Case $l_j \leq a_j \leq h_j$ where $l_j \leq 0 < h_j$: z^β -rule

In this case, the root point is on the line segment $(l_j \mathbb{1}\{w_{jk} \geq 0\} + h_j \mathbb{1}\{w_{jk} \leq 0\}, a_j)$. One can show that

$$R_k = \max \left(\sum_j a_j w_{jk} + b_k, 0 \right) \quad (3.46)$$

$$= \max \left(\sum_j (l_j \mathbb{1}\{w_{jk} \geq 0\} + h_j \mathbb{1}\{w_{jk} \leq 0\}) w_{jk} + b_k, 0 \right) \quad (3.47)$$

$$= \max \left(\sum_j l_j w_{jk}^+ + h_j w_{jk}^- + b_k, 0 \right) \quad (3.48)$$

$$= 0 \quad (3.49)$$

Hence, the search direction is

$$\begin{aligned} v_j &= a_j - \tilde{a}_j \\ &= a_j - l_j \mathbb{1}\{w_{jk} \geq 0\} - h_j \mathbb{1}\{w_{jk} \leq 0\} \end{aligned}$$

Figure 3.12 illustrates details of the search direction. Therefore,

$$\begin{aligned} R_j &= \sum_k \frac{w_{jk}(a_j - l_j \mathbb{1}\{w_{jk} \geq 0\} - h_j \mathbb{1}\{w_{jk} \leq 0\})}{\sum_j w_{jk}(a_j - l_j \mathbb{1}\{w_{jk} \geq 0\} - h_j \mathbb{1}\{w_{jk} \leq 0\})} R_k \\ &= \sum_k \frac{a_j w_{jk} - l_j w_{jk}^- - h_j w_{jk}^+}{\sum_j a_j w_{jk} - l_j w_{jk}^- - h_j w_{jk}^+} R_k \end{aligned}$$

In summary, DTD is a theory for explaining nonlinear computations of neural network through decomposing relevance score between successive layers. Its propagation

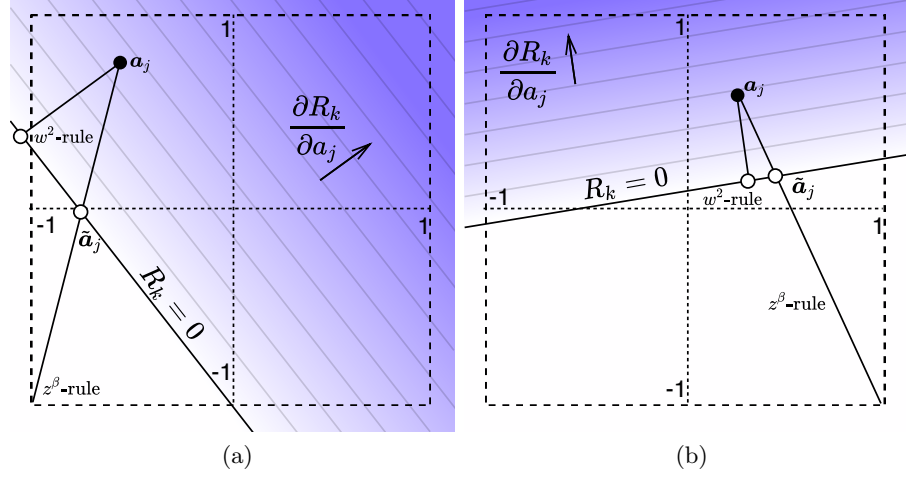


Figure 3.12: R_k functional view and root points from z^β -rule with $-1 < a_j < 1$

rules ensures conservation property. Given the rules above, the relevance scores can be propagated using LRP Algorithm 1.

Lastly, as DTD provides more general propagation rules than the $\alpha\beta$ -rule from LRP, I will use DTD and LRP interchangeably throughout the thesis. In particular, it will be mentioned explicitly if $\alpha\beta$ is being used, otherwise the rule is a DTD's rule. Table 3.1 concludes the details when such DTD rules should be used.

Input Domain	LRP Propagation Rule
w^2 -rule : Real values, $a_j \in R$	$R_j = \sum_k \frac{w_{jk}^2}{\sum_j w_{jk}^2} R_k$
z^+ -rule : ReLU activations, $a_j \geq 0$	$R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k$
z^β -rule : Pixel Intensities, $l_j \leq a_j \leq h_j, l_j \leq 0 \leq h_j$	$R_j = \sum_k \frac{a_j w_{jk} - l_j w_{jk}^- - h_j w_{jk}^+}{\sum_j a_j w_{jk} - l_j w_{jk}^- - h_j w_{jk}^+} R_k$

Table 3.1: LRP rules from Deep Taylor Decomposition

4 Experiments

4.1 Dataset

4.1.1 MNIST

MNIST[LeCun and Cortes, 2010] is one of the most popular dataset that machine learning partitioners use to benchmark machine learning algorithms. The dataset consists of 60,000 training and 10,000 testing samples. Each sample is a grayscale 28x28 image of a digit between from 0 to 9.

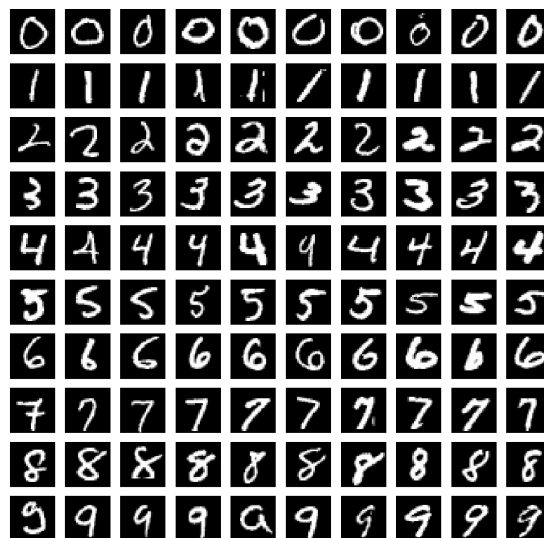


Figure 4.1: MNIST Dataset

State-of-the-art algorithms can classify MNIST with accuracy higher than 0.99, while classical ones, such as SVC or RandomForest, are able to achieve around 0.97[Xiao et al., 2017].

4.1.2 Fashion-MNIST

Xiao et. al.[Xiao et al., 2017] propose a novel dataset, called Fashion-MNIST dataset, as a replacement of MNIST dataset for benchmarking machine learning algorithms. According to [Xiao et al., 2017], Fashion-MNIST brings more challenging to the problem and more representative to modern computer vision tasks. It contains images of fashion products from 10 categories. Fashion-MNIST is comparable to MNIST in every aspects, such as the size of training and testing set, image dimension and data format, hence one

can easily apply existing algorithms that work with MNIST to Fashion-MNIST without any change.



Figure 4.2: Fashion-MNIST Dataset TODO : Figure use same size as MNIST figure

[Xiao et al., 2017] also reports benchmarking results of classical machine learning algorithms on Fashion-MNIST. On average, they achieve accuracy between 0.85 to 0.89. According to Fashion-MNIST's page¹, A. Brock reports the state-of-the-art result with 0.97 accuracy using Wide Residual Network(WRN)[Zagoruyko and Komodakis, 2016] and standard data preprocessing and augmentation.

4.2 RNN Cell Architectures

In this section, I will describe architectures or RNN cell evaluated in this thesis. To make the descriptions concise, I denote notations as follows:

- $\mathbf{a}_t^{(l)}$: activation vector of layer l at step t

¹<https://github.com/zalandoresearch/fashion-mnist>

- \mathbf{x}_t : subset of $x_i \in \mathbf{x}$ corresponding to step t and assume to be reshaped into a column vector

4.2.1 Shallow Cell

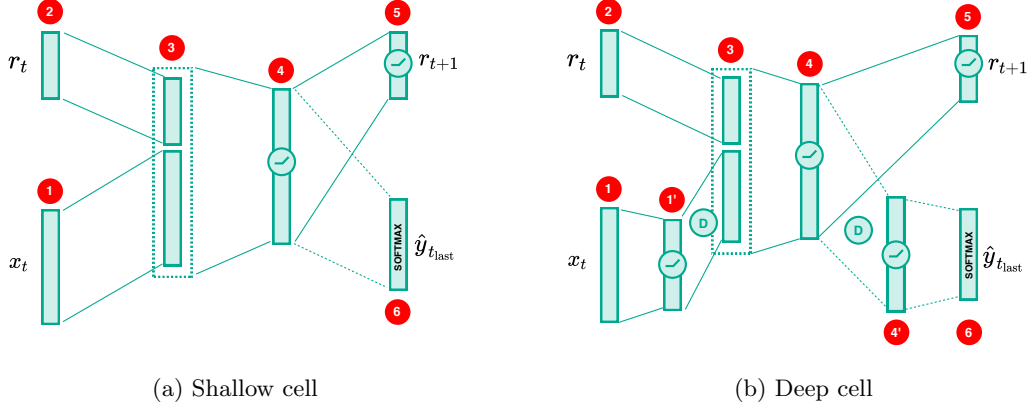


Figure 4.3: Shallow and Deep Cell Architecture

As shown in Figure 4.3a, Shallow cell first concatenates input \mathbf{x}_t (1) and recurrent input \mathbf{r}_t (2) at layer (3) as one vector before computing $\mathbf{a}_t^{(4)}$ of layer (4). Then, the next recurrent input \mathbf{r}_{t+1} (5) is derived from $\mathbf{a}_t^{(4)}$. In the last step t_{last} , the raw output \mathbf{h} is computed from $\mathbf{a}_{t_{\text{last}}}^{(4)}$ and applied to softmax function to compute class probabilities $\hat{\mathbf{y}}$ (6).

4.2.2 Deep Cell Architecture

Figure 4.3b illustrates the architecture of Deep cell. It can be viewed as an extension of the Shallow cell with 2 additional layers, namely (1') and (4'). The ideas of introducing the layers are to let (1') learn representations of the problem, while (4) can focus on combining information from past and current step, which enables (4') to compute more fine-grained decision probabilities. Dropout is applied at $\mathbf{a}^{(1')}$, and $\mathbf{a}_{t_{\text{last}}}^{(4)}$ for computing $\mathbf{a}^{(4')}$.

Two variations of Deep cell are also experimented, namely DeepV2 and ConvDeep, shown on Figure 4.4. The former has one additional layer (1'') with dropout regularization between (1'). On the other hand, the latter replaces fully connected layers between (1) and (3) with 2 convolutional and max pooling layers, $[\text{C1}, \text{P1}]$ and $[\text{C2}, \text{P2}]$.

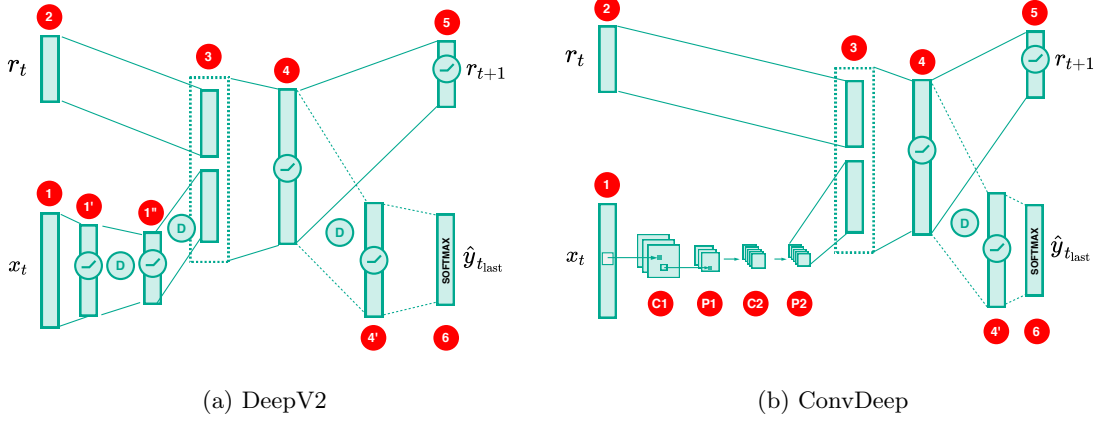


Figure 4.4: DeepV2 and ConvDeep Cell Architecture

4.3 Setup

I implemented experiments using TensorFlow². weights $w_{ij} \in \mathbf{W}$ and biases $b_j \in \mathbf{b}$ are initialized as follows:

$$w_{ij} \sim \Psi(\mu = 0, \sigma = 0.1, [-2\sigma, 2\sigma])$$

$$b_j = \ln(e^{0.01} - 1)$$

where $\Psi(\cdot)$ denotes Truncated Normal Distribution with valid region between $[-2\sigma, 2\sigma]$.

Activations $\mathbf{a}^{(l)}$ of layer l is calculated using :

$$\mathbf{h}^{(l)} = \mathbf{W}^T \mathbf{a}^{(l-1)} - \sigma_s(\mathbf{b})$$

$$\mathbf{a}^{(l)} = \sigma_r(\mathbf{h}^{(l)})$$

where

$$\sigma_r(h_j) = \max(0, h_j) \quad (\text{ReLU function})$$

$$\sigma_s(b_j) = \log(1 + \exp b_j) \quad (\text{Softplus function})$$

The reason of using $\sigma_s(b_j)$ for bias term is due to the non positive bias assumption of DTD. Moreover, $\sigma'_s(b_j)$ is $(0, \infty)$, as a result the network has more flexibility to adjust decision through back-propagation rather than using $\sigma_r(b_j)$. With this setting, the initial value of bias term $\sigma_s(b_j)$ is then 0.01.

Speaking about training methodology, I use Adam[Kingma and Ba, 2014] optimizer to train networks. Preliminary study shows that relevance heatmaps from networks trained by Adam are less noisy than the ones from other optimizers, such as Stochastic Gradient Descent(SGD). Number of epochs is set to 200, while dropout probability is

²<http://tensorflow.org/>

0.5 and batch size is 50. Learning rate is not fixed as it is likely that different network architectures requires different value to achieve good performance, hence left adjustable. Table 4.1 summaries the setting of hyperparameters.

Hyperparameter	Value
Optimizer	Adam
Epoch	200
Dropout Probability	0.5
Batch size	50

Table 4.1: Hyperparameter Summary

Traditionally, number of neurons in each layer ($n^{(l)}$) is another hyperparameter that we can adjust. However, as the goal is to compare relevance heatmaps from different architectures, those numbers are fixed and chosen in such a way that total number of variables in each architecture are equivalent. Figure 4.5 illustrates the details of the settings.

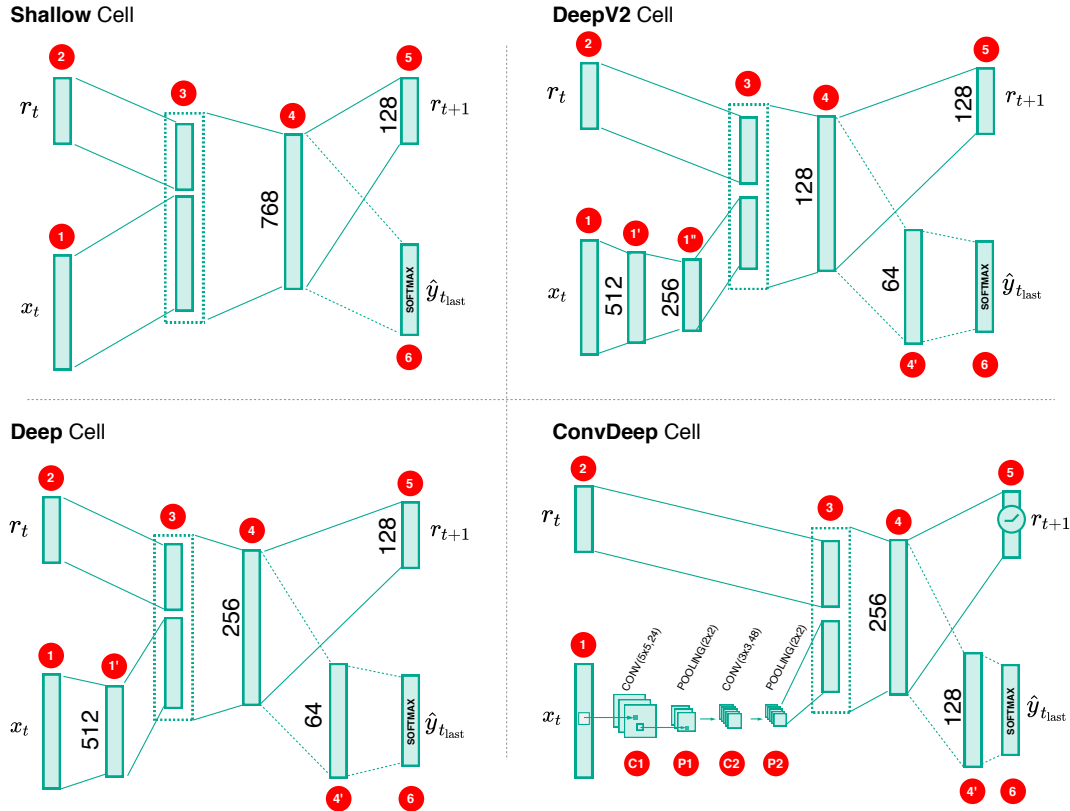


Figure 4.5: Number of neurons in each layer for each cell architecture

- **Shallow Cell**

$$\{n^{(4)}\} = \{768\}$$

- **Deep Cell**

$$\{n^{(1')}, n^{(4)}, n^{(4')}\} = \{512, 256, 64\}$$

- **DeepV2 Cell**

$$\{n^{(1')}, n^{(1'')}, n^{(4)}, n^{(4')}\} = \{512, 256, 128, 64\}$$

- **ConvDeep Cell :**

$$\{n^{(C1)}, n^{(P1)}\} = \{CONV(5 \times 5, 24), POOL(2 \times 2)\}$$

$$\{n^{(C2)}, n^{(P2)}\} = \{CONV(3 \times 3, 48), POOL(2 \times 2)\}$$

$$\{n^{(4)}, n^{(4')}\} = \{256, 128\}$$

where $CONV(x, y)$ is a convolutional operator with y filters whose kernel size is \mathbb{R}^x . Similarly, $POOL(x)$ is a pooling operator with kernel size \mathbb{R}^x .

Noting that, $n^{(5)}$ is set at 128 for all architectures and 0 when the sequence length of the problem is 1. $n^{(6)}$ is equal to the number of categories of a problem, for example $n^{(6)} = 10$ MNIST. Table 4.2 shows the total numbers of variables in details.

Cell Architecture	Sequence Length			
	1	4	7	14
Shallow	610570	355722	291210	248202
Deep	550346	314954	271946	243274
DeepV2	575050	306890	263882	235210
ConvDeep	647594	283178	197162	197162

Table 4.2: Total variables in each architecture and sequence length

Lastly, as the quality of relevance heatmap depending on performance of the model, the minimum classification accuracy is set as in Table 4.3.

Dataset	Minimum Accuracy
MNIST	0.98
Fashion-MNIST	0.85

Table 4.3: Classification Accuracy Criteria

4.4 Experiment 1 : Sequence Classification

4.4.1 Problem Formulation

To demonstrate how well RNNs can distribute relevant quantities to input space, I formulated an artificial classification problem in which each image sample \mathbf{x} is column-wise split into non-overlapping $(\mathbf{x}_t)_{t=1}^T$. The RNN classifier needs to summarize information from the sequence $(\mathbf{x}_t)_{t=1}^T$ to answer what is the class of \mathbf{x} .

Figure 4.6 illustrates the setting. Here, a MNIST sample $\mathbf{x} \in \mathbb{R}^{28,28}$ is divided to a sequence of $(\mathbf{x}_t \in \mathbb{R}^{28,7})_{t=1}^4$. At time step t , \mathbf{x}_t is presented to the RNN classifier which yields recurrent input \mathbf{r}_{t+1} for the next step. For the last step T , in this example $T = 4$, the RNN classifier computes $f(\mathbf{x}) \in \mathbb{R}^{10}$ and the class that \mathbf{x} belongs to.

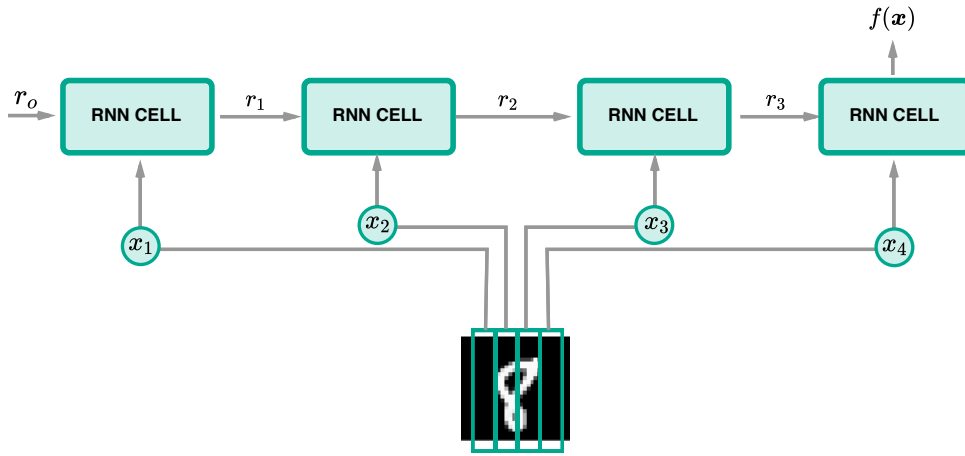


Figure 4.6: General setting of RNN classifiers in this thesis

Denote g_r and g_f function that the RNN with parameters θ uses to compute \mathbf{r}_{t+1} and $f(\mathbf{x})$ respectively. The overall computation can be summarized as follows:

$$\begin{aligned} \mathbf{r}_{t+1} &= g_r(\theta, \mathbf{x}_t, \mathbf{r}_t) \\ &\vdots \\ f(\mathbf{x}) &= g_f(\theta, \mathbf{x}_{t_{\text{last}}}, \mathbf{r}_{t_{\text{last}}}) \\ \hat{\mathbf{y}} &= \text{softmax}(f(\mathbf{x})), \end{aligned}$$

where $\hat{\mathbf{y}}$ is the class probabilities and $\mathbf{r}_0 = \mathbf{0}$. To explain the model, $R(\mathbf{x})$ is set to the value of $f(\mathbf{x})$ that is corresponding to the true target class.

TODO : Figure figure show how $R(\mathbf{x})$ is set?

4.4.2 Result

I began this preliminary experiment with Shallow and Deep architecture. They were trained on MNIST and FashionMNIST with sequence length $T = \{1, 4, 7\}$. Table 4.4

Dataset	Sequence Length		
	1	4	7
MNIST / FashionMNIST	$\mathbb{R}^{28,28}$	$\mathbb{R}^{28,7}$	$\mathbb{R}^{28,4}$

Table 4.4: Dimensions of \mathbf{x}_t for each dataset and sequence length

	Shallow	Deep
SEQ-1	xx.xx%	xx.xx%
SEQ-4	xx.xx%	xx.xx%
SEQ-7	xx.xx%	xx.xx%

Table 4.5: Model Accuracy

shows dimensions of \mathbf{x}_t for different sequence length and Table 4.5 summaries accuracy of the trained models. To simplify the manuscript, I am going to use *ARCHITECTURE-T* convention to denote a RNN with *ARCHITECTURE* trained on sequence length T . For example, Deep-7 refers to the Deep RNN architecture trained on $(\mathbf{x}_t \in \mathbb{R}^{28,4})_{t=1}^7$.

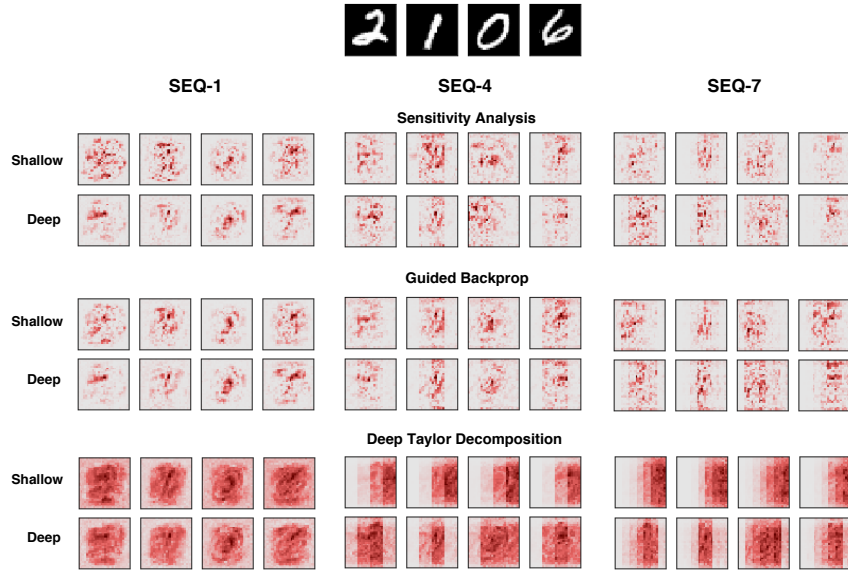


Figure 4.7: Relevance heatmaps from Shallow and Deep Cell trained on MNIST with different sequence lengths

Figure 4.7 shows relevance heatmaps from Shallow and Deep architecture trained on MNIST. We can see general characteristics of each explanation technique. In particular, sensitivity analysis(SA) and guided backprop(GB) heatmaps are sparse, while the ones from deep Taylor decomposition(DTD) are more diffuse throughout \mathbf{x} . When applying

these techniques to Shallow-1 and Deep-1, the relevance heatmaps look similar regardless of the architectures. As the sequence length is increased, SA and GB heatmaps are still almost identical for Shallow-4,7 and Deep-4,7. However, this is not the case for DTD. From the figure, we can see that Shallow-4,7 and Deep-4,7 return significantly different relevance heatmaps from DTD method. In particular, Shallow-4,7's heatmaps are mainly concentrated on the right part of \mathbf{x} associating to last time steps, while Deep-4,7's ones are appropriately highlighted at the actual content area of \mathbf{x} .

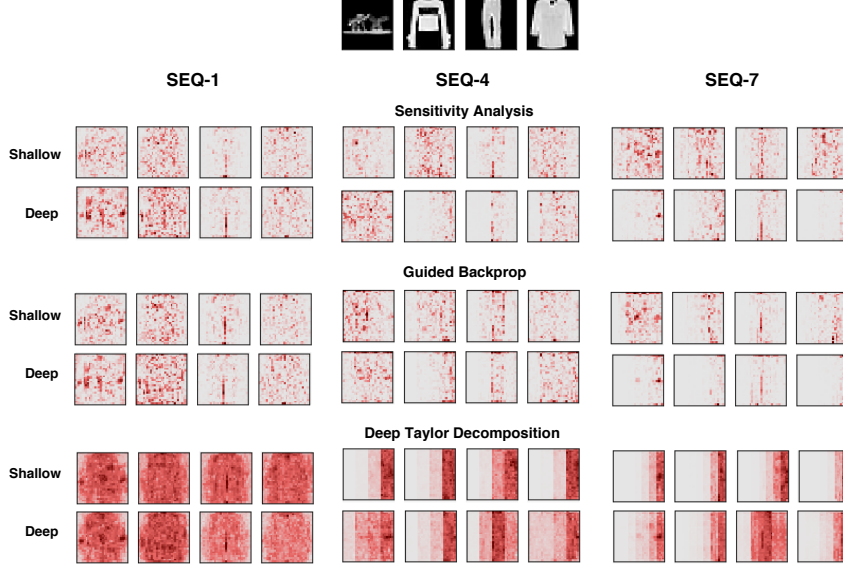


Figure 4.8: Relevance heatmaps from Shallow and Deep Cell trained on FashionMNIST with different sequence lengths

Relevance heatmaps of Shallow and Deep architecture trained on FashionMNIST are shown on Figure 4.8. Similar to MNIST, we do not see any remarkable difference on SA and GB heatmaps of the two architectures : only that Deep-4,7 produces slightly more sparse heatmaps. Although the wrong concentration issue of DTD seems to appear on both Shallow-4,7's and Deep-4,7's heatmaps, we still can observe proper highlight from Deep architecture on some samples. For example, the trouser sample, we can see that Deep-4,7 architecture manage to distribute high relevance scores to area of the trouser. Latent features of FashionMNIST might be one of the reasons why Deep architecture does not distribute relevance scores to early steps for the FashionMNIST samples. Consider *Shoe* and *Ankle Boot* samples in Figure 4.2. One can see that their front part are similar and only the heel part that determines the differences between the two categories.

Figure 4.9 presents relevance heatmaps of MNIST *Class 1* and FashionMNIST *Class Trouser* samples. These samples were chosen to emphasize the impact of RNN architecture on DTD explanation. In particular, these samples have $\mathbf{x}_{t'}$ containing features primarily locating at the center, or middle of the sequence, hence appropriate relevance heatmaps should be highlighted at $\mathbf{x}_{t'}$ and possibly its neighbors. As expected, we

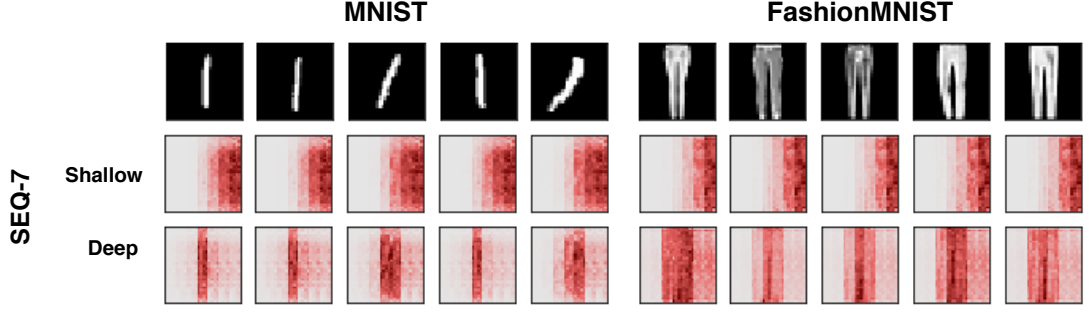


Figure 4.9: DTD relevance heatmaps of MNIST *Class 1* and FashionMNIST *Class Trouser* samples from Shallow-7 and Deep-7.

can see Deep-7 produces sound results in which the heatmaps have high intensity value where $\mathbf{x}_{t'}$ approximately locate, while Shallow-7 mainly assigns relevance quantities to \mathbf{x}_t for $t \approx T$. Figure 4.10 further shows a quantitative evidence of this problem. Here, distribution of relevance score from Shallow-7 and Deep-7 are plotted across time step $t = \{1, \dots, 7\}$. The distributions are computed from all test samples in MNIST *Class 1* and FashionMNIST *Class Trouser* respectively as well as the data distributions that are derived from pixel intensity values. We can see that Deep-7's relevance distributions are similar to the data distributions, while Shallow-7's ones diverge completely. Particularly, one can see that Shallow-7 distributes more than 90% of relevance scores to the last 3 steps, namely \mathbf{x}_5 , \mathbf{x}_6 and \mathbf{x}_7 .

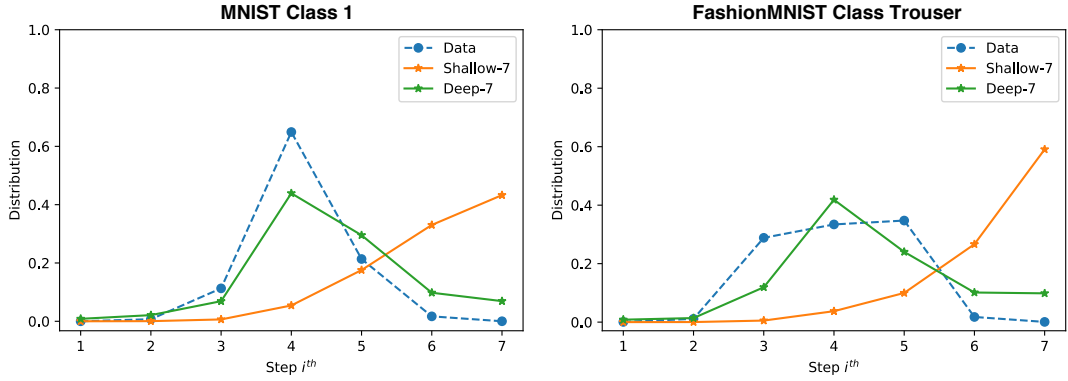


Figure 4.10: Pixel intensity and DTD relevance distribution from Shallow-7 and Deep-7 averaged over MNIST *Class 1* and FashionMNIST *Class Trouser* test population.

4.4.3 Summary

Results from this first experiment seem to suggest that choice of RNN architectures has an impact on quality of its explanation. In particular, as presented in Figure 4.9 and Figure 4.10, quality of deep Taylor decomposition (DTD) explanation is significantly influenced by the architecture. In contrast, we do see such notable effect from sensitivity analysis (SA) and guided backprop (GB) method. In the following experiment, I am going to present a methodical evaluation of this impact in detail.

4.5 Experiment 2 : Majority Sample Classification

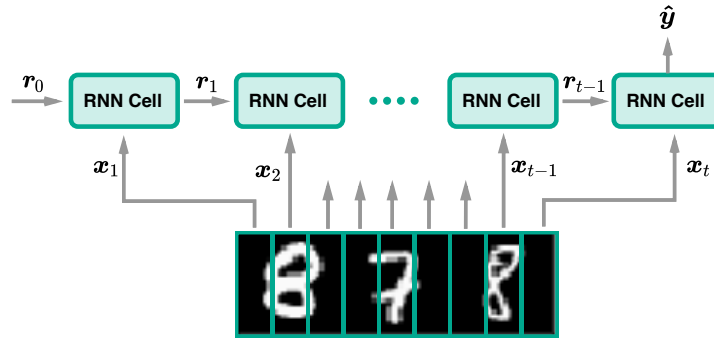


Figure 4.11: Majority Sample Classification Problem

4.5.1 Problem Formulation

When neural network are trained, one can apply explantation techniques to the models to get relevance heatmaps of samples. The heatmap of sample \mathbf{x} illustrates important features in \mathbf{x} that the trained network utilizes to perform its objective prediction, such as classification. Therefore, one needs to know the ground truth of these latent features in order to methodologically evaluate how well the model distributes relevance scores to \mathbf{x} 's input space. However, this knowledge is not trivial to find as it is an incident from the optimization process in high-dimensional space that we in fact seek to understand.

To alleviate this challenge, I constructed two artificial datasets using samples from MNIST and FashionMNIST respectively. Let's consider MNIST. The new dataset was constructed as follows: each original sample $\tilde{\mathbf{x}} \in \mathbb{R}^{28,28}$, I randomly selected 2 additional samples : one from the same class of $\tilde{\mathbf{x}}$ and the other one from a different class. Then, these 3 samples are concatenated in random order yielding a sample $\mathbf{x} \in \mathbb{R}^{28,84}$ of the new dataset.

With this construction procedure, one possible objective is to train RNNs to predict the majority class in each sequence $(\mathbf{x}_t)_{t=1}^T$. For example, consider $\mathbf{x} = \{8, 7, 8\}$ shown in Figure 4.11, the majority class here is *Class 8*. Because we already know time step t' that are corresponding to original samples $\tilde{\mathbf{x}}$ belonging to the majority group in each sequence \mathbf{x} , we can simply compute percentage of relevance scores assigned to t' to quantitatively evaluate explainability of a RNN architecture.

I also introduces two additional variation of Deep architecture, namely DeepV2 and ConvDeep. DeepV2 has one more layer after the first fully-connected layer, while ConvDeep instead employs a sequence of convolutional and pooling operation after the input layer. Figure X shows details of the architectures.

4.5.2 Result

Cell Architecture	No. variables	Accuracy	
		MNIST	FashionMNIST
Shallow	184330	98.39%	90.15%
Deep	153578	98.42%	90.60%
DeepV2	161386	98.38%	90.43%
ConvDeep	151802	99.07%	92.43%

Table 4.6: Number of trainable variables and model accuracy for Majority Sample Classification experiment.

In this experiment, I chose $T = 12$, hence $(\mathbf{x}_t \in \mathbb{R}^{28,7})_{t=1}^{12}$. Table 4.6 shows number of trainable variables and accuracy of the trained models. These trained models have equivalent number of variables and accuracy, hence comparing heatmaps of these models is a fair evaluation.

As can be seen from Figure 4.12, the deeper architecture, the fewer relevant scores distributed to irrelevant region. This effect happens across all explanation methods. This result further supports what has shown in Experiment 1. Moreover, although relevance heatmaps from Shallow-12, Deep-12, and DeepV2-12 generally look noisy, increasing the depth of architecture seems to reduce the noise in the heatmaps. On the other hand, ConvDeep-12 does not only properly assign relevance quantities to the right time steps, but its heatmaps are also sound. In particular, one can easily perceive features of \mathbf{x} from ConvDeep-12's GB and DTD heatmaps.

Figure 4.13 presents a quantitive evaluation of the impact from the depth of architecture to the explanation. Here, the measurement is the percentage of relevance quantities assigned to regions of majority class : these regions are indicated by the red triangles in Figure 4.12, averaged over test sample population. Results from Figure 4.13 indicate that the depth of architecture indeed improves quality of the explanations. In particular, the percentage of correct relevance assignment of each explanation technique increases as more layers introduced. This enhancement can be seen clearly from the result of Fash-

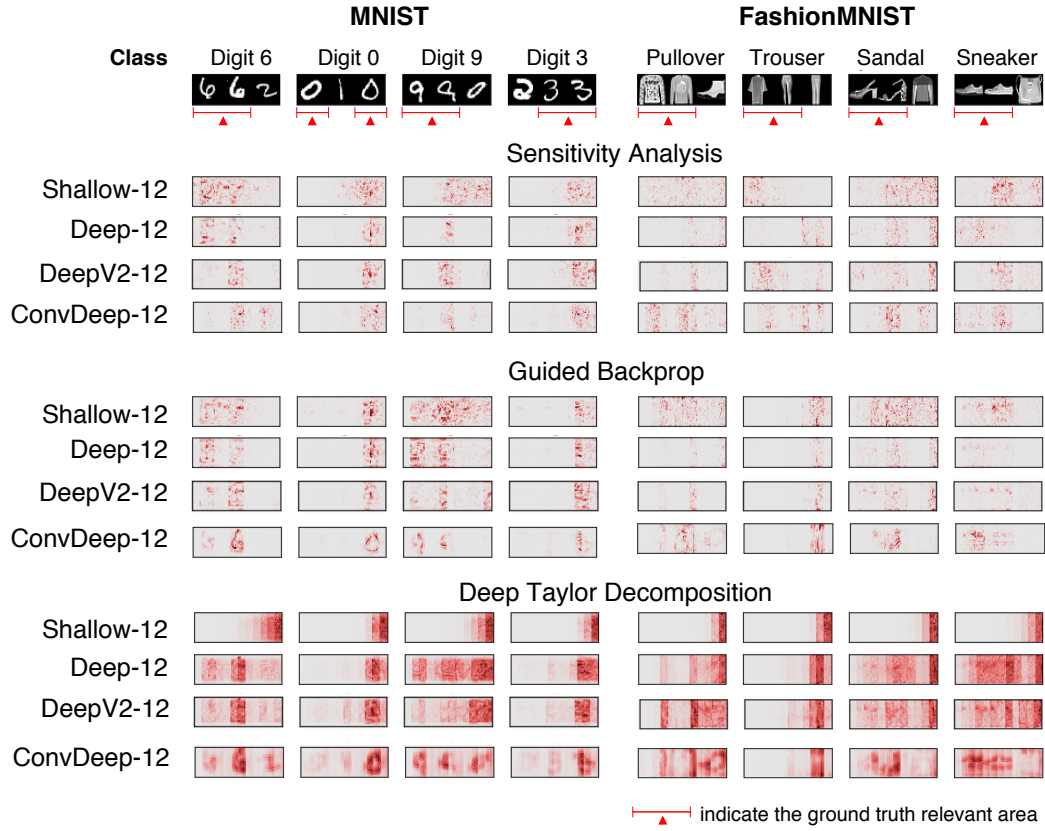


Figure 4.12: Relevance heatmaps for MSC problem

ionMNIST. Moreover, as expected, ConvDeep-12 improves the percentage even more. Additionally, we can observe that the difference between the percentage of the baseline, Shallow-12, and the deep architectures changes with different proposition across methods. In particular, we see the difference of SA and DTD are slightly larger than the difference of GB. This implies that some explanation methods get more benefit from the depth of architecture.

4.5.3 Summary

The outcome of this experiment quantitatively confirms that the depth of architecture has impacts on explanation of the model. It also shows that the depth of architecture affects explanation in different level on different methods. More precisely, comparing to guided backprop(GB), quality of sensitivity analysis(SA) and deep Taylor decomposition(DTD) explanation is more depend on the depth.

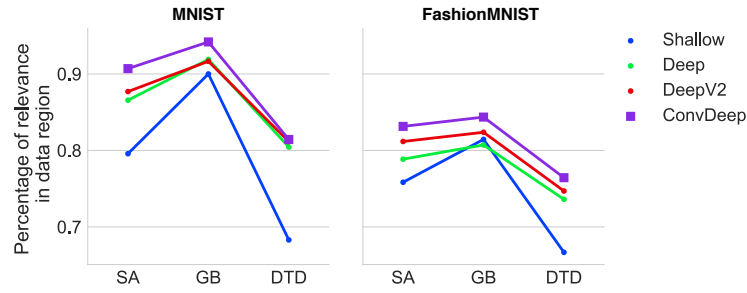


Figure 4.13: Percentage of relevance in data region

4.6 Experiment 3 : Improve Relevance Distribution

4.6.1 Proposal 1 : Gating units

TODO : Figure : R-LSTM

4.6.2 Proposal 2 : Stationary Dropout

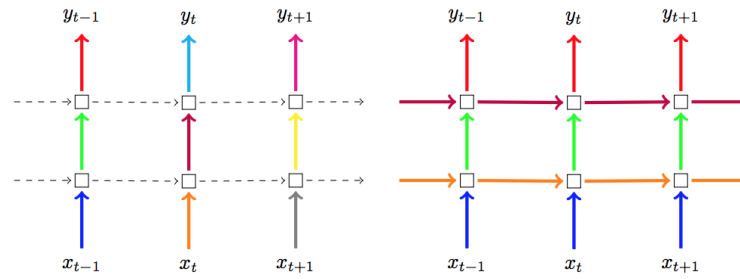


Figure 4.14: Variational Dropout[Gal and Ghahramani, 2016]

4.6.3 Proposal 3 : Convolutional layers with literal connections

4.6.4 Result

4.6.5 Summary

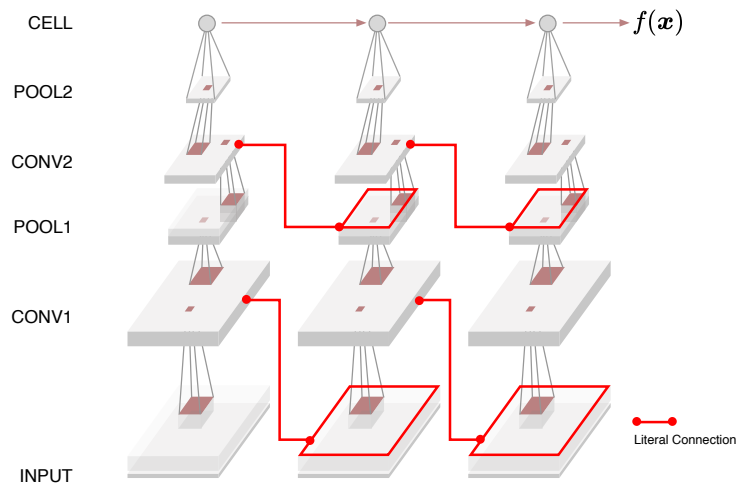


Figure 4.15: ConvDeep with Literal connections

5 Conclusion

5.1 Summary

5.2 Challenges

5.3 Future work

References

- [Abadi et al., 2016] Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mane, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viegas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2016). TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems.
- [Bach et al., 2016] Bach, S., Binder, A., Montavon, G., Muller, K.-R., and Samek, W. (2016). Analyzing classifiers: Fisher vectors and deep neural networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 2912–2920.
- [Binder et al., 0906] Binder, A., Montavon, G., Lapuschkin, S., Müller, K.-R., and Samek, W. (2016/09/06). Layer-Wise Relevance Propagation for Neural Networks with Local Renormalization Layers. In Artificial Neural Networks and Machine Learning – ICANN 2016, Lecture Notes in Computer Science, pages 63–71. Springer, Cham.
- [Cho et al., 2014] Cho, K., van Merriënboer, B., Gülçehre, Ç., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1724–1734. Association for Computational Linguistics.
- [Erhan et al., 2010] Erhan, D., Courville, A., and Bengio, Y. (October 2010). Understanding Representations Learned in Deep Architectures.
- [Gal and Ghahramani, 2016] Gal, Y. and Ghahramani, Z. (2016). A Theoretically Grounded Application of Dropout in Recurrent Neural Networks. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R., editors, Advances in Neural Information Processing Systems 29, pages 1019–1027. Curran Associates, Inc.
- [Greff et al., 2017] Greff, K., Srivastava, R. K., Koutník, J., Steunebrink, B. R., and Schmidhuber, J. (2017). LSTM: A search space odyssey. IEEE transactions on neural networks and learning systems.
- [He et al., 2015] He, K., Zhang, X., Ren, S., and Sun, J. (2015). Deep Residual Learning for Image Recognition. CoRR, abs/1512.03385.

- [Hochreiter and Schmidhuber, 1997] Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. Neural computation, 9(8):1735–1780.
- [Jia et al., 2014] Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., and Darrell, T. (2014). Caffe: Convolutional Architecture for Fast Feature Embedding.
- [Jozefowicz et al., 2015] Jozefowicz, R., Zaremba, W., and Sutskever, I. (2015). An empirical exploration of recurrent network architectures. Journal of Machine Learning Research.
- [Kingma and Ba, 2014] Kingma, D. P. and Ba, J. (2014). Adam: A Method for Stochastic Optimization. In Proceedings of the 3rd International Conference on Learning Representations (ICLR).
- [Krizhevsky et al., 2012] Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). ImageNet Classification with Deep Convolutional Neural Networks. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, Advances in Neural Information Processing Systems 25, pages 1097–1105. Curran Associates, Inc.
- [LeCun et al., 2001] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (2001). Gradient-Based Learning Applied to Document Recognition. In Intelligent Signal Processing, pages 306–351. IEEE Press.
- [LeCun and Cortes, 2010] LeCun, Y. and Cortes, C. (2010). MNIST handwritten digit database.
- [Montavon et al., 2017a] Montavon, G., Lapuschkin, S., Binder, A., Samek, W., and Müller, K.-R. (May 1, 2017a). Explaining nonlinear classification decisions with deep Taylor decomposition. Pattern Recognition, 65:211–222.
- [Montavon et al., 2017b] Montavon, G., Samek, W., and Müller, K.-R. (2017b). Methods for Interpreting and Understanding Deep Neural Networks.
- [Pascanu et al., 2012] Pascanu, R., Mikolov, T., and Bengio, Y. (2012). Understanding the exploding gradient problem. CoRR, abs/1211.5063.
- [Simonyan et al., 2013] Simonyan, K., Vedaldi, A., and Zisserman, A. (2013). Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. CoRR, abs/1312.6034.
- [Simonyan and Zisserman, 2014] Simonyan, K. and Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition. CoRR, abs/1409.1556.
- [Springenberg et al., 2014] Springenberg, J. T., Dosovitskiy, A., Brox, T., and Riedmiller, M. A. (2014). Striving for Simplicity: The All Convolutional Net. CoRR, abs/1412.6806.

- [Szegedy et al., 2014] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. (2014). Going Deeper with Convolutions.
- [Tieleman and Hinton, 2012] Tieleman, T. and Hinton, G. (2012). Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude.
- [Xiao et al., 2017] Xiao, H., Rasul, K., and Vollgraf, R. (2017). Fashion-MNIST: A Novel Image Dataset for Benchmarking Machine Learning Algorithms.
- [Zagoruyko and Komodakis, 2016] Zagoruyko, S. and Komodakis, N. (2016). Wide Residual Networks.
- [Zeiler, 2012] Zeiler, M. D. (2012). ADADELTA: An Adaptive Learning Rate Method. CoRR, abs/1212.5701.
- [Zeiler and Fergus, 2013] Zeiler, M. D. and Fergus, R. (2013). Visualizing and Understanding Convolutional Networks. CoRR, abs/1311.2901.