

(Draft : January 13, 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

Thesis Supervisor
Prof. Dr. Klaus-Robert Müller

Thesis Advisor
Dr. Grégoire Montavon

~~First of all I would like to thank Prof. Dr. Thomas Magedanz at the Fraunhofer Institute FOKUS for giving me the opportunity to carry out state of the art research in this field.~~

~~Special thanks to Mister X and Mister Y for their guidance. — Some personal words... —~~

~~Furthermore I would like to thank ...~~

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	xii
List of Figures	xv
List of Tables	xvii
1 Introduction	1
2 Related Work	3
3 Background	5
3.1 Neural Networks	5
3.1.1 Learning Algorithms	5
3.1.2 CNNs	5
3.1.3 RNNs	5
3.2 Explainability of Neural Networks	5
3.2.1 Sensitivity Analysis	7
3.2.2 Simple Taylor Decomposition	7
3.2.3 Layer-wise Relevance Propagation	8
3.3 Deep Taylor Decomposition	10
4 Experiments	17
4.1 Dataset	17
4.1.1 MNIST	17
4.1.2 Fashion-MNIST	17
4.2 The Problem	18
4.3 RNN Cell Architectures	20
4.3.1 Shallow Cell	20
4.3.2 Deep Cell Architecture	21
4.4 Setup	21
4.5 Results	24
4.5.1 MNIST	24
4.5.2 Fashion-MNIST	26
References	27

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	RNN Structure	5
3.2	Comparison between Global and Local Analysis	6
3.3	The LOF caption	9
3.4	An illustration of R_k functional view and root point candidates	12
3.5	R_k functional view and root points from z^+ -rule	13
3.6	R_k functional view and root points from z^β -rule with $-1 < a_j < 1$	15
4.1	MNIST Dataset	17
4.2	Fashion-MNIST Dataset	18
4.3	General setting of RNN classifiers in this thesis	19
4.4	Shallow and Deep Cell Architecture	20
4.5	DeepV2 and ConvDeep Cell Architecture	21
4.6	Number of neurons in each layer for each cell architecture	23
4.7	Comparison of relevance heatmaps from LRP between Shallow and Deep Cell	25
4.8	Relevance heatmaps Shallow and Deep Cell (The samples are randomly chosen from MNIST's testing set)	25
4.9	Comparison of pixel intensity distribution from MNIST Class 1 testing population and relevance distribution explained by the models	26

List of Tables

3.1	LRP rules from Deep Taylor Decomposition	15
4.1	Dimensions of \mathbf{x}_t for each dataset and sequence length	20
4.2	Hyperparameter Summary	22
4.3	Total variables in each architecture and sequence length	24
4.4	Classification Accuracy Criteria	24
4.5	Model Accuracy	24

1 Introduction

2 Related Work

3 Background

3.1 Neural Networks

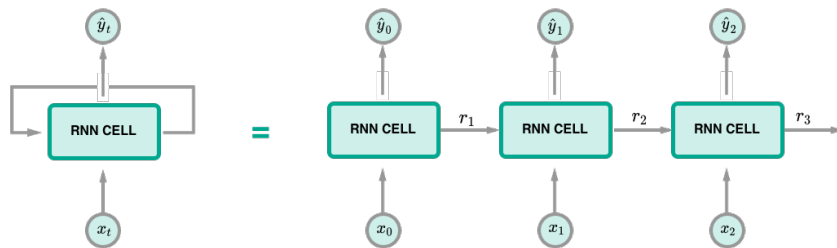


Figure 3.1: RNN Structure

3.1.1 Learning Algorithms

Loss functions

SGD

ADAM

Backpropagation Through Time

Vanishing and exploding gradient

3.1.2 CNNs

3.1.3 RNNs

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.[?] 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[?]” 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?



Figure 3.2: Comparison between Global and Local Analysis

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

Definition 3.2.2. Positivity Property

Definition 3.2.3. Consistency

3.2.1 Sensitivity Analysis

Sensitivity analysis 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,

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

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[?], the derivation 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 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.[?] 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

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

Hence, the decomposition can be simplified to :

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

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.

3.2.3 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.[?] 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 3.3. R_j and R_k are relevance score of neurons j, k in successive layers. [?] formulates the general form of relevance propagation as :

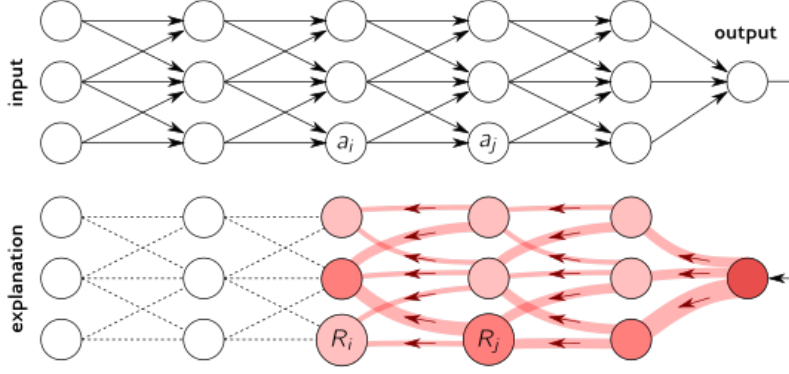
$$R_j = \sum_k \delta_{j \leftarrow k} R_k, \quad (3.1)$$

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. [?] suggests

$$\delta_{j \leftarrow k} = \alpha \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} - \beta \frac{a_j w_{jk}^-}{\sum_j a_j w_{jk}^-}, \quad (3.2)$$

Figure 3.3: LRP Framework¹

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 3.1, [?] calls $\alpha\beta$ -rule.

$$R_j = \sum_k \left(\alpha \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} - \beta \frac{a_j w_{jk}^-}{\sum_j a_j w_{jk}^-} \right) R_k \quad (3.3)$$

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 rewrite $\alpha\beta$ -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

“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 [?]

¹Source: <http://heatmapping.org>

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

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

```

3.3 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.[?] proposes the method to explain decisions of neural networks 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 :

$$R_k = R_k \Big|_{\tilde{a}_j} + \sum_j \frac{\partial R_k}{\partial a_j} \Big|_{a_j=\tilde{a}_j} (a_j - \tilde{a}_j) + \zeta_k \quad (3.4)$$

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 3.4 is simplified to

$$R_k = \sum_j \underbrace{\frac{\partial R_k}{\partial a_j} \Big|_{a_j=\tilde{a}_j} (a_j - \tilde{a}_j)}_{R_{j \leftarrow k}} \quad (3.5)$$

Moreover, as the relevance propagated back, [?] 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

$$R_j = \sum_k R_{j \leftarrow k} \quad (3.6)$$

Combining Equation 3.5 and 3.6 yields :

$$\begin{aligned}
R_j &= \sum_k R_{j \leftarrow k} \\
\sum_j R_j &= \sum_j \sum_k R_{j \leftarrow k} \\
\sum_j R_j &= \sum_k \sum_j R_{j \leftarrow k} \\
\sum_j R_j &= \sum_k R_k
\end{aligned} \tag{3.7}$$

Demonstrated by [?], Equation 3.7 holds for all j, k and all subsequent layers. Hence, this results in conservation property which guarantee that no relevance loss during the propagations.

$$\sum_i R_i = \dots = \sum_j R_j = \sum_k R_k = \dots = f(\mathbf{x}) \tag{3.8}$$

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

$$R_k = \max \left(0, \sum_j a_j w_{jk} + b_k \right), \tag{3.9}$$

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 .

$$\tilde{a}_j = a_j - tv_j, \tag{3.10}$$

The root point is then the intersection point between Equation 3.10 and 3.9. Hence,

$$0 = \sum_j (a_j - tv_j) w_{jk} + b_k \tag{3.11}$$

$$t \sum_j v_j w_{jk} = R_k \tag{3.12}$$

$$t = \frac{R_k}{\sum_j v_j w_{jk}} \tag{3.13}$$

$$\tag{3.14}$$

Therefore, R_j can be computed by :

$$R_j = \sum_k \left. \frac{\partial R_k}{\partial a_j} \right|_{a_j - \tilde{a}_j} (a_j - \tilde{a}_j) \quad (3.15)$$

$$= \sum_k w_{jk} t v_j \quad (3.16)$$

$$= \sum_k \frac{v_j w_{jk}}{\sum_j v_j w_{jk}} R_k \quad (3.17)$$

Noticing here is that \tilde{a}_j is not necessary the closest point to the line $R_k = 0$ in Euclidean distance, because \tilde{a}_j might not be in the same domain as a_j , hence v_j needs to be chosen according to the domain of a_j . Consider an example on Figure 3.4, if $a_j \in \mathbb{R}^+$, then \tilde{a}_j must be also in \mathbb{R}^+ . In the following, I will summarize how a_j can be computed for each domain of a_j .

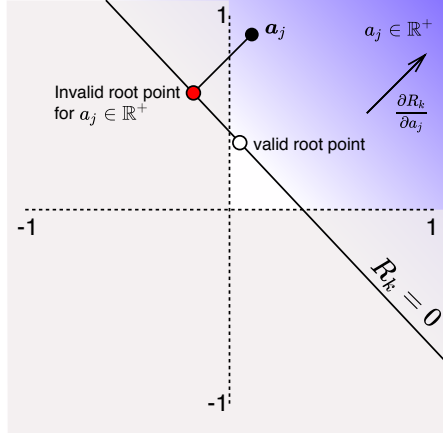


Figure 3.4: An illustration of R_k functional view and root point candidates

Case $a_j \in \mathbb{R} : w^2$ -rule

Trivially, the search direction v_j is just the direction of gradient $\frac{\partial R_k^{(l+1)}}{\partial a_j^{(l)}}$:

$$v_j = w_{jk}$$

Hence,

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

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 3.5, R_k has a root at $\mathbf{a}_j \mathbb{1}\{w_{jk} < 0\}$, because of:

$$R_k = \max \left(\sum_j a_j w_{jk} + b_k, 0 \right) \quad (3.18)$$

$$= \max \left(\sum_j a_j \mathbb{1}\{w_{jk} < 0\} w_{jk} + b_k, 0 \right) \quad (3.19)$$

$$= \max \left(\sum_j a_j w_{jk}^- + b_k, 0 \right) \quad (3.20)$$

$$= 0 \quad (3.21)$$

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

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

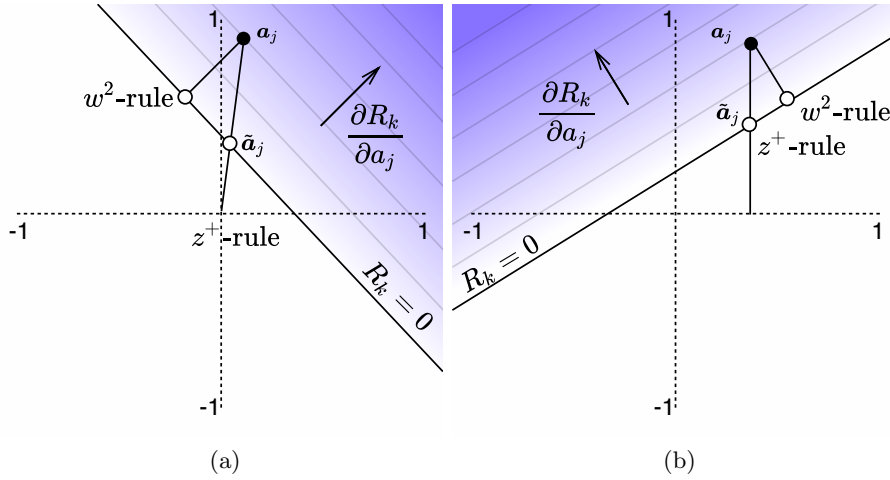


Figure 3.5: R_k functional view and root points from z^+ -rule

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 $\alpha = 1, \beta = 0$.

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

$$R_k = \max \left(\sum_j a_j w_{jk} + b_k, 0 \right) \quad (3.22)$$

$$= \max \left(\sum_j (l_j \mathbb{1}\{w_{jk} \geq 0\} + h_j \mathbb{1}\{w_{jk} \leq 0\}) w_{jk} + b_k, 0 \right) \quad (3.23)$$

$$= \max \left(\sum_j l_j w_{jk}^+ + h_j w_{jk}^- + b_k, 0 \right) \quad (3.24)$$

$$= 0 \quad (3.25)$$

Hence, the search direction is

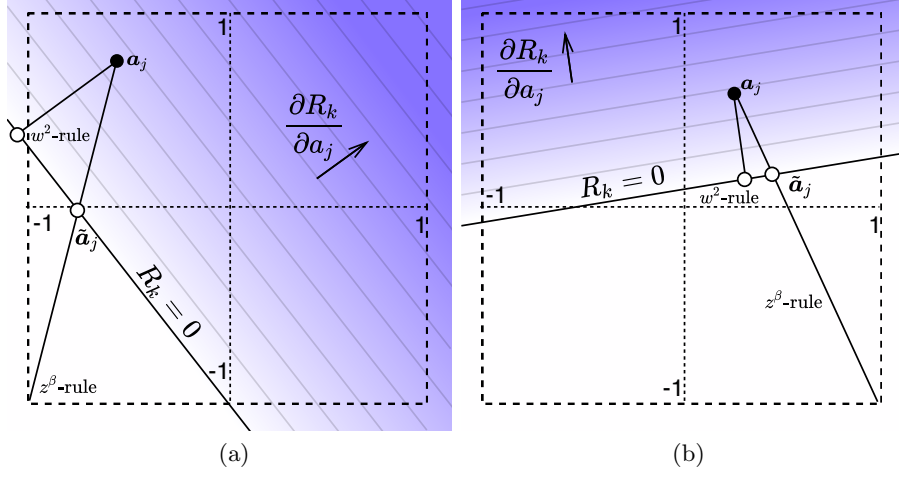
$$\begin{aligned} v_j &= a_j - \tilde{a}_j \\ &= a_j - l_j \mathbb{1}\{w_{jk} \geq 0\} - h_j \mathbb{1}\{w_{jk} \leq 0\} \end{aligned}$$

Figure 3.6 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 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.

Figure 3.6: R_k functional view and root points from z^β -rule with $-1 < a_j < 1$

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[?] 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.

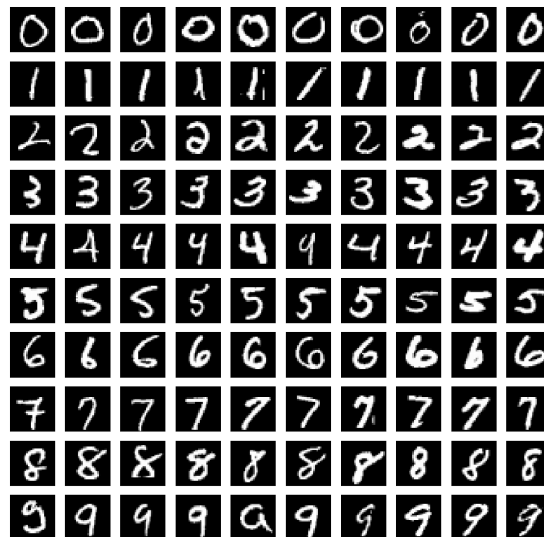


Figure 4.1: MNIST Dataset

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[?].

4.1.2 Fashion-MNIST

Xiao et. al.[?] propose a novel dataset, called Fashion-MNIST dataset, as a replacement of MNIST dataset for benchmarking machine learning algorithms. According to [?], 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.

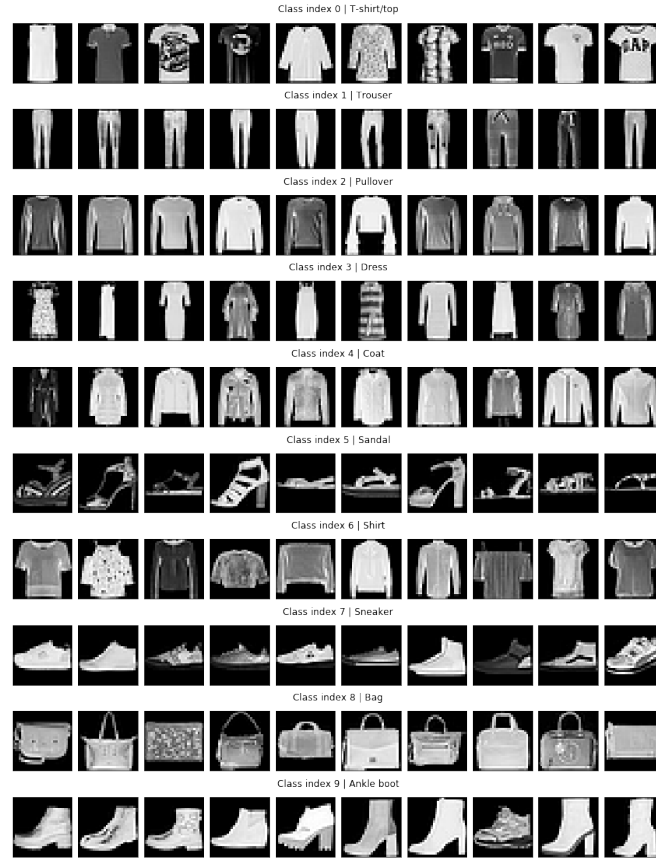


Figure 4.2: Fashion-MNIST Dataset

[?] 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)[?] and standard data preprocessing and augmentation.

4.2 The Problem

Because the main objective of the thesis is to demonstrate how well vanilla RNNs can propagate relevant quantity back to input space, we propose an artificial classification problem in which each sample is split into chunks and fed sequentially to the classifier.

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

The classifier needs to summarize information along the sequence to make the final decision at the end.

Figure 4.3 illustrates a general setting to the proposed problem. Here, a MNIST sample $\mathbf{x} \in \mathbb{R}^{28,28}$ is column-wise split into a sequence of $\{\mathbf{x}_t\}_{t=0}^3, \mathbf{x}_t \in \mathbb{R}^{28,7}$. For a step t , \mathbf{x}_t is fed to a RNN cell, whose parameters denoted by θ , yielding recurrent input \mathbf{r}_{t+1} for the next step. For the last step $t_{\text{last}} = 3$, we compute the final decision to the classification problem $\hat{\mathbf{y}} \in \mathbb{R}^{10}$.

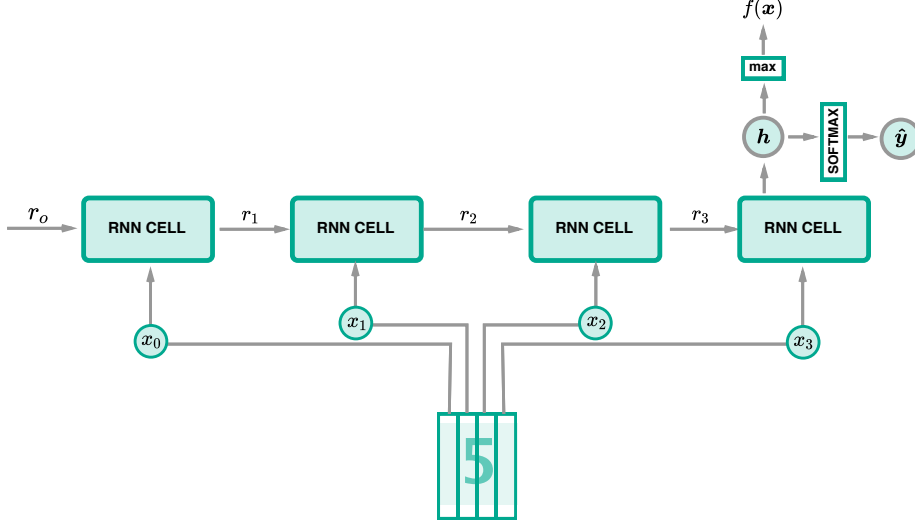


Figure 4.3: General setting of RNN classifiers in this thesis

Denote g_r and g_h a function that a RNN cell uses to compute \mathbf{r}_{t+1} and raw output $\mathbf{h} \in \mathbb{R}^{10}$ at the last step respectively. The whole computation can be summarized as follows:

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

where $\hat{\mathbf{y}}$ is the class prediction probabilities and $\mathbf{r}_0 = \mathbf{0}$. In this case, the relevance score $f(\mathbf{x})$ that will be propagated is:

$$f(\mathbf{x}) = \max_i h_i$$

In this thesis, I performed experiments using sequence length $SEQ = \{1, 4, 7, 14\}$. Table 4.1 shows dimensions of \mathbf{x}_t for different sequence length.

Dataset	Sequence Length			
	1	4	7	14
MNIST / Fashion-MNIST	$\mathbb{R}^{28,28}$	$\mathbb{R}^{28,7}$	$\mathbb{R}^{28,4}$	$\mathbb{R}^{28,2}$

Table 4.1: Dimensions of \mathbf{x}_t for each dataset and sequence length

4.3 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
- \mathbf{x}_t : subset of $x_i \in \mathbf{x}$ corresponding to step t and assume to be reshaped into a column vector

4.3.1 Shallow Cell

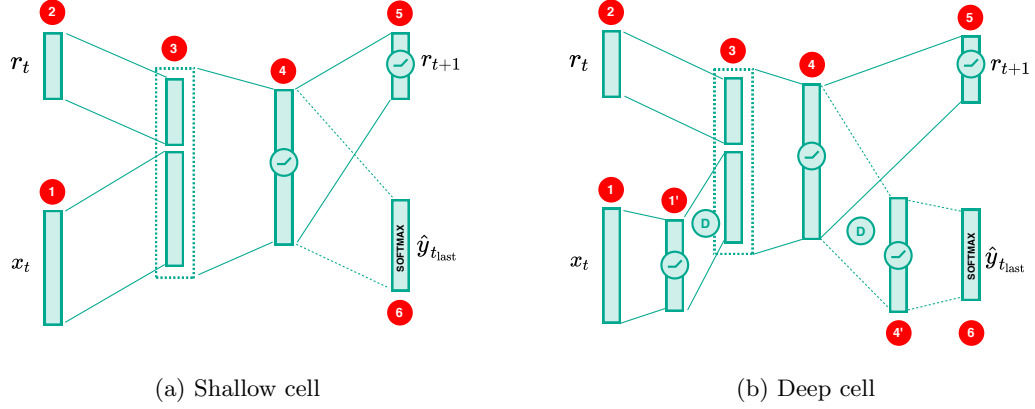


Figure 4.4: Shallow and Deep Cell Architecture

As shown in Figure 4.4a, 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).

4.3.2 Deep Cell Architecture

Figure 4.4b 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 $\mathbf{a}^{(4')}$.

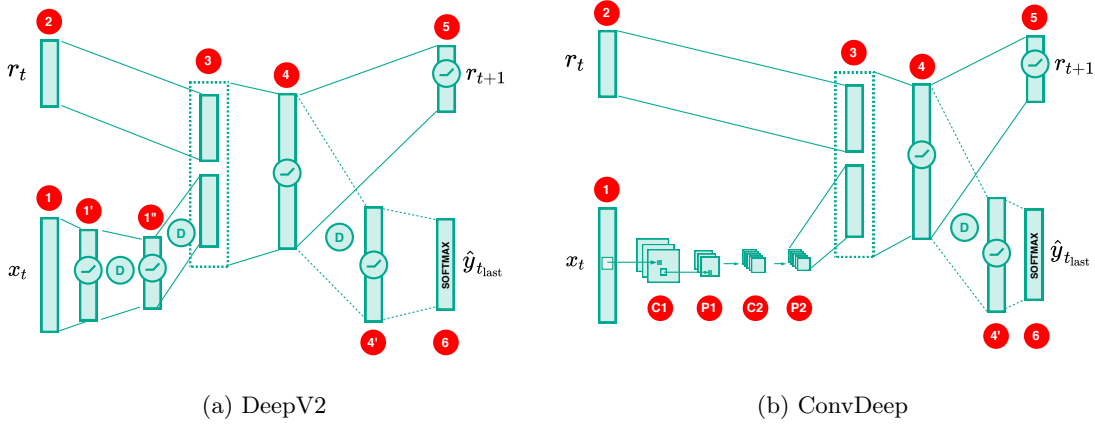


Figure 4.5: DeepV2 and ConvDeep Cell Architecture

Two variations of Deep cell are also experimented, namely DeepV2 and ConvDeep, shown on Figure 4.5. 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, $[\mathbf{C1}, \mathbf{P1}]$ and $[\mathbf{C2}, \mathbf{P2}]$.

4.4 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 $[-2\sigma, 2\sigma]$.

²<http://tensorflow.org/>

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

$$\begin{aligned}\mathbf{h}^{(l)} &= \mathbf{W}^T \mathbf{a}^{(l-1)} - \sigma_s(\mathbf{b}) \\ \mathbf{a}^{(l)} &= \sigma_r(\mathbf{h}^{(l)})\end{aligned}$$

where

$$\begin{aligned}\sigma_r(h_j) &= \max(0, h_j) && \text{(ReLU function)} \\ \sigma_s(b_j) &= \log(1 + \exp b_j) && \text{(Softplus function)}\end{aligned}$$

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[?] 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 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 4.2 summaries the setting of hyperparameters.

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

Table 4.2: 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 4.6 illustrates the details of the settings.

- **Shallow Cell**

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

- **Deep Cell**

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

- **DeepV2 Cell**

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

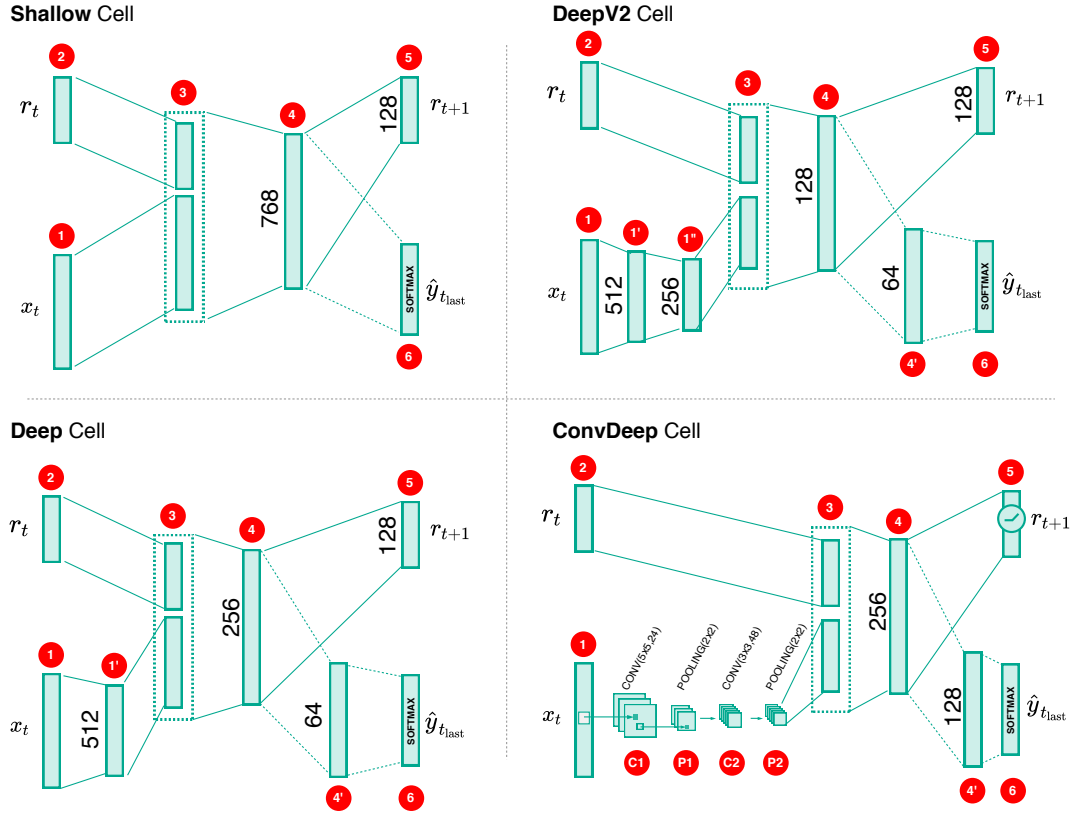


Figure 4.6: Number of neurons in each layer for each cell architecture

- **ConvDeep Cell :**

$$\begin{aligned} \{n^{(C1)}, n^{(P1)}\} &= \{CONV(5 \times 5, 24), POOL(2 \times 2)\} \\ \{n^{(C2)}, n^{(P2)}\} &= \{CONV(3 \times 3, 48), POOL(2 \times 2)\} \\ \{n^{(4)}, n^{(4')}\} &= \{256, 128\} \end{aligned}$$

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 4.3 shows the total numbers of variables in details.

Lastly, as the quality of relevance heatmap depending on performance of the model, the minimum classification accuracy is set as in Table 4.4.

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 4.3: Total variables in each architecture and sequence length

Dataset	Minimum Accuracy
MNIST	0.98
Fashion-MNIST	0.85

Table 4.4: Classification Accuracy Criteria

4.5 Results

In this section, I am going to present what I have found from experiments. First, I am going to discuss results from MNIST and then move on to ones from Fashion-MNIST. Moreover, I will use CELL_NAME-SEQ convention to denote a RNN network with CELL_NAME cell trained on a sequence length SEQ. For example, Deep-7 means a RNN with Deep cell architecture trained on data whose $\mathbf{x}^{(\alpha)}$ is a sequence of length 7.

4.5.1 MNIST

As a starting point, I first trained classifiers using Shallow and Deep cell with $SEQ = \{1, 4, 7\}$ on MNSIT using setting described in Section 4.4. As shown in Table 4.5, these models are equivalent in term of accuracy performance, hence their heatmaps can be compared.

Figure 4.7 demonstrates relevance heatmaps from LRP and the trained models. Both Shallow-1 and Deep-1 produce sound heatmaps, noting here that $\mathbf{x}_t \in \mathbb{R}^{28,28}$. For $SEQ = \{4, 7\}$, their heatmaps are significantly different. Deep-4,7’s heatmaps contains high intensity values around the center of the images, while Shallow-4,7’s ones mostly concentrate on the right of the images. This implies that Shallow-4,7 are not able to

	Shallow	Deep
SEQ-1	0.9813	0.9788
SEQ-4	0.9830	0.9861
SEQ-7	0.9807	0.9852

Table 4.5: Model Accuracy

propagate relevance quantities through recurrent channels.

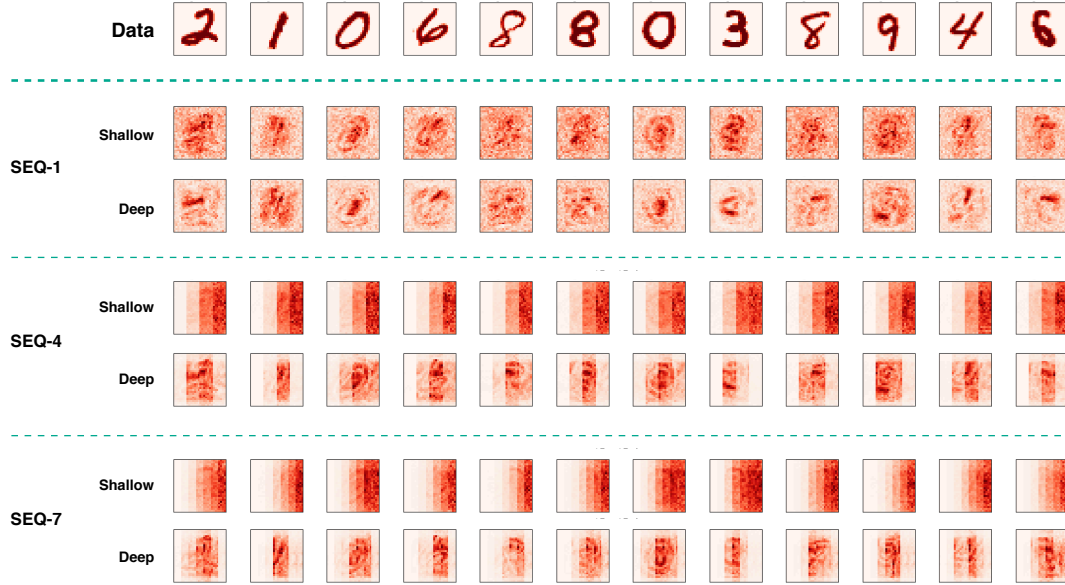


Figure 4.7: Comparison of relevance heatmaps from LRP between Shallow and Deep Cell

Moreover, we could observe that Deep-4,7 are not only able to propagate relevance back through the input sequence, it also appropriately captures features of the input as the corresponding shapes of the samples are present on the heatmaps. On the other hand, Shallow-4,7 do not show such capabilities.

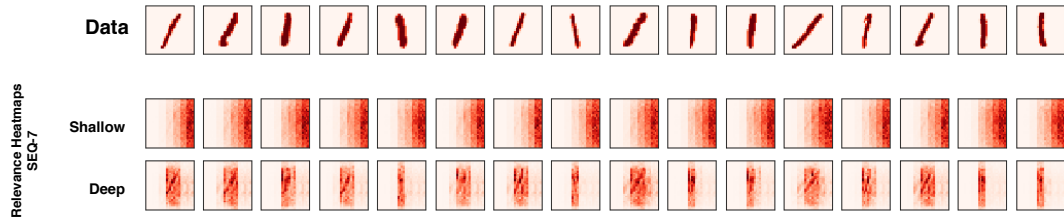


Figure 4.8: Relevance heatmaps Shallow and Deep Cell
(The samples are randomly chosen from MNIST's testing set)

Figure 4.8 shows more concrete evidence. Here, 16 samples from *Class 1* were randomly chosen from the test set and fed to Shallow and Deep as $SEQ = 7$ meaning $\mathbf{x}_t \in \mathbb{R}^{28,4}$. Intuitively, the heatmaps from Deep-7 align nicely to the data, while Shallow-7's ones do not seem to present any characteristic of the inputs. In fact, Deep-7 manages to perfectly distribute relevance scores to region that the data actually lie, for example the last 2 images.

Figure 4.9 also conveys the same insight, but from the aspect of *Class 1* population in testing set. The figure is a comparison of the distribution of pixel intensity from the *Class 1* population, represented with dashed line, and the distributions of relevance quantities from the heatmaps explained by the models. In the plot, value of the distribution for each step- i^{th} is aggregated from corresponding pixels of the data and relevance heatmaps. We can clearly see that the distribution of relevance produced by Deep-7 looks similar to the input distribution, while the one from Shallow-7 diverges significantly.

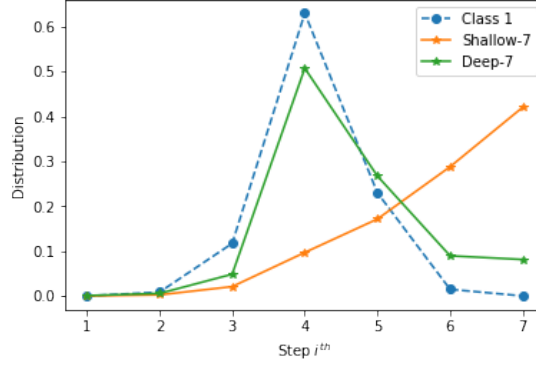


Figure 4.9: Comparison of pixel intensity distribution from MNIST Class 1 testing population and relevance distribution explained by the models

It is worth mentioning that samples in MNIST Class 1 are relatively simple, hence comparing such distributions can give some insight of the explainability of the models. However, in real world setting, qualifying explainability of a neural network via the distributions might not give meaningful intuition, because input tends to have complex structure and the network is likely to produce relevance heatmap that has different structure from the original image. We shall see this behavior in Fashion-MNIST experiments.

However, one interesting observation I have found is that sensitivity analysis ...

4.5.2 Fashion-MNIST

References