

CS 4501/6501 Interpretable Machine Learning

Neural Networks and Deep Learning

Yangfeng Ji

Department of Computer Science
University of Virginia



Overview

1. What is a white box looks like?
2. What is a neural network?
3. Why we need neural network models?
4. Why neural network is a black box?

What is a white box looks like?

Linear Models

Directly modeling a linear classifier as

$$h_y(\boldsymbol{x}) = \boldsymbol{w}_y^\top \boldsymbol{x} + b_y \quad (1)$$

- ▶ $\boldsymbol{x} \in \mathbb{N}^V$: vector, bag-of-words representation
- ▶ $\boldsymbol{w}_y \in \mathbb{R}^V$: vector, classification weights associated with label y
- ▶ $b_y \in \mathbb{R}$: scalar, label bias in the training set y

Logistic Regression

Rewrite the linear decision function in the log probabilistic form

$$\log P(y | x) \propto \underbrace{\mathbf{w}_y^\top \mathbf{x} + b_y}_{h_y(x)} \quad (2)$$

Logistic Regression

Rewrite the linear decision function in the log probabilistic form

$$\log P(y | x) \propto \underbrace{w_y^\top x + b_y}_{h_y(x)} \quad (2)$$

or, the probabilistic form is

$$P(y | x) \propto \exp(w_y^\top x + b_y) \quad (3)$$

Logistic Regression

Rewrite the linear decision function in the log probabilistic form

$$\log P(y | x) \propto \underbrace{w_y^\top x + b_y}_{h_y(x)} \quad (2)$$

or, the probabilistic form is

$$P(y | x) \propto \exp(w_y^\top x + b_y) \quad (3)$$

To make sure $P(y | x)$ is a valid definition of probability, we need to make sure $\sum_y P(y | x) = 1$,

$$P(\textcolor{red}{y} | x) = \frac{\exp(w_{\textcolor{red}{y}}^\top x + b_{\textcolor{red}{y}})}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x + b_{y'})} \quad (4)$$

Alternative Form

Rewriting x and w as

- ▶ $x^\top = [x_1, x_2, \dots, x_V, \mathbf{1}]$
- ▶ $w_y^\top = [w_1, w_2, \dots, w_V, \mathbf{b}_y]$

allows us to have a more concise form

$$P(y \mid x) = \frac{\exp(w_y^\top x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x)} \quad (5)$$

Alternative Form

Rewriting x and w as

- ▶ $x^\top = [x_1, x_2, \dots, x_V, \mathbf{1}]$
- ▶ $w_y^\top = [w_1, w_2, \dots, w_V, \mathbf{b}_y]$

allows us to have a more concise form

$$P(y \mid x) = \frac{\exp(w_y^\top x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x)} \quad (5)$$

Comments:

- ▶ $\frac{\exp(a)}{\sum_{a'} \exp(a')}$ is the **softmax** function
- ▶ This form works with any size of \mathcal{Y} — it does not have to be a binary classification problem.

Binary Classifier

Assume $\mathcal{Y} = \{\text{NEG}, \text{POS}\}$, then the corresponding logistic regression classifier with $Y = \text{Pos}$ is

$$P(Y = \text{Pos} \mid x) = \frac{1}{1 + \exp(-w^\top x)} \quad (6)$$

where w is the only parameter.

Binary Classifier

Assume $\mathcal{Y} = \{\text{NEG}, \text{POS}\}$, then the corresponding logistic regression classifier with $Y = \text{Pos}$ is

$$P(Y = \text{Pos} \mid \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \quad (6)$$

where \mathbf{w} is the only parameter.

- ▶ $P(Y = \text{NEG} \mid \mathbf{x}) = 1 - P(Y = \text{Pos} \mid \mathbf{x})$

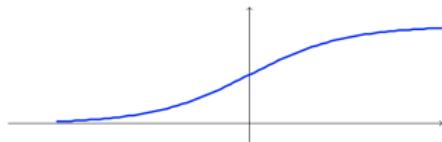
Binary Classifier

Assume $\mathcal{Y} = \{\text{NEG}, \text{POS}\}$, then the corresponding logistic regression classifier with $Y = \text{Pos}$ is

$$P(Y = \text{Pos} | \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \quad (6)$$

where \mathbf{w} is the only parameter.

- ▶ $P(Y = \text{NEG} | \mathbf{x}) = 1 - P(Y = \text{Pos} | \mathbf{x})$
- ▶ $\frac{1}{1+\exp(-z)}$ is the Sigmoid function



Two Questions on Building LR Models

... of building a logistic regression classifier

$$P(y | x) = \frac{\exp(w_y^\top x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x)} \quad (7)$$

- ▶ How to learn the parameters $W = \{w_y\}_{y \in \mathcal{Y}}$?

Two Questions on Building LR Models

... of building a logistic regression classifier

$$P(y | x) = \frac{\exp(w_y^\top x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x)} \quad (7)$$

- ▶ How to learn the parameters $W = \{w_y\}_{y \in \mathcal{Y}}$?
- ▶ Can x be better than the bag-of-words representations?

Review: (Log)-likelihood Function

With a collection of training examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$, the likelihood function of $\{w_y\}_{y \in \mathcal{Y}}$ is

$$L(W) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}) \quad (8)$$

and the **log**-likelihood function is

$$\ell(\{w_y\}) = \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}) \quad (9)$$

Log-likelihood Function of a LR Model

With the definition of a LR model

$$P(y | x) = \frac{\exp(w_y^\top x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x)} \quad (10)$$

the log-likelihood function is

$$\ell(W) = \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}) \quad (11)$$

$$= \sum_{i=1}^m \left\{ w_{y^{(i)}}^\top x^{(i)} - \log \sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x^{(i)}) \right\} \quad (12)$$

Given the training examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$, $\ell(W)$ is a **function of** $W = \{w_y\}$.

Optimization with Gradient

MLE is equivalent to **minimize** the Negative Log-Likelihood (NLL) as

$$\begin{aligned}\text{NLL}(W) &= -\ell(W) \\ &= \sum_{i=1}^m \left\{ -w_{y^{(i)}}^\top x^{(i)} + \log \sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x) \right\}\end{aligned}$$

then, the parameter w_y associated with label y can be updated as

$$w_y \leftarrow w_y - \eta \cdot \frac{\partial \text{NLL}(\{w_y\})}{\partial w_y}, \quad \forall y \in \mathcal{Y} \quad (13)$$

where η is called **learning rate**.

Optimization with Gradient (II)

Two questions answered by the update equation

- (1) which direction?
- (2) how far it should go?

Optimization with Gradient (II)

Two questions answered by the update equation

- (1) which direction?
- (2) how far it should go?

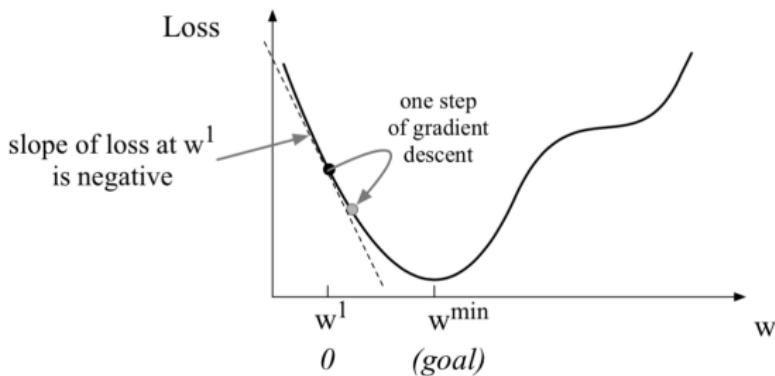
$$w_y \leftarrow w_y - \underbrace{\eta}_{(2)} \cdot \underbrace{\frac{\partial \text{NLL}(\{w_y\})}{\partial w_y}}_{(1)} \quad (14)$$

Optimization with Gradient (II)

Two questions answered by the update equation

- (1) which direction?
- (2) how far it should go?

$$w_y \leftarrow w_y - \underbrace{\eta}_{(2)} \cdot \underbrace{\frac{\partial \text{NLL}(\{w_y\})}{\partial w_y}}_{(1)} \quad (14)$$



Training Procedure

Steps for parameter estimation, given the current parameter $\{w_y\}$

1. Compute the derivative

$$\frac{\partial \text{NLL}(\{w_y\})}{\partial w_y}, \quad \forall y \in \mathcal{Y}$$

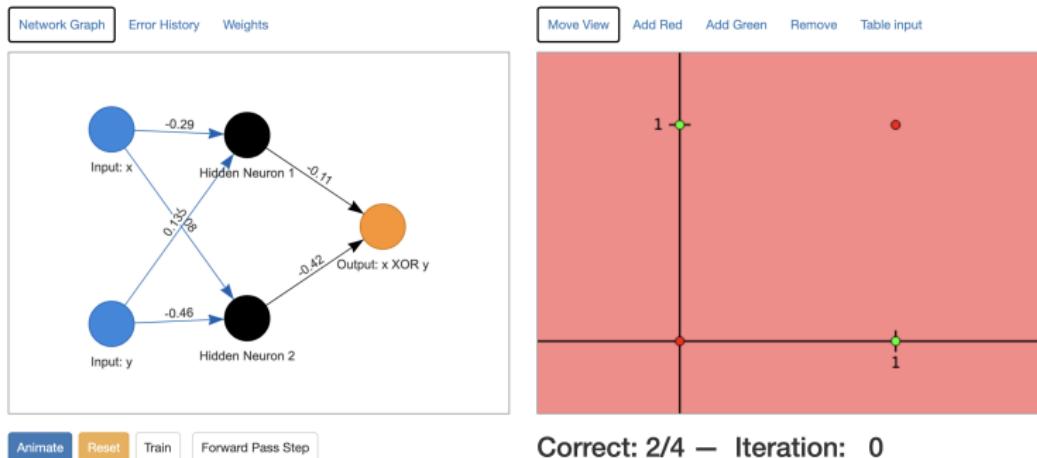
2. Update parameters with

$$w_y \leftarrow w_y - \eta \cdot \frac{\partial \text{NLL}(\{w_y\})}{\partial w_y}, \quad \forall y \in \mathcal{Y}$$

3. If not **done**, return to step 1

Neural Network Demo

A simple demo with 2-dimensional inputs



<https://phiresky.github.io/neural-network-demo/>

Sentiment Classification

A few training examples from Yelp Review

| Label | Text |
|-------|---|
| 5 | Love the staff, love the meat, love the place. Prepare for a long line around lunch or dinner hours ... |
| 5 | Super simple place but amazing nonetheless. It's been around since the 30's and they still serve the same thing ... |

Sentiment Classification

A few training examples from Yelp Review

| Label | Text |
|-------|---|
| 5 | Love the staff, love the meat, love the place. Prepare for a long line around lunch or dinner hours ... |
| 5 | Super simple place but amazing nonetheless. It's been around since the 30's and they still serve the same thing ... |
| 1 | Actually I would like to give them a big fat zero. Any vet's office that would tell ... |

Sentiment Classification

A few training examples from Yelp Review

| Label | Text |
|-------|---|
| 5 | Love the staff, love the meat, love the place. Prepare for a long line around lunch or dinner hours ... |
| 5 | Super simple place but amazing nonetheless. It's been around since the 30's and they still serve the same thing ... |
| 1 | Actually I would like to give them a big fat zero. Any vet's office that would tell ... |
| 2 | OK so first off the the burger was great as far as the taste. But I got super sick after eating it ... |

Classification Configuration and Performance

| | |
|--------------------------|---------------------|
| Vocabulary size | 17,490 |
| Training size | 40K |
| Development size | 5K |
| Classifier | Logistic regression |
| Regularization parameter | $C = 1$ |
| Training accuracy | 88.48% |
| Development accuracy | 61.22% |

You can find the demo code via the [link](#)

Interpretability: Global

| RATING | FEATURES |
|--------|--|
| 1 | worst awful horrible disgusting disgusted joke terrible zero luck pathetic |
| 2 | meh mediocre tacky nope renovations cheddar sand- which underwhelmed passable tasteless |

Interpretability: Global

| RATING | FEATURES |
|--------|--|
| 1 | worst awful horrible disgusting disgusted joke terrible zero luck pathetic |
| 2 | meh mediocre tacky nope renovations cheddar sand- which underwhelmed passable tasteless |
| 3 | feelings mains healthier remind hearty bleu overrated unsure smelling rules |

Interpretability: Global

| RATING | FEATURES |
|--------|---|
| 1 | worst awful horrible disgusting disgusted joke terrible zero luck pathetic |
| 2 | meh mediocre tacky nope renovations cheddar sand- which underwhelmed passable tasteless |
| 3 | feelings mains healthier remind hearty bleu overrated unsure smelling rules |
| 4 | default hankering drawback bojangles pleasantly hazel- nut customize gratuity excellent tremendously |
| 5 | phenomenal incredible amazing gem excellent pleas- antly hesitate master magnificent spotless |

Interpretability: Local

The prediction on the following is 5

I love the service here , they ' re on it !

After pre-processing, we remove the high-frequency word I and the punctuation

Interpretability: Local

The prediction on the following is 5

I love the service here , they ' re on it !

After pre-processing, we remove the high-frequency word I and the punctuation

| FEATURE | CLASSIFICATION WEIGHT |
|---------|-----------------------|
| love | 0.85 |
| service | 0.04 |
| on | 0.01 |
| they | -0.00 |
| it | -0.02 |
| the | -0.06 |
| re | -0.13 |
| here | -0.15 |

What is a neural network?

Logistic Regression

- ▶ An unified form for $y \in \{-1, +1\}$

$$p(Y = +1 \mid \mathbf{x}) = \frac{1}{1 + \exp(-\langle \mathbf{w}, \mathbf{x} \rangle)} \quad (15)$$

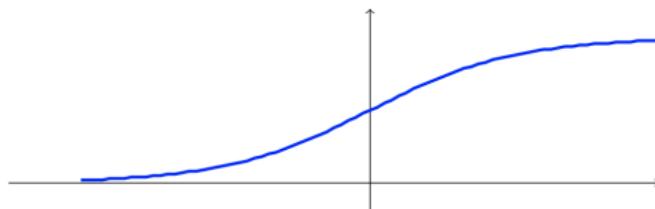
Logistic Regression

- ▶ An unified form for $y \in \{-1, +1\}$

$$p(Y = +1 \mid \mathbf{x}) = \frac{1}{1 + \exp(-\langle \mathbf{w}, \mathbf{x} \rangle)} \quad (15)$$

- ▶ The sigmoid function $\sigma(a)$ with $a \in \mathbb{R}$

$$\sigma(a) = \frac{1}{1 + \exp(-a)} \quad (16)$$

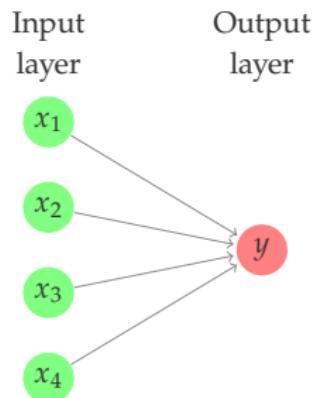


Graphical Representation

- ▶ A specific example of LR

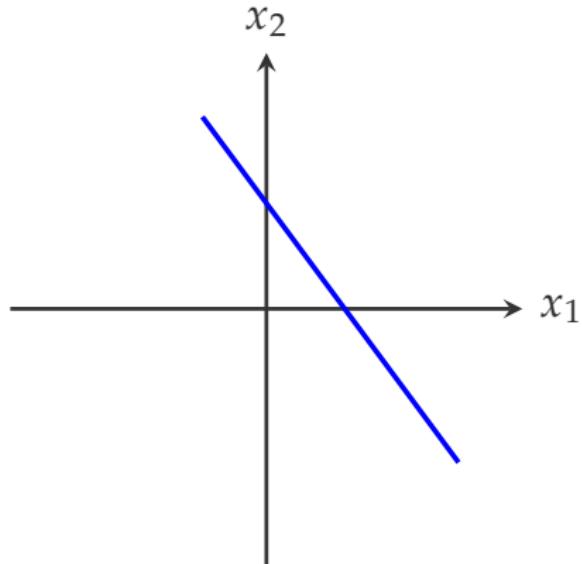
$$p(Y = 1 \mid \boldsymbol{x}) = \sigma\left(\sum_{j=1}^4 w_j x_{\cdot j}\right) \quad (17)$$

- ▶ The graphical representation of this LR model is



Capacity of a LR

Logistic regression gives a linear decision boundary



From LR to Neural Networks

Build upon logistic regression, a simple neural network can be constructed as

$$z_k = \sigma\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (18)$$

$$P(y = 1 \mid x) = \sigma\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (19)$$

- ▶ $x \in \mathbb{R}^d$: d -dimensional input
- ▶ $y \in \{-1, +1\}$ (binary classification problem)

From LR to Neural Networks

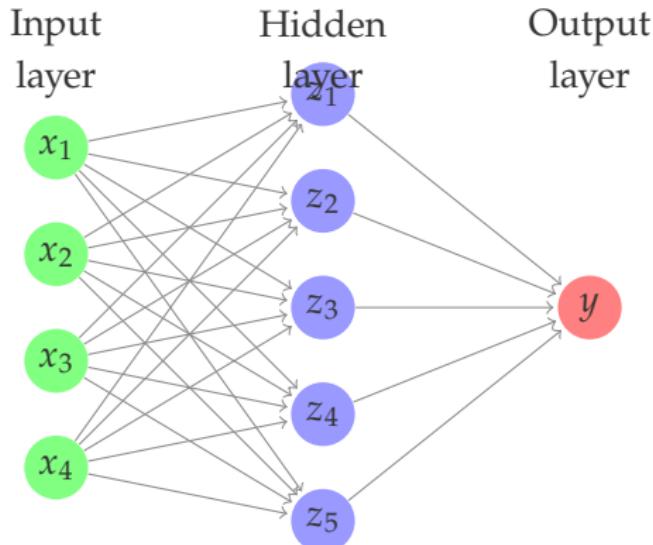
Build upon logistic regression, a simple neural network can be constructed as

$$z_k = \sigma\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (18)$$

$$P(y = 1 \mid x) = \sigma\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (19)$$

- ▶ $x \in \mathbb{R}^d$: d -dimensional input
- ▶ $y \in \{-1, +1\}$ (binary classification problem)
- ▶ $\{w_{k,i}^{(1)}\}$ and $\{w_k^{(o)}\}$ are two sets of the parameters, and
- ▶ K is the number of hidden units, each of them has the same form as a LR.

Graphical Representation

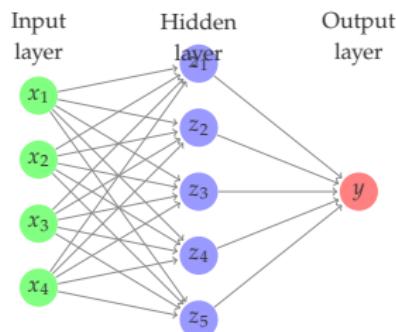


- ▶ Depth: 2 (two-layer neural network)
- ▶ Width: 5 (the maximal number of units in each layer)

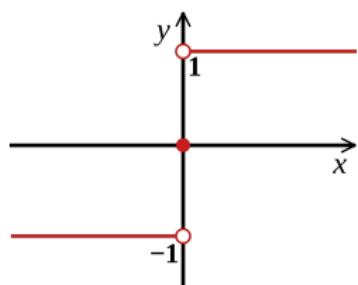
Hypothesis Space

The hypothesis space of neural networks is usually defined by the **architecture** of the network, which includes

- ▶ the nodes in the network,
- ▶ the connections in the network, and
- ▶ the **activation function** (e.g., σ)

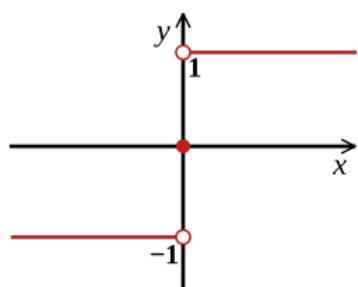


Other Activation Functions

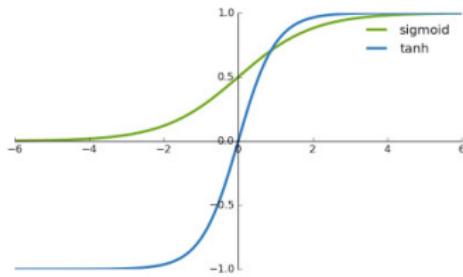


(a) Sign function

Other Activation Functions

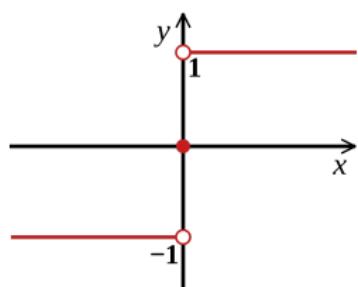


(a) Sign function

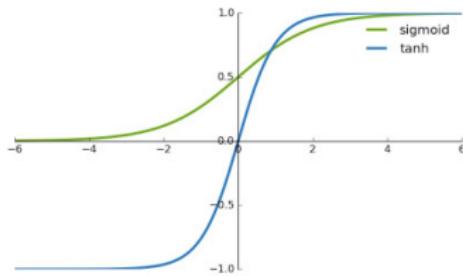


(b) Tanh function

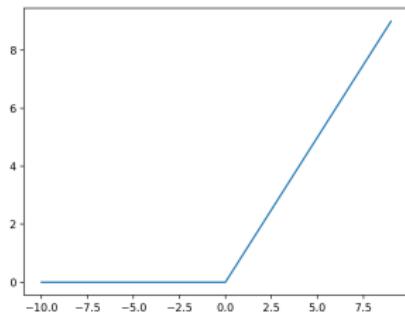
Other Activation Functions



(a) Sign function



(b) Tanh function



(c) ReLU function
[Jarrett et al., 2009]

Mathematical Formulation

- ▶ Element-wise formulation

$$z_k = \sigma\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (20)$$

$$P(y = +1 \mid x) = \sigma\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (21)$$

Mathematical Formulation

- ▶ Element-wise formulation

$$z_k = \sigma\left(\sum_{j=1}^d w_{k,j}^{(1)} x_{\cdot,j}\right) \quad k \in [K] \quad (20)$$

$$P(y = +1 \mid \boldsymbol{x}) = \sigma\left(\sum_{k=1}^K w_k^{(o)} z_k\right) \quad (21)$$

- ▶ Matrix-vector formulation

$$\boldsymbol{z} = \sigma(\mathbf{W}^{(1)} \boldsymbol{x}) \quad (22)$$

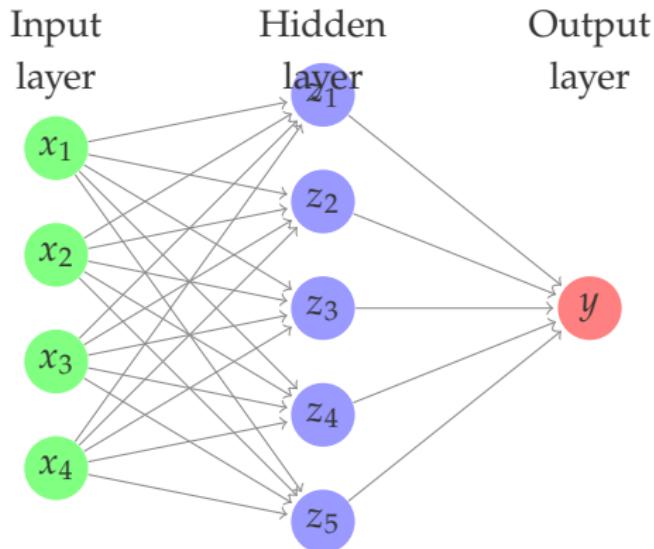
$$P(y = +1 \mid \boldsymbol{x}) = \sigma((\mathbf{w}^{(o)})^\top \boldsymbol{z}) \quad (23)$$

where $\mathbf{W}^{(1)} \in \mathbb{R}^{K \times d}$ and $\mathbf{w}^{(o)} \in \mathbb{R}^K$

Network Architecture

We are going to build a simple neural network for text classification.
It includes three layers as the previous example

- ▶ Input layer
- ▶ Hidden layer
- ▶ Output layer



Example

Consider the following special case, where we have a 4-dimensional BoW representation $x \in \mathbb{R}^4$ and a weight matrix $W \in \mathbb{R}^{5 \times 4}$

$$Wx = \begin{bmatrix} 0.1 & 0.3 & 0.7 & 0.9 \\ 0.2 & 0.8 & 0.3 & 0.5 \\ 0.4 & 0.8 & 0.6 & 0.1 \\ 0.7 & 0.2 & 0.9 & 0.2 \\ 0.4 & 0.5 & 0.8 & 0.9 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} \quad (24)$$

(25)

Example

Consider the following special case, where we have a 4-dimensional BoW representation $x \in \mathbb{R}^4$ and a weight matrix $W \in \mathbb{R}^{5 \times 4}$

$$Wx = \begin{bmatrix} 0.1 & 0.3 & 0.7 & 0.9 \\ 0.2 & 0.8 & 0.3 & 0.5 \\ 0.4 & 0.8 & 0.6 & 0.1 \\ 0.7 & 0.2 & 0.9 & 0.2 \\ 0.4 & 0.5 & 0.8 & 0.9 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} \quad (24)$$

$$= \begin{bmatrix} 0.3 \\ 0.8 \\ 0.8 \\ 0.2 \\ 0.5 \end{bmatrix} + \begin{bmatrix} 0.9 \\ 0.5 \\ 0.1 \\ 0.2 \\ 0.9 \end{bmatrix} \quad (25)$$

Example

Consider the following special case, where we have a 4-dimensional BoW representation $x \in \mathbb{R}^4$ and a weight matrix $W \in \mathbb{R}^{5 \times 4}$

$$Wx = \begin{bmatrix} 0.1 & 0.3 & 0.7 & 0.9 \\ 0.2 & 0.8 & 0.3 & 0.5 \\ 0.4 & 0.8 & 0.6 & 0.1 \\ 0.7 & 0.2 & 0.9 & 0.2 \\ 0.4 & 0.5 & 0.8 & 0.9 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} \quad (24)$$

$$= \begin{bmatrix} 0.3 \\ 0.8 \\ 0.8 \\ 0.2 \\ 0.5 \end{bmatrix} + \begin{bmatrix} 0.9 \\ 0.5 \\ 0.1 \\ 0.2 \\ 0.9 \end{bmatrix} \quad (25)$$

- ▶ Each column vector in W corresponds one word in the BoW representation
- ▶ The column vectors can be considered as representations of words, in other words, *word embeddings*

Output Layer

For the same text classification task, we use the following configuration:

- ▶ Input dimension: $x \in \mathbb{R}^{17K}$
- ▶ Hidden layer: $h \in \mathbb{R}^{32}$
- ▶ Output layer: $y \in \{1, \dots, 5\}$

You can find an extremely simple implementation via the same link, the dev accuracy is 65%

Why we need neural network models?

Distributed Representation

- ▶ Bag-of-words representations

| VOCAB | coffee | love | like | ... | tea | you |
|-------|--------|------|------|-----|-----|-----|
| love | (0 | 1 | 0 | ... | 0 | 0) |

¹I made up those numbers for illustration purpose

Distributed Representation

- ▶ Bag-of-words representations

| VOCAB | coffee | love | like | ... | tea | you |
|-------|--------|------|------|-----|-----|-----|
| love | (0 | 1 | 0 | ... | 0 | 0) |

- ▶ Distributed representations¹

| | | | | |
|------|------|-----|------|------|
| love | (0.1 | 0.8 | -1.0 | 0.3) |
|------|------|-----|------|------|

¹I made up those numbers for illustration purpose

Distributed Representation

- ▶ Bag-of-words representations

| VOCAB | coffee | love | like | ... | tea | you |
|-------|--------|------|------|-----|-----|-----|
| love | (0 | 1 | 0 | ... | 0 | 0) |
| like | (0 | 0 | 1 | ... | 0 | 0) |

- ▶ Distributed representations¹

love (0.1 0.8 -1.0 0.3)

¹I made up those numbers for illustration purpose

Distributed Representation

- ▶ Bag-of-words representations

| VOCAB | coffee | love | like | ... | tea | you |
|-------|--------|------|------|-----|-----|-----|
| love | (0 | 1 | 0 | ... | 0 | 0) |
| like | (0 | 0 | 1 | ... | 0 | 0) |

- ▶ Distributed representations¹

| | | | | |
|------|------|-----|------|------|
| love | (0.1 | 0.8 | -1.0 | 0.3) |
| like | (0.2 | 0.7 | -0.9 | 0.3) |

¹I made up those numbers for illustration purpose

Distributed Representation

- ▶ Bag-of-words representations

| VOCAB | coffee | love | like | ... | tea | you |
|-------|--------|------|------|-----|-----|-----|
| love | (0 | 1 | 0 | ... | 0 | 0) |
| like | (0 | 0 | 1 | ... | 0 | 0) |

- ▶ Distributed representations¹

| | | | | |
|------|------|-----|------|------|
| love | (0.1 | 0.8 | -1.0 | 0.3) |
| like | (0.2 | 0.7 | -0.9 | 0.3) |

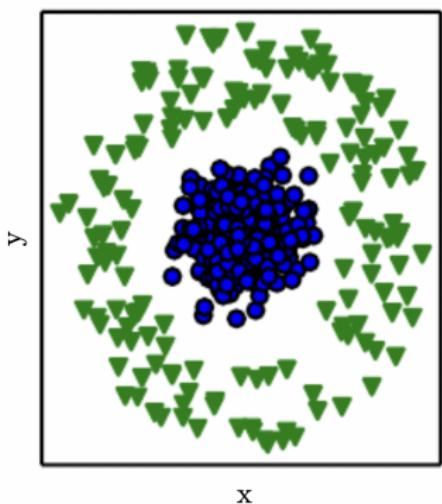
- ▶ Distributed representations allow simple algebraic operations for semantic meanings
 - ▶ E.g., the cosine value between two word embeddings measures their semantic similarity

¹I made up those numbers for illustration purpose

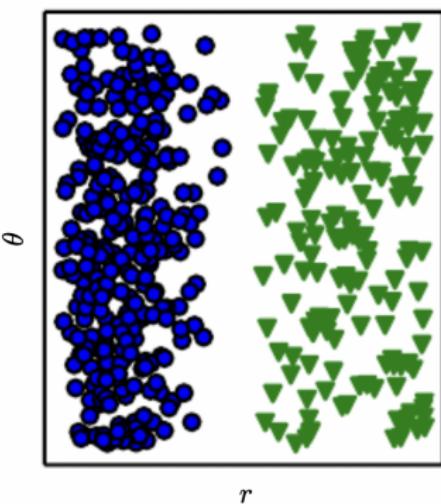
The Value of Different Representations

Different representations lead to different levels of challenge in machine learning

Cartesian coordinates

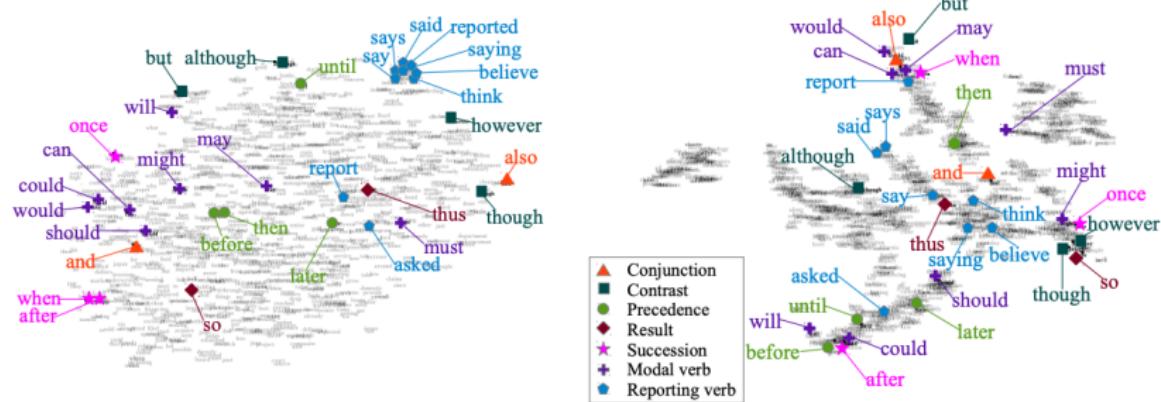


Polar coordinates



Advantage of Representation Learning

Driven by supervision signals, the model can learn some task-specific information and encode it in word embeddings



Similar advantage exists in any other supervised learning tasks [Bengio et al., 2013]

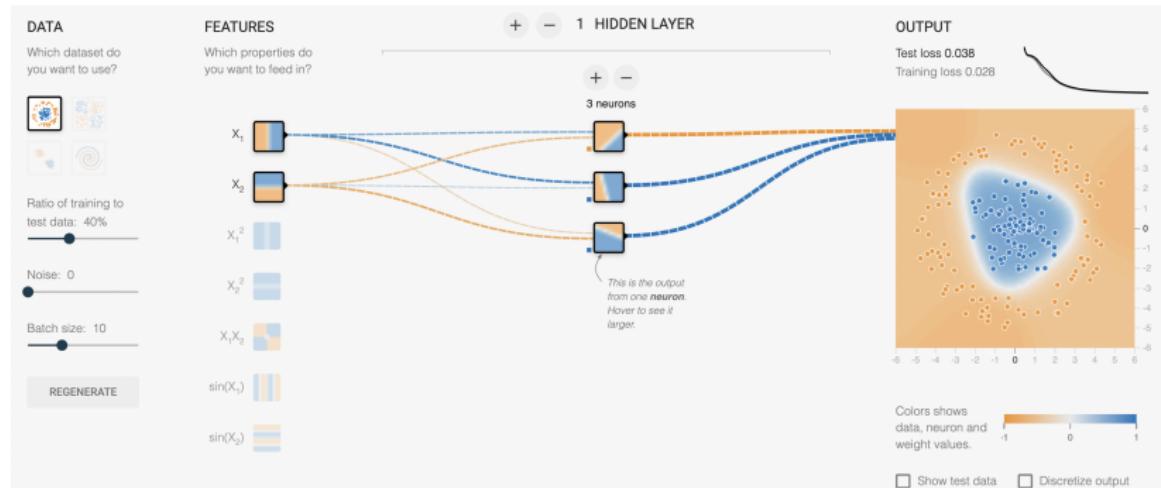
Neural Network as Universal Approximators

Neural network as universal approximators

- ▶ With arbitrary width and bounded depth [Cybenko, 1989]
- ▶ With arbitrary depth and limited width [Kidger and Lyons, 2020]

Model Capacity via Function Composition

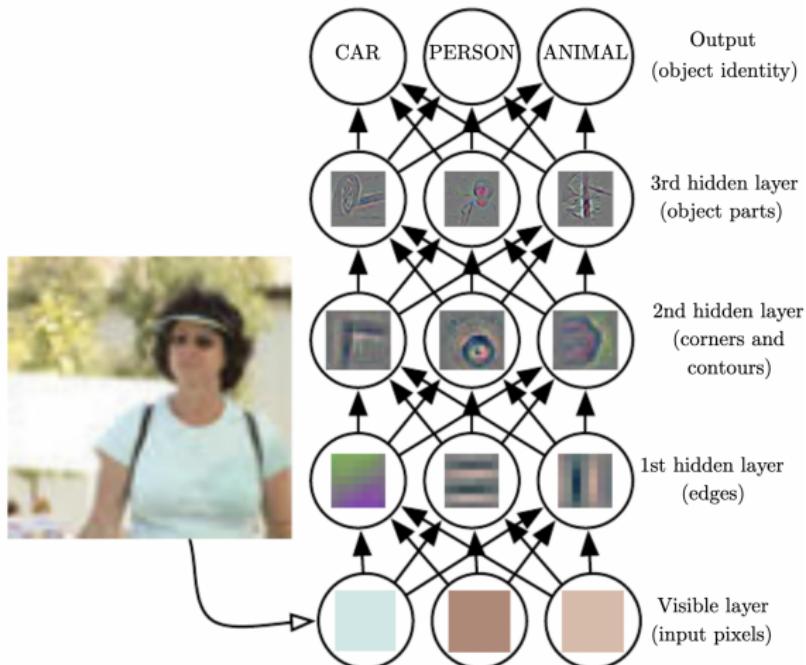
A Toy Example about Function Composition:



<https://playground.tensorflow.org/>

Model Capacity via Function Composition

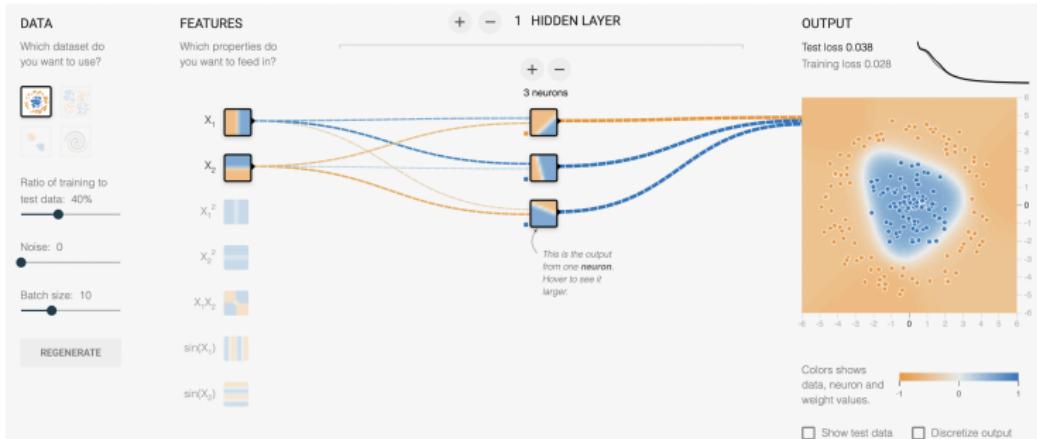
An example of function composition to extract high-level features



Why neural network is a black
box?

Is Neural Network a Black Box?

Not exactly, we can analyze what it learns when the model is small



Example

Model interpretability: model predictions can be interpreted as *certain* rules associated with inputs²

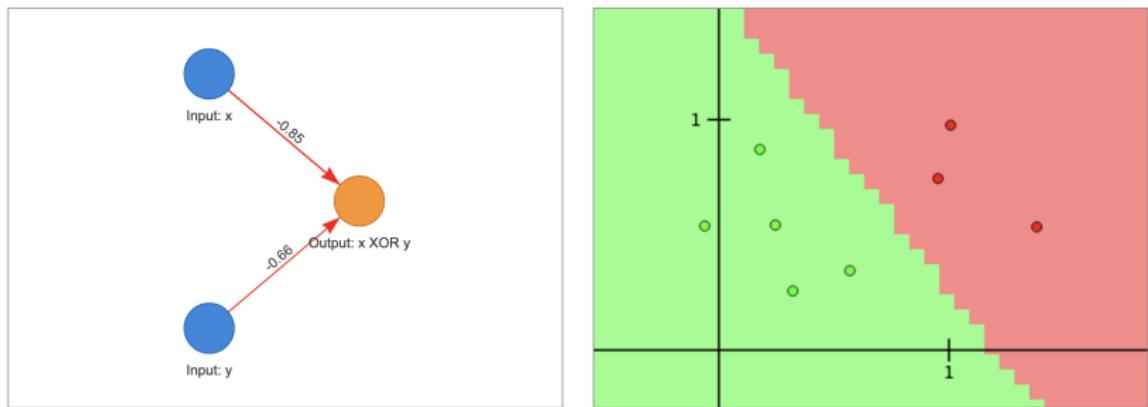


This is equivalent to a logistic regression model

²This is by no means a formal definition of interpretability

Example

Model interpretability: model predictions can be interpreted as *certain rules associated with inputs*²

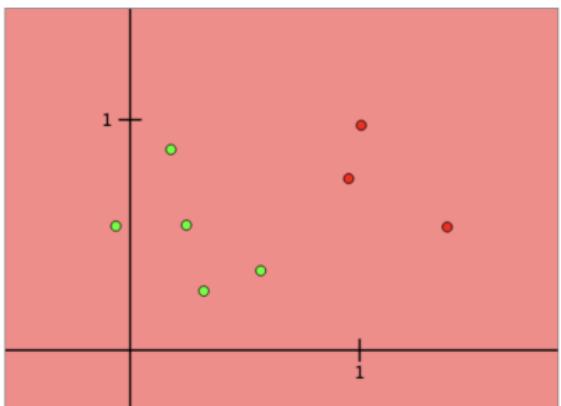
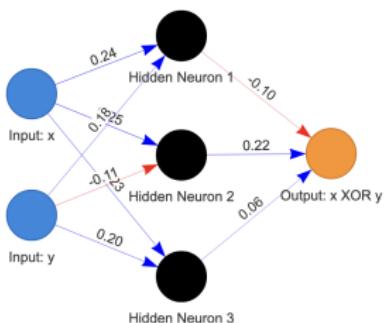


This is equivalent to a logistic regression model

²This is by no means a formal definition of interpretability

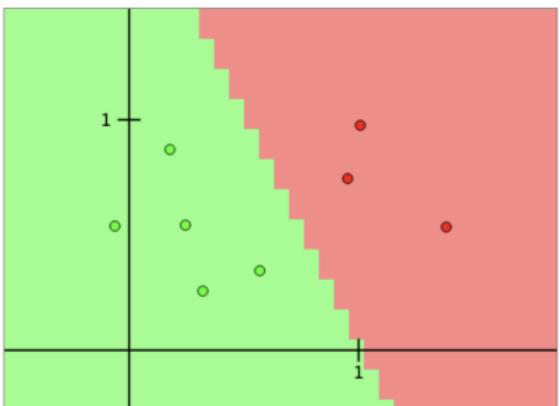
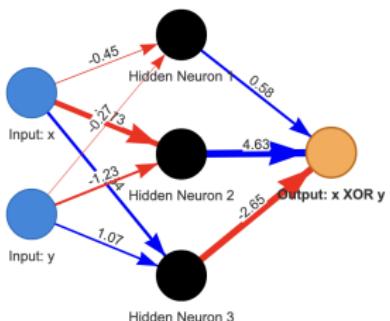
With a Simple Neural Network

With randomly initialized weights (classifying all examples as negative), the neural network (with one hidden layer) can easily learn a classifier with 100%



With a Simple Neural Network

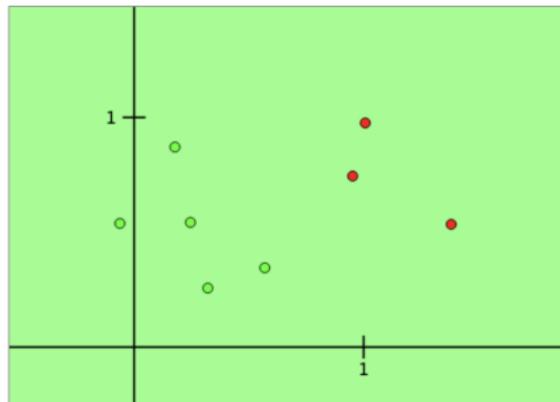
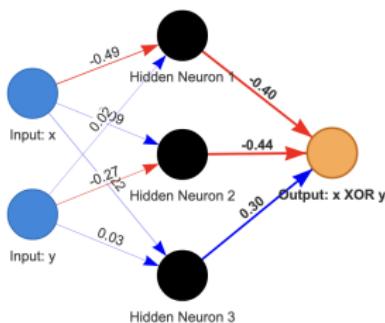
With randomly initialized weights (classifying all examples as negative), the neural network (with one hidden layer) can easily learn a classifier with 100%



With the visualization, it is not difficult to identify the second and the third neurons are important.

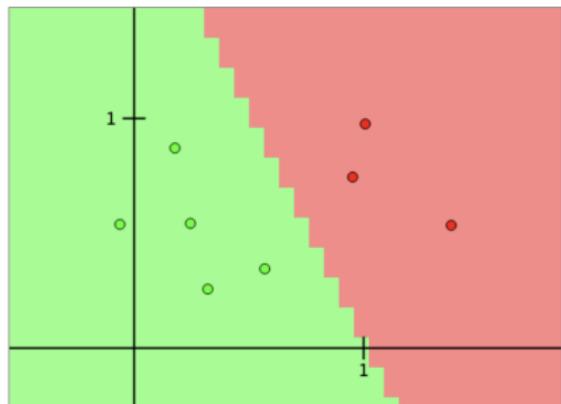
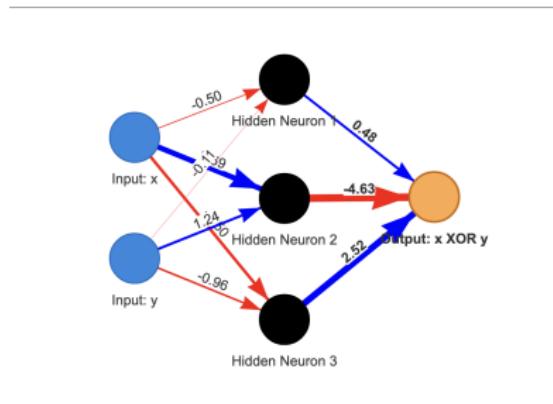
Distributed Representations (II)

With another set of randomly initialized weights (classifying all examples as negative), the learned classifier gives a very similar decision boundary



Distributed Representations (II)

With another set of randomly initialized weights (classifying all examples as negative), the learned classifier gives a very similar decision boundary



With the visualization, it is not difficult to identify the second and the third neurons are important.

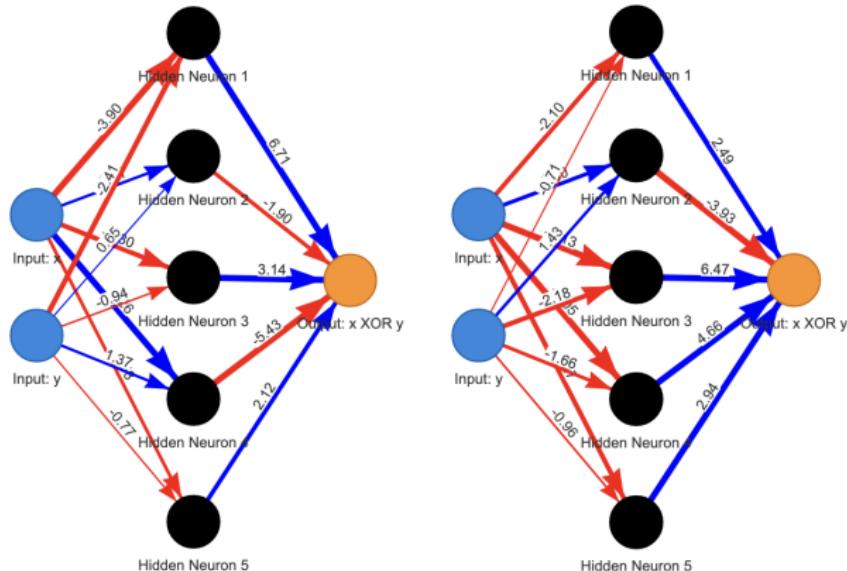
Distributed Representation (III)

- ▶ The contributions of the second and third neurons from the these two neural networks are contradicted with each other

Distributed Representation (III)

- ▶ The contributions of the second and third neurons from the these two neural networks are contradicted with each other
- ▶ Actually, it is still explainable, if we also consider the contribution from the previous layer

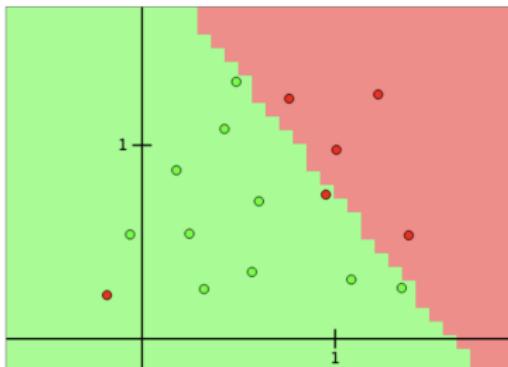
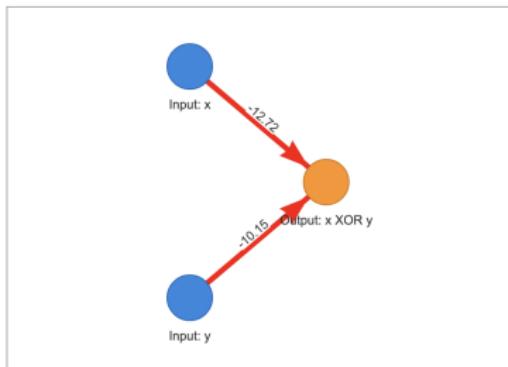
Try to Explain the Following Two Models?



- ▶ The neural network has more parameters than the task actually needs
- ▶ The contributions of hidden neurons are *randomly* distributed

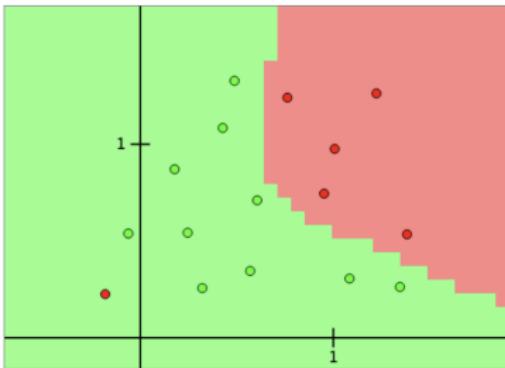
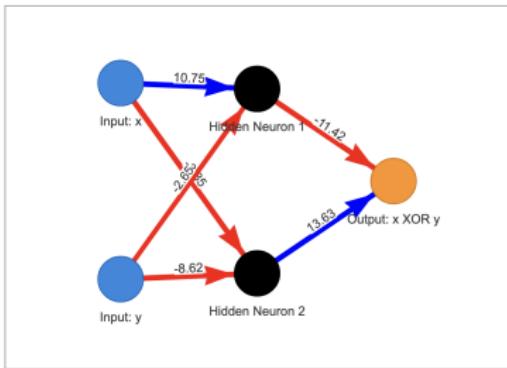
Number of Layers

Similar challenge for interpreting the contributions when we increase the number of layers



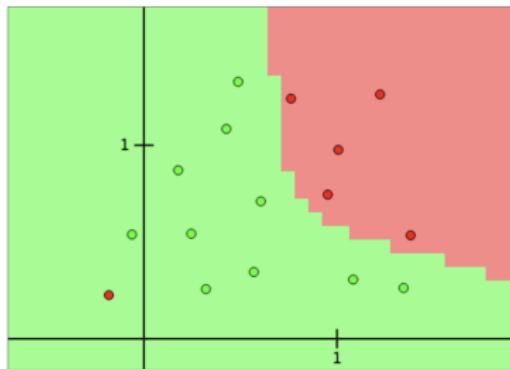
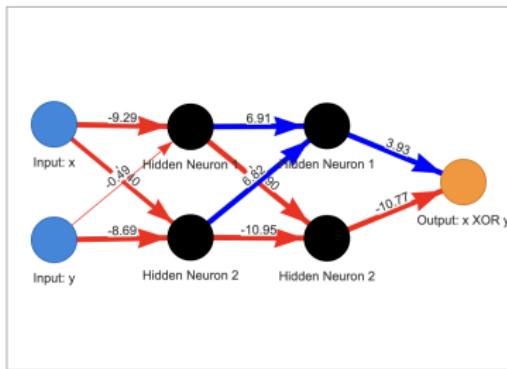
Number of Layers

Similar challenge for interpreting the contributions when we increase the number of layers



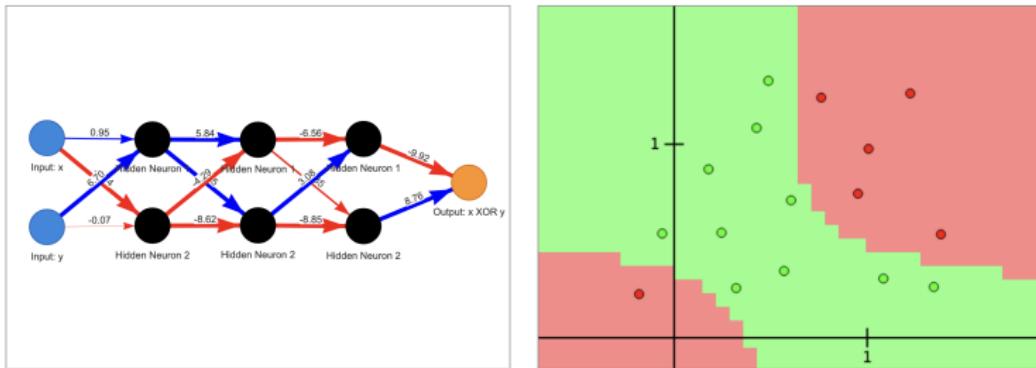
Number of Layers

Similar challenge for interpreting the contributions when we increase the number of layers



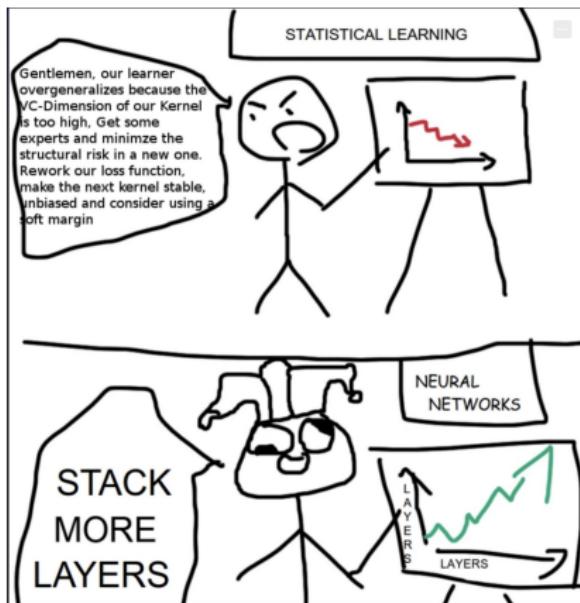
Number of Layers

Similar challenge for interpreting the contributions when we increase the number of layers



Number of Layers

Similar challenge for interpreting the contributions when we increase the number of layers



Question

Can we train a neural network that maintains good performance and
is also interpretable?

Reference

-  Bengio, Y., Courville, A., and Vincent, P. (2013).
Representation learning: A review and new perspectives.
IEEE transactions on pattern analysis and machine intelligence, 35(8):1798–1828.
-  Cybenko, G. (1989).
Approximation by superpositions of a sigmoidal function.
Mathematics of control, signals and systems, 2(4):303–314.
-  Goodfellow, I., Bengio, Y., and Courville, A. (2016).
Deep Learning.
MIT Press.
<http://www.deeplearningbook.org>.
-  Jarrett, K., Kavukcuoglu, K., Ranzato, M., and LeCun, Y. (2009).
What is the best multi-stage architecture for object recognition?
In *Proceedings of the 12th International Conference on Computer Vision*, pages 2146–2153. IEEE.
-  Jurafsky, D. and Martin, J. H. (2022).
Speech and Language Processing.
Online.
-  Kidger, P. and Lyons, T. (2020).
Universal approximation with deep narrow networks.
In *Conference on learning theory*, pages 2306–2327. PMLR.