Advanced ICT approaches Machine learning

Hands-on deep learning

prof. dr. Bojan Cestnik

IPS, 20. 11. 2024

Slides from http://udlbook.com, by Simon J.D. Prince

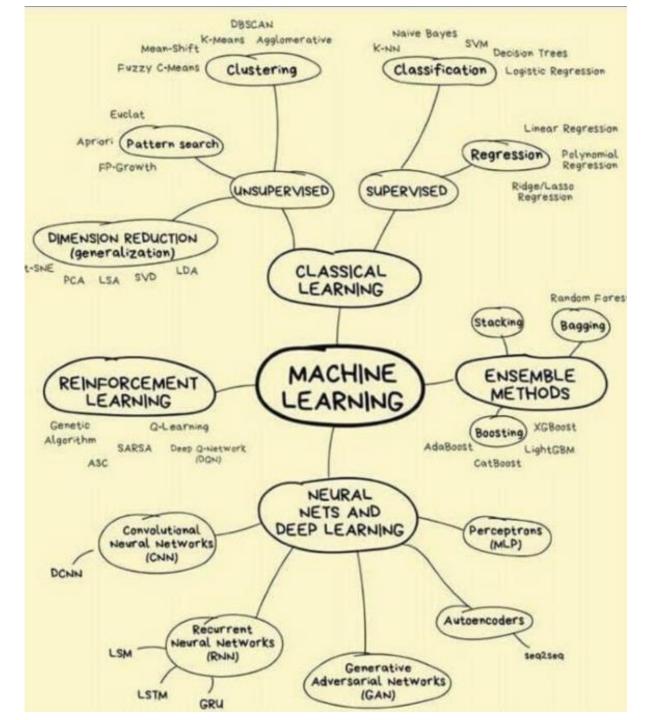
Contents

Machine learning – introduction

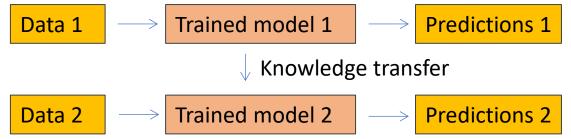
- 1. (Linear) regression
- 2. Approximation with shallow neural network
- 3. Approximation with deep neural network
- 4. Neural Networks: Playground Exercises
- 5. nanoGPT demo

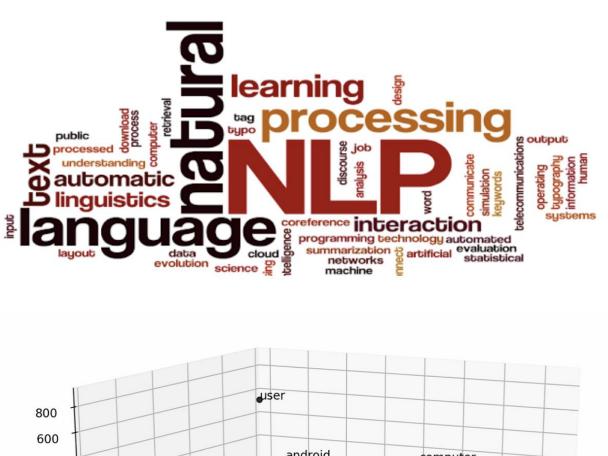
Machine Learning (ML)

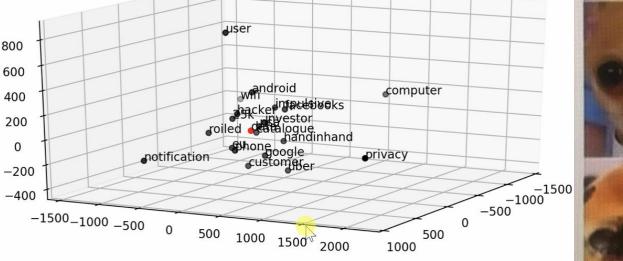
- Developing algorithms and models that enable machines to learn from data
- Key concepts:
 - **Supervised learning**: ML algorithms learn from labeled data (input-output pairs) to make predictions about new, unseen data
 - **Unsupervised learning**: ML algorithms learn from unlabeled data to identify patterns, relationships, or structures in the data
 - **Reinforcement learning**: ML algorithms learn by interacting with an environment and receiving feedback as rewards or penalties
 - **Transfer learning**: ML algorithms focus on storing knowledge gained while solving one problem and applying it to a different but related problem

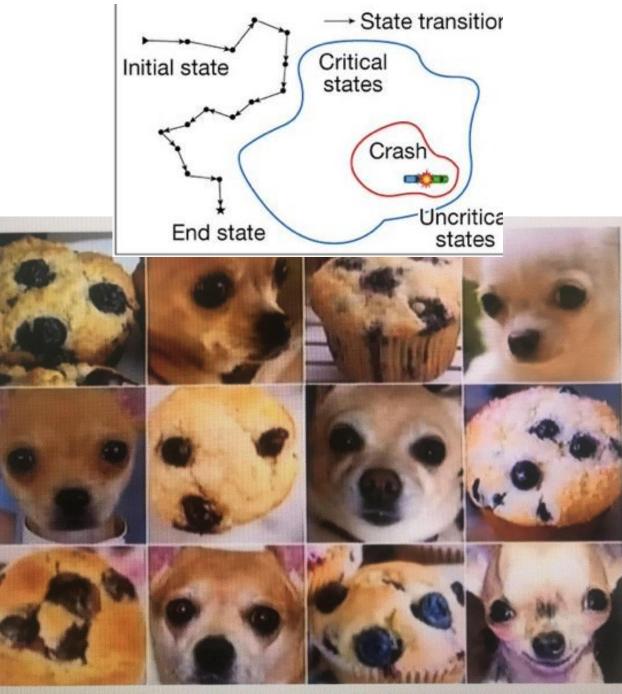


Transfer learning



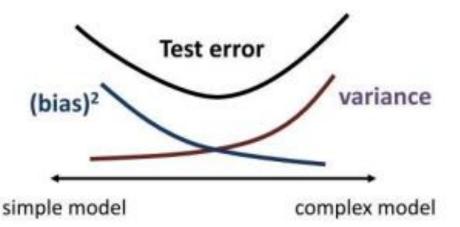


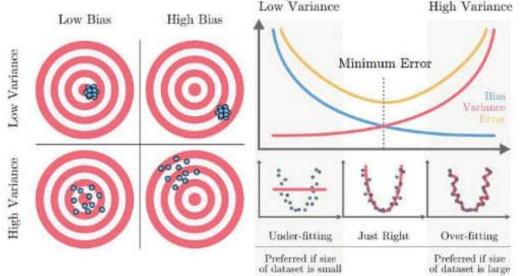




Machine Learning (ML)

- Learning models from training data (generalization)
- Measuring performance on unseen test data
- Model overfitting: too specialized on its training data
- Model underfitting: overly simplistic, inadequate capture of underlying patterns in the data
- Result: poor performance on unseen data

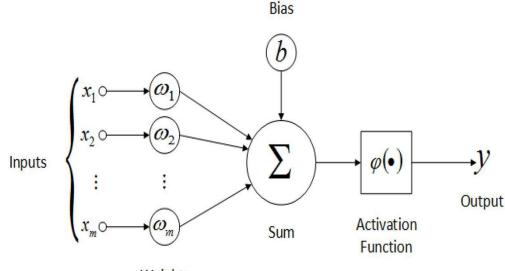




Low Bias

Neural Networks

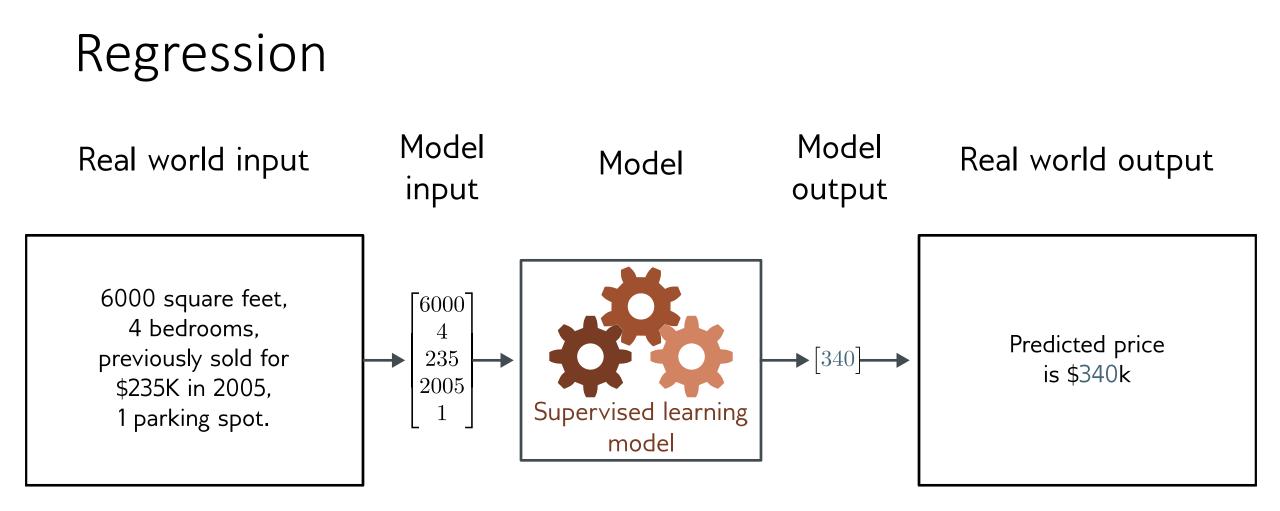
 Inspired by the structure and function of biological neurons and neural systems



- Designed to process, learn, and represent complex patterns in data
- Key concepts:
 - **Neurons**: the basic building blocks of neural networks, which receive input signals, process them, and produce output signals
 - Activation functions: nonlinear functions applied to the weighted sum of a neuron's inputs to determine its output signal
 - **Backpropagation**: an algorithm used to train neural networks by minimizing the error between predicted outputs and actual outputs through gradient descent and weight adjustments

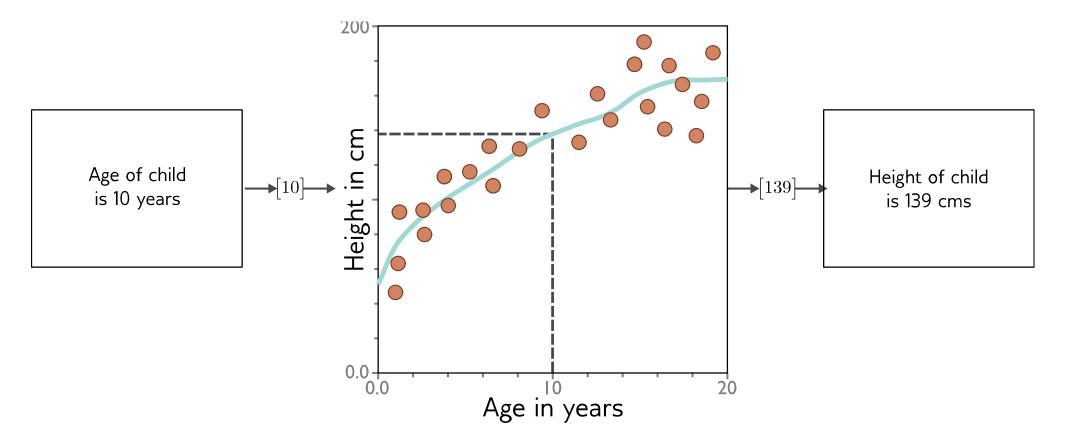
1. Linear regression

- Data points (x_1, y_1) , (x_2, y_2) , (x_3, y_3) , ... (x_{n-1}, y_{n-1}) , (x_n, y_n)
- y = k * x + n
- Loss function
- Estimation of k and n from data

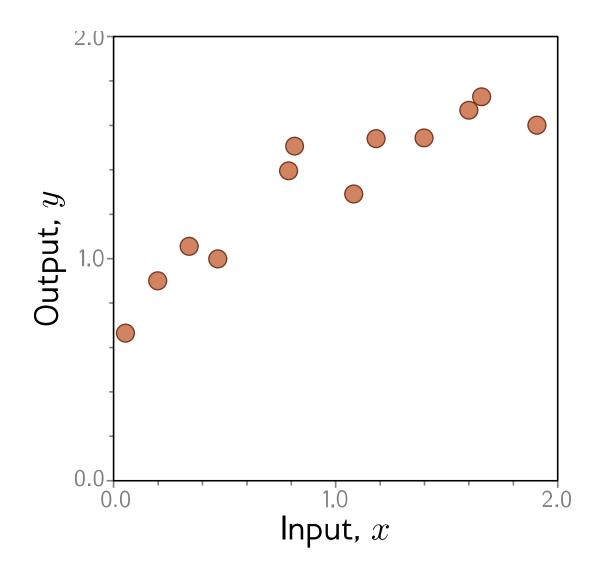


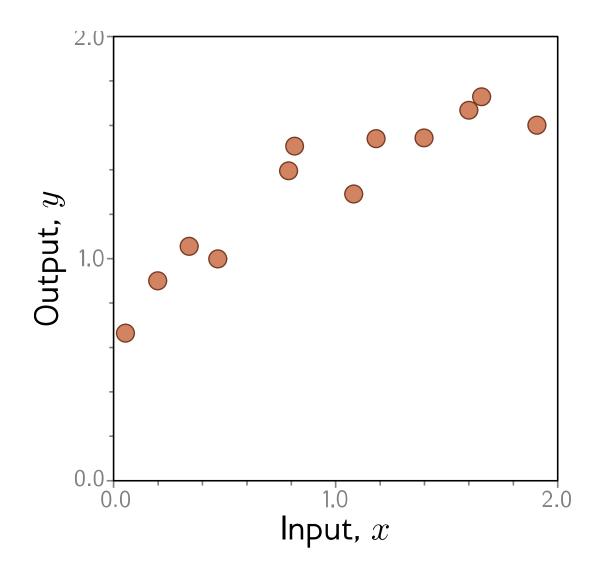
- Univariate regression problem (one output, real value)
- Fully connected network

What is a supervised learning model?



- An equation relating input (age) to output (height)
- Search through family of possible equations to find one that fits training data well



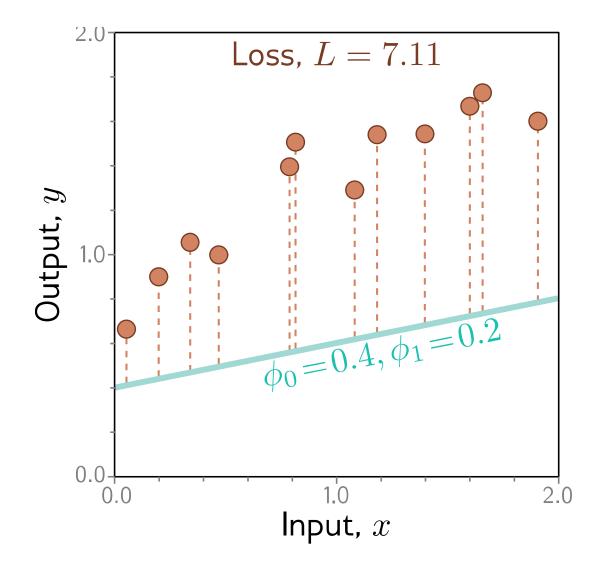


Loss function:

$$L[\phi] = \sum_{i=1}^{I} (f[x_i, \phi] - y_i)^2$$
$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

Example: 1D Linear regression loss function



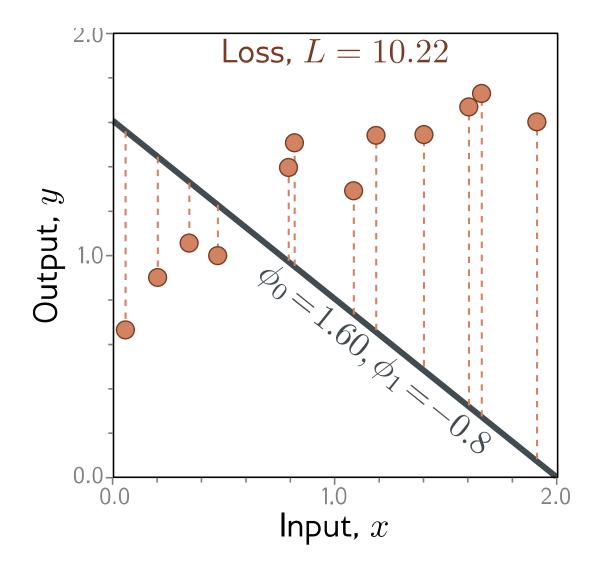
Loss function:

L

$$[\phi] = \sum_{i=1}^{I} (\mathbf{f}[x_i, \phi] - y_i)^2$$
$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

Example: 1D Linear regression loss function

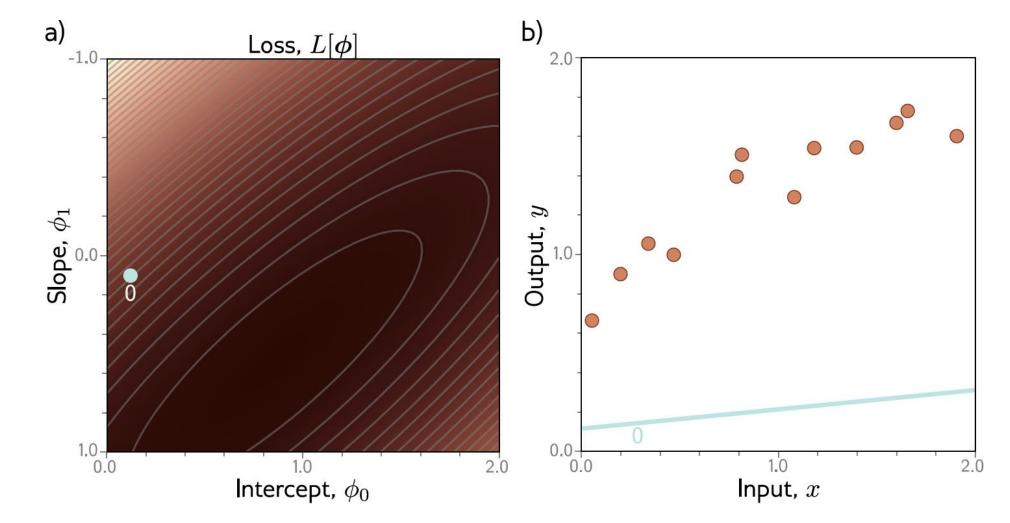


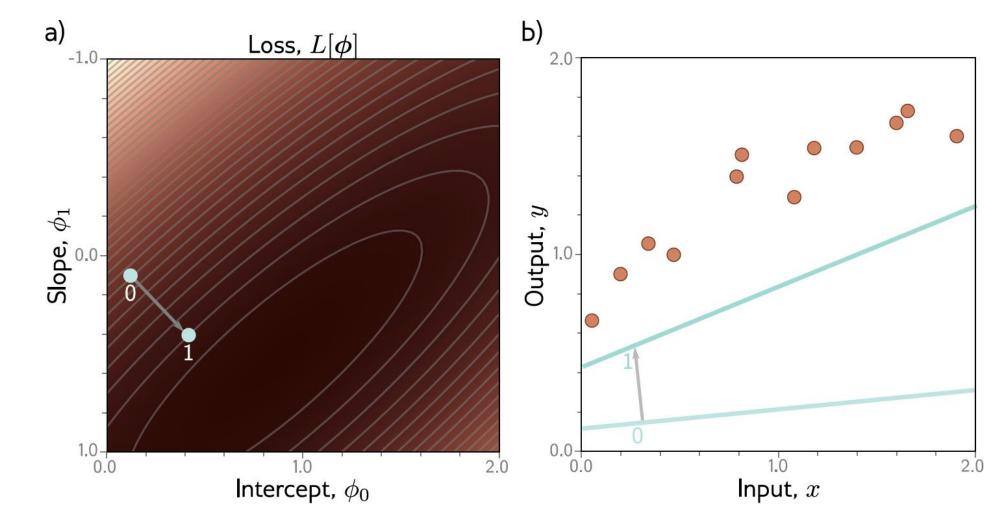
Loss function:

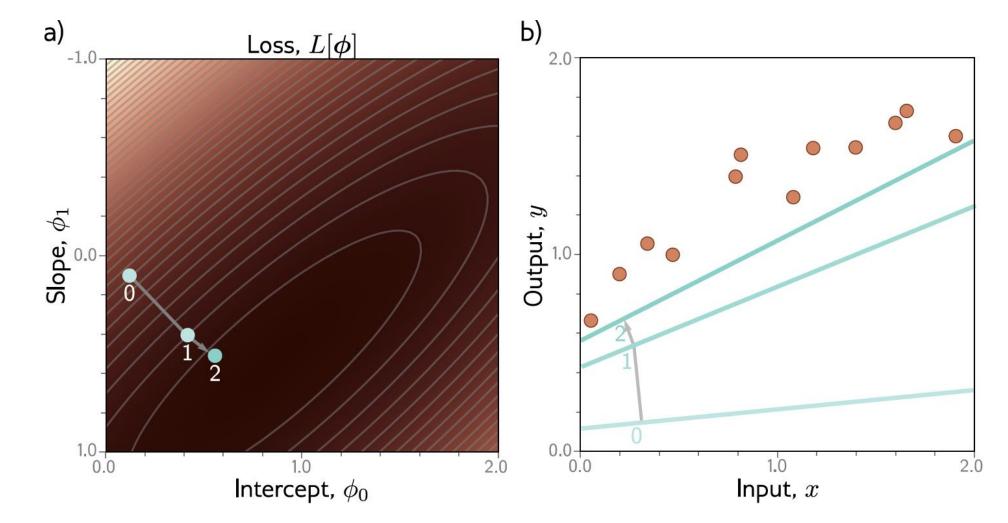
L

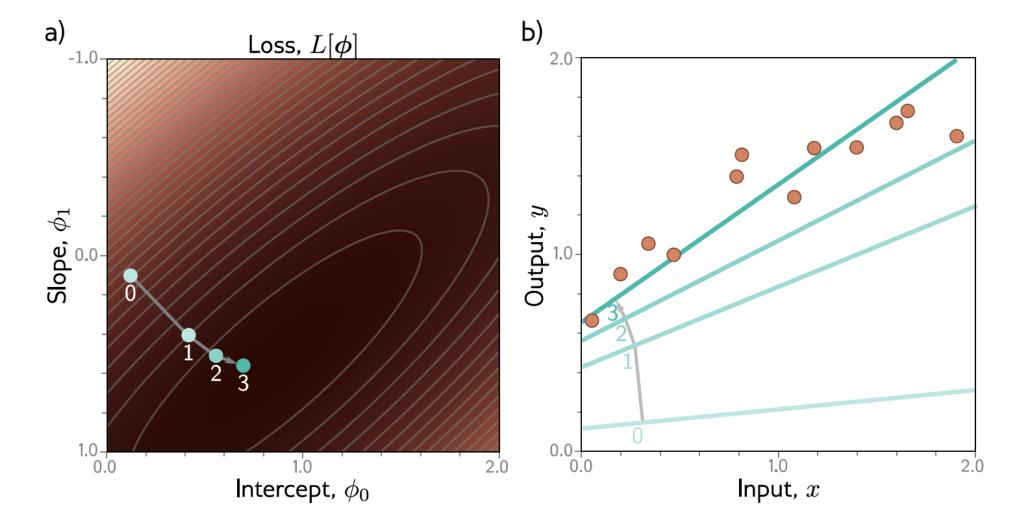
$$[\phi] = \sum_{i=1}^{I} (\mathbf{f}[x_i, \phi] - y_i)^2$$
$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

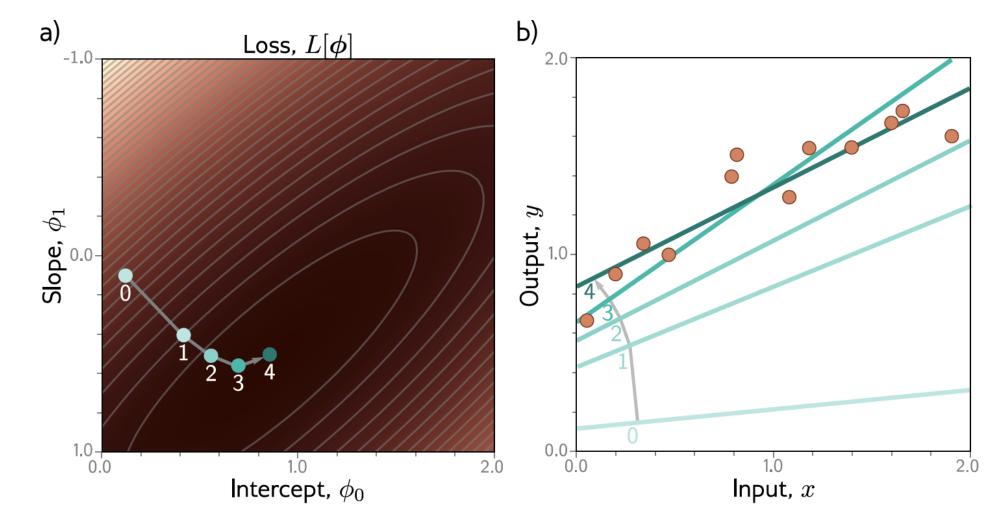
"Least squares loss function"











This technique is known as gradient descent

Possible objections

- But you can fit the line model in closed form!
 - Yes but we won't be able to do this for more complex models
- But we could exhaustively try every slope and intercept combo!
 - Yes but we won't be able to do this when there are a million parameters

Example in python notebook

• Notebook 2.1 Supervised Learning

2. Approximation with shallow neural network

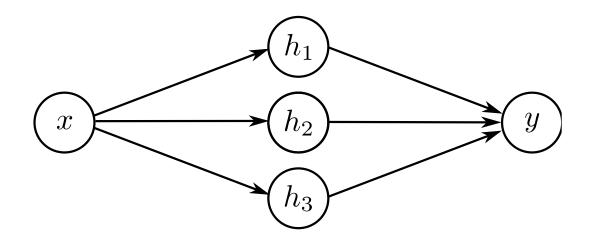
- Data points (x_1, y_1) , (x_2, y_2) , (x_3, y_3) , ... (x_{n-1}, y_{n-1}) , (x_n, y_n)
- Loss function
- Estimation of parameters (from training data)

Depicting neural networks

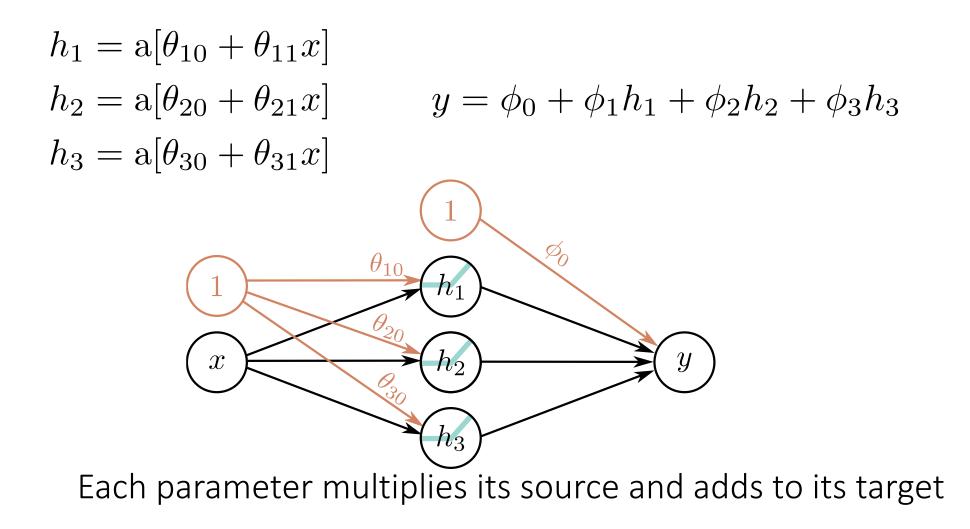
$$h_{1} = a[\theta_{10} + \theta_{11}x]$$

$$h_{2} = a[\theta_{20} + \theta_{21}x] \qquad y = \phi_{0} + \phi_{1}h_{1} + \phi_{2}h_{2} + \phi_{3}h_{3}$$

$$h_{3} = a[\theta_{30} + \theta_{31}x]$$



Depicting neural networks



1D Linear Regression

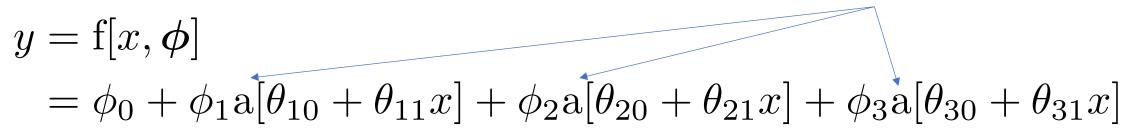
$$y = \mathbf{f}[x, \phi]$$
$$= \phi_0 + \phi_1 x$$

Example shallow network

 $y = f[x, \phi]$ = $\phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$

 $y = f[x, \phi]$ = $\phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$

Activation function



Activation function

$$y = f[x, \phi]$$

= $\phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$

$$\mathbf{a}[z] = \operatorname{ReLU}[z] = \begin{cases} 0 & z < 0 \\ z & z \ge 0 \end{cases}.$$

Rectified Linear Unit

(particular kind of activation function)

Activation function

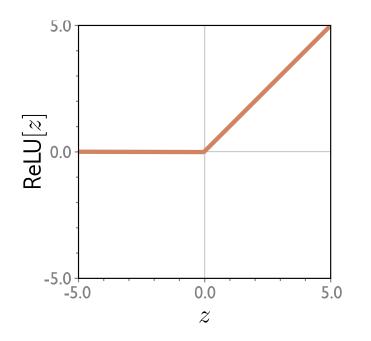
$$y = f[x, \phi]$$

= $\phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$

$$\mathbf{a}[z] = \operatorname{ReLU}[z] = \begin{cases} 0 & z < 0 \\ z & z \ge 0 \end{cases}.$$

Rectified Linear Unit

(particular kind of activation function)



$$y = f[x, \phi]$$

= $\phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$

This model has 10 parameters:

$$\boldsymbol{\phi} = \{\phi_0, \phi_1, \phi_2, \phi_3, \theta_{10}, \theta_{11}, \theta_{20}, \theta_{21}, \theta_{30}, \theta_{31}\}$$

- Represents a family of functions
- Parameters determine particular function
- Given parameters can perform inference (run equation)
- Given training dataset $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}$
- Define loss function $L[\phi]$ (least squares)
- Change parameters to minimize loss function

 $y = \phi_0 + \phi_1 \mathbf{a}[\theta_{10} + \theta_{11}x] + \phi_2 \mathbf{a}[\theta_{20} + \theta_{21}x] + \phi_3 \mathbf{a}[\theta_{30} + \theta_{31}x].$

Hidden units

$y = \phi_0 + \phi_1 \mathbf{a}[\theta_{10} + \theta_{11}x] + \phi_2 \mathbf{a}[\theta_{20} + \theta_{21}x] + \phi_3 \mathbf{a}[\theta_{30} + \theta_{31}x].$

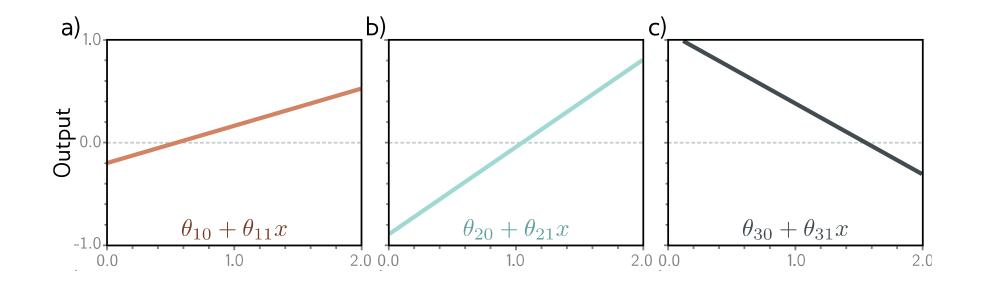
Break down into two parts:

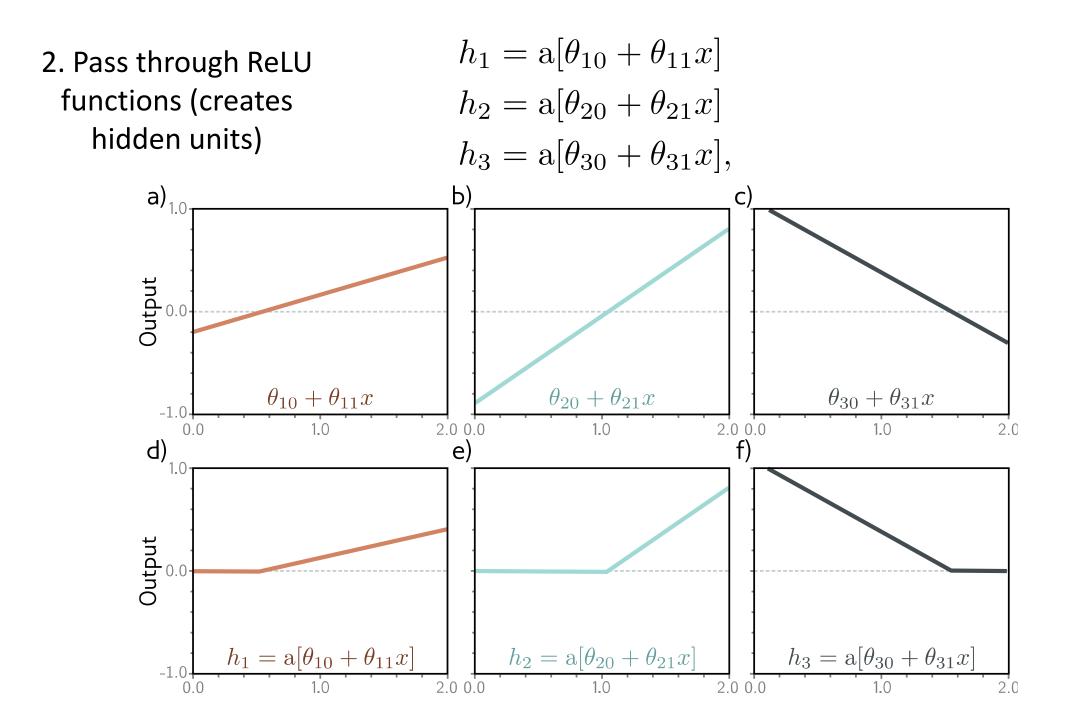
$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

where:

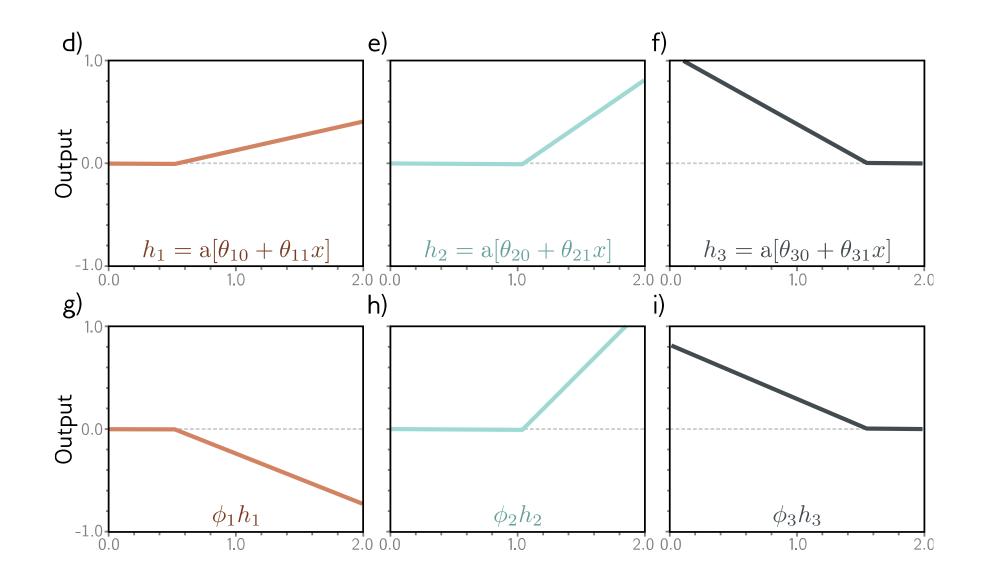
Hidden units
$$\begin{cases} h_1 = a[\theta_{10} + \theta_{11}x] \\ h_2 = a[\theta_{20} + \theta_{21}x] \\ h_3 = a[\theta_{30} + \theta_{31}x] \end{cases}$$

1. compute three linear functions

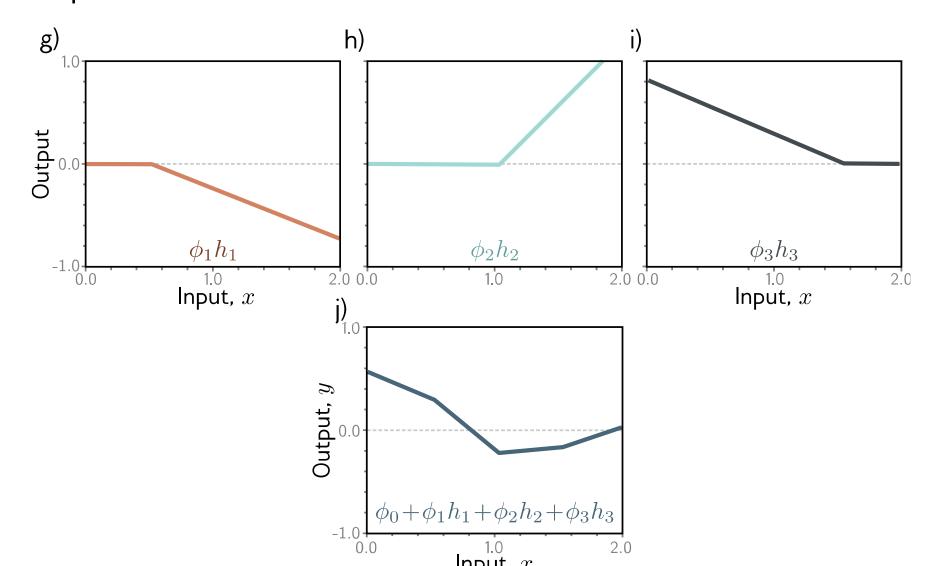




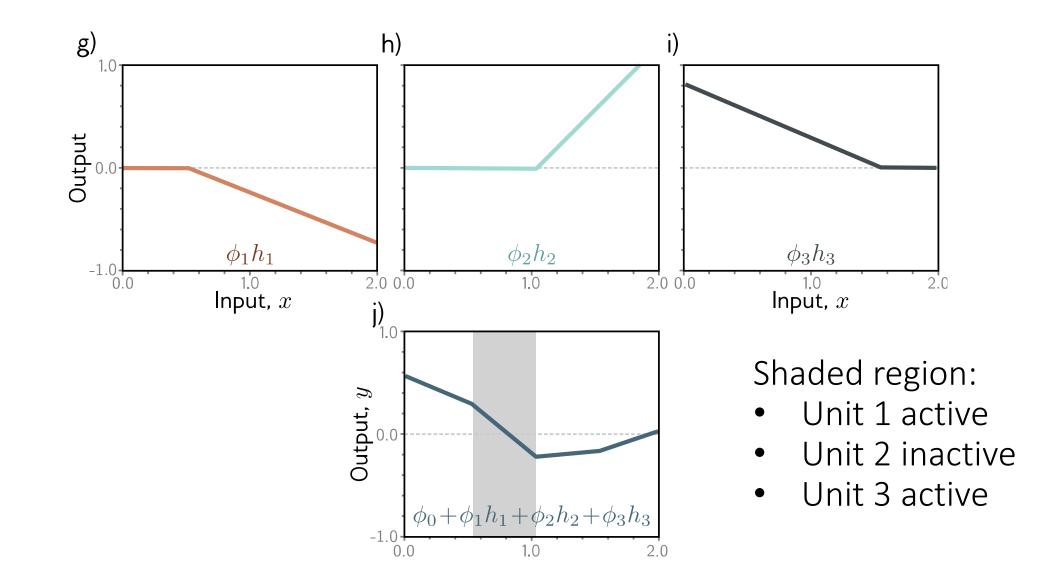
2. Weight the hidden units



4. Sum the weighted hidden units to create $y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$ output



Activation pattern = which hidden units are activated



Example in python notebook

• Notebook 3.1 Shallow neural networks I

3. Approximation with deep neural network

- Data points (x_1, y_1) , (x_2, y_2) , (x_3, y_3) , ... (x_{n-1}, y_{n-1}) , (x_n, y_n)
- Loss function
- Estimation of parameters

Deep neural networks

- Networks with more than one hidden layer
- Intuition becomes more difficult

Composing two networks.

Network 1:

 $h_1 = \mathbf{a}[\theta_{10} + \theta_{11}x]$ $h_2 = \mathbf{a}[\theta_{20} + \theta_{21}x]$ $h_3 = \mathbf{a}[\theta_{30} + \theta_{31}x]$

$$\begin{bmatrix} x \\ 1 \\ x \end{bmatrix} \qquad y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3 \\ \begin{bmatrix} 1 \\ 1 \\ x \end{bmatrix}$$

Network 2:

 $h'_{1} = a[\theta'_{10} + \theta'_{11}y]$ $h'_{2} = a[\theta'_{20} + \theta'_{21}y] \qquad y' = \phi'_{0} + \phi'_{1}h'_{1} + \phi'_{2}h'_{2} + \phi'_{3}h'_{3}$ $h'_{3} = a[\theta'_{30} + \theta'_{31}y]$

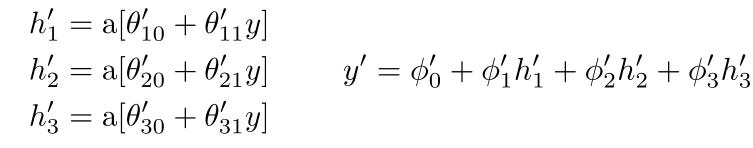
Composing two networks.

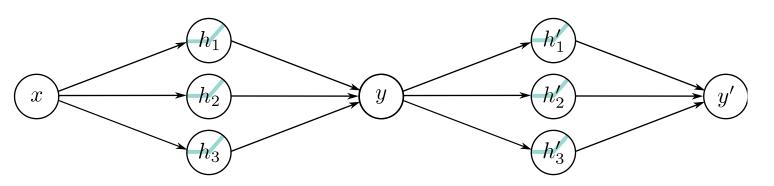
Network 1:

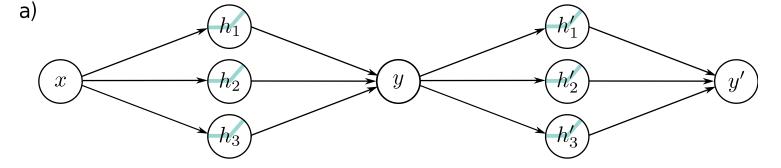
 $h_1 = \mathbf{a}[\theta_{10} + \theta_{11}x]$ $h_2 = \mathbf{a}[\theta_{20} + \theta_{21}x]$ $h_3 = \mathbf{a}[\theta_{30} + \theta_{31}x]$

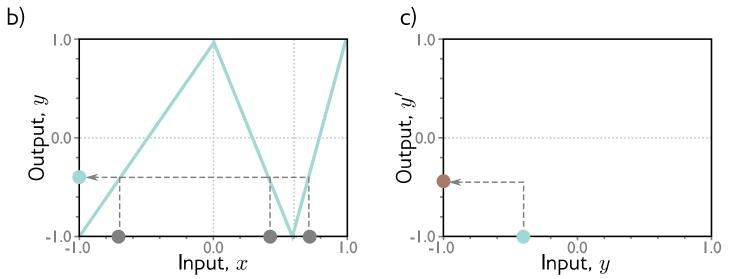
$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

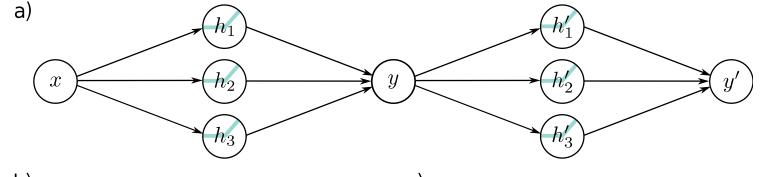
Network 2:

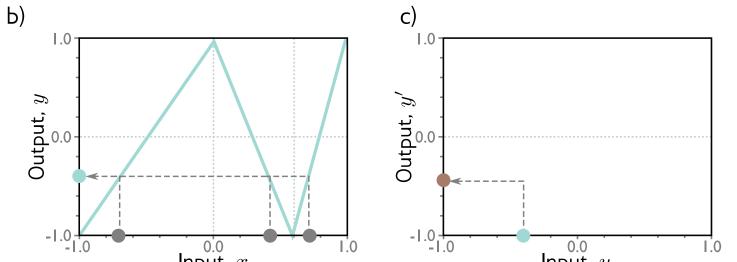




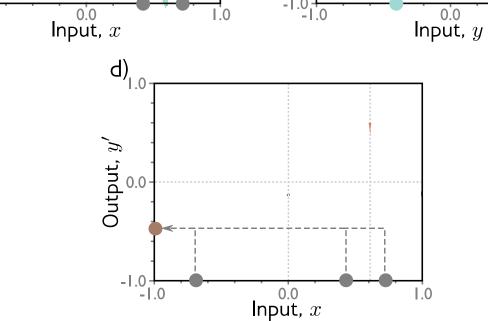


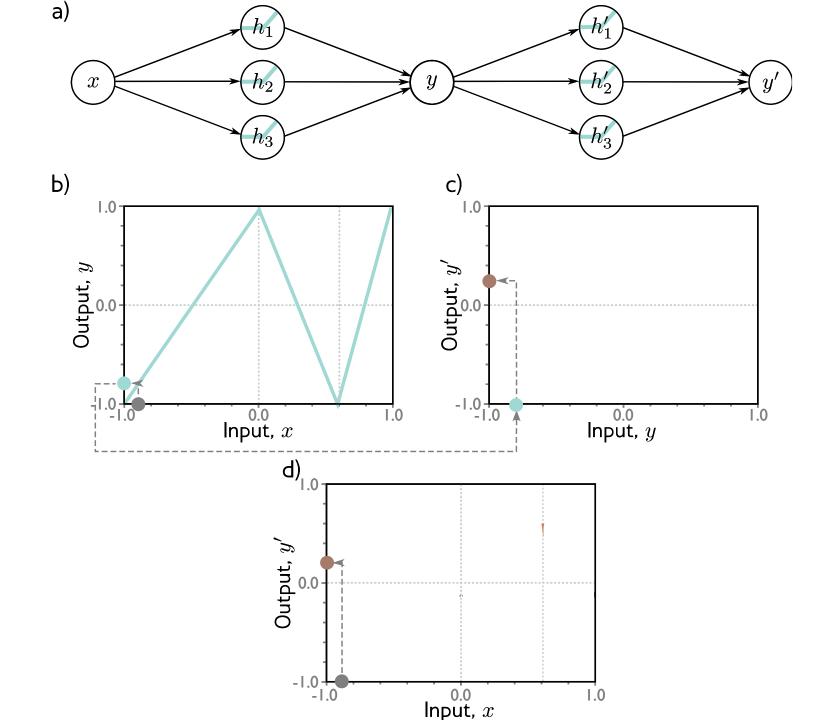


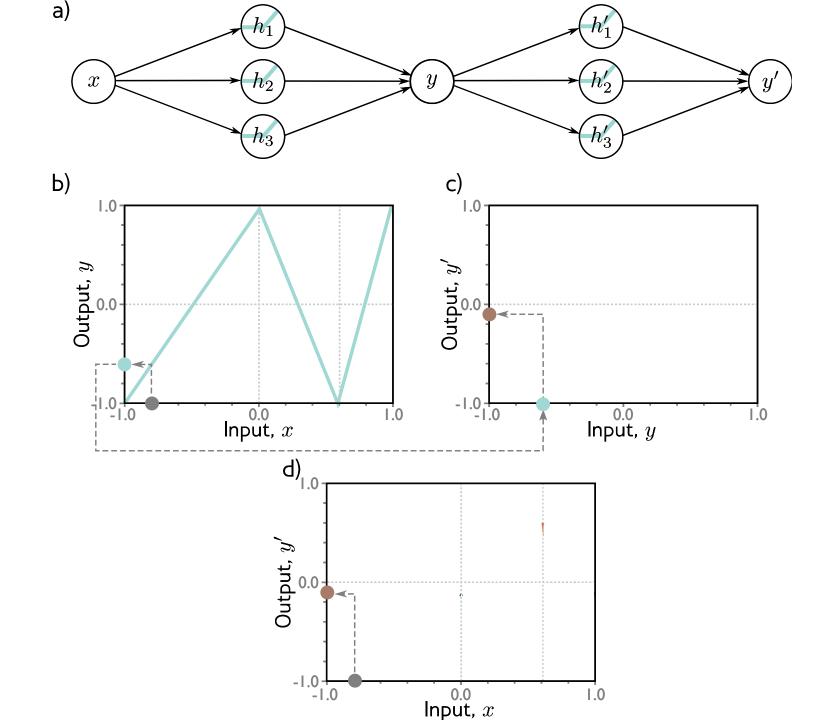


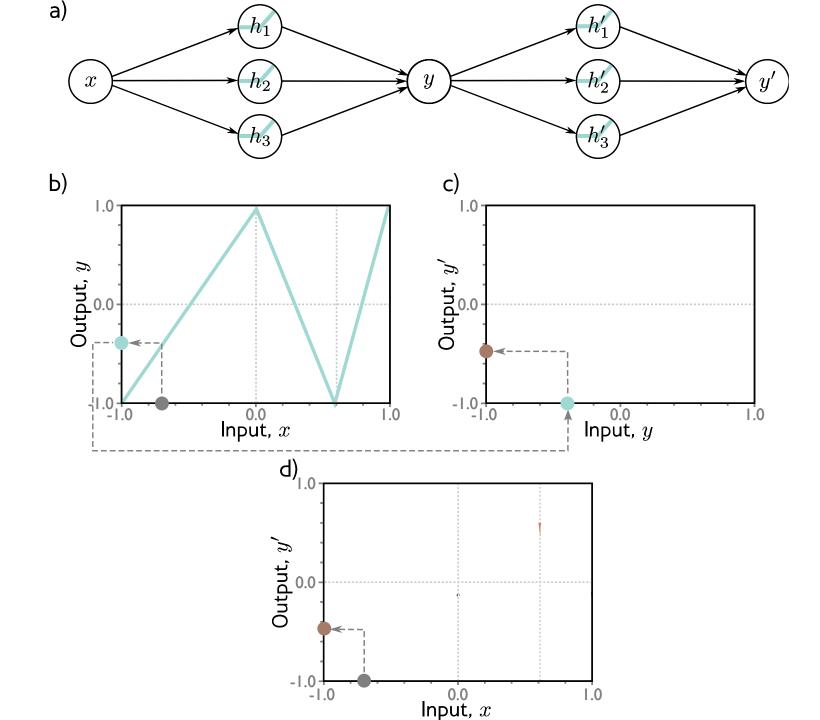


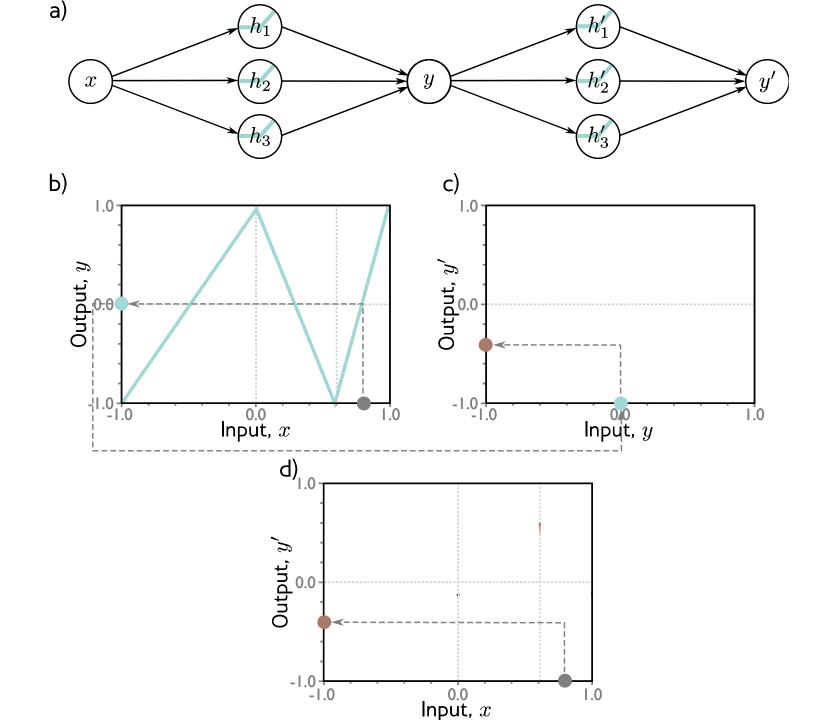
1.0

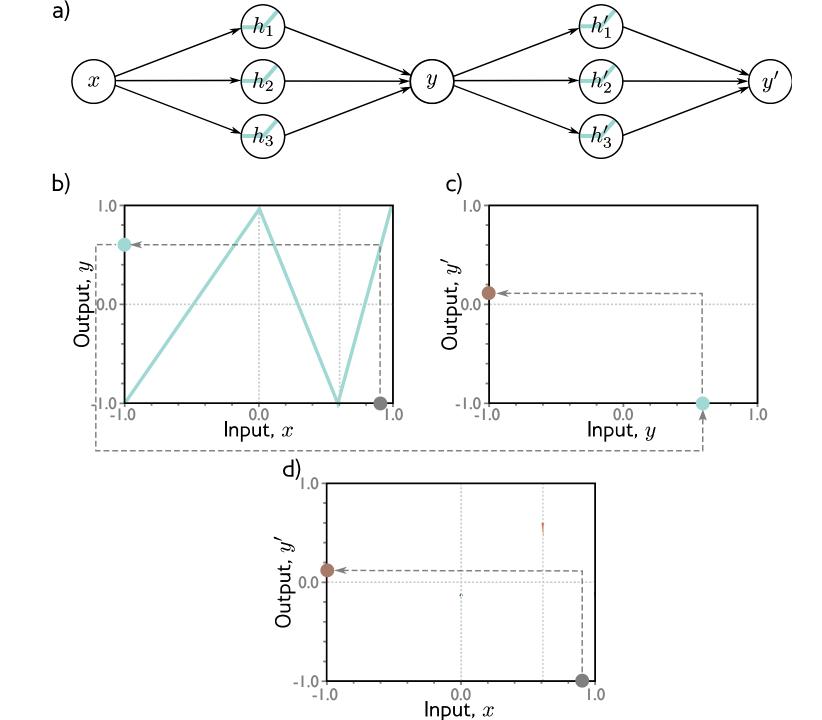




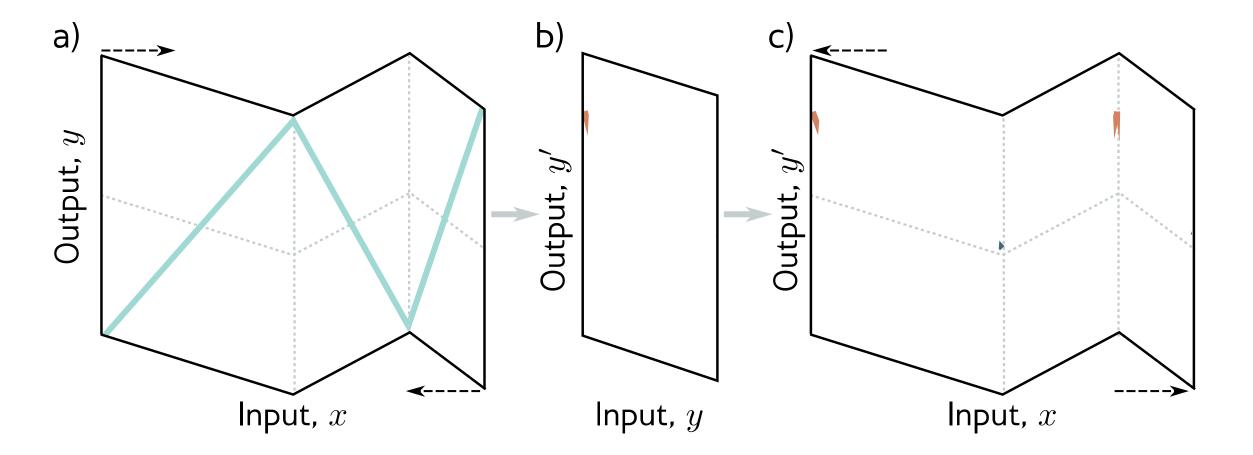




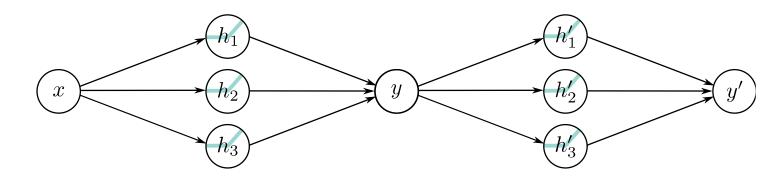




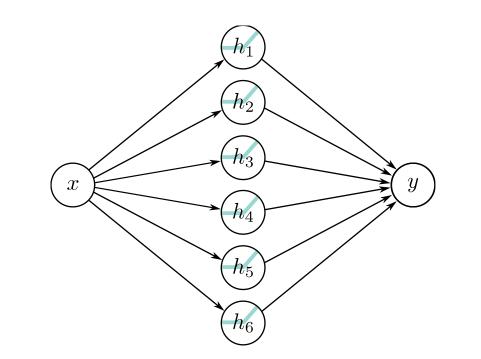




Comparing to shallow with six hidden units



- 20 parameters
- (at least) 9 regions



- 19 parameters
- Max 7 regions

Combine two networks into one

Network 1:

 $h_1 = \mathbf{a}[\theta_{10} + \theta_{11}x]$ $h_2 = \mathbf{a}[\theta_{20} + \theta_{21}x]$ $h_3 = \mathbf{a}[\theta_{30} + \theta_{31}x]$

 $h'_{1} = a[\theta'_{10} + \theta'_{11}y]$

x]
$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

x]

Network 2:

$$h'_{2} = a[\theta'_{20} + \theta'_{21}y] \qquad y' = \phi'_{0} + \phi'_{1}h'_{1} + \phi'_{2}h'_{2} + \phi'_{3}h'_{3}$$
$$h'_{3} = a[\theta'_{30} + \theta'_{31}y]$$

Hidden units of second network in terms of first:

$$\begin{aligned} h'_1 &= a[\theta'_{10} + \theta'_{11}y] &= a[\theta'_{10} + \theta'_{11}\phi_0 + \theta'_{11}\phi_1h_1 + \theta'_{11}\phi_2h_2 + \theta'_{11}\phi_3h_3] \\ h'_2 &= a[\theta'_{20} + \theta'_{21}y] &= a[\theta'_{20} + \theta'_{21}\phi_0 + \theta'_{21}\phi_1h_1 + \theta'_{21}\phi_2h_2 + \theta'_{21}\phi_3h_3] \\ h'_3 &= a[\theta'_{30} + \theta'_{31}y] &= a[\theta'_{30} + \theta'_{31}\phi_0 + \theta'_{31}\phi_1h_1 + \theta'_{31}\phi_2h_2 + \theta'_{31}\phi_3h_3] \end{aligned}$$

Create new variables

$$\begin{aligned} h_1' &= a[\theta_{10}' + \theta_{11}'y] &= a[\theta_{10}' + \theta_{11}'\phi_0 + \theta_{11}'\phi_1h_1 + \theta_{11}'\phi_2h_2 + \theta_{11}'\phi_3h_3] \\ h_2' &= a[\theta_{20}' + \theta_{21}'y] &= a[\theta_{20}' + \theta_{21}'\phi_0 + \theta_{21}'\phi_1h_1 + \theta_{21}'\phi_2h_2 + \theta_{21}'\phi_3h_3] \\ h_3' &= a[\theta_{30}' + \theta_{31}'y] &= a[\theta_{30}' + \theta_{31}'\phi_0 + \theta_{31}'\phi_1h_1 + \theta_{31}'\phi_2h_2 + \theta_{31}'\phi_3h_3] \end{aligned}$$

$$h_{1}' = a[\psi_{10} + \psi_{11}h_{1} + \psi_{12}h_{2} + \psi_{13}h_{3}]$$

$$h_{2}' = a[\psi_{20} + \psi_{21}h_{1} + \psi_{22}h_{2} + \psi_{23}h_{3}]$$

$$h_{3}' = a[\psi_{30} + \psi_{31}h_{1} + \psi_{32}h_{2} + \psi_{33}h_{3}]$$

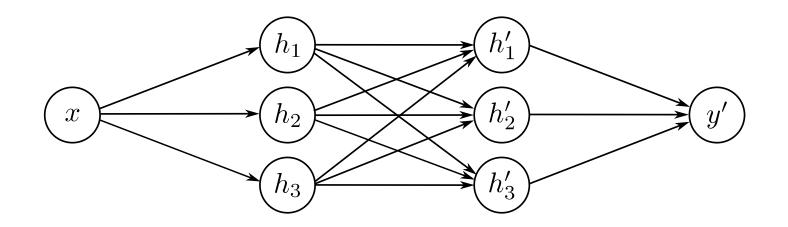
Two-layer network

$$h_{1} = a[\theta_{10} + \theta_{11}x] \qquad h'_{1} = a[\psi_{10} + \psi_{11}h_{1} + \psi_{12}h_{2} + \psi_{13}h_{3}]$$

$$h_{2} = a[\theta_{20} + \theta_{21}x] \qquad h'_{2} = a[\psi_{20} + \psi_{21}h_{2} + \psi_{22}h_{2} + \psi_{23}h_{3}]$$

$$h_{3} = a[\theta_{30} + \theta_{31}x] \qquad h'_{3} = a[\psi_{30} + \psi_{31}h_{2} + \psi_{32}h_{2} + \psi_{33}h_{3}]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$



Two-layer network as one equation

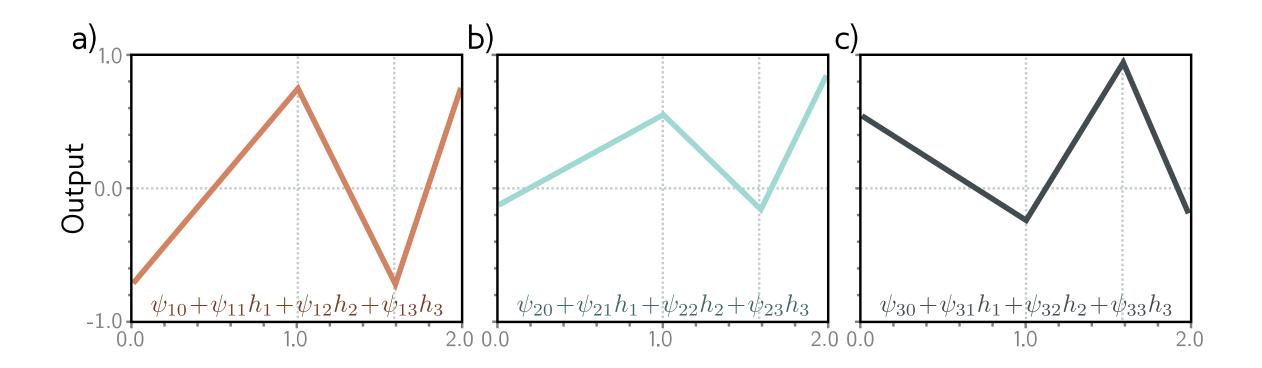
$$h_{1} = a[\theta_{10} + \theta_{11}x] \qquad h_{1}' = a[\psi_{10} + \psi_{11}h_{1} + \psi_{12}h_{2} + \psi_{13}h_{3}]$$

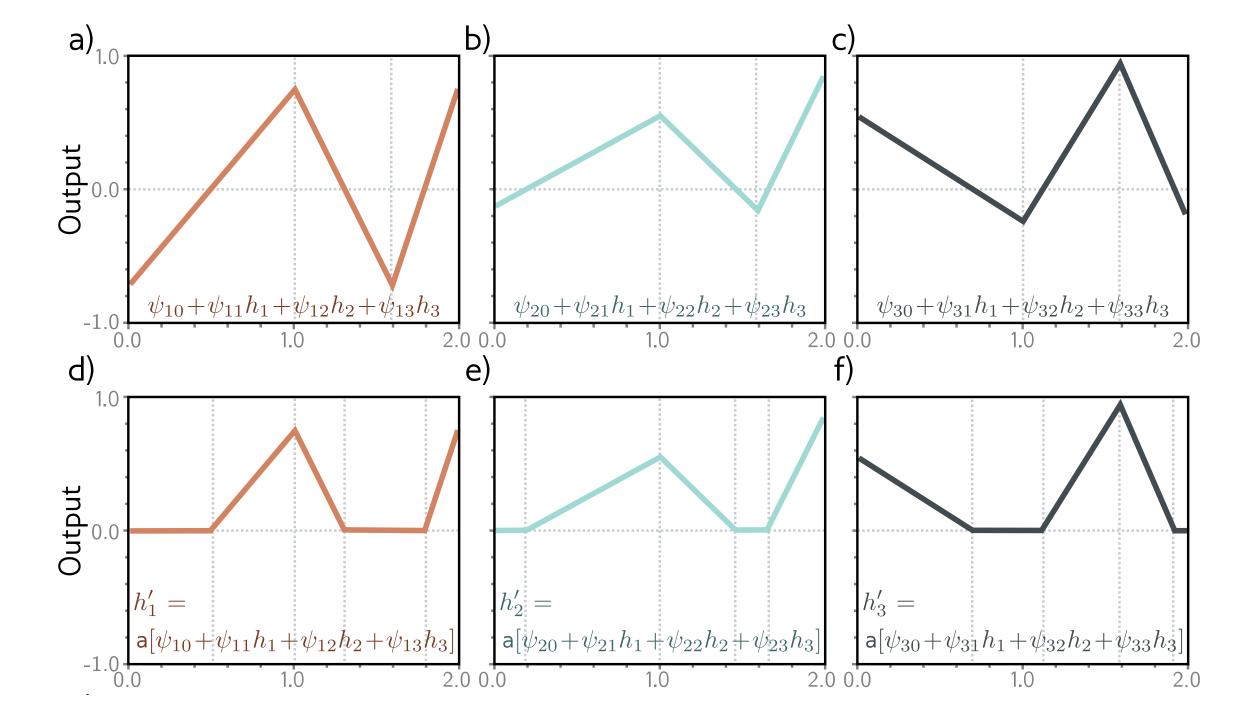
$$h_{2} = a[\theta_{20} + \theta_{21}x] \qquad h_{2}' = a[\psi_{20} + \psi_{21}h_{1} + \psi_{22}h_{2} + \psi_{23}h_{3}]$$

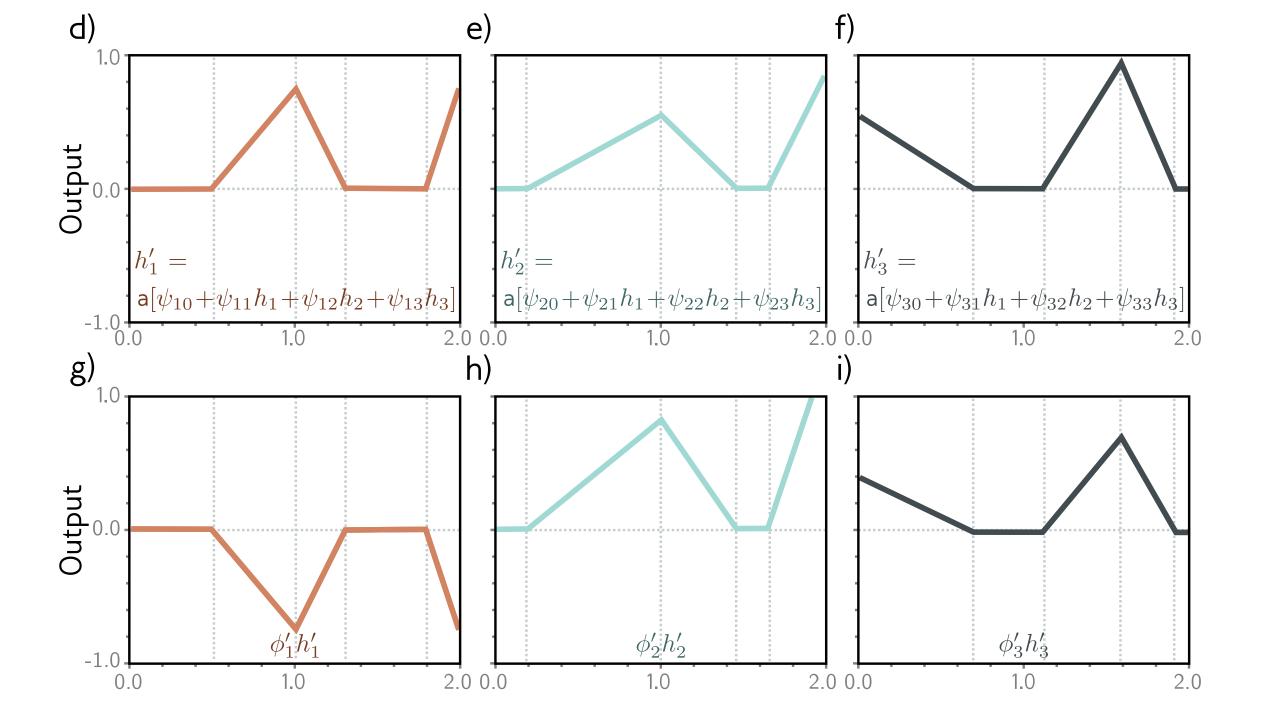
$$h_{3} = a[\theta_{30} + \theta_{31}x] \qquad h_{3}' = a[\psi_{30} + \psi_{31}h_{1} + \psi_{32}h_{2} + \psi_{33}h_{3}]$$

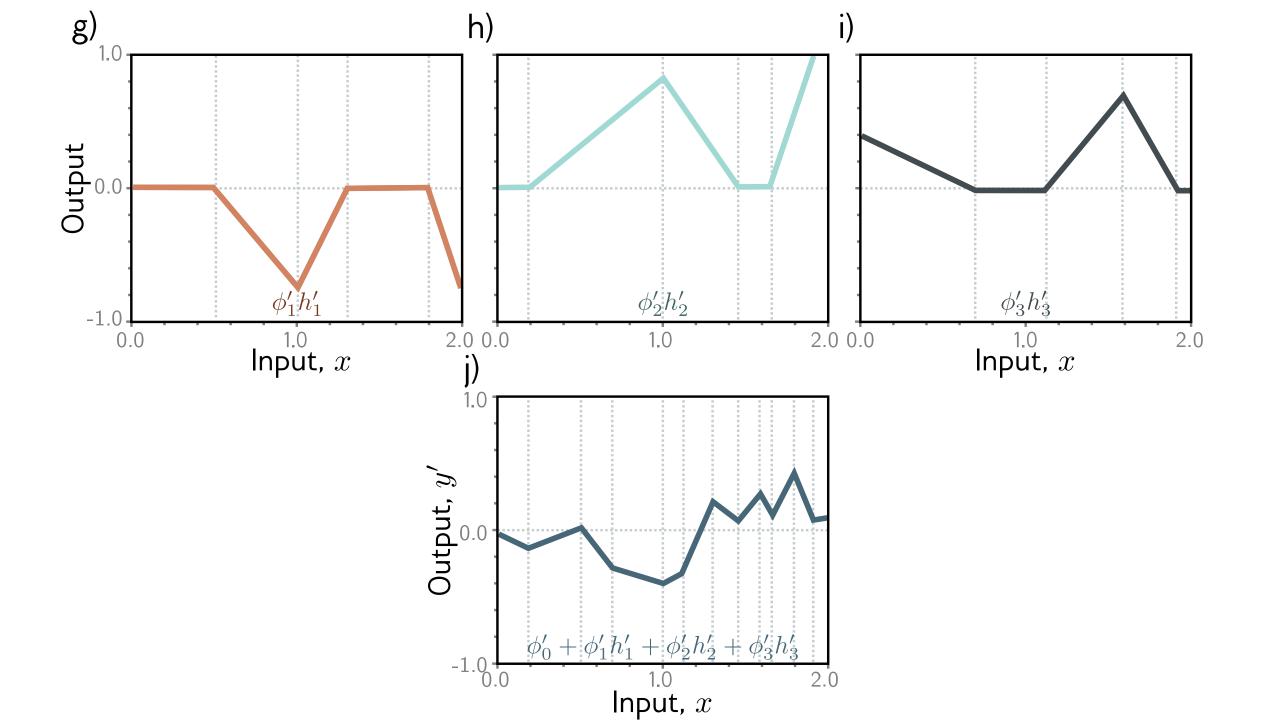
$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

 $y' = \phi'_{0} + \phi'_{1}a \left[\psi_{10} + \psi_{11}a[\theta_{10} + \theta_{11}x] + \psi_{12}a[\theta_{20} + \theta_{21}x] + \psi_{13}a[\theta_{30} + \theta_{31}x]\right]$ $+ \phi'_{2}a[\psi_{20} + \psi_{21}a[\theta_{10} + \theta_{11}x] + \psi_{22}a[\theta_{20} + \theta_{21}x] + \psi_{23}a[\theta_{30} + \theta_{31}x]]$ $+ \phi'_{3}a[\psi_{30} + \psi_{31}a[\theta_{10} + \theta_{11}x] + \psi_{32}a[\theta_{20} + \theta_{21}x] + \psi_{33}a[\theta_{30} + \theta_{31}x]]$









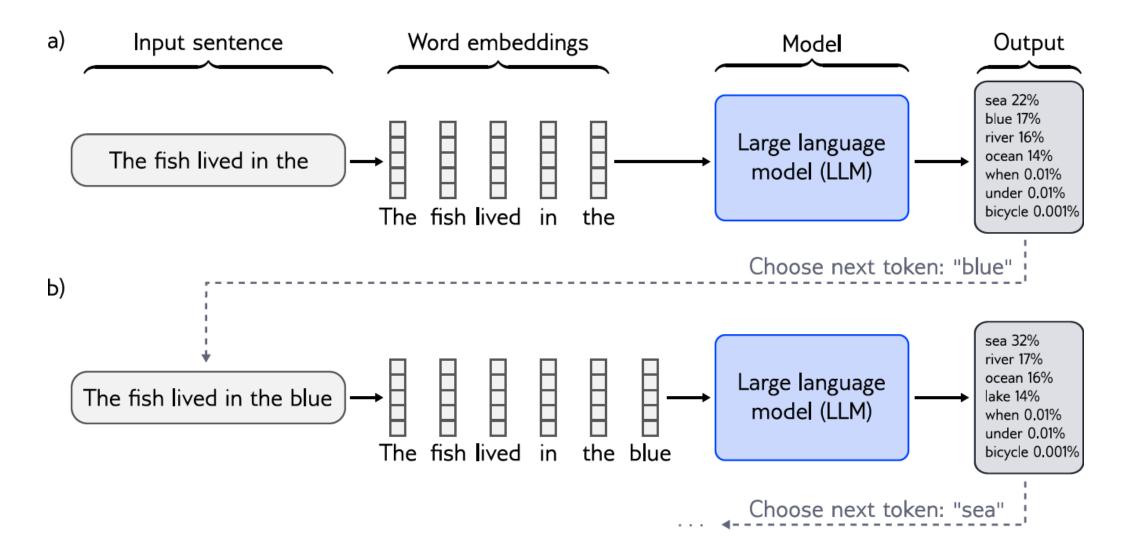
4. Neural Networks: Playground Exercises

 <u>https://developers.google.com/machine-learning/crash-</u> <u>course/introduction-to-neural-networks/playground-exercises</u>

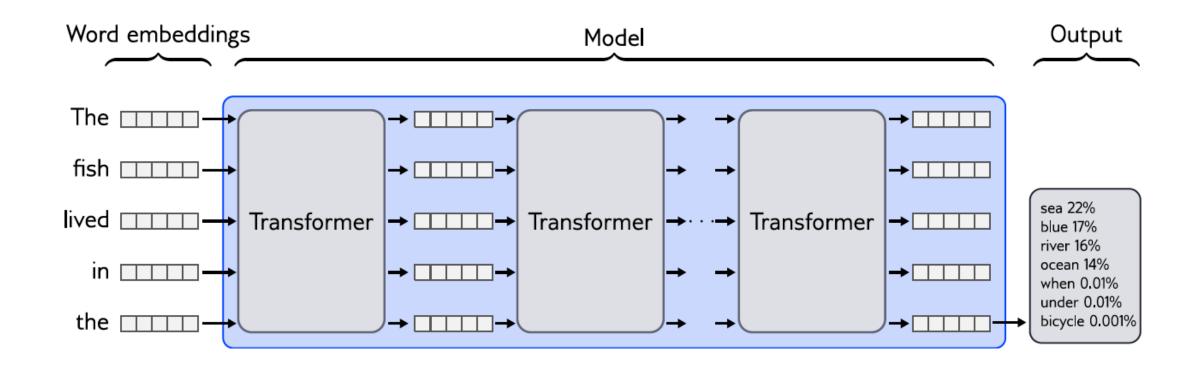
5. nanoGPT demo

- by Andrei Karpathy
- <u>https://www.youtube.com/watch?v=kCc8FmEb1nY</u>
- <u>https://www.youtube.com/watch?v=6jTQ61tBeoQ</u>

Decoder model

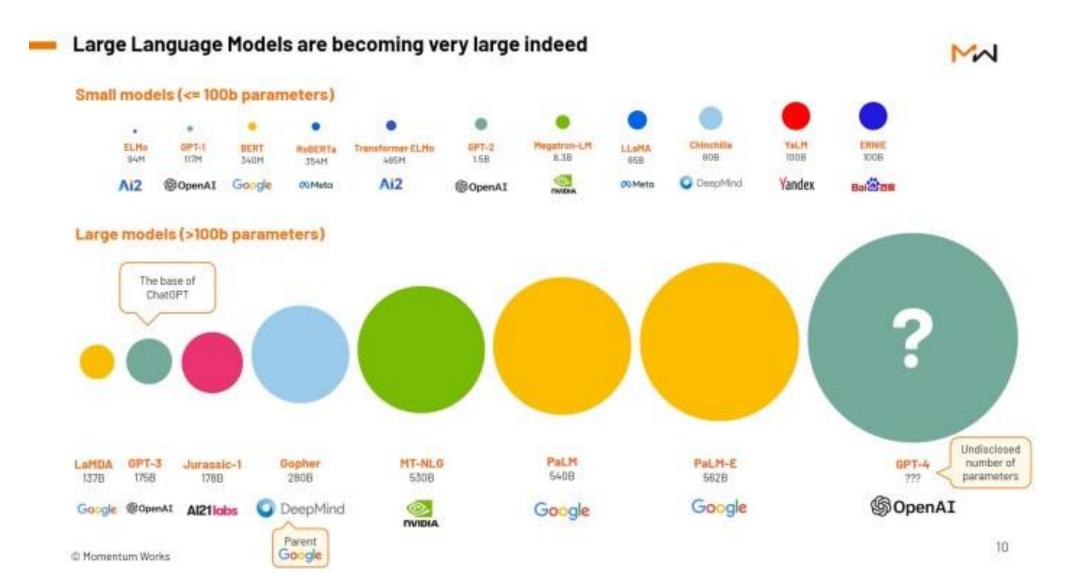


Predicting next



GPT3 (Brown et al. 2020)

- Sequence lengths are 2048 tokens long
- Batch size is 3.2 million tokens.
- 96 transformer layers (some of which implement a sparse version of attention), each of which processes a word embedding of size 12288.
- 96 heads in the self-attention layers and the value, query, and key dimension is 128.
- 300 billion tokens
- 175 billion parameters



Conclusion

- Demos of linear regression, shallow NN and deep NN
- NN training parameters
- Generative LLM