



ARISTOTLE UNIVERSITY OF THESSALONIKI



FACULTY OF ENGINEERING

Pattern Recognition & Machine Learning

Neural Networks

Panagiotis C. Petrantonakis

Assistant Professor

Dept. of Electrical and Computer Engineering

ppetrant@ece.auth.gr

Fall Semester

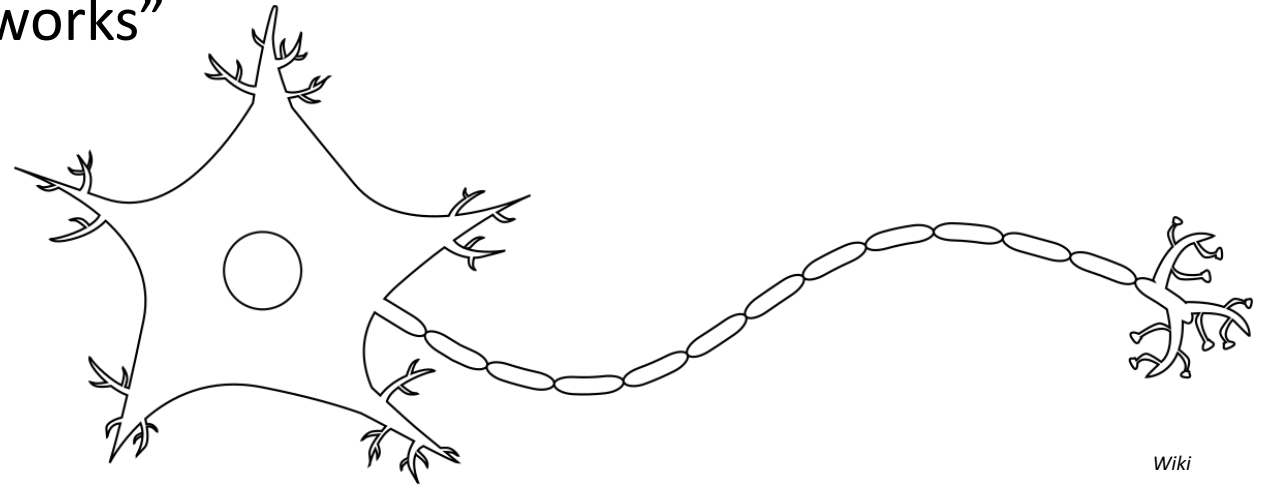
Until now...

- Dataset: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ and targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D$
- Fixed features: $\boldsymbol{\varphi}(\mathbf{x}) = (\varphi_0(\mathbf{x}), \dots, \varphi_M(\mathbf{x}))$, where $\varphi_0 = 1$
- Classification: $y(\mathbf{x}, \mathbf{w}) = f(\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}))$
 - f : nonlinear activation function (e.g., logistic sigmoid)
- In the SVM case we had basis functions centered on training samples and then selecting a number of them to define the decision rule (support vectors)
 - convex objective function
 - again fixed features $\boldsymbol{\varphi}(\mathbf{x})$
- What if we want adaptive basis functions?
 - Neural Networks!

Neural Network

- Artificial “Neuron” and “Neural Networks”

- Mathematical interpretation of information processing in the brain.
- McCulloch and Pitts (1943)
- Widrow and Hoff (1960)
- Rosenblatt (1962)
- Rumelhart et al. (1986)
- etc. etc.



- A vast variety of different implementations and algorithms regarding the above mentioned mathematical interpretation.
 - Biological plausibility
- Here we will talk for a class of NN that are of great practical value:
 - Multilayer Perceptrons (MLP)!

Perceptron

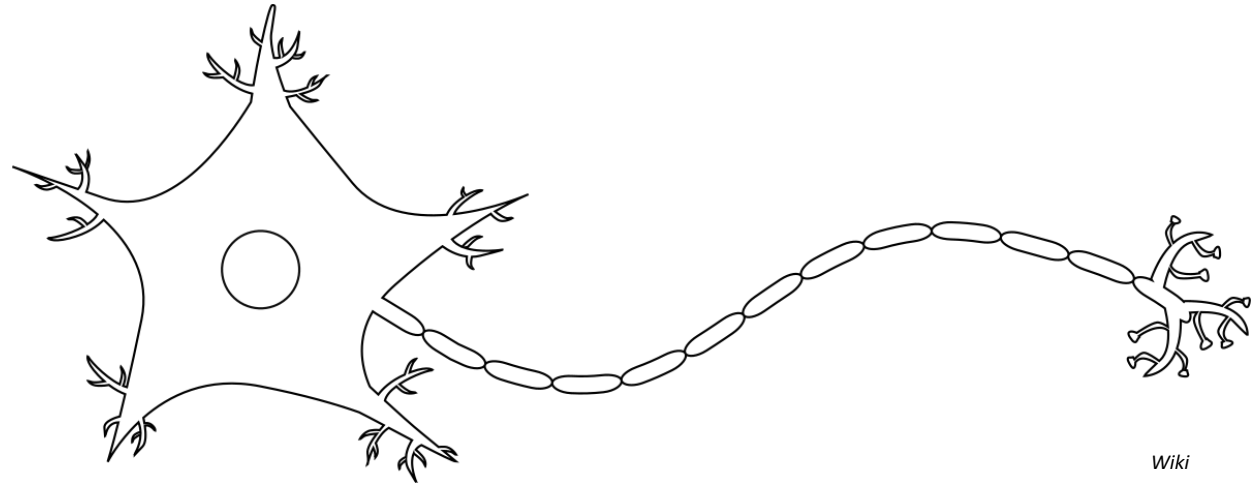
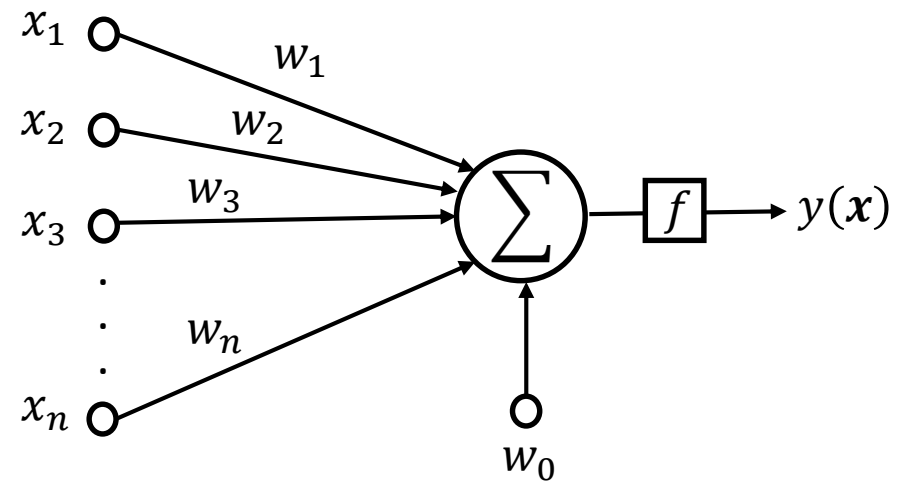
$$y(\mathbf{x}) = f(\mathbf{w}^t \mathbf{x} + w_0)$$

If we incorporate w_0 :

$$y(\mathbf{x}) = f(\mathbf{w}^t \mathbf{x})$$

In the general case:

$$y(\mathbf{x}) = f(\mathbf{w}^t \boldsymbol{\varphi}(\mathbf{x}))$$



Wiki

From a single perceptron to a MLP

x_1 ○
 x_2 ○
·
·
·
 x_D ○

From a single perceptron to a MLP

x_0 ●

x_1 ○

x_2 ○

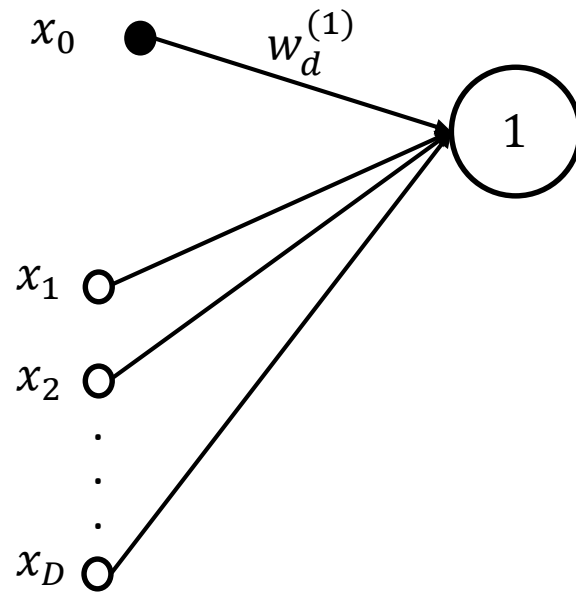
⋮

⋮

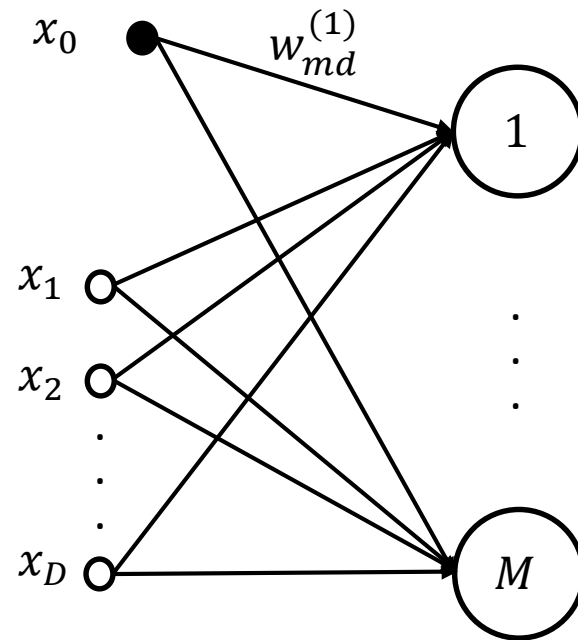
⋮

x_D ○

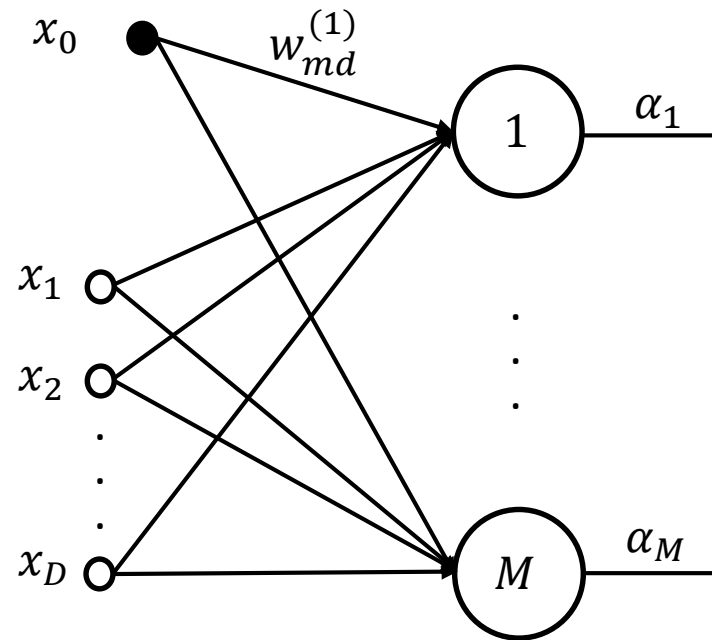
From a single perceptron to a MLP



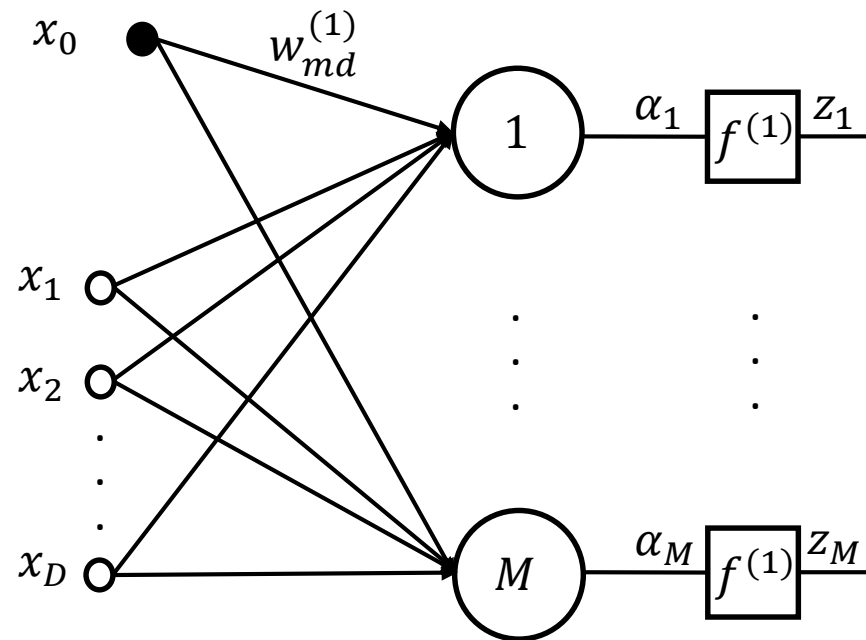
From a single perceptron to a MLP



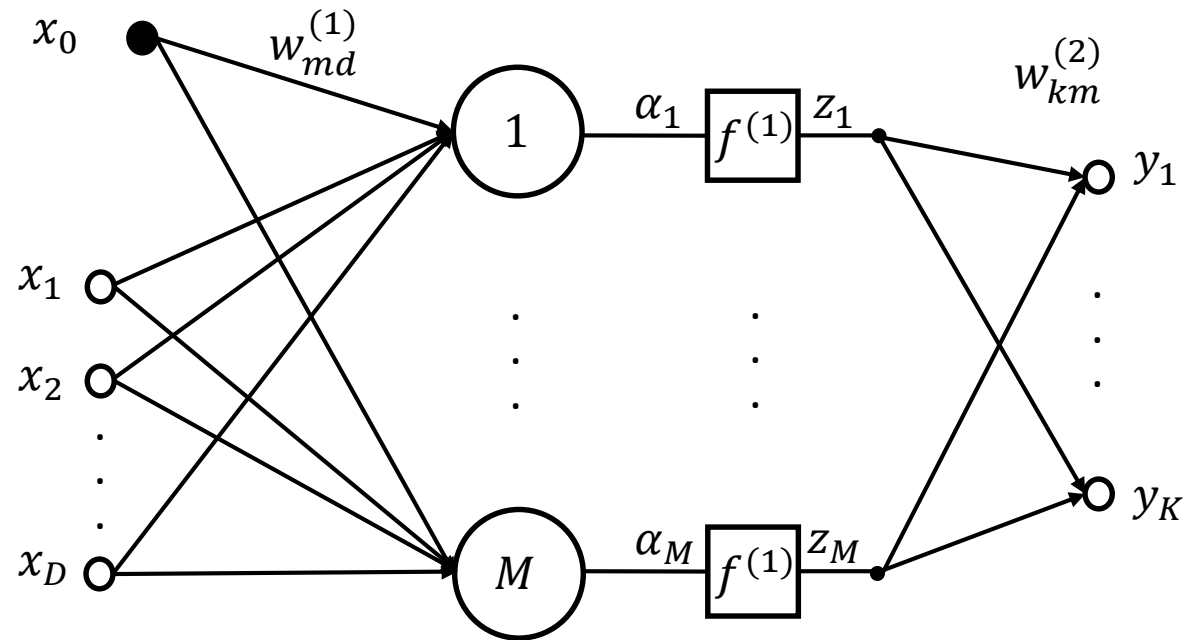
From a single perceptron to a MLP



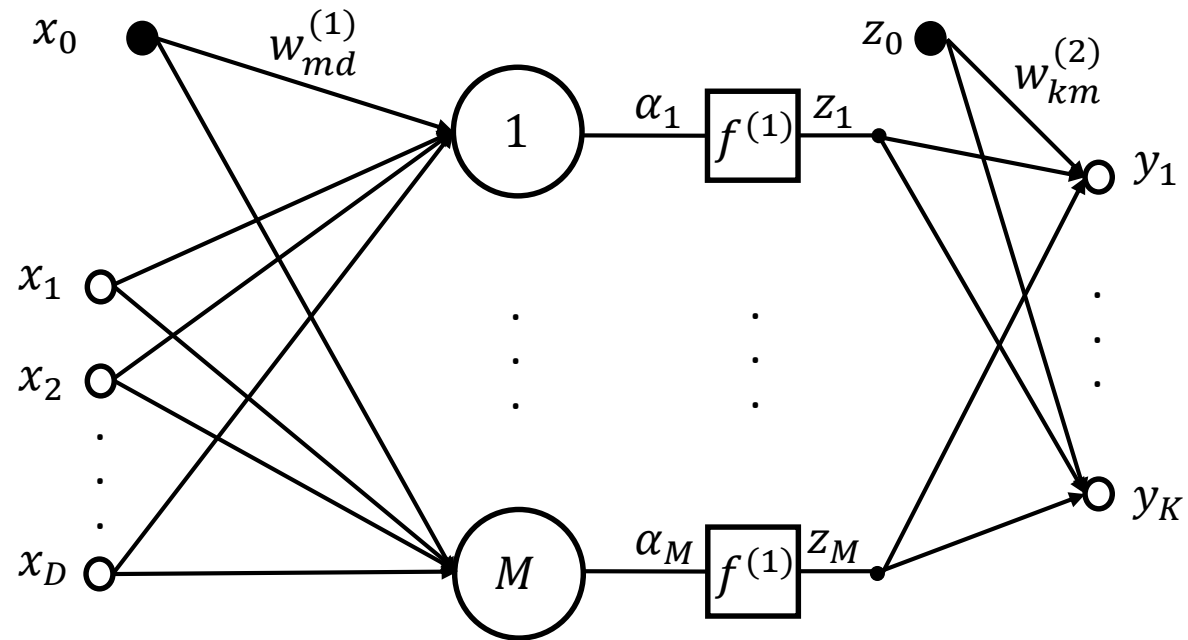
From a single perceptron to a MLP



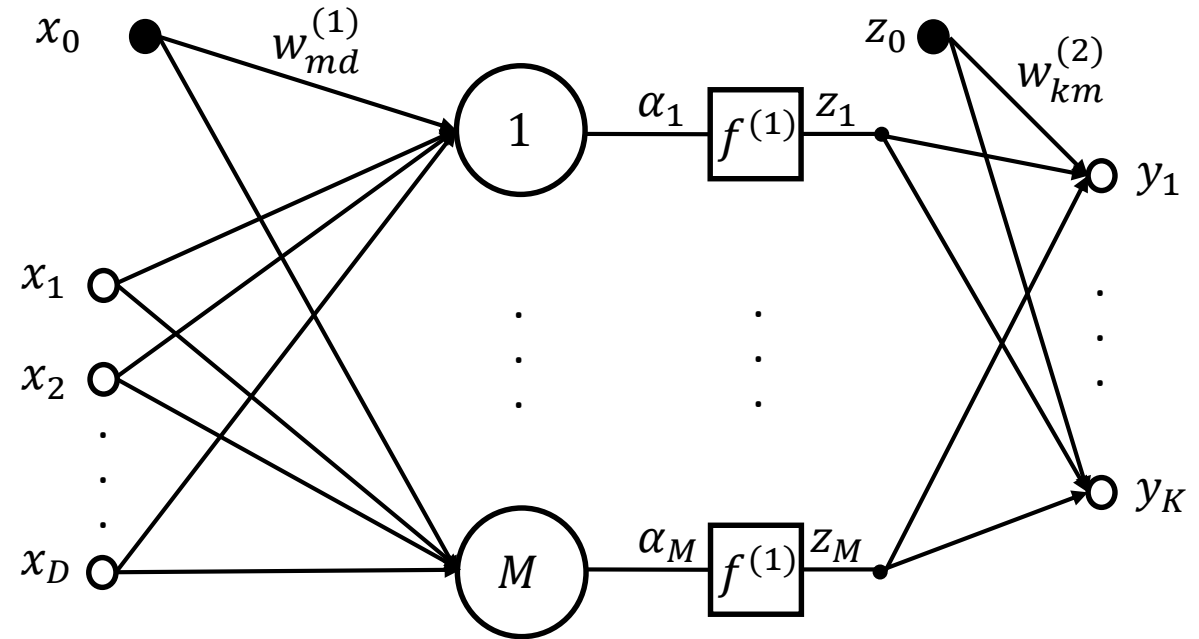
Two layer-Perceptron



Two layer-Perceptron

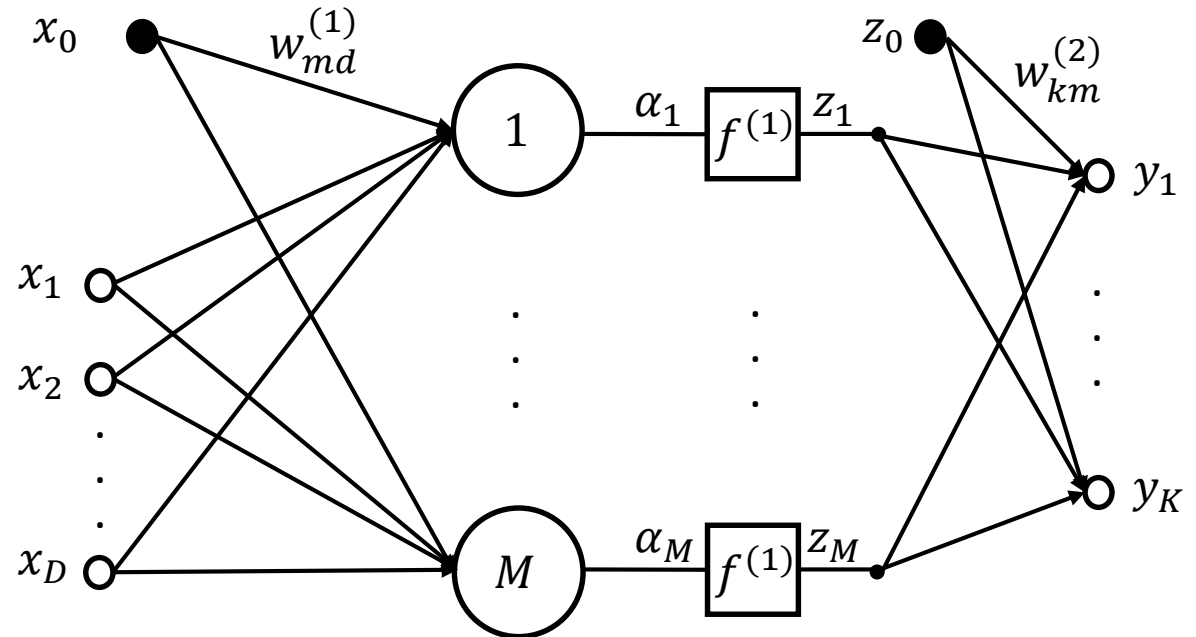


Two layer-Perceptron



Two layer-Perceptron

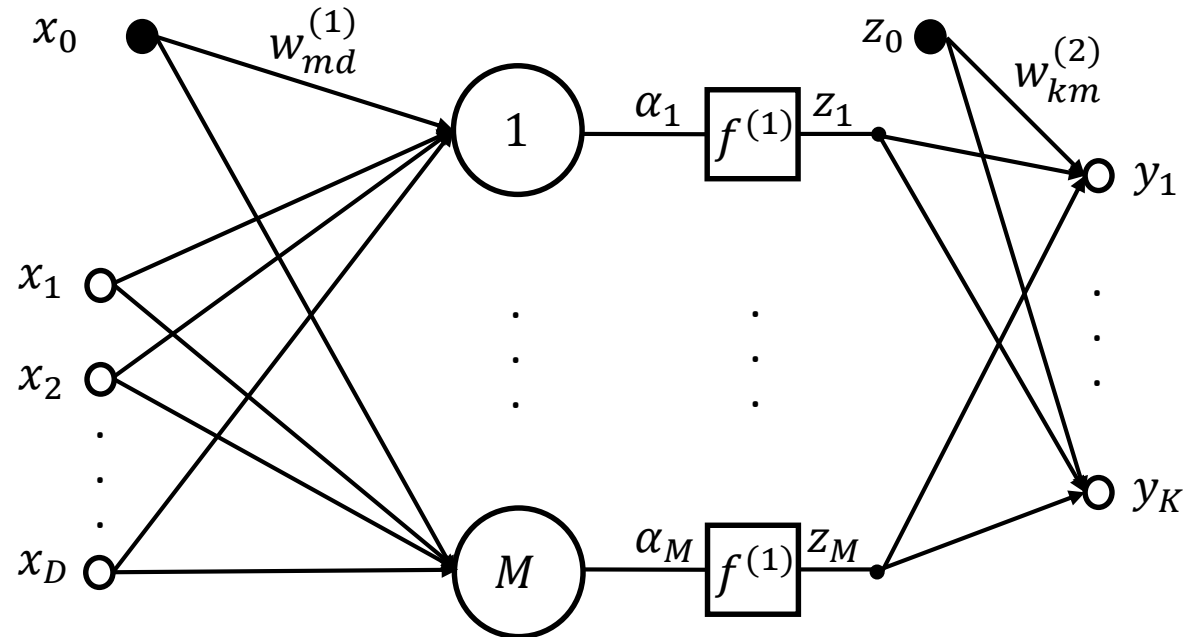
$$y_k = \sum_{m=0}^M w_{km}^{(2)} z_m$$



Two layer-Perceptron

$$y_k = \sum_{m=0}^M w_{km}^{(2)} z_m$$

$$y_k = \sum_{m=0}^M w_{km}^{(2)} f^{(1)}(\alpha_m)$$

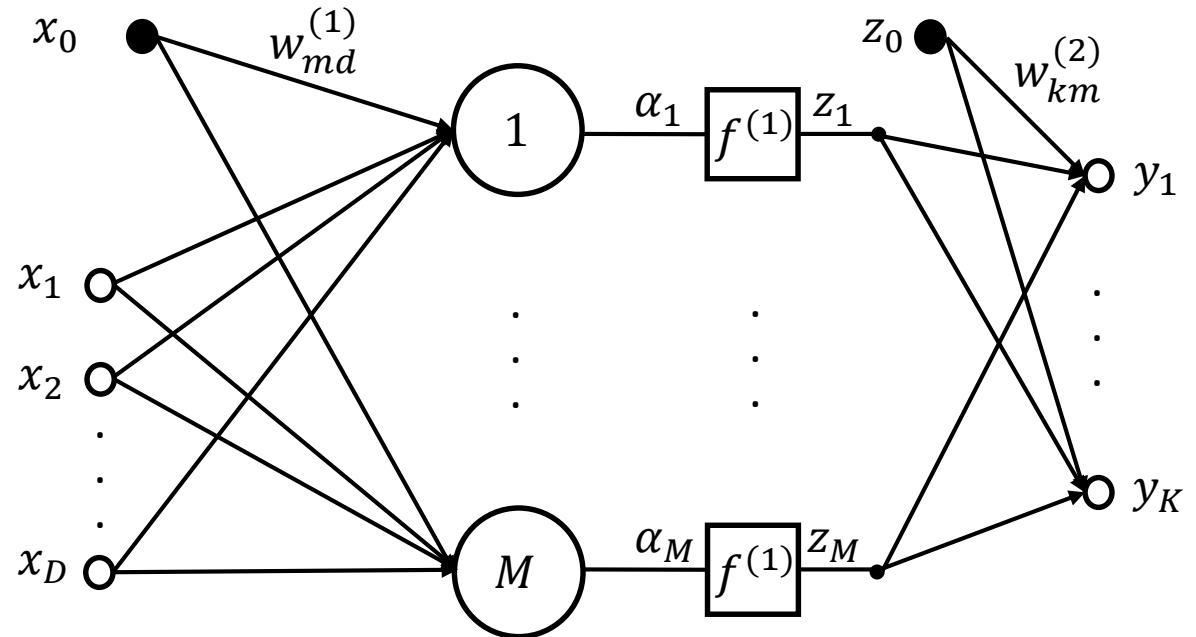


Two layer-Perceptron

$$y_k = \sum_{m=0}^M w_{km}^{(2)} z_m$$

$$y_k = \sum_{m=0}^M w_{km}^{(2)} f^{(1)}(\alpha_m)$$

$$y_k = \sum_{m=0}^M w_{km}^{(2)} f^{(1)}\left(\sum_{d=0}^D w_{md}^{(1)} x_d\right)$$

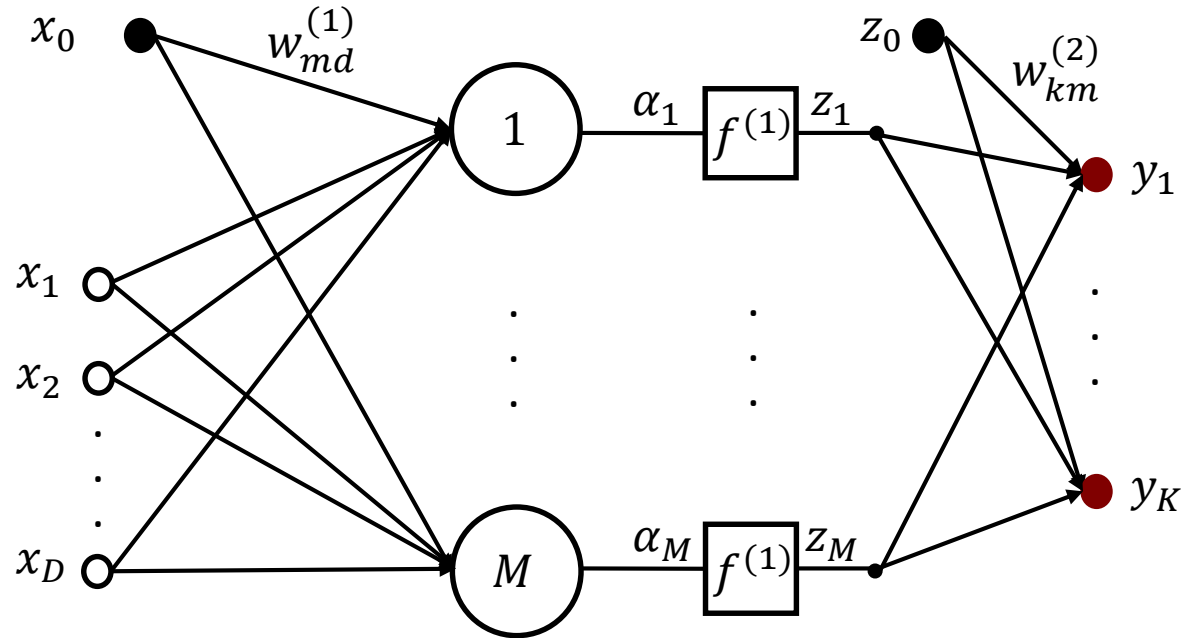


Two layer-Perceptron

$$y_k = \sum_{m=0}^M w_{km}^{(2)} z_m$$

$$y_k = \sum_{m=0}^M w_{km}^{(2)} f^{(1)}(\alpha_m)$$

$$y_k = \sum_{m=0}^M w_{km}^{(2)} f^{(1)}\left(\sum_{d=0}^D w_{md}^{(1)} x_d\right)$$



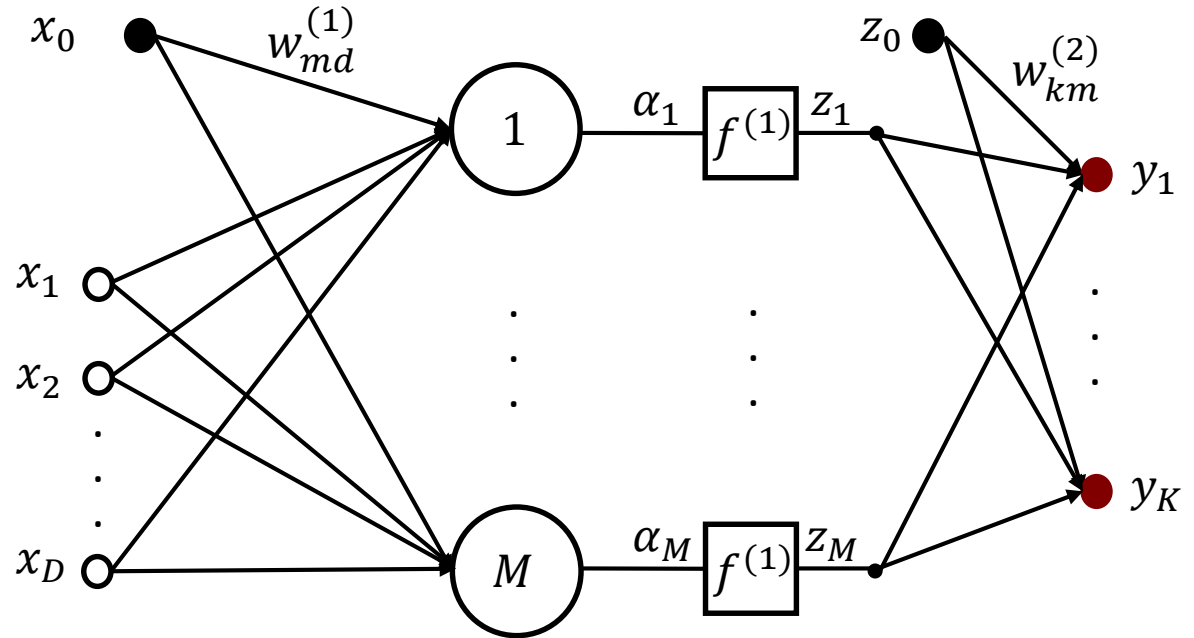
Usually: $y_k = f^{(2)}\left(\sum_{m=0}^M w_{km}^{(2)} f^{(1)}\left(\sum_{d=0}^D w_{md}^{(1)} x_d\right)\right)$

Two layer-Perceptron

$$y_k = \sum_{m=0}^M w_{km}^{(2)} z_m$$

$$y_k = \sum_{m=0}^M w_{km}^{(2)} f^{(1)}(\alpha_m)$$

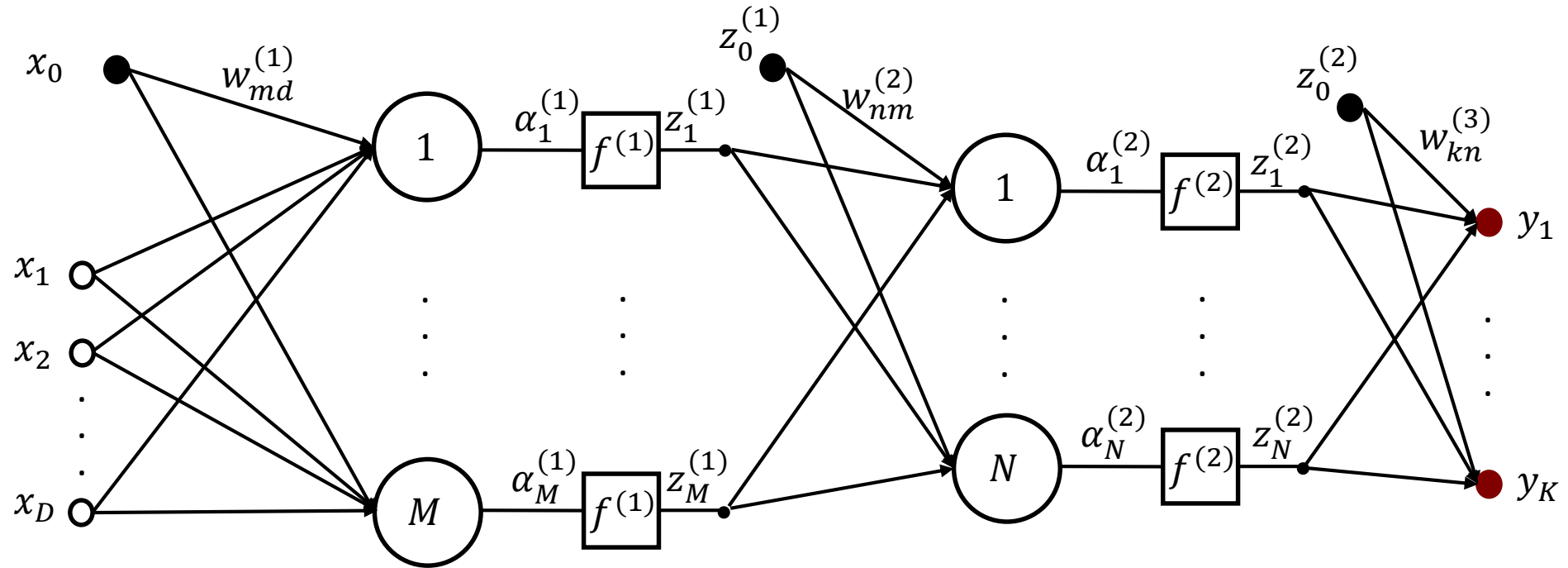
$$y_k = \sum_{m=0}^M w_{km}^{(2)} f^{(1)}\left(\sum_{d=0}^D w_{md}^{(1)} x_d\right)$$



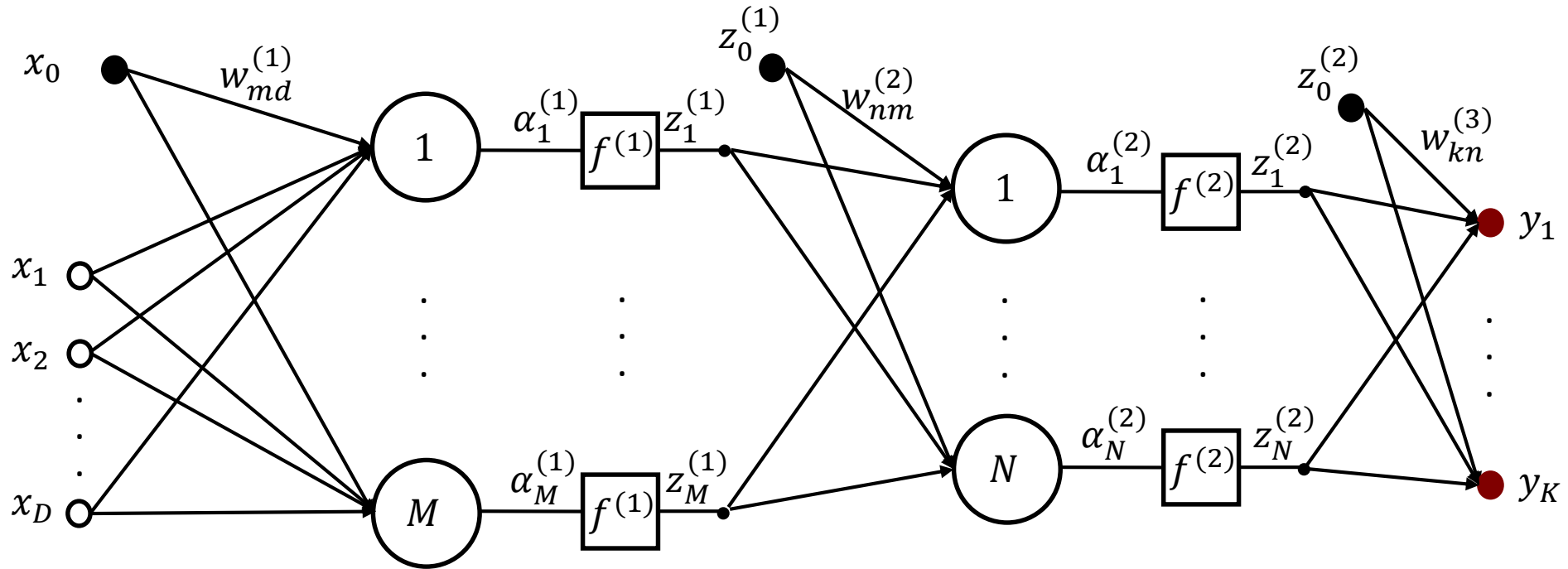
Usually: $y_k = f^{(2)}\left(\sum_{m=0}^M w_{km}^{(2)} f^{(1)}\left(\sum_{d=0}^D w_{md}^{(1)} x_d\right)\right) = f^{(2)}\left(\sum_{m=0}^M w_{km}^{(2)} \varphi(\mathbf{x})\right)$

- Thus we use parametric forms of basis functions in advance and allow them to be *adaptive*!
- Adaptive: learn through training!

Three-layer Perceptron

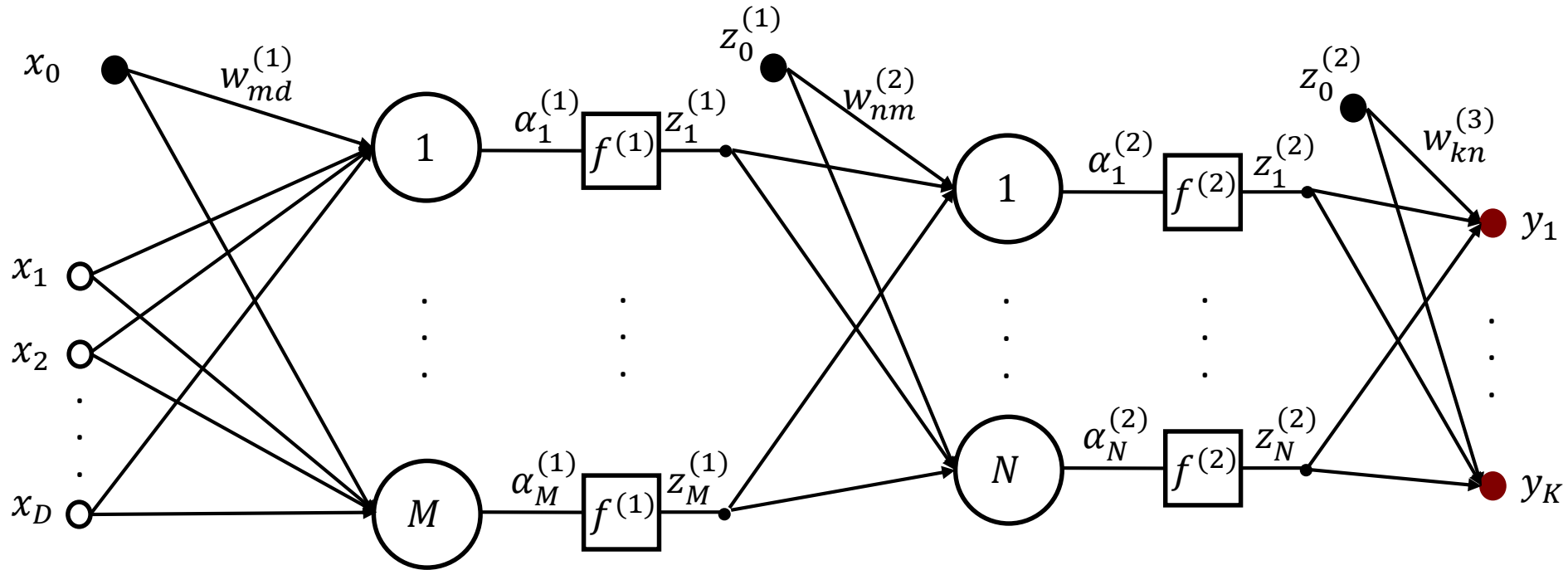


Three-layer Perceptron



$$y_k = f^{(3)} \left(\sum_{n=0}^N w_{kn}^{(3)} f^{(2)} \left(\sum_{m=0}^M w_{nm}^{(2)} f^{(1)} \left(\sum_{d=0}^D w_{md} x_d \right) \right) \right)$$

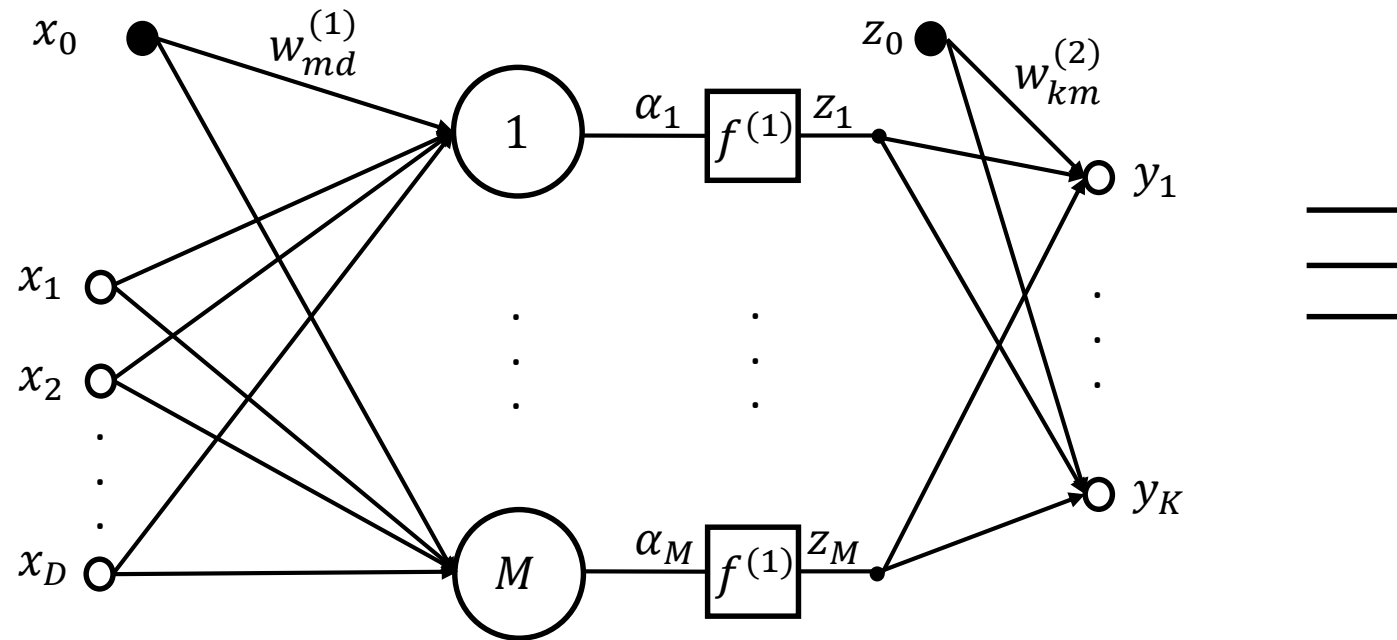
Three-layer Perceptron



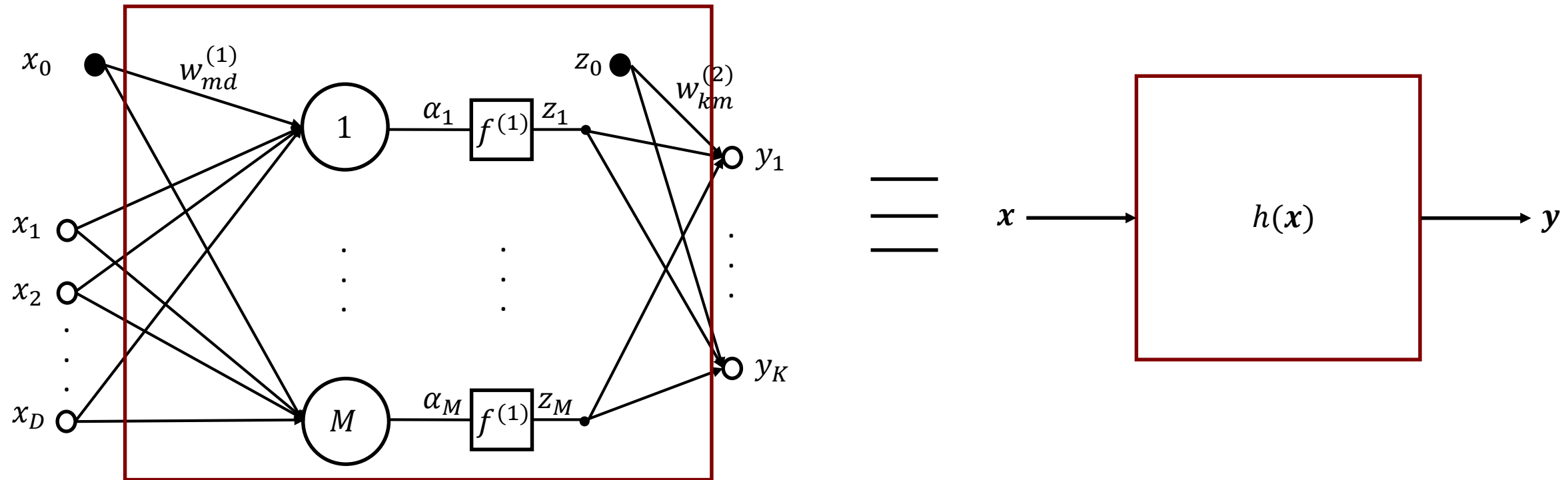
$$y_k = f^{(3)} \left(\sum_{n=0}^N w_{kn}^{(3)} f^{(2)} \left(\sum_{m=0}^M w_{nm}^{(2)} f^{(1)} \left(\sum_{d=0}^D w_{md} x_d \right) \right) \right) = f^{(3)} \left(\alpha^{(3)} \left(f^{(2)} \left(\alpha^{(2)} \left(f^{(1)} \left(\alpha^{(1)}(x) \right) \right) \right) \right) \right)$$

- Very complicated basis functions as functions of x (sequential mapping from layer to layer)

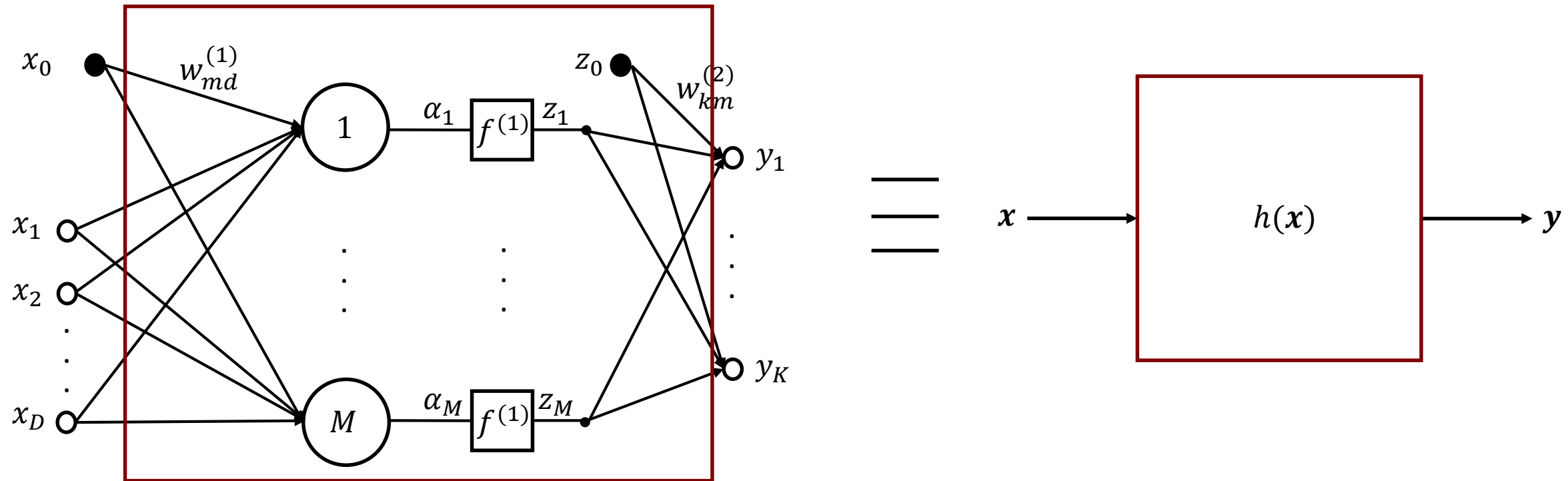
Function Approximators



Function Approximators

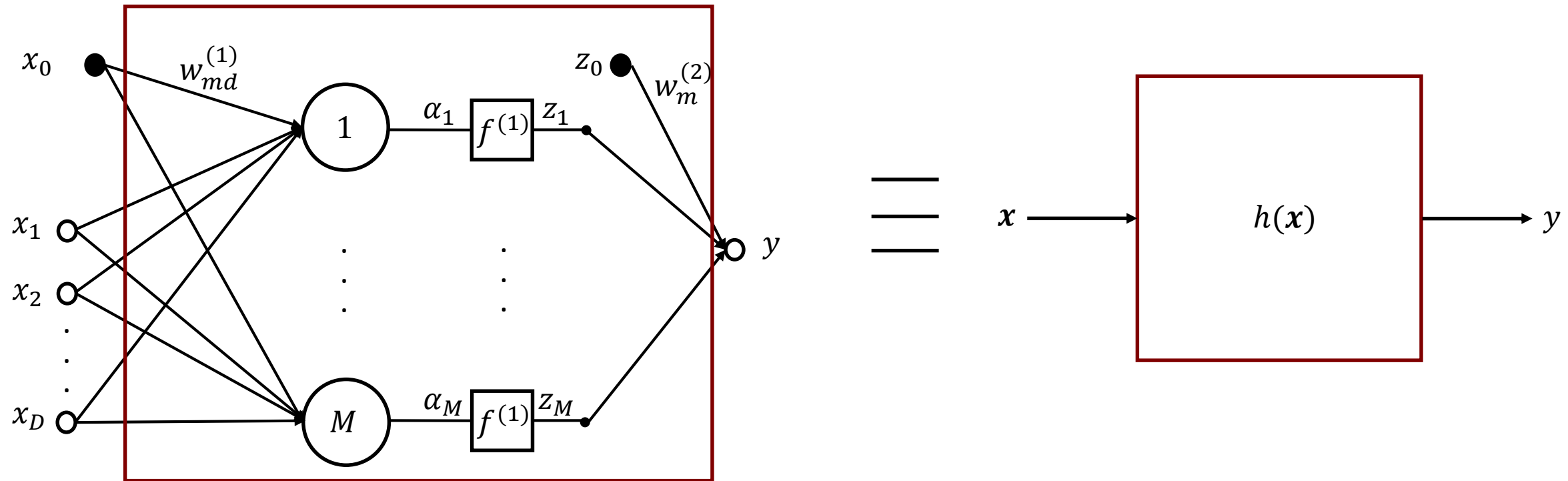


Function Approximators



$$\mathbf{y}(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}) = \sum_{m=0}^M w_{km}^{(2)} f^{(1)} \left(\sum_{d=0}^D w_{md}^{(1)} x_d \right) = \sum_{m=0}^M w_{km}^{(2)} \varphi_m(\mathbf{x})$$

Function Approximators



$$y(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{w}^{(2)}) = \sum_{m=0}^M w_m^{(2)} f^{(1)} \left(\sum_{d=0}^D w_{md}^{(1)} x_d \right) = \sum_{m=0}^M w_m^{(2)} \varphi_m(\mathbf{x})$$

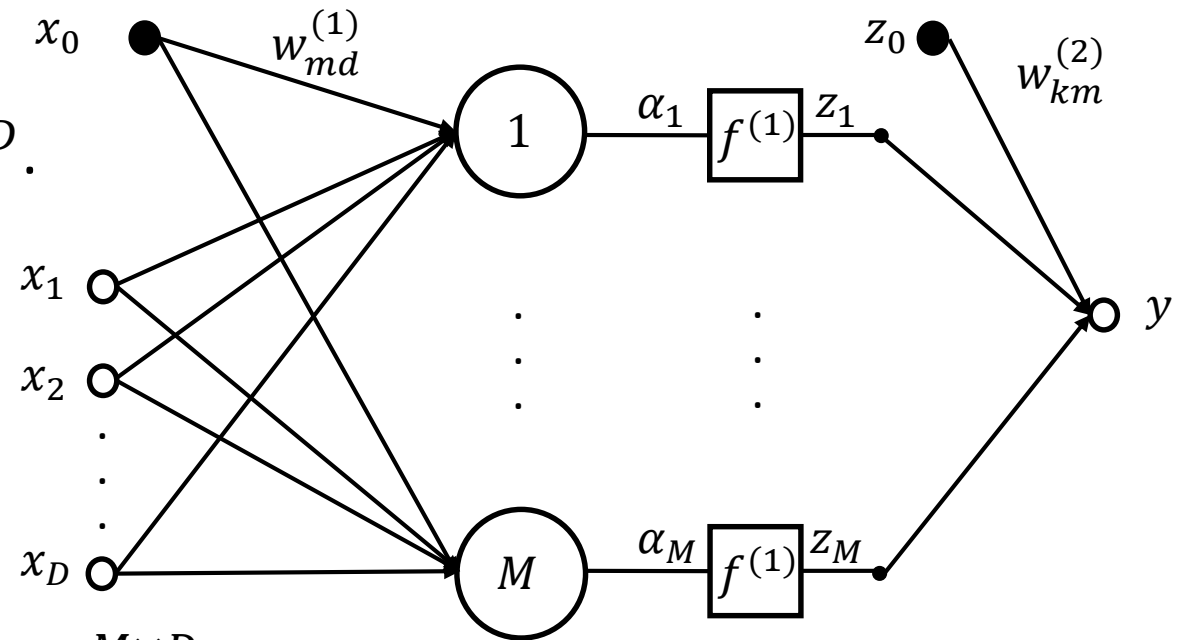
Universal approximators

Theorem:

- Let $h(\cdot)$ be a continuous function on \mathbb{R}^D .
- Let $f(\cdot)$ be a fixed analytic function which is **not** polynomial.
- Given any small number $\varepsilon > 0$, we can find a number M , $\mathbf{w}^{(2)} \in \mathbb{R}^M$ and $\mathbf{W}^{(1)} \in \mathbb{R}^{M \times D}$ such that:

$$|h(x) - y(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{w}^{(2)})| < \varepsilon$$

Note: For smaller ε we usually need larger M .



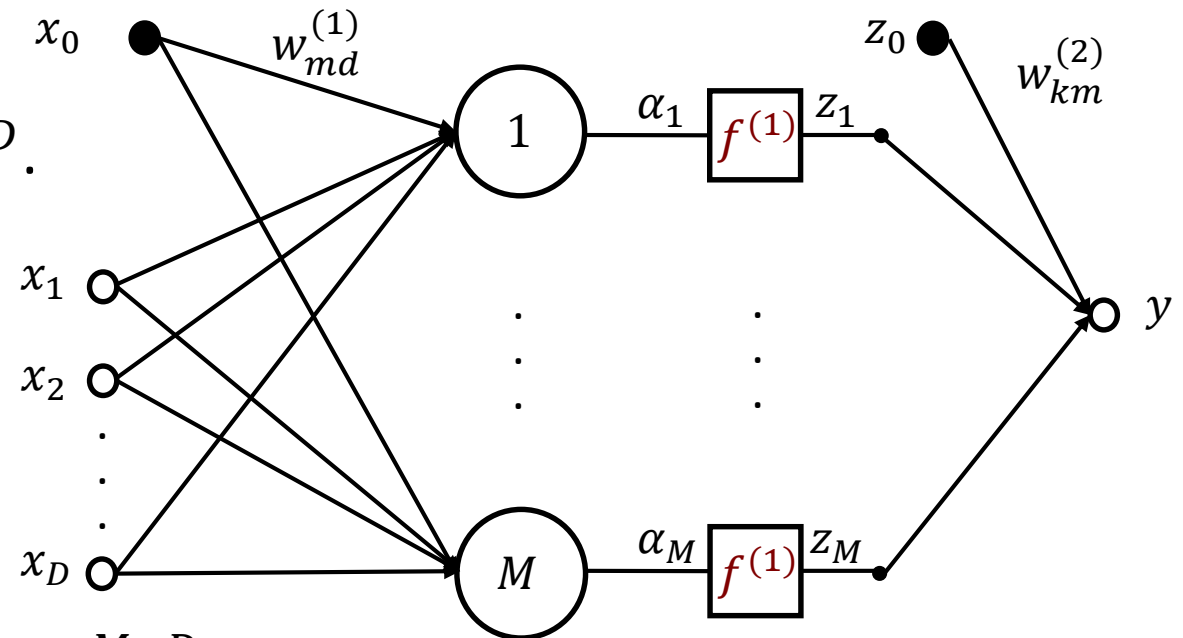
Universal approximators

Theorem:

- Let $h(\cdot)$ be a continuous function on \mathbb{R}^D .
- Let $f(\cdot)$ be a fixed analytic function which is **not** polynomial.
- Given any small number $\varepsilon > 0$, we can find a number M , $\mathbf{w}^{(2)} \in \mathbb{R}^M$ and $\mathbf{W}^{(1)} \in \mathbb{R}^{M \times D}$ such that:

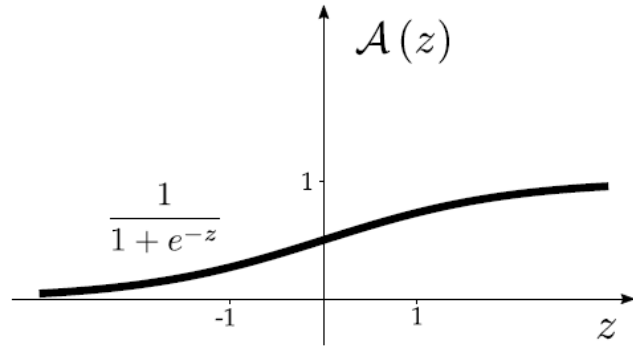
$$|h(x) - y(\mathbf{x}, \mathbf{W}^{(1)}, \mathbf{w}^{(2)})| < \varepsilon$$

Note: For smaller ε we usually need larger M .

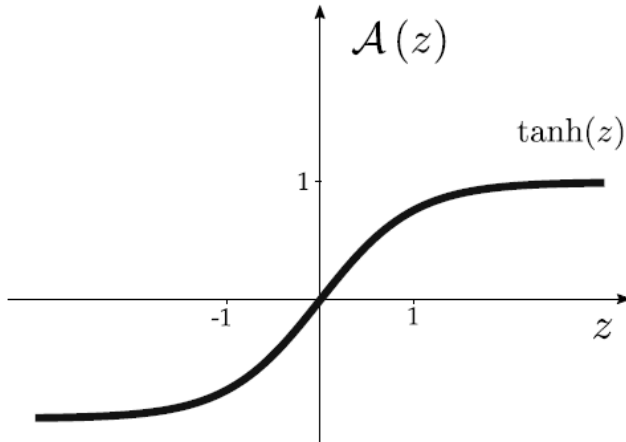


Activation Functions

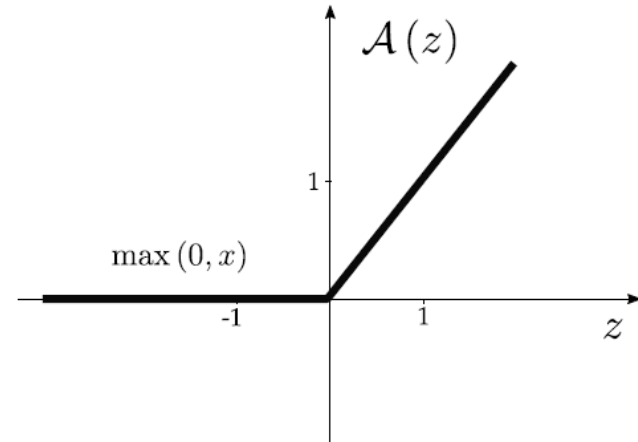
Logistic sigmoid



Tanh

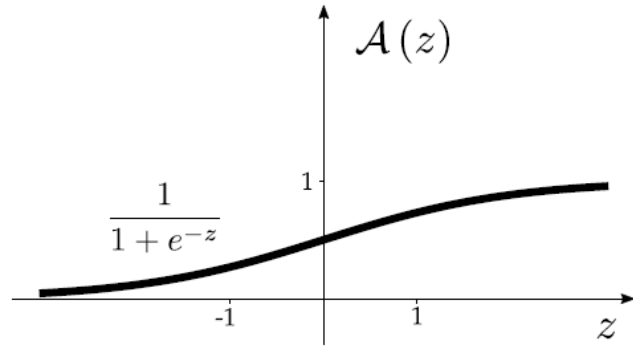


Rectified Linear Unit (RELU)

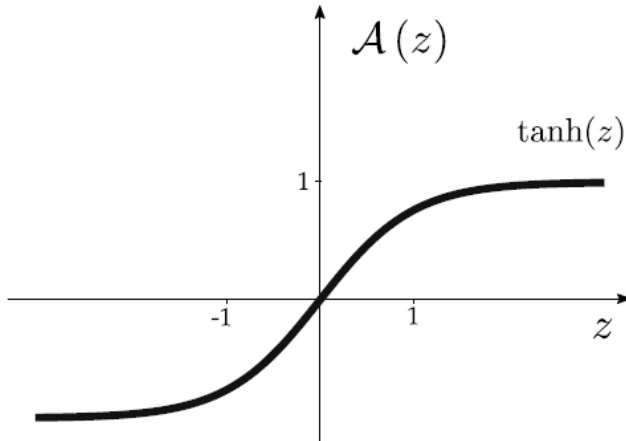


Activation Functions

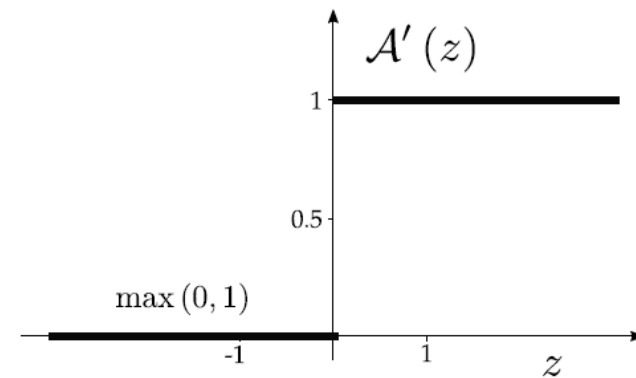
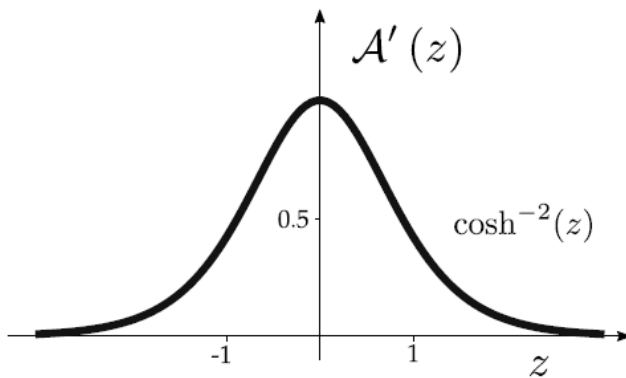
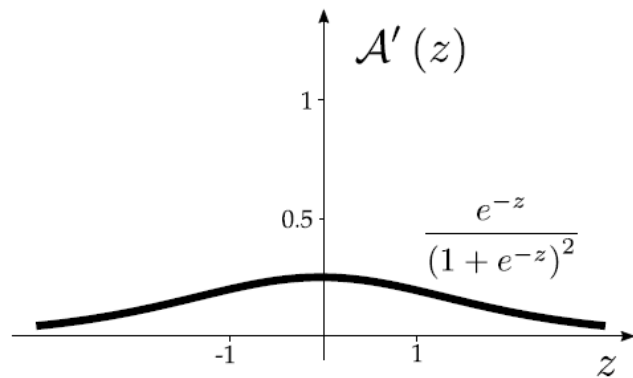
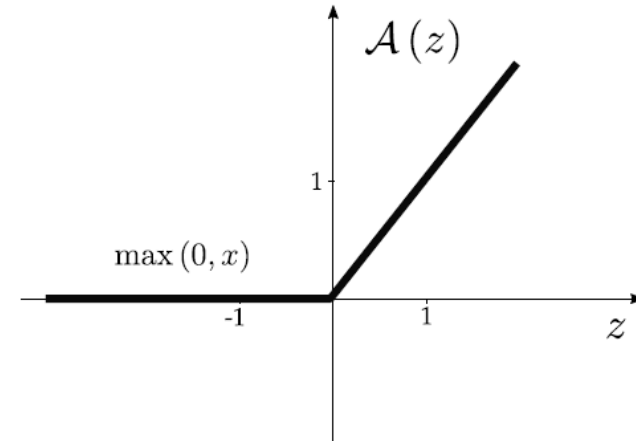
Logistic sigmoid



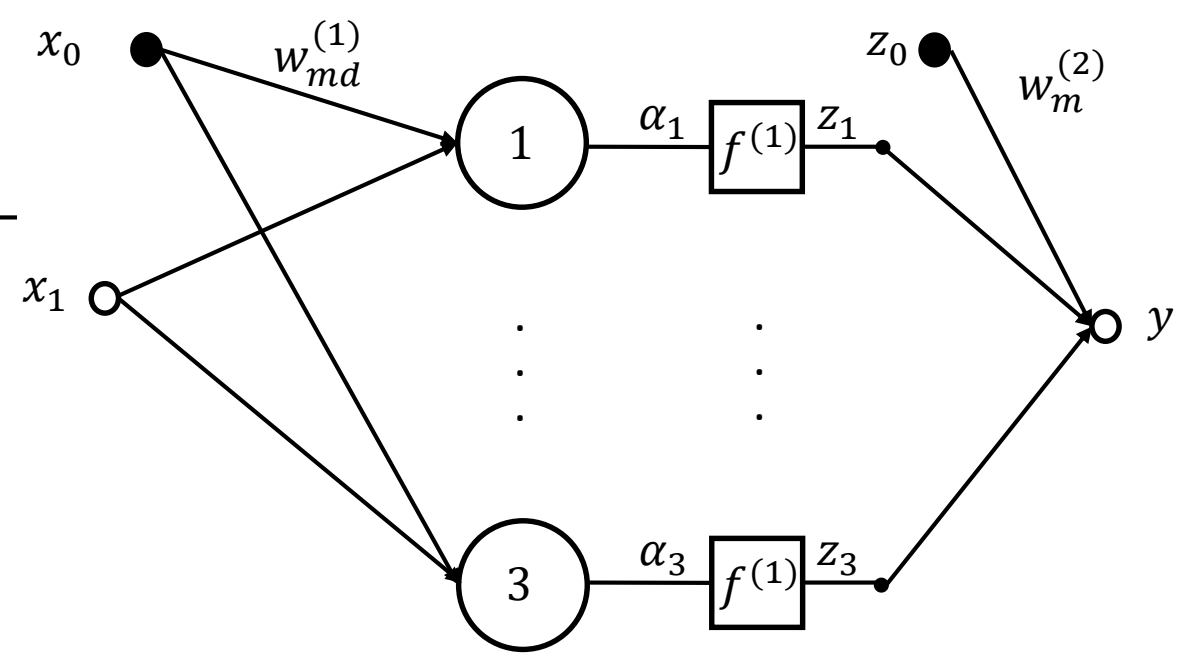
Tanh



Rectified Linear Unit (RELU)

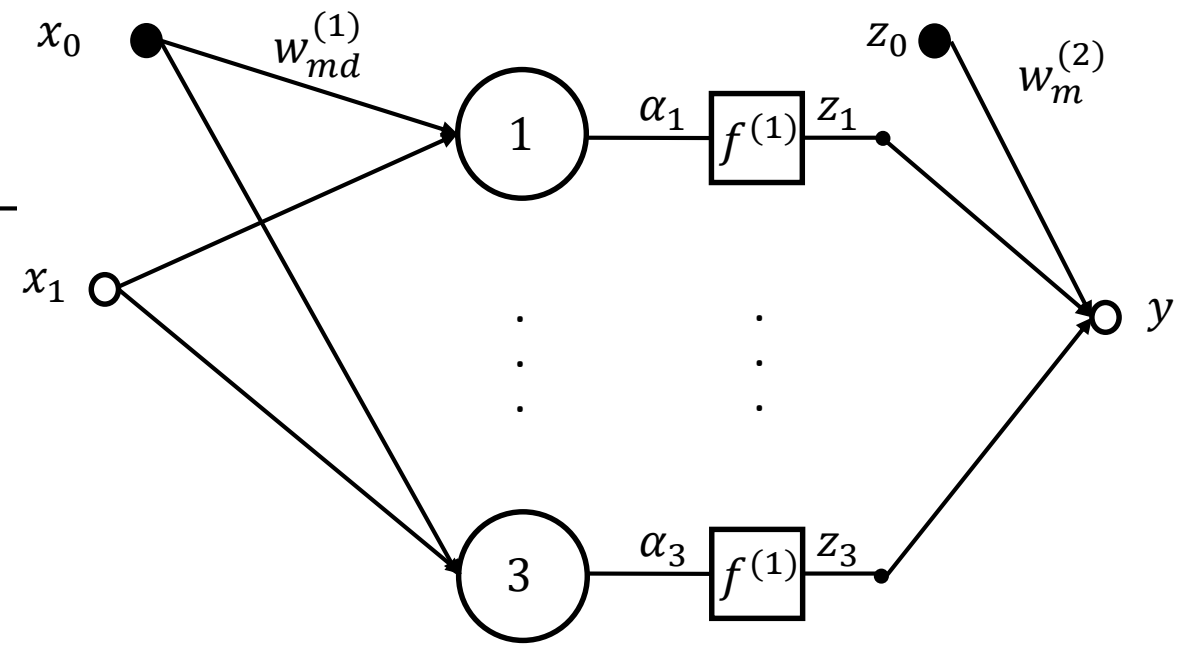
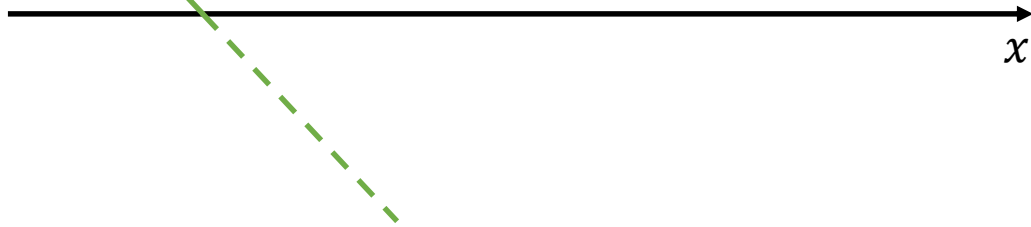


NN with ReLU

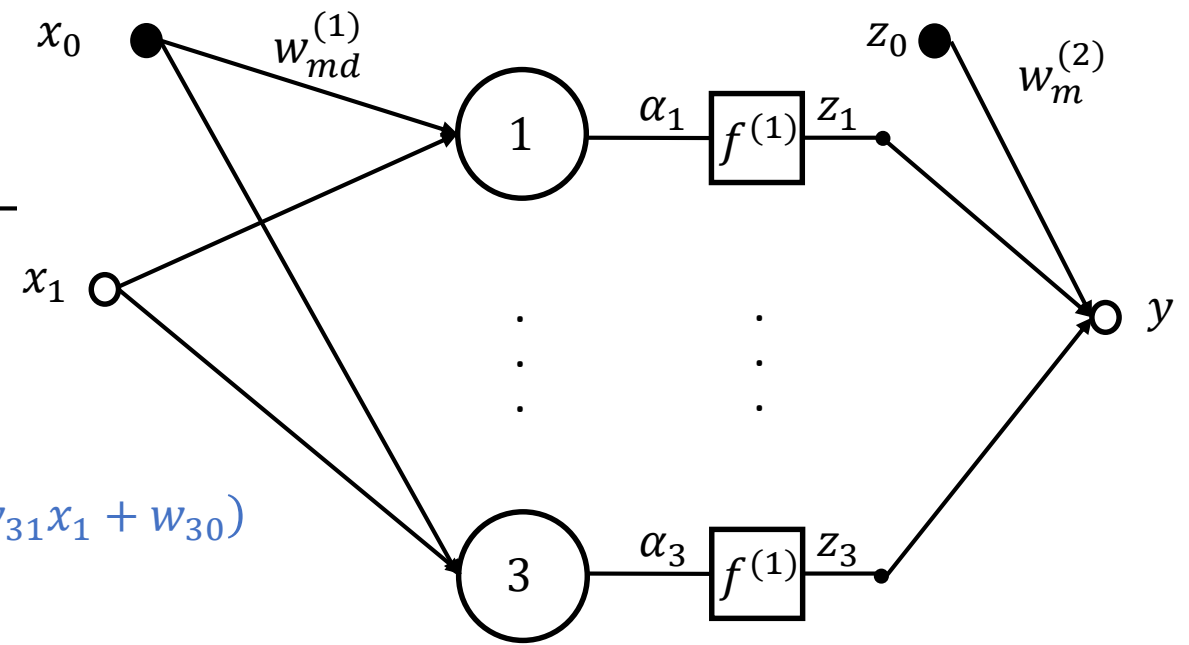
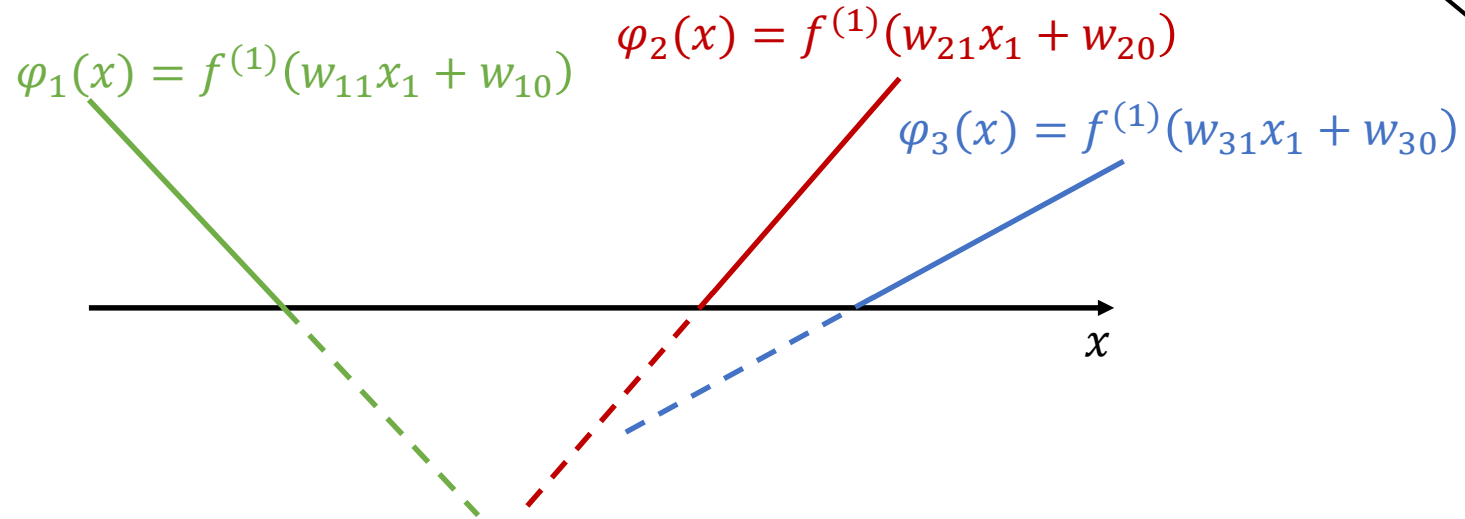


NN with ReLU

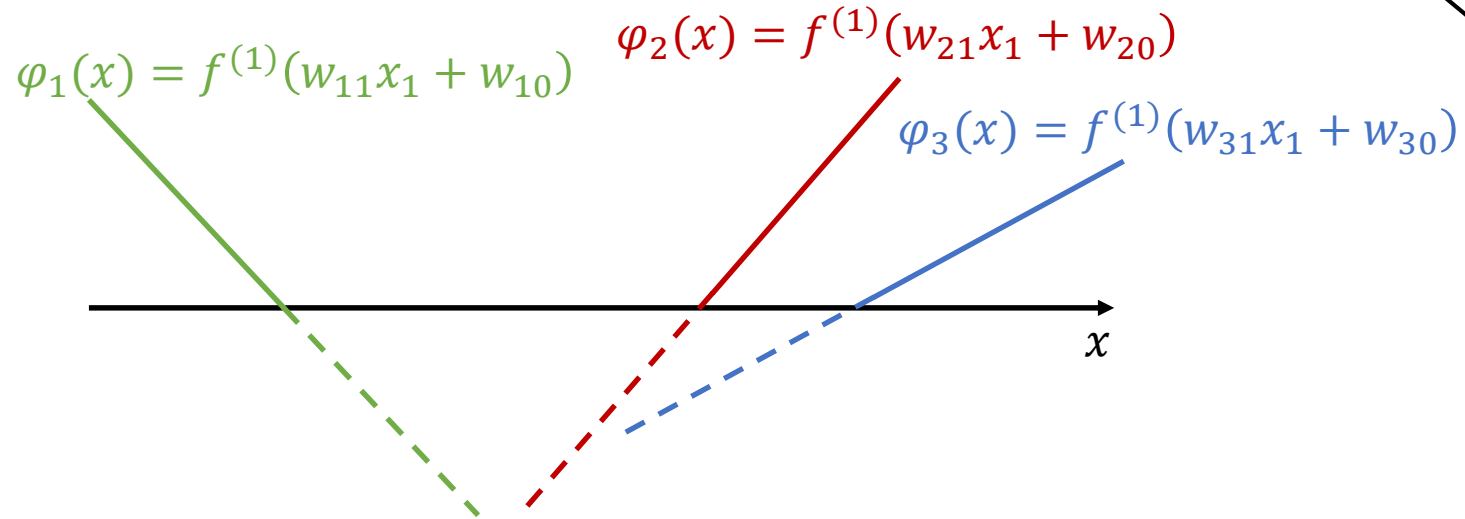
$$z_1 = \varphi_1(x) = f^{(1)}(w_{11}x_1 + w_{10})$$



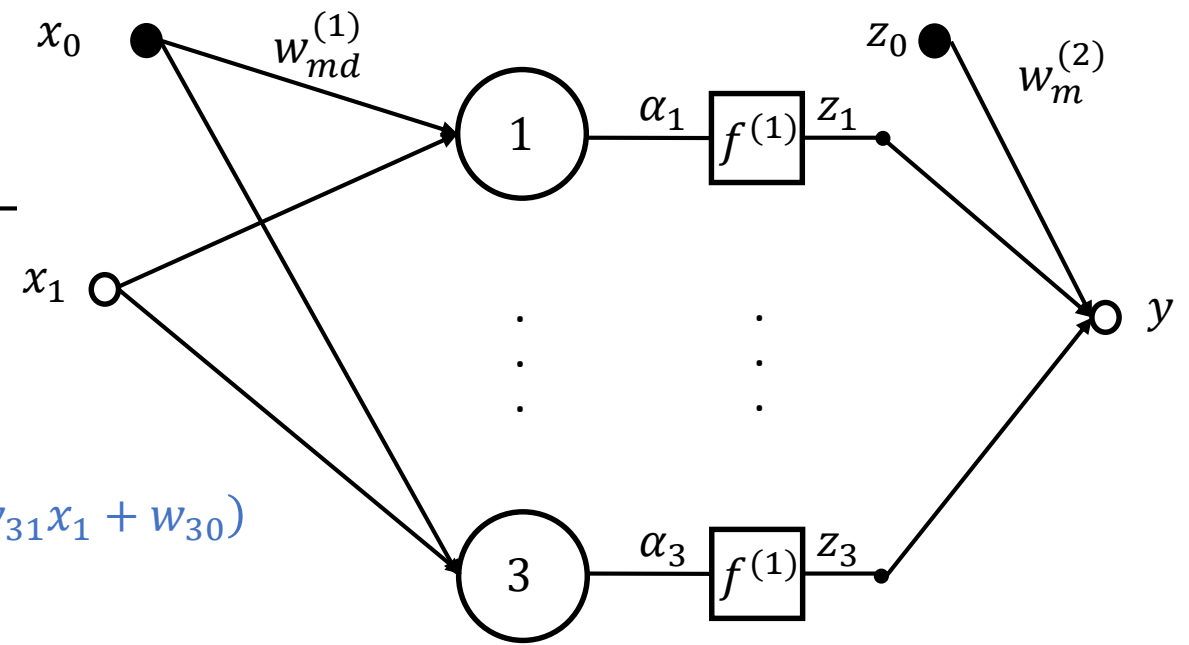
NN with ReLU



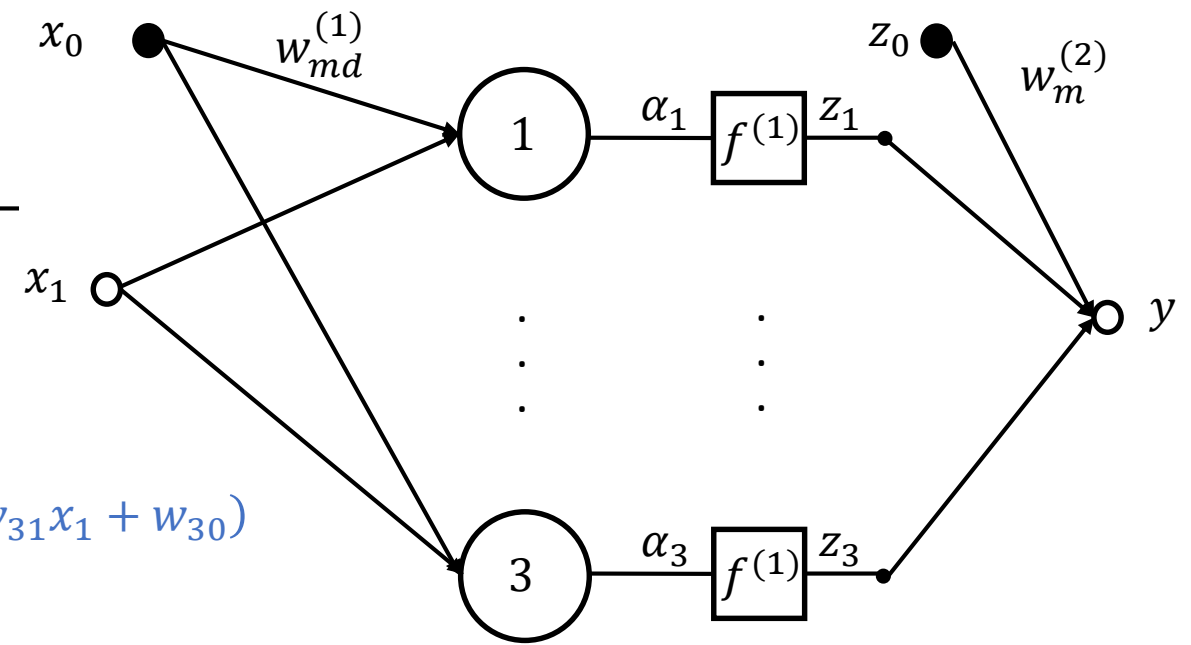
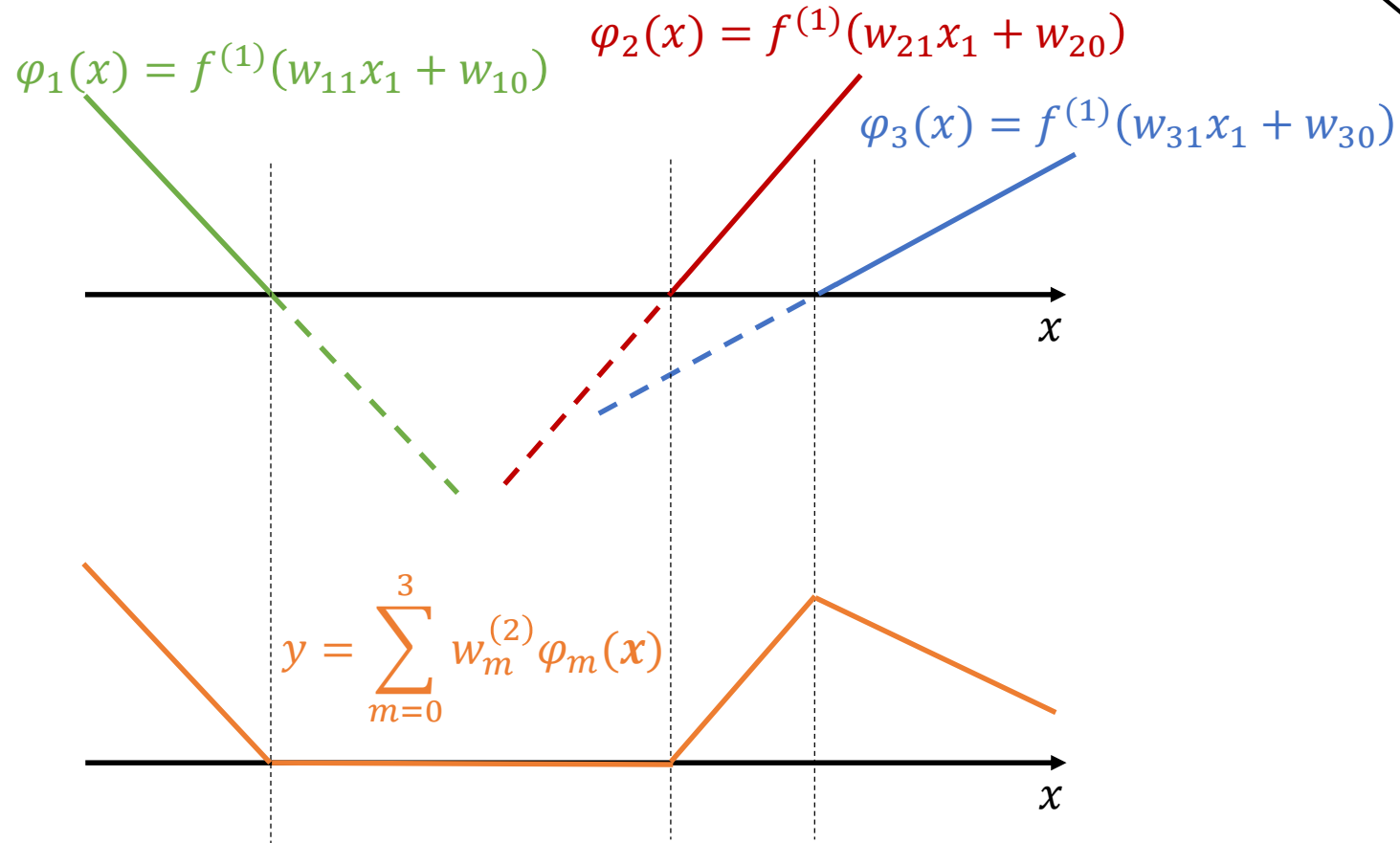
NN with ReLU



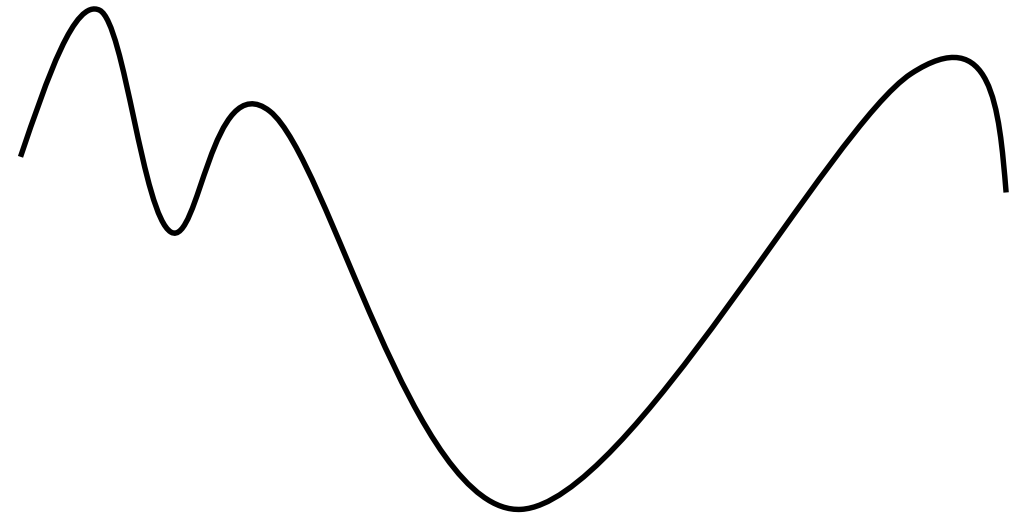
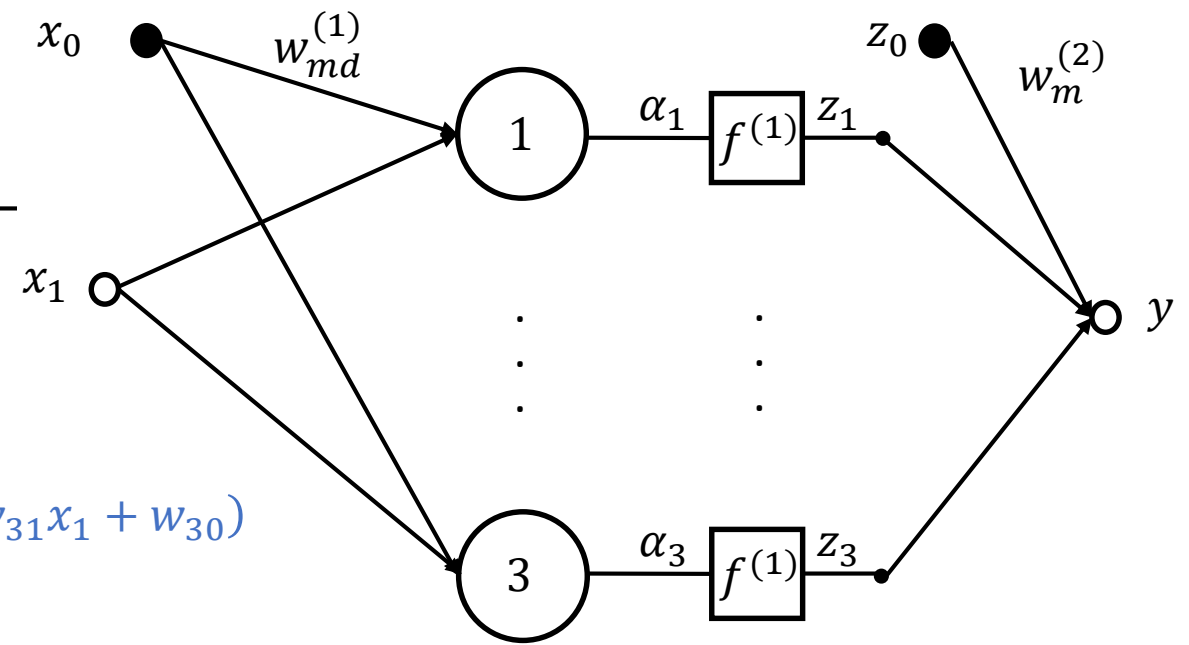
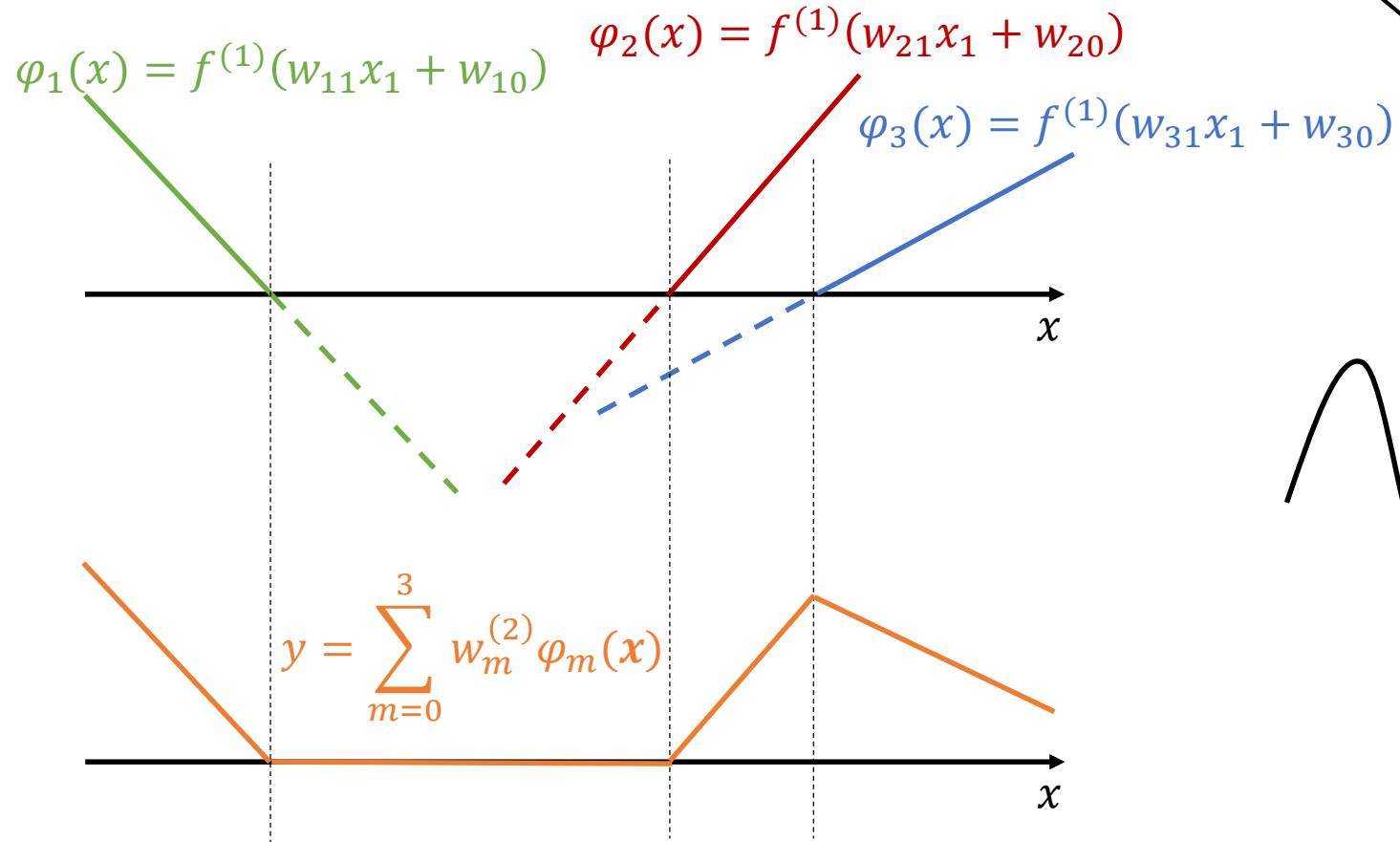
$$y = \sum_{m=0}^3 w_m^{(2)} \varphi_m(x)$$



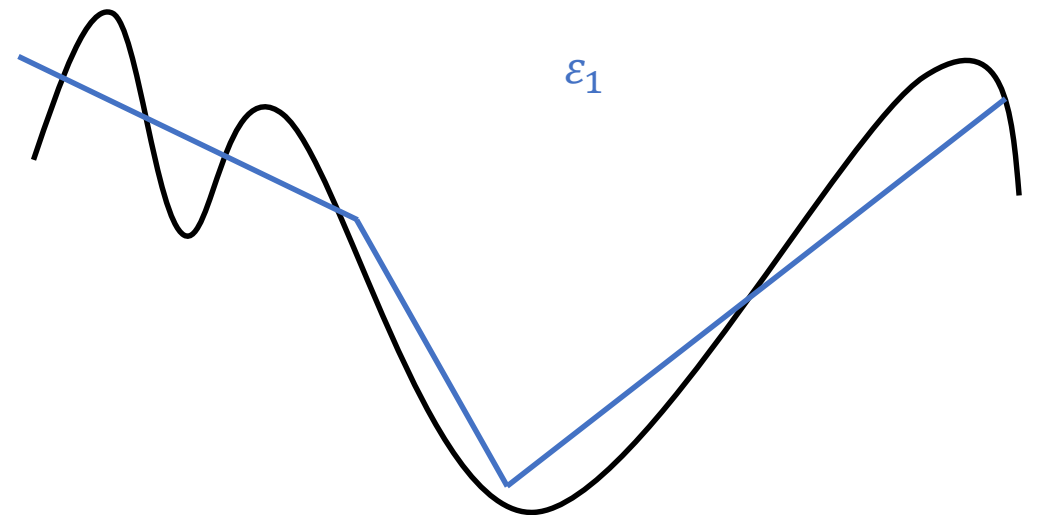
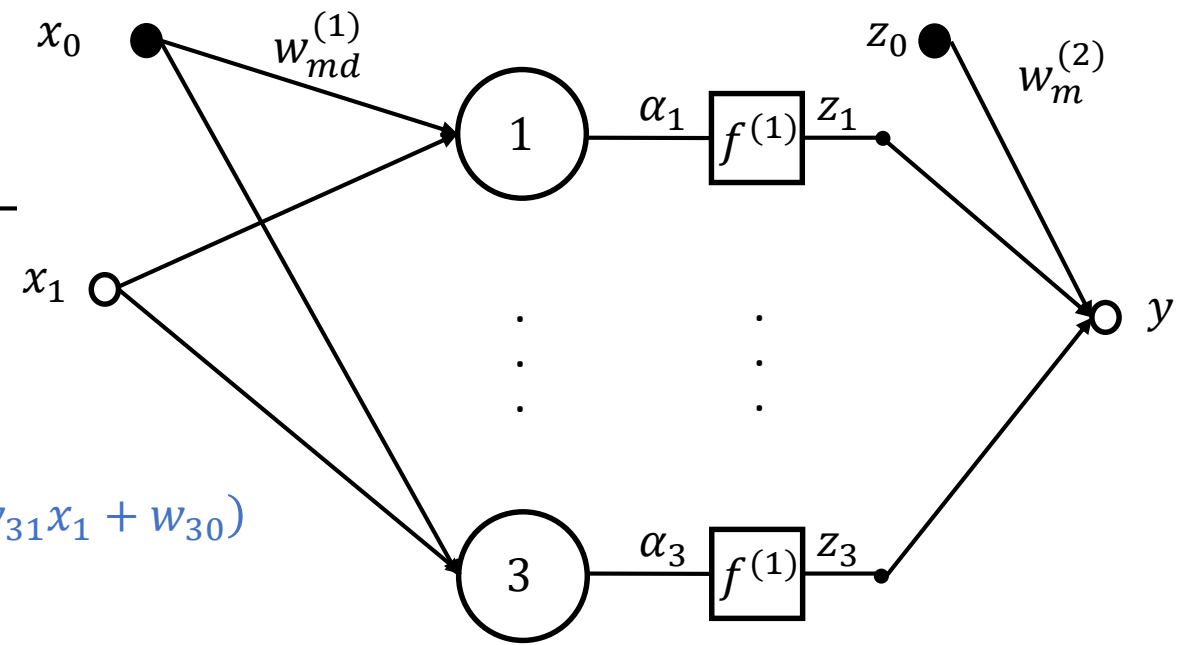
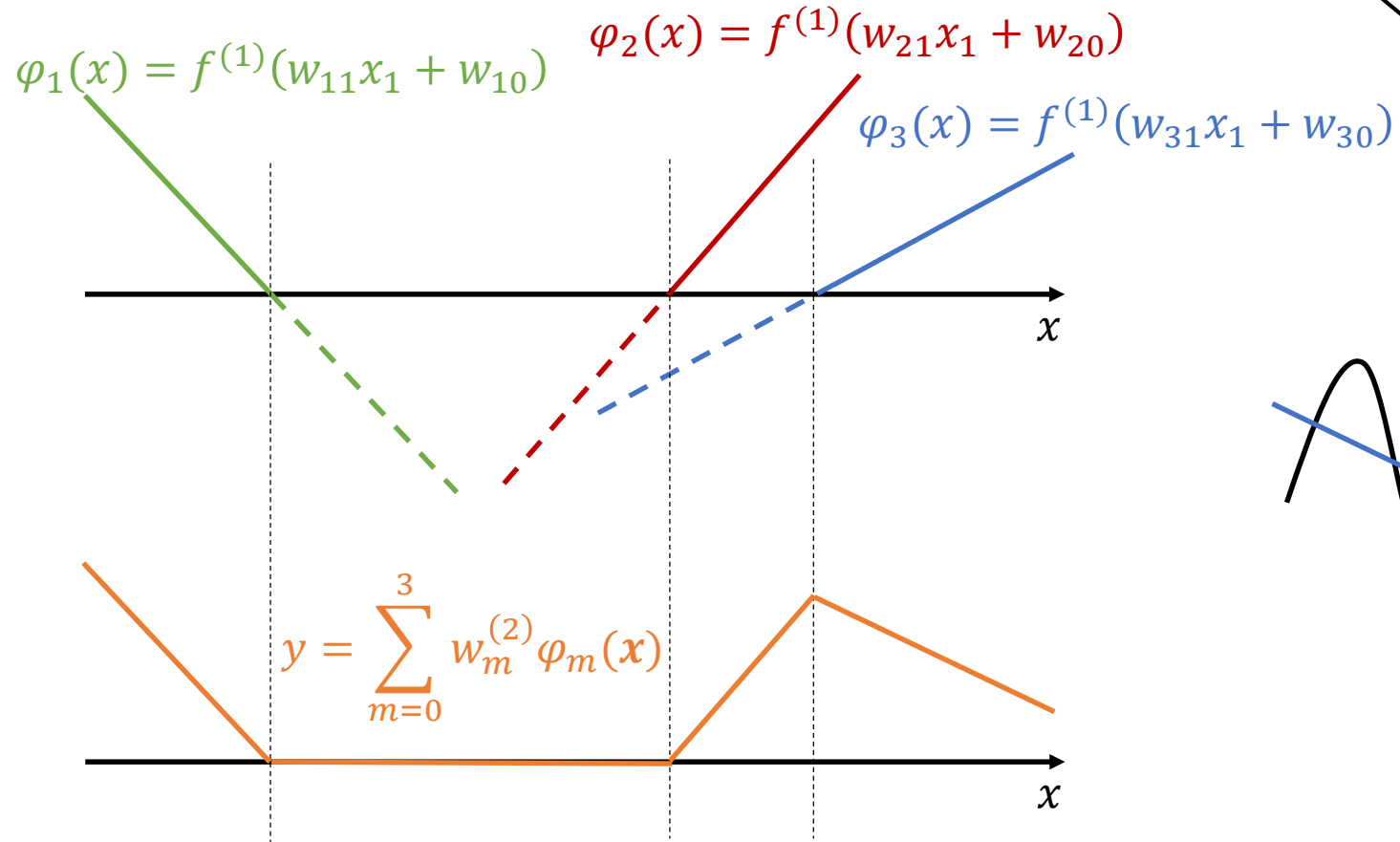
NN with ReLU



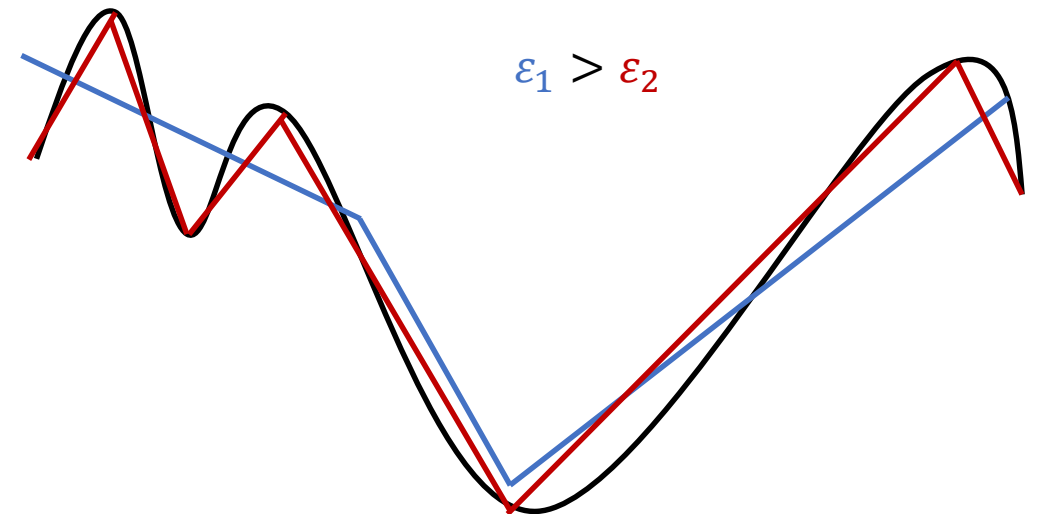
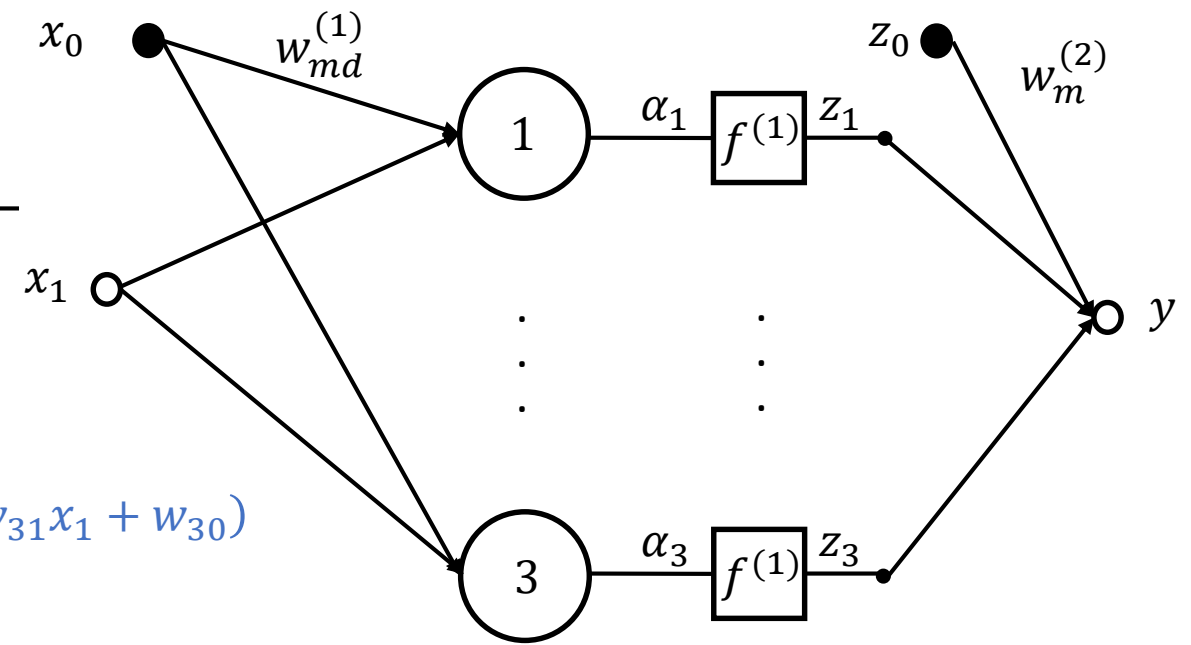
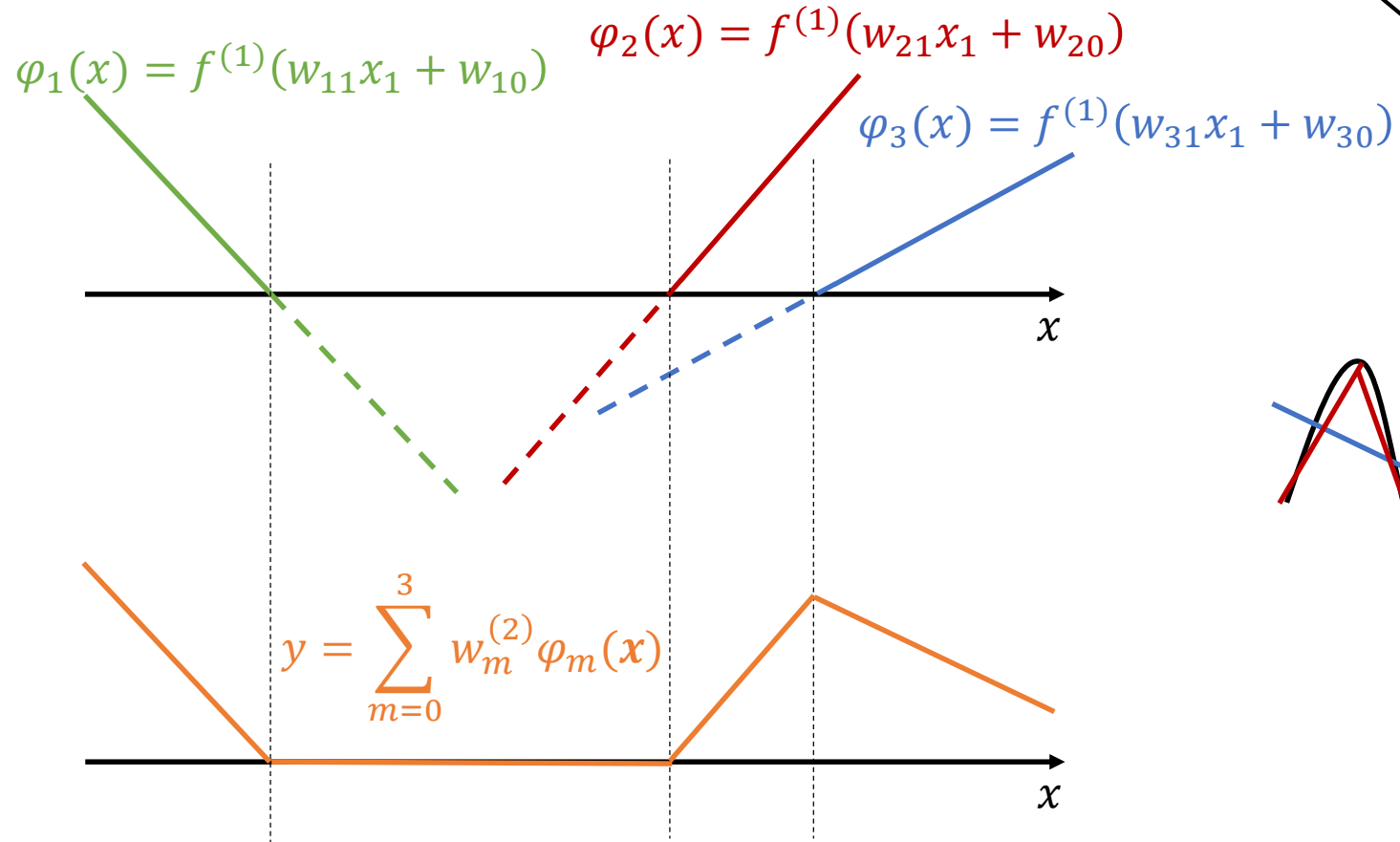
NN with ReLU



NN with ReLU



NN with ReLU



Expressive power of ReLU Networks

- The expressive power of ReLU-DNN networks is represented by the **number of linear regions** (see the example in the previous slide).
- In particular:

$$\#linear\ regions = width^{depth \times D}$$

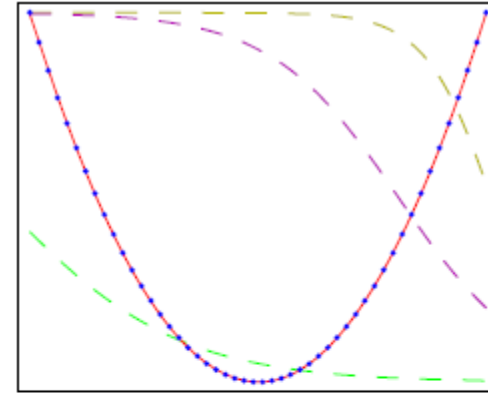
- Thus the **expressive power** is polynomial in width but **exponential in depth**.
- With fixed network capacity:

$$\#parameters = width^2 \times depth$$

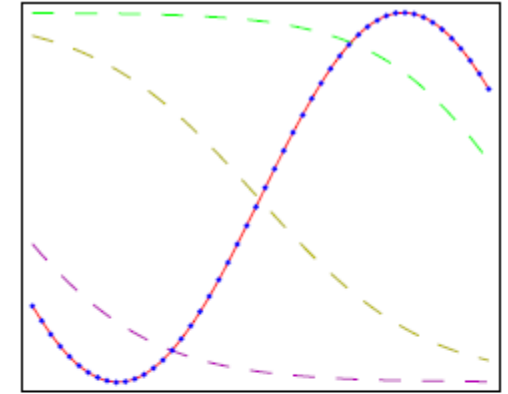
- Most expressive power is gained by **going deeper instead of making wider** networks (per layer) and staying shallow!

Examples: Function Approximations

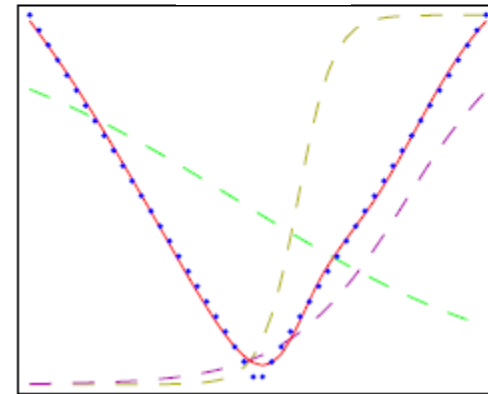
- $N = 50$
- $\#layers = 2$
- $M = 3$ (hidden units)
- Activation function: $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- One linear output unit
- Dashed curves: Hidden units output



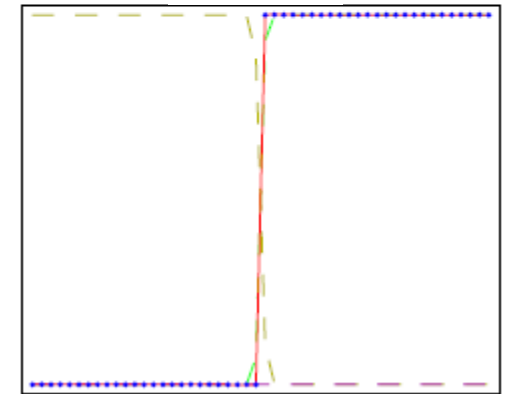
x^2



$\sin(x)$



$|x|$

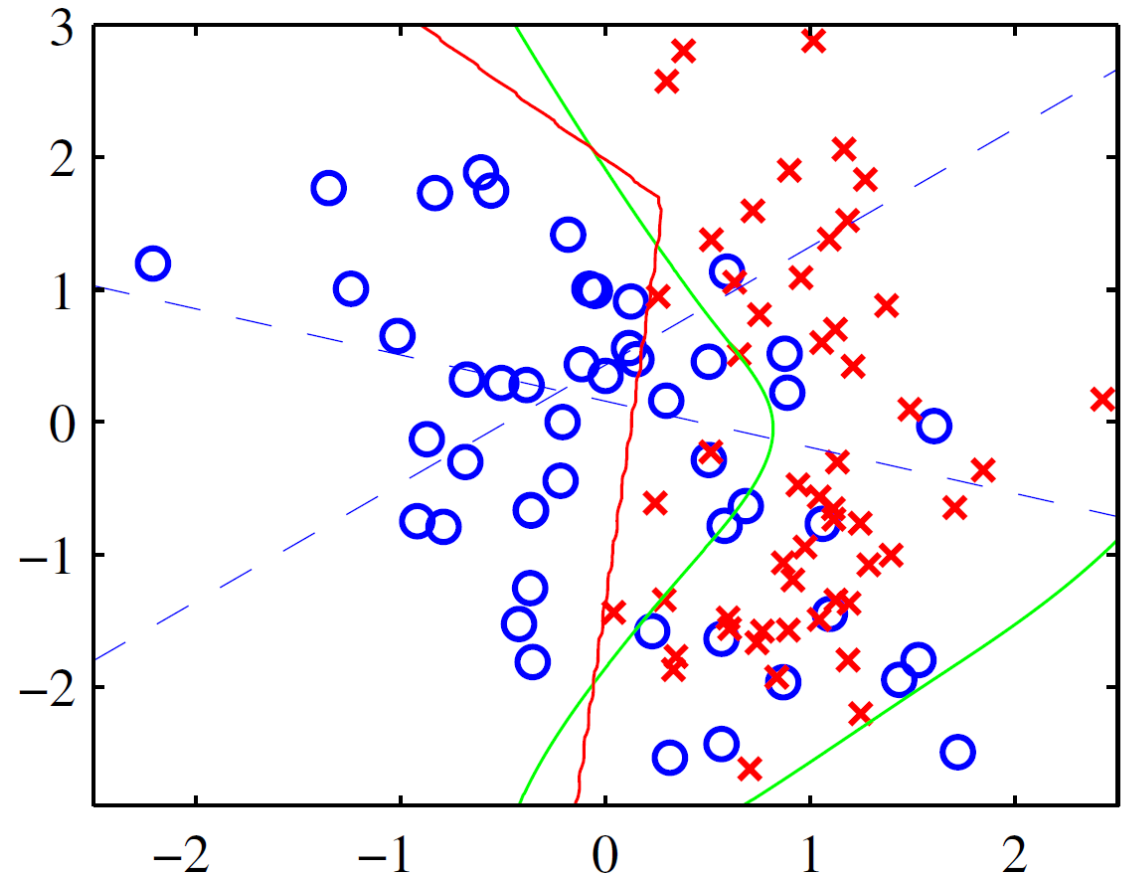


$H(x)$

Bishop

Example: Classification

- $\#layers = 2$
- $M = 2$ (hidden units)
- Activation function: $\tanh(x)$
- One output unit with sigmoid activation
- Red line: 2-LP decision boundary
- Green line: optimal decision boundary



Bishop

Thus far...

- We have defined:
 - Number of inputs
 - Hidden layers
 - Output layers
 - Activation functions (of hidden layers)
- We will use now a **probabilistic interpretation** of the network outputs to choose:
 - Number of outputs
 - Output activation function
 - Loss function!

Binary classification

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t \in [0,1]$.

Binary classification

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t \in [0,1]$.
- For binary classification we define one output with target distribution:

$$p(t = 1|\mathbf{x}) = y(\mathbf{x}, \mathbf{w}) \text{ (output of the NN)}$$

Binary classification

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t \in [0,1]$.
- For binary classification we define one output with target distribution:

$$p(t = 1|\mathbf{x}) = y(\mathbf{x}, \mathbf{w}) \text{ (output of the NN)}$$

- In general:

$$p(t|\mathbf{x}, \mathbf{w}) = y(\mathbf{x}, \mathbf{w})^t (1 - y(\mathbf{x}, \mathbf{w}))^{1-t} \text{ (Bernoulli distribution)}$$

Binary classification

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t \in [0,1]$.
- For binary classification we define one output with target distribution:

$$p(t = 1|\mathbf{x}) = y(\mathbf{x}, \mathbf{w}) \text{ (output of the NN)}$$

- In general:

$$p(t|\mathbf{x}, \mathbf{w}) = y(\mathbf{x}, \mathbf{w})^t (1 - y(\mathbf{x}, \mathbf{w}))^{1-t} \text{ (Bernoulli distribution)}$$

- Thus, we have a **single output unit**: $y(\mathbf{x}, \mathbf{w}) = f^L(a^{out}) = \sigma(a^{out}) \in [0,1]$.
- Maximum likelihood or Minimum Negative log likelihood:

Binary classification

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t \in [0,1]$.
- For binary classification we define one output with target distribution:

$$p(t = 1|\mathbf{x}) = y(\mathbf{x}, \mathbf{w}) \text{ (output of the NN)}$$

- In general:

$$p(t|\mathbf{x}, \mathbf{w}) = y(\mathbf{x}, \mathbf{w})^t (1 - y(\mathbf{x}, \mathbf{w}))^{1-t} \text{ (Bernoulli distribution)}$$

- Thus, we have a **single output unit**: $y(\mathbf{x}, \mathbf{w}) = f^L(a^{out}) = \sigma(a^{out}) \in [0,1]$.
- Maximum likelihood or Minimum Negative log likelihood:

$$E(\mathbf{w}) = -\sum_{n=1}^N t_n \ln(y(\mathbf{x}_n, \mathbf{w})) + (1 - t_n) \ln(1 - y(\mathbf{x}_n, \mathbf{w})) \text{ (Cross Entropy Loss)}$$

Classification with c classes

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_N)$, where $\mathbf{x}_n \in \mathbb{R}^D$, and $\mathbf{t}_n = [t_{n1}, \dots, t_{nc}]$ (one hot encoding).

Classification with c classes

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_N)$, where $\mathbf{x}_n \in \mathbb{R}^D$, and $\mathbf{t}_n = [t_{n1}, \dots, t_{nc}]$ (one hot encoding).
- For multi-class classification we define c output units and we assume target distribution:

$$p(\mathbf{t}_n | \mathbf{x}_n, \mathbf{w}) = \prod_{i=1}^c y_i(\mathbf{x}_n, \mathbf{w})^{t_{ni}} \text{ (Generalized Bernoulli)}$$

where

$$y_i(\mathbf{x}_n, \mathbf{w}) = p(\omega_i | \mathbf{x}_n, \mathbf{w}) \text{ (output of the NN)}$$

Classification with c classes

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_N)$, where $\mathbf{x}_n \in \mathbb{R}^D$, and $\mathbf{t}_n = [t_{n1}, \dots, t_{nc}]$ (one hot encoding).
- For multi-class classification we define c output units and we assume target distribution:

$$p(\mathbf{t}_n | \mathbf{x}_n, \mathbf{w}) = \prod_{i=1}^c y_i(\mathbf{x}_n, \mathbf{w})^{t_{ni}} \text{ (Generalized Bernoulli)}$$

where

$$y_i(\mathbf{x}_n, \mathbf{w}) = p(\omega_i | \mathbf{x}_n, \mathbf{w}) \text{ (output of the NN)}$$

- For each output unit: $y_i(\mathbf{x}_n, \mathbf{w}) = f^L(a_i^{out}) = \frac{\exp(a_i^{out})}{\sum_{j=1}^c \exp(a_j^{out})}$ (softmax activation function) with $\sum p(\omega_i | \mathbf{x}_n, \mathbf{w}) = \sum y_i(\mathbf{x}_n, \mathbf{w}) = 1$

Classification with c classes

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_N)$, where $\mathbf{x}_n \in \mathbb{R}^D$, and $\mathbf{t}_n = [t_{n1}, \dots, t_{nc}]$ (one hot encoding).
- For multi-class classification we define c output units and we assume target distribution:

$$p(\mathbf{t}_n | \mathbf{x}_n, \mathbf{w}) = \prod_{i=1}^c y_i(\mathbf{x}_n, \mathbf{w})^{t_{ni}} \text{ (Generalized Bernoulli)}$$

where

$$y_i(\mathbf{x}_n, \mathbf{w}) = p(\omega_i | \mathbf{x}_n, \mathbf{w}) \text{ (output of the NN)}$$

- For each output unit: $y_i(\mathbf{x}_n, \mathbf{w}) = f^L(a_i^{out}) = \frac{\exp(a_i^{out})}{\sum_{j=1}^c \exp(a_j^{out})}$ (softmax activation function) with $\sum p(\omega_i | \mathbf{x}_n, \mathbf{w}) = \sum y_i(\mathbf{x}_n, \mathbf{w}) = 1$
- Maximum likelihood or Minimum Negative log likelihood:

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{i=1}^c t_{ni} \ln(y_i(\mathbf{x}_n, \mathbf{w})) \text{ (Cross Entropy Loss)}$$

Regression

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t_n \in \mathbb{R}$.

Regression

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t_n \in \mathbb{R}$.
- We assume target distribution: $p(t|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}_n, \mathbf{w}), \beta^{-1})$. Here the mean of the gaussian is \mathbf{x} –dependent. Actually, it is the output of the NN.

Regression

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t_n \in \mathbb{R}$.
- We assume target distribution: $p(t|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}_n, \mathbf{w}), \beta^{-1})$. Here the mean of the gaussian is \mathbf{x} –dependent. Actually, it is the output of the NN.
- We adopt a single output NN. In particular, $y(\mathbf{x}_n, \mathbf{w}) = f^L(a^{out})$.
- Because the outputs are real valued we assume an **identity output activation** function: $y(\mathbf{x}_n, \mathbf{w}) = f^L(a^{out}) = a^{out}$.

Regression

- Data: inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, targets $\mathbf{t} = (t_1, \dots, t_N)$, where $\mathbf{x}_n \in \mathbb{R}^D, t_n \in \mathbb{R}$.
- We assume target distribution: $p(t|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}_n, \mathbf{w}), \beta^{-1})$. Here the mean of the gaussian is \mathbf{x} –dependent. Actually, it is the output of the NN.
- We adopt a single output NN. In particular, $y(\mathbf{x}_n, \mathbf{w}) = f^L(a^{out})$.
- Because the outputs are real valued we assume an identity output activation function: $y(\mathbf{x}_n, \mathbf{w}) = f^L(a^{out}) = a^{out}$.
- Maximum likelihood or Minimum Negative log likelihood:

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi$$

- In essence: $E(\mathbf{w}) = \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2$

Losses overview

- Binary classification:
 - **Assumed** Bernoulli target distribution
 - NN makes prediction of the probability of class ω_1
 - Output activation is logistic sigmoid
 - Minimize cross entropy loss
- Multi-class classification
 - **Assumed** generalized Bernoulli target distribution
 - NN makes prediction of the probability of each class
 - Output activation is softmax function
 - Minimize (multi-class) cross entropy loss.
- Regression
 - **Assumed** gaussian target distribution
 - NN makes prediction for the mean
 - Output activation function is the identity function
 - Minimize least squares.

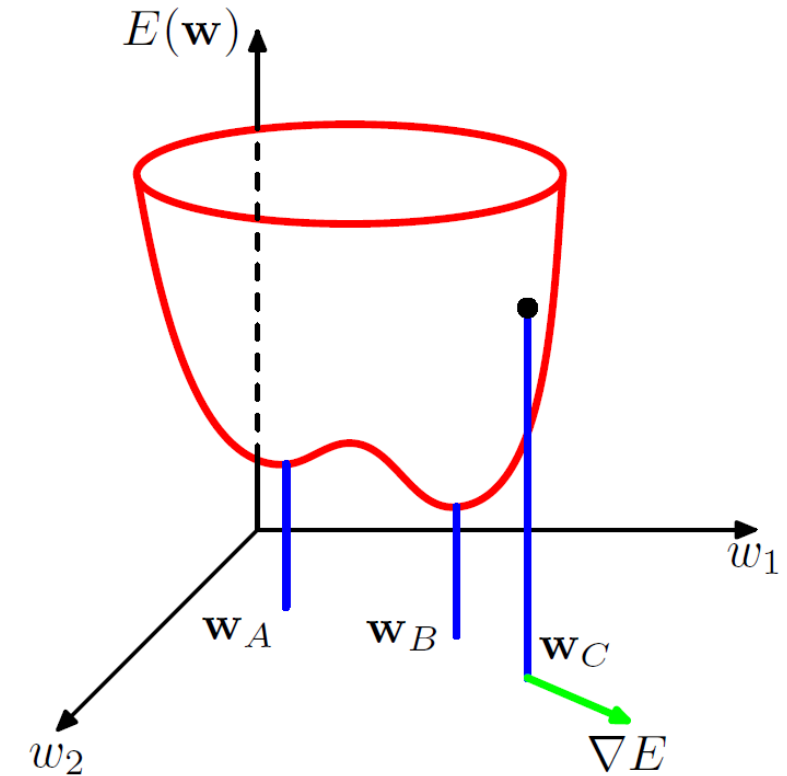
Training

- Thus far:
 - Structure
 - different components of the NNs
 - activation functions
 - loss functions
- Next:

Training: How can I adjust \mathbf{w} in order to approximate the function I am interested in.

Optimization

- For each task I have a different loss function $E(\mathbf{w})$ that I want to minimize.
- Parameter optimization: $\mathbf{w}_{opt} = \operatorname{argmin}_{\mathbf{w}} E(\mathbf{w})$
- A big problem: $E(\mathbf{w})$ is not convex in \mathbf{w} , thus several local optimal exist!
- How to reach a global minimum?

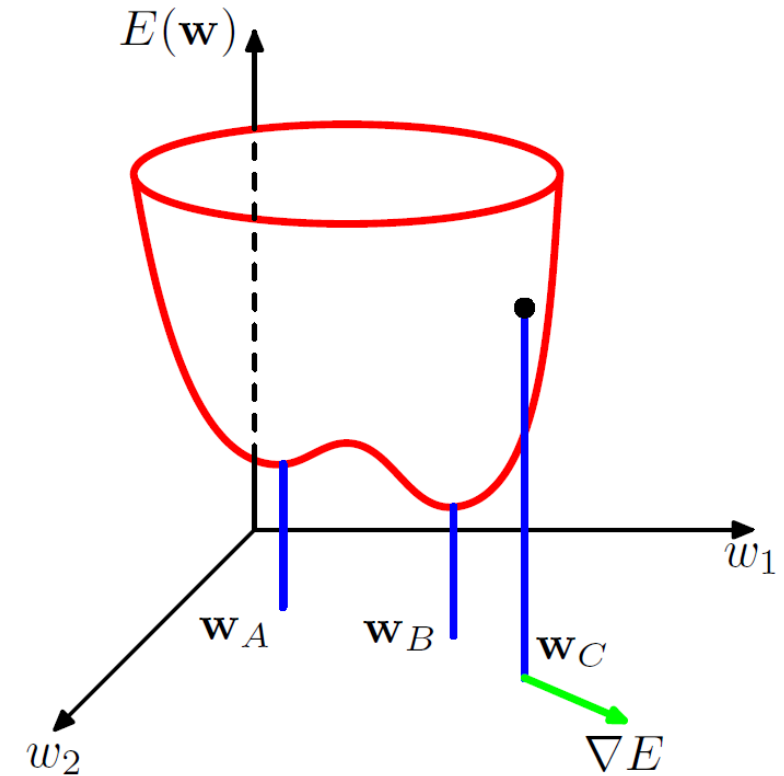


Bishop

Optimization

- For each task I have a different loss function $E(\mathbf{w})$ that I want to minimize.
- Parameter optimization: $\mathbf{w}_{opt} = \operatorname{argmin}_{\mathbf{w}} E(\mathbf{w})$
- A big problem: $E(\mathbf{w})$ is not convex in \mathbf{w} , thus several local optima exist!
- How to reach a global minimum?

Stochastic Gradient Descent



Bishop

Stochastic gradient descent

- Gradient descent: $\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla E(\mathbf{w}^i)$

where $E(\mathbf{w}^i) = \sum_{n=1}^N E_n(\mathbf{w}^i)$

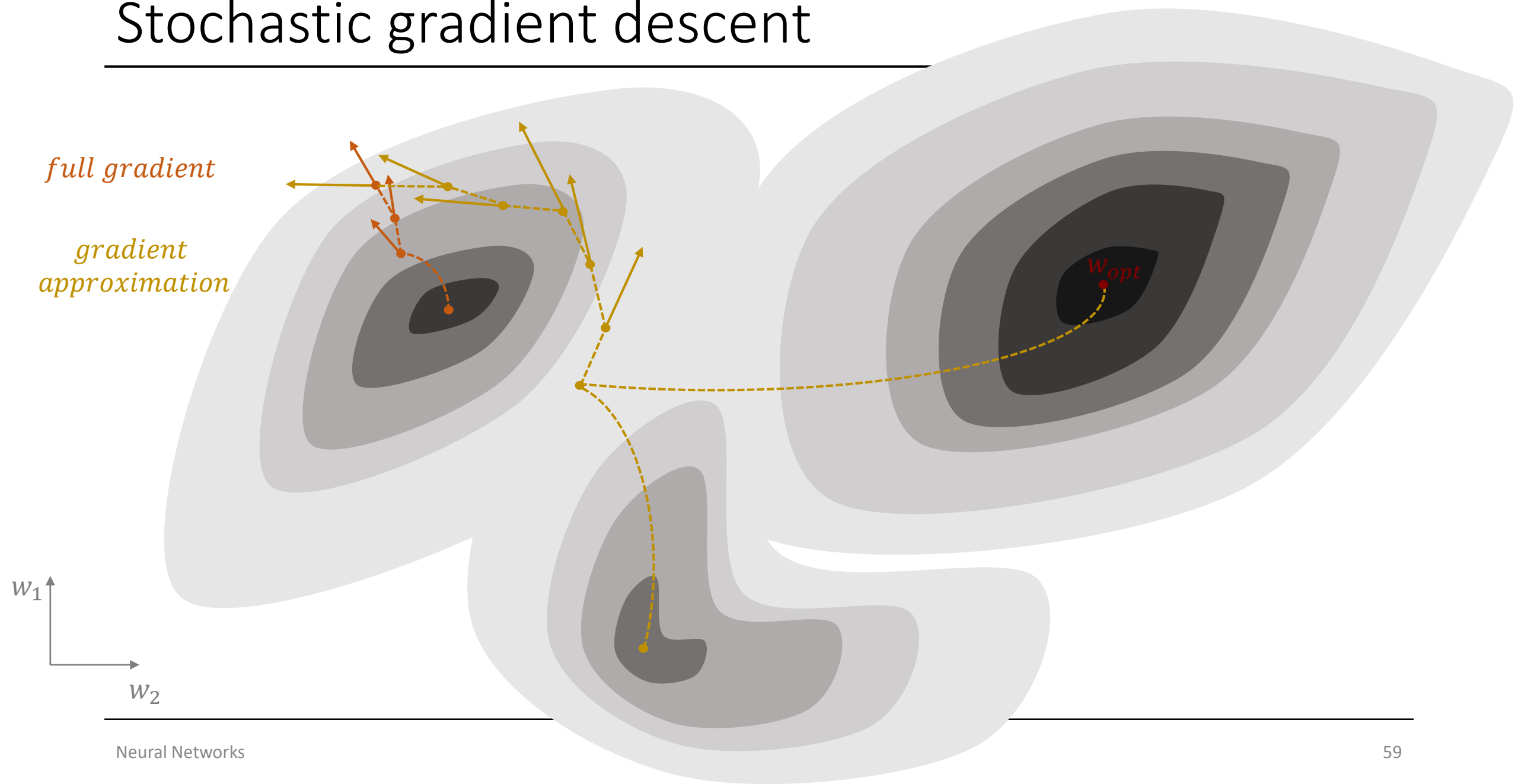
- By incorporating Stochasticity:

$$\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla E_n(\mathbf{w}^i)$$

or

$$\mathbf{w}^{i+1} = \mathbf{w}^i - \eta \nabla \sum_{l=1}^L E_l(\mathbf{w}^i) \text{ (minibatch)}$$

Stochastic gradient descent



Comments on SGD

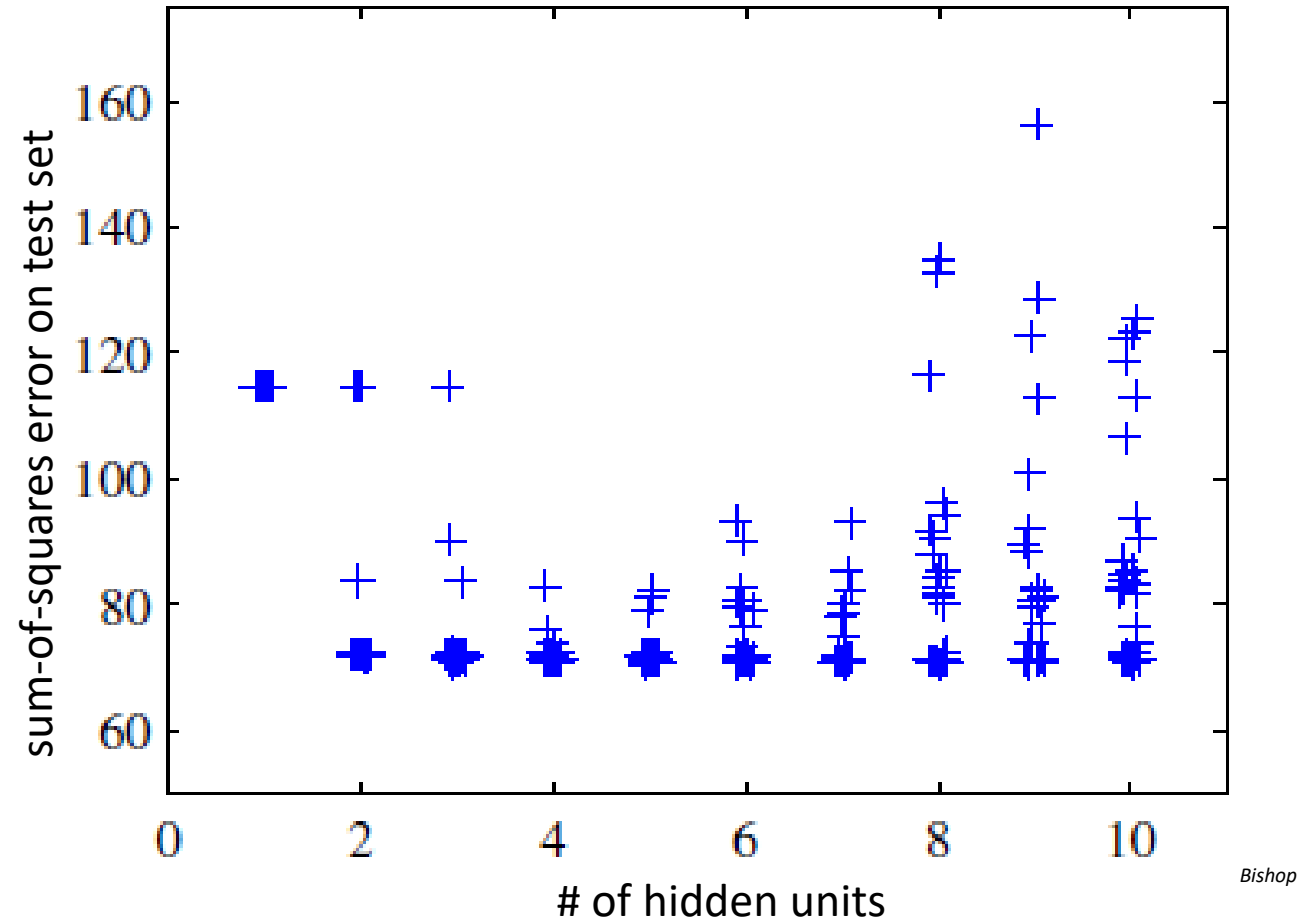
- Learning rate η too small: slow convergence
- Learning rate η too large: oscillations (no convergence)
- We most often use rate scheduling with smaller and smaller rates over time.
- SGD is more likely to escape a local minimum as:

$$\nabla E(\mathbf{w}) = \mathbf{0} \text{ does not necessarily imply } \nabla E_n(\mathbf{w}) = \mathbf{0}!$$

- As the **number of hidden units increases** the landscape becomes more **complex** and thus **more initializations** are needed to prove the robustness of the network

Comments on SGD

- Try several initialization of \mathbf{w} to test uncertainty on performance





ARISTOTLE UNIVERSITY OF THESSALONIKI



FACULTY OF ENGINEERING

Questions?

Pattern Recognition & Machine Learning
Neural Networks