

GR5260

Programming for Quantitative & Computational Finance

Instructor: Ka Yi Ng

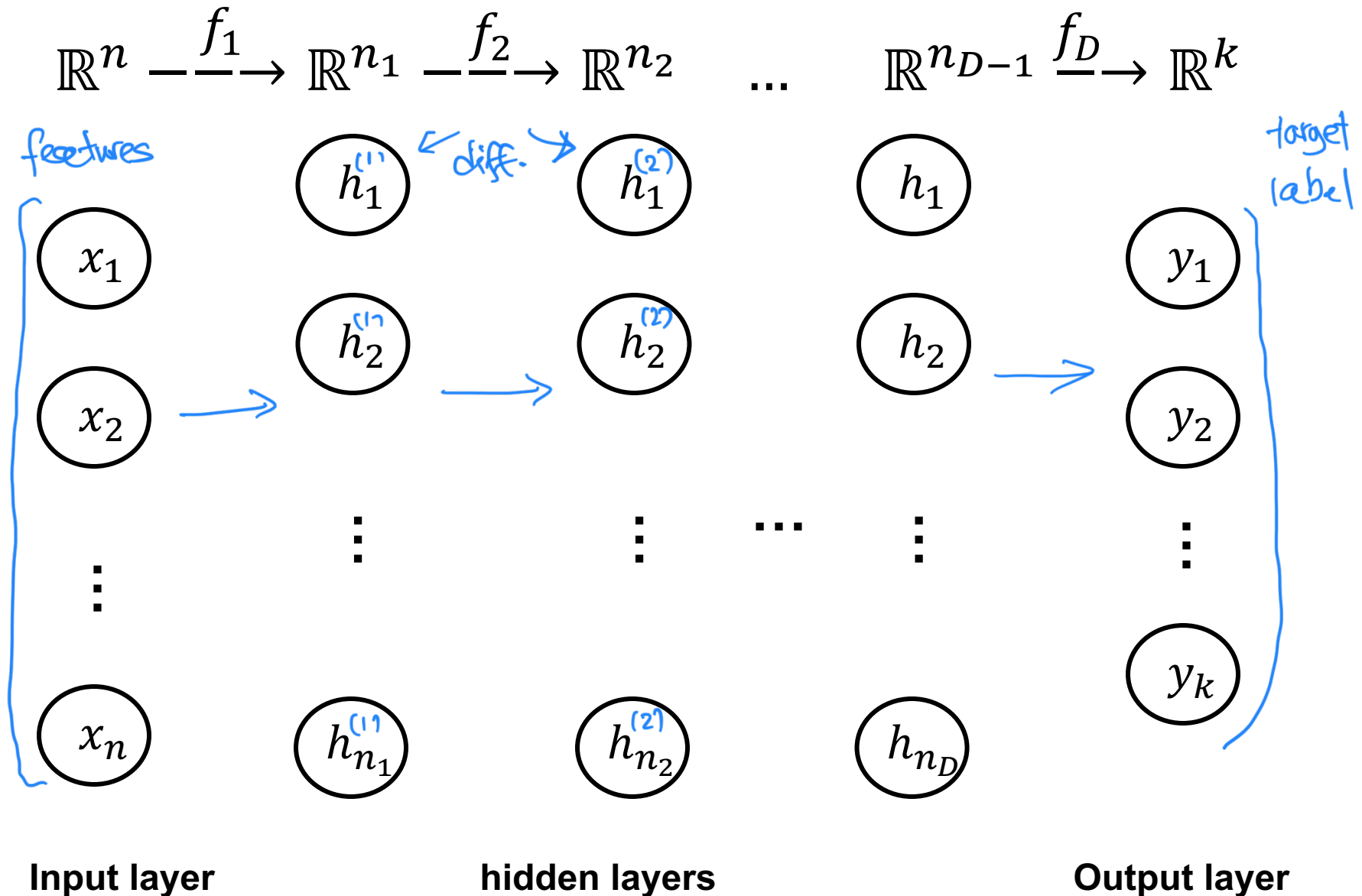
April 28, 2023

Neural networks

Feedforward Neural Networks

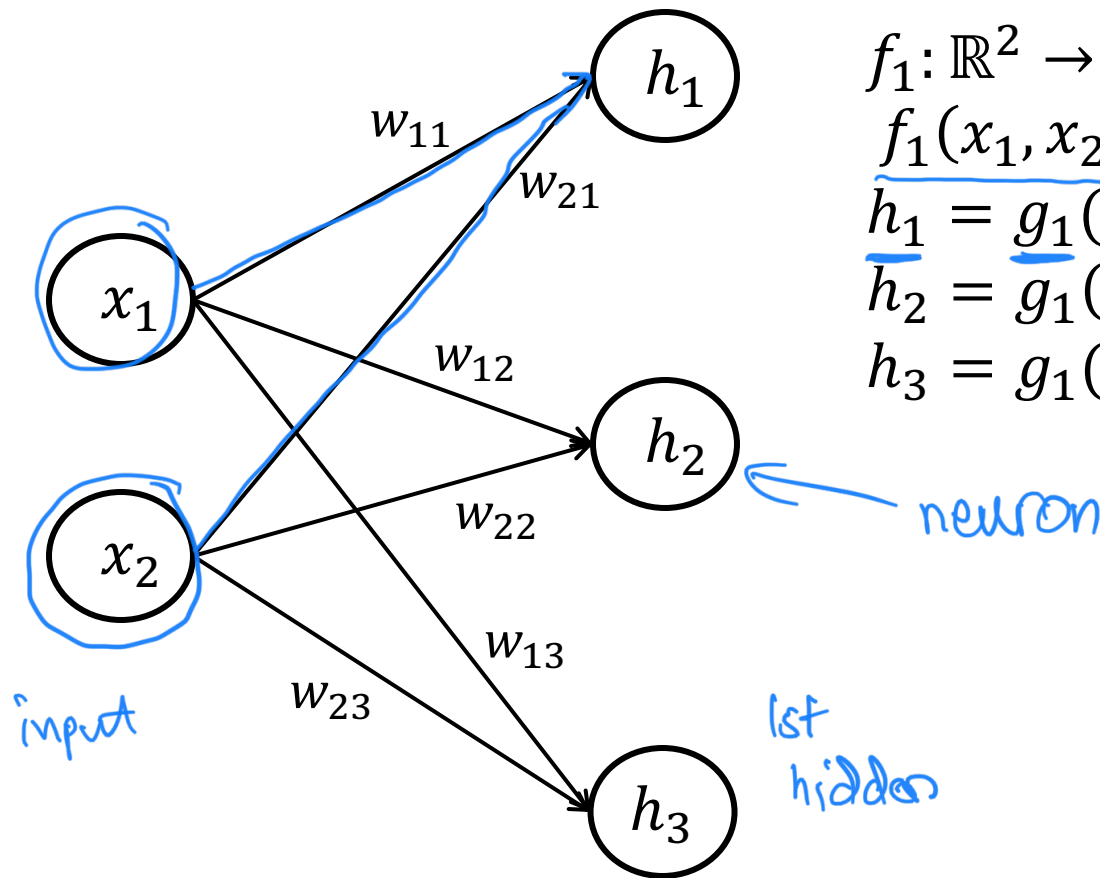
- Estimate the input-output mapping $f: \mathbb{R}^n \rightarrow \mathbb{R}^k$ using a parametrized function $\hat{f}: \mathbb{R}^n \rightarrow \mathbb{R}^k$ of the following form:
- $\hat{f} = f_D \circ f_{D-1} \circ \dots \circ f_2 \circ f_1$
where $f_1: \mathbb{R}^n \rightarrow \mathbb{R}^{n_1}, f_2: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}, \dots,$
 $f_{D-1}: \mathbb{R}^{n_{D-2}} \rightarrow \mathbb{R}^{n_{D-1}}, f_D: \mathbb{R}^{n_{D-1}} \rightarrow \mathbb{R}^k$
- $f_i(\mathbf{x}) = (g_i(W_{i1}\mathbf{x} + b_{i1}), g_i(W_{i2}\mathbf{x} + b_{i2}), \dots, g_i(W_{in_i}\mathbf{x} + b_{in_i}))$
- W_{ij} is a weight vector of same dimension as \mathbf{x}
- b_{ij} is a scalar, called bias
- $g_i: \mathbb{R} \rightarrow \mathbb{R}$ is an activation function (eg. $\max(0, z)$, $\tanh(z)$)

Feedforward Neural Networks



Feedforward Neural Networks

$$f_i(\mathbf{x}) = (g_i(W_{i1}\mathbf{x} + b_{i1}), g_i(W_{i2}\mathbf{x} + b_{i2}), \dots, g_i(W_{in_i}\mathbf{x} + b_{in_i}))$$



$$f_1: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$f_1(x_1, x_2) = (h_1, h_2, h_3)$$

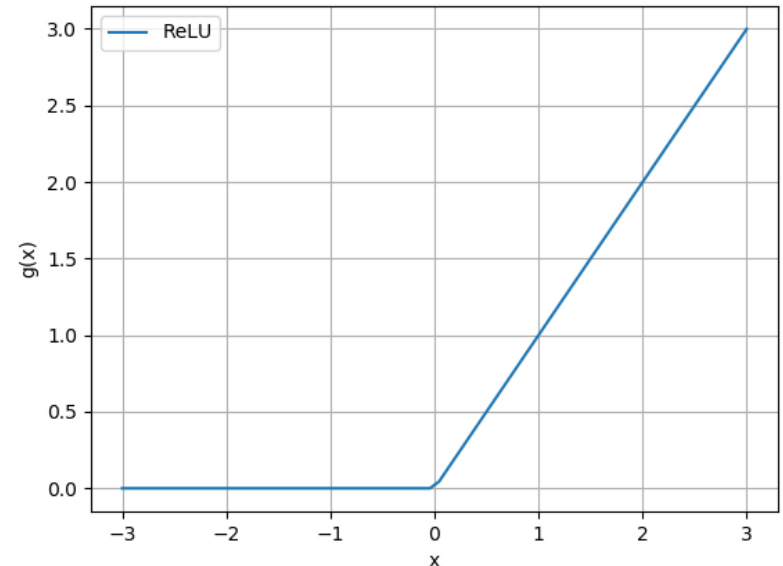
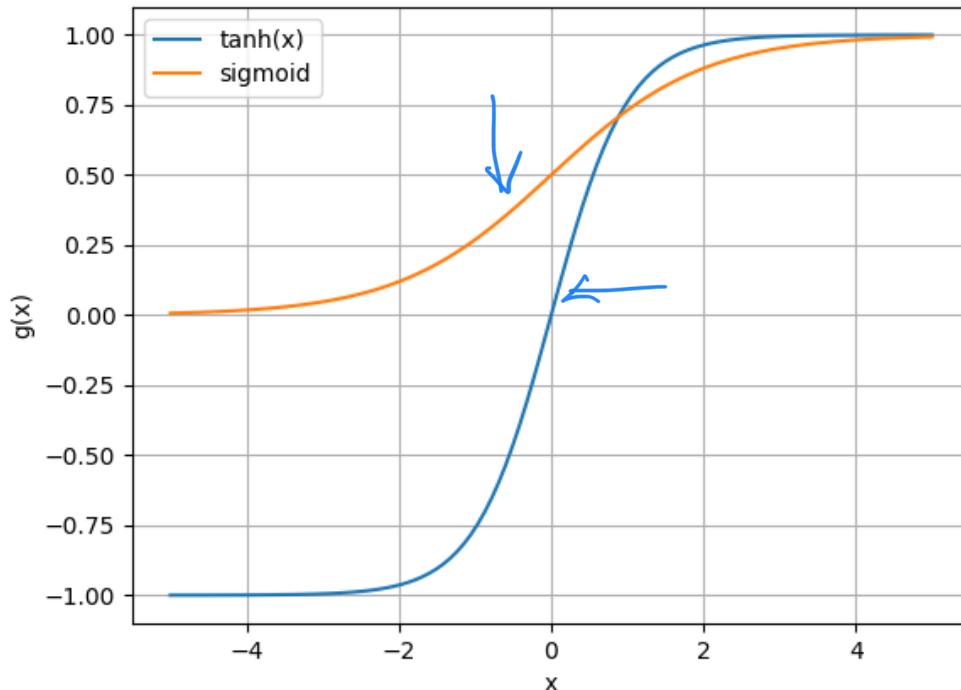
$$h_1 = g_1(w_{11}x_1 + w_{21}x_2 + b_{11})$$

$$h_2 = g_1(w_{12}x_1 + w_{22}x_2 + b_{12})$$

$$h_3 = g_1(w_{13}x_1 + w_{23}x_2 + b_{13})$$

Activation functions

- Rectified linear unit (ReLU): $g(z) = \max(0, z)$
- Hyperbolic tangent $g(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$
- Logistic sigmoid $g(z) = \frac{1}{1 + \exp(-z)}$



More examples in https://en.wikipedia.org/wiki/Activation_function

Feedforward Neural Networks

- Hyper-parameters:
 - L : depth of neural network
 - n_1, \dots, n_{L-1} : width of the hidden layers
 - g_1, \dots, g_L : activation functions
- Model parameters of $\hat{f}(\mathbf{x}; \{W_{ij}, b_{ij}\}_{i,j})$: W_{ij} 's, b_i 's
- Training set: $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$

Goal: solve for W_{ij} 's and b_{ij} 's that minimize the average cost function $J(\mathbf{w})$ over the training set

$$J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m L(\hat{f}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

- Stochastic Gradient Descent or Adam is used
- Weight initialization: Glorot uniform $w_{ij} \sim \text{Unif}(-a, a)$
- Back propagation to compute gradients

Universal Approximation

- Let \mathcal{N} be a feedforward network with one hidden layer and linear output units (ie. $\hat{f} = f_2 \circ f_1$ and f_2 is linear)
- (Hornik 1991) *If the activation function is continuous, bounded and non-constant, \mathcal{N} can approximate any continuous function on a compact subset $X \subset \mathbb{R}^n$ to any degree of accuracy provided that sufficiently many hidden units are available.*
- *Activation function example: $\tanh(z)$*
- (Leshno 1993) *\mathcal{N} with a locally bounded piecewise continuous activation can approximate any continuous function to any degree of accuracy if and only if the activation function is not a polynomial.*
- *Activation function example: $\max(0, z)$*

FNN: Regression


- Input layer: feature vector $\mathbf{x} \in \mathbb{R}^n$
- Output layer: label value y , $g_D = \text{identity}$
- Hyper-parameters:
 - D : depth of neural network
 - n_1, \dots, n_{D-1} : width of the hidden layers
 - g_1, \dots, g_{D-1} : activation functions
- Cost function: $L(y, \hat{y}) = \|y - \hat{y}\|^2$
- Training set: $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$
- Goal: solve for W_{ij} 's and b_i 's that minimize the average cost function $J(\mathbf{w})$ over the training set

$$J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m (\hat{f}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

FNN: Binary classification

- Feature space: $X \subset \mathbb{R}^n$, output: $Y = \{0, 1\}$
- Training set: $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$
- Cost function: take a probabilistic viewpoint
- What is this likelihood $P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$?
- Consider $P(y|\mathbf{x})$ conditional probability of y given \mathbf{x}
- Assume that it follows a Bernoulli distribution with parameter $p(\mathbf{x})$
- $P(y = 1|\mathbf{x}) = p(\mathbf{x})$, $P(y = 0|\mathbf{x}) = 1 - p(\mathbf{x})$
- $P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$ can be expressed as:
- $\prod_{i=1}^m P(y = \mathbf{y}^{(i)} | \mathbf{x}^{(i)}) = \prod_{i=1}^m p(\mathbf{x})^{y^{(i)}} (1 - p(\mathbf{x}))^{1-y^{(i)}}$
- Find $p(\mathbf{x})$ that maximizes the likelihood of observing the training dataset

FNN: Binary classification

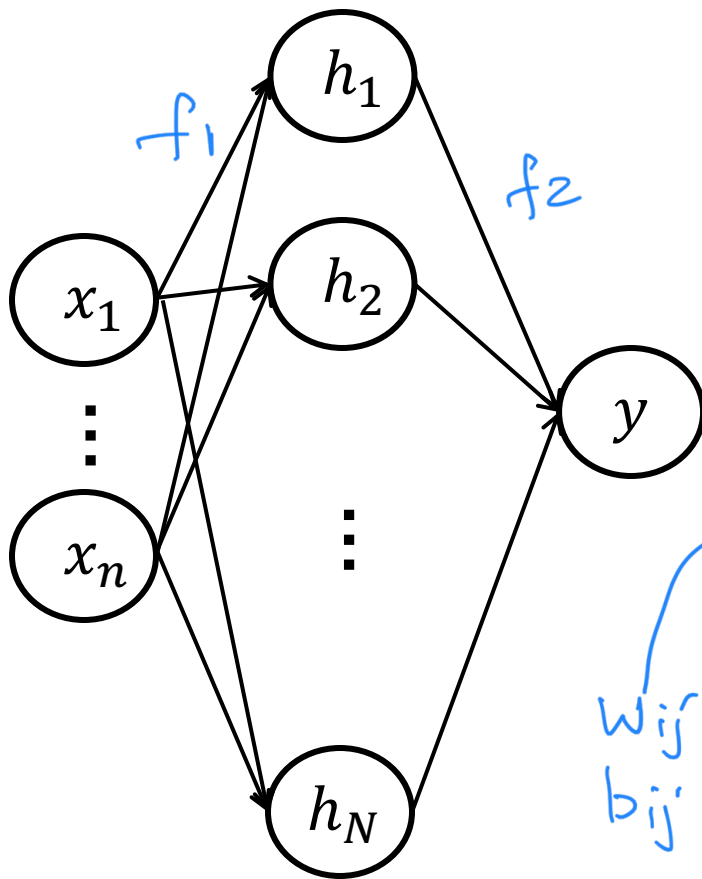
- Write $p(x)$ as a parametrized function $p(x; \theta)$ 
- Likelihood function $h(\theta) = \prod_{i=1}^m p(x; \theta)^{y^{(i)}} (1 - p(x; \theta))^{1-y^{(i)}}$
- Goal: Find $\hat{\theta}$ that maximizes the likelihood function
- Equivalently: Find $\hat{\theta}$ that minimizes the negative log-likelihood function $J(\theta) = -\log(h(\theta))$
- $J(\theta) = -\frac{1}{m} \sum_{i=1}^m \underbrace{y^{(i)} \log(p(x; \theta)) + (1 - y^{(i)}) \log(1 - p(x; \theta))}_{\text{Cost function}}$

Cost function

- Use feedforward neural network $\hat{f}(x; \theta)$ to approximate $p(x; \theta)$ where θ is the set of W_{ij} 's and b_{ij} 's in the network

FNN: Binary classification

$$\hat{f}: \mathbb{R}^n \rightarrow \mathbb{R} \text{ with } \hat{f} = f_2 \circ f_1$$



$f_1: \mathbb{R}^n \rightarrow \mathbb{R}^N$ with $f_1(\mathbf{x}) = (h_1, \dots, h_N)$
and $h_i = g(W_{1i}\mathbf{x} + b_1)$
 g is an activation function (eg. ReLU)

$f_2: \mathbb{R}^N \rightarrow \mathbb{R}$ with
 $f_2(\mathbf{x}) = \sigma(W_2\mathbf{x} + b_2)$
 $\sigma(z) = \frac{1}{1+\exp(-z)}$ logistic sigmoid

$\hat{f}(\mathbf{x}; \boldsymbol{\theta})$ is an estimate of $p(\mathbf{x}; \boldsymbol{\theta})$

Prediction: \vec{x}

if $\hat{f}(\mathbf{x}) \geq 0.5$, $\hat{y} = 1$; otherwise, $\hat{y} = 0$

FNN: K-class classification

$$\hat{f}: \mathbb{R}^n \rightarrow \mathbb{R}^K, \hat{f} = f_2 \circ f_1$$

$f_1: \mathbb{R}^n \rightarrow \mathbb{R}^N$ with $f_1(\mathbf{x}) = (h_1, \dots, h_N)$
and $h_i = g(W_{1i}\mathbf{x} + b_1)$

g is an activation function (eg. ReLU)

$f_2: \mathbb{R}^N \rightarrow \mathbb{R}^3$ with $f_2(\mathbf{h}) = s(W_2\mathbf{h} + b_2)$

$s: \mathbb{R}^3 \rightarrow \mathbb{R}^3$ softmax function

$s(z_1, z_2, z_3) = (s_1, s_2, s_3)$ where

$$s_i = s_i(\mathbf{z}) = \frac{\exp(z_i)}{\sum_{j=1}^3 \exp(z_j)}$$

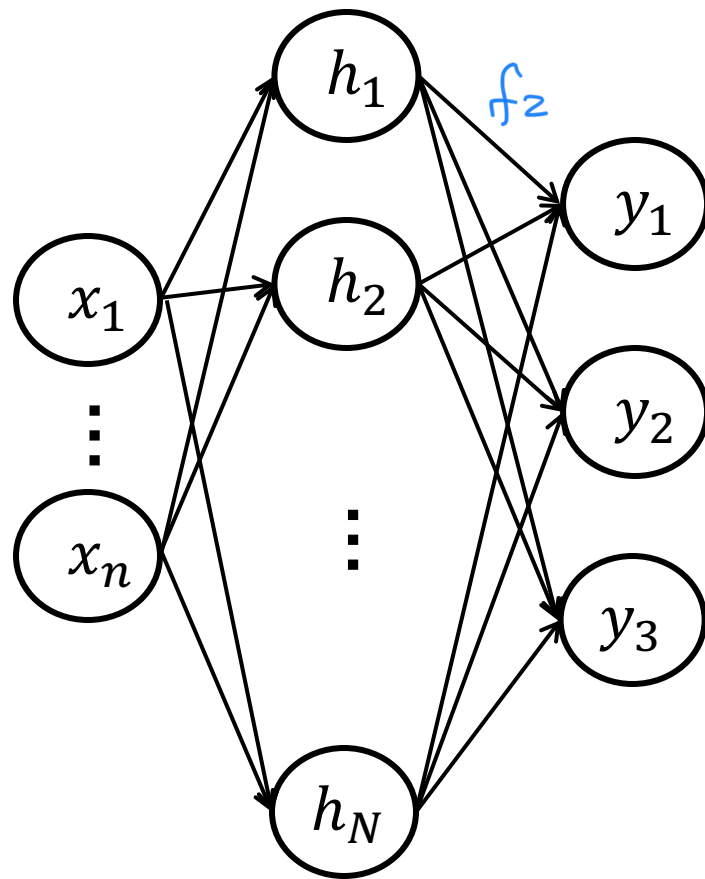
$\hat{f}(\mathbf{x}; \boldsymbol{\theta})$: estimate of $p_i(\mathbf{x}; \boldsymbol{\theta})$'s
ie. estimate of $P(y = i | \mathbf{x})$

Cost function:

- negative log-likelihood

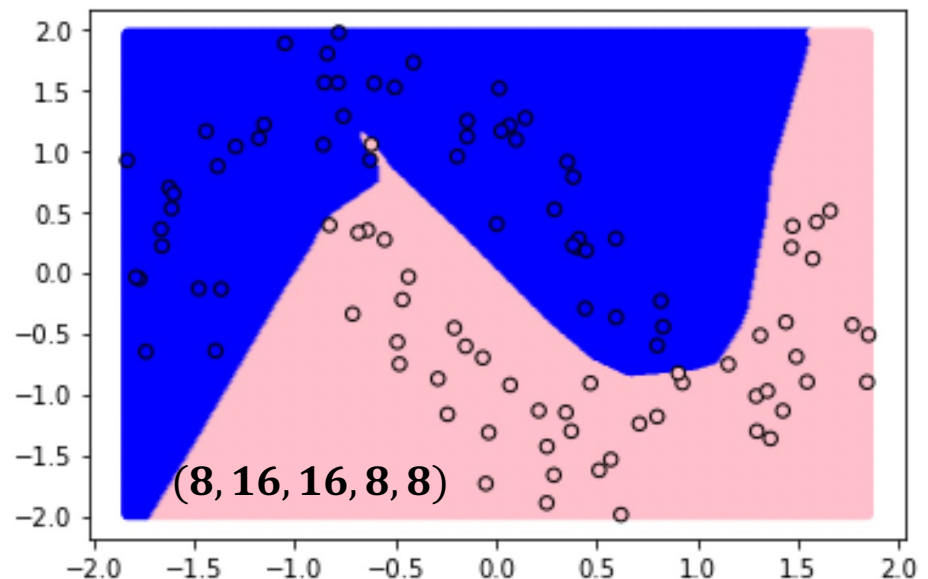
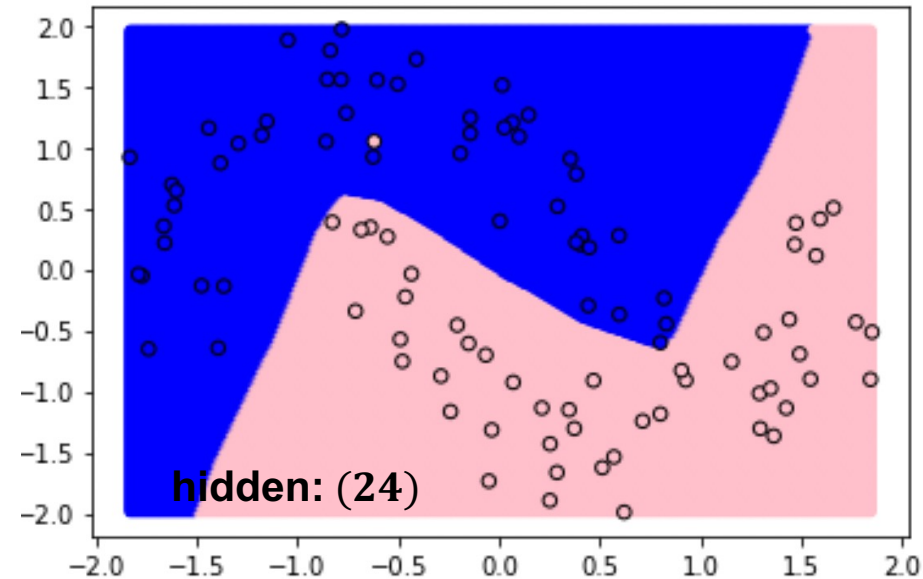
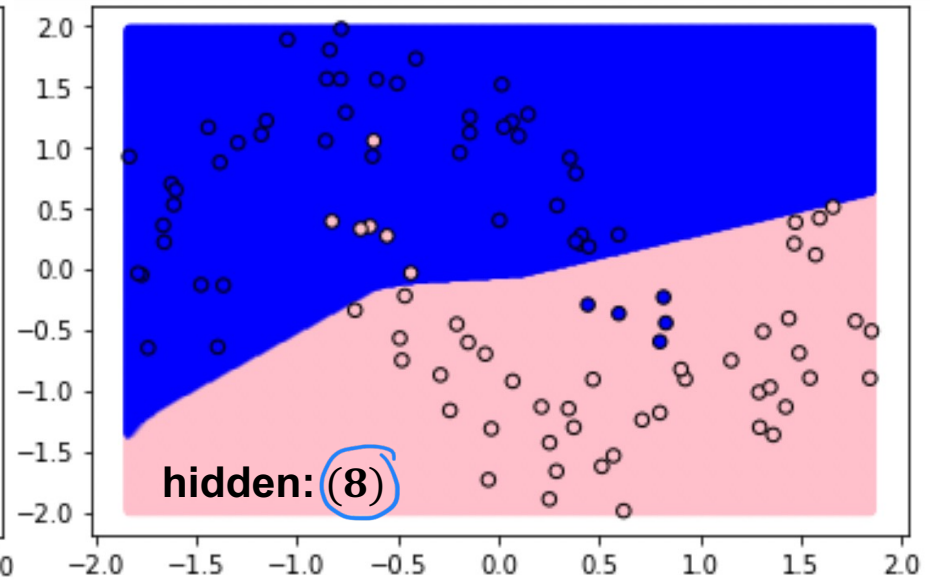
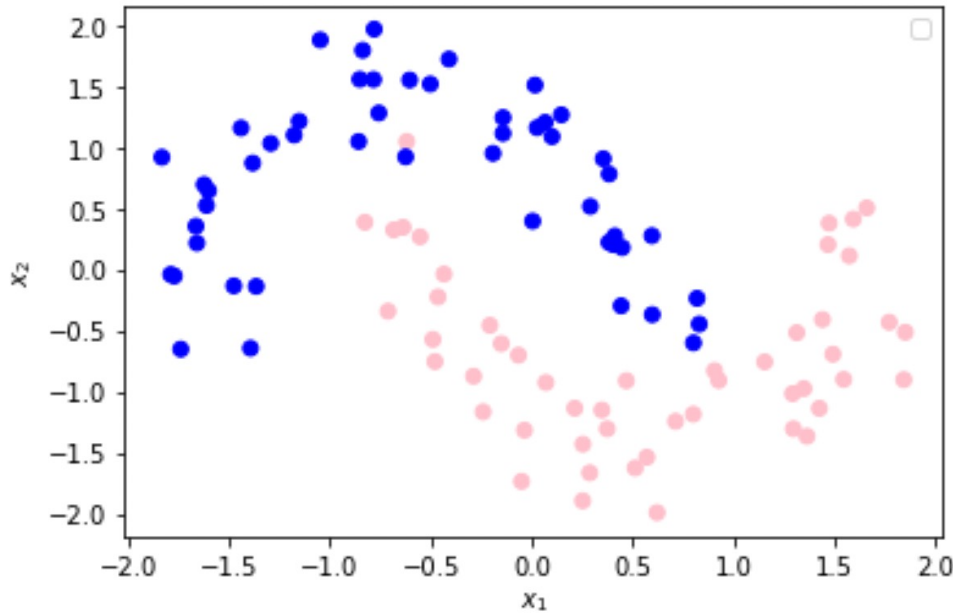
Prediction:

- class with highest predicted probability



MLP: width & depth

ReLU



Time series forecasting

- Classical linear time series models AR(p)
- $X(t) = \phi_0 + \phi_1 X(t-1) + \dots + \phi_p X(t-p) + \varepsilon$
- Classical ML models and ensembles
 - Features: lagged returns
- Recurrent neural networks (RNN)
 - Hidden state $h(t)$: information up to time t
 - $h(t) = f(X(t), h(t-1))$
 - Recurrent layer as a parametrization of function f
 - $h(t) = g(X(t)W_{xh} + h(t-1)W_{hh} + b_h)$
 - W_{xh} and W_{hh} are weight matrices, b_h bias vector
 - g : an activation function
 - Output layer: $y = s(h(t)W + b)$
- LSTM networks
 - Different parametrization of f

logistic function

