

Neural networks

COMS 4771 Fall 2025

Feature maps revisited

Justification for simple statistical models (e.g., logistic regression):

- ▶ They are reasonable with a judicious choice of features or feature map
- ▶ In linear models, best prediction of Y given $X = x$ is based entirely on

$$w^\top \varphi(x)$$

where φ is the feature map

1 / 34

Weierstrass approximation theorem: For any continuous function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, any bounded region $B \subset \mathbb{R}^d$, and any $\varepsilon > 0$, there exists a polynomial $g: \mathbb{R}^d \rightarrow \mathbb{R}$ such that

$$\max_{x \in B} |f(x) - g(x)| \leq \varepsilon$$

- ▶ Polynomials give good approximations uniformly over an interval (Cf. Taylor's theorem: only guarantees local approximations)
- ▶ Universal justification of polynomial expansion + linear functions
- ▶ Caveat: Degree of g may be large (e.g., growing with d and $1/\varepsilon$)
 - ▶ Somewhat ameliorated by kernel methods

2 / 34

Kernel machine: function learned by kernel method

$$g(x) = \sum_{i=1}^n \alpha_i k(x, x^{(i)})$$

where $k(\cdot, \cdot)$ is the kernel function, and $x^{(1)}, \dots, x^{(n)}$ are from training data

3 / 34

Stone-Weierstrass approximation theorem: For any continuous function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, any bounded region $B \subset \mathbb{R}^d$, and any $\varepsilon > 0$, there exists a function $g: \mathbb{R}^d \rightarrow \mathbb{R}$ of the form

$$g(x) = \sum_{i=1}^p \alpha_i \exp(x^\top w^{(i)})$$

such that

$$\max_{x \in B} |f(x) - g(x)| \leq \varepsilon$$

- ▶ Can replace “exp” with other “activation functions”
- ▶ Caveat: p may be large
- ▶ Another interpretation: linear function $\alpha^\top \varphi(x)$ with feature map

$$\varphi(x) = (\exp(x^\top w^{(1)}), \dots, \exp(x^\top w^{(p)}))$$

Except the $w^{(i)}$'s may need to depend on f

- ▶ This kind of function is called a (two-layer) neural network

4 / 34

Kernel machine

$$g(x) = \sum_{i=1}^n \alpha_i k(x, x^{(i)})$$

- ▶ Only α_i 's are learned using data

(Two-layer) neural network

$$g(x) = \sum_{i=1}^p \alpha_i \exp(x^\top w^{(i)})$$

- ▶ Both α_i 's and $w^{(i)}$'s are learned
- ▶ Can use $p > n$

Neural networks as straight-line programs

Very abbreviated history:

- ▶ McCulloch and Pitts (early 1940s):
Neural networks as computational model for brain
- ▶ Arnold and Kolmogorov (late 1950s):
Solved Hilbert's 13th problem (about polynomial roots) using neural networks
- ▶ **Modern use of neural networks with Linnainmaa's autodiff (early 1970s) started with Werbos (early 1980s)**
- ▶ Many other researchers have since discovered other approximation-theoretic properties and practical uses of neural networks (e.g., Cybenko, Rumelhart and Hinton, LeCun)

Today, for machine learning purposes: a [neural network](#) is (almost) any function f such that $f(x)$ can be computed by a [straight-line program](#)

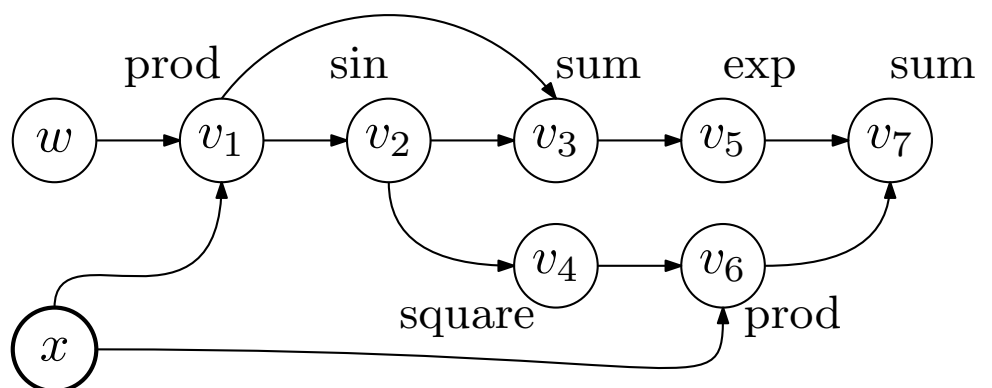
6 / 34

[Straight-line program](#): each line declares a new variable as a function of inputs (e.g., x), numerical parameters (e.g., w), or previously defined variables

Example:

$$f(x) = \exp(xw + \sin(xw)) + \sin^2(xw)x$$

$v_1 := \text{prod}(x, w)$
 $v_2 := \sin(v_1)$
 $v_3 := \text{sum}(v_1, v_2)$
 $v_4 := \text{square}(v_2)$
 $v_5 := \exp(v_3)$
 $v_6 := \text{prod}(v_4, x)$
 $v_7 := \text{sum}(v_5, v_6)$



Computation directed acyclic graph (DAG) $G = (V, E)$

7 / 34

Example (allowing slightly more advanced functions):

$$f(x) = \alpha_0 + \sum_{i=1}^p \alpha_i \text{logistic}(x^\top w^{(i)} + b^{(i)})$$

$$v_1 := \text{logistic}(x^\top w^{(1)} + b^{(1)})$$

$$v_2 := \text{logistic}(x^\top w^{(2)} + b^{(2)})$$

$$\vdots$$

$$v_p := \text{logistic}(x^\top w^{(p)} + b^{(p)})$$

$$\hat{y} := \alpha_0 + \alpha_1 \cdot v_1 + \alpha_2 \cdot v_2 + \cdots + \alpha_p \cdot v_p$$

- ▶ v_1, \dots, v_p called [hidden units](#) (antiquated terminology)
- ▶ A single-line using modern numerical software (e.g., pytorch):

$$\hat{y} := \alpha_0 + \alpha^\top \text{logistic}(Wx + b)$$

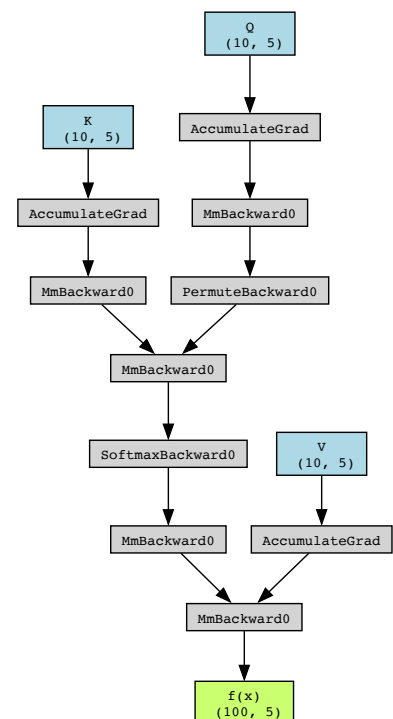
Parameters: $W \in \mathbb{R}^{p \times d}$, $b, \alpha \in \mathbb{R}^p$, $\alpha_0 \in \mathbb{R}$
 $\text{logistic}: \mathbb{R} \rightarrow \mathbb{R}$ is applied component-wise

8 / 34

Example (in pytorch):

```
K = torch.randn(d, p, requires_grad=True)
Q = torch.randn(d, p, requires_grad=True)
V = torch.randn(d, p, requires_grad=True)
```

```
def f(x):
    k = x @ K
    q = x @ Q
    a = torch.softmax(k @ q.T, dim=1)
    return a @ x @ V
```



9 / 34

In practice, neural network “architectures” (i.e., program “templates”) are built by using/composing component modules

Simplest module is fully-connected layer:

$$h \mapsto \sigma(Wh + b)$$

Affine transformation followed by non-linear transformation

Some examples of non-linear transformations σ :

- ▶ Rectified linear unit: $\text{relu}(t) = [t]_+ = \max\{0, t\}$
- ▶ Hyperbolic tangent: $\tanh(t) = 2 \text{logistic}(t) - 1$
- ▶ Softmax: $\text{softmax}: \mathbb{R}^k \rightarrow \mathbb{R}^k$, where $\text{softmax}(u)_i = \frac{\exp(u_i)}{\sum_{j=1}^k \exp(u_j)}$

10 / 34

Q: How to choose the “architecture”?

A: Primary constraint is that functions should be “differentiable” (explained later)

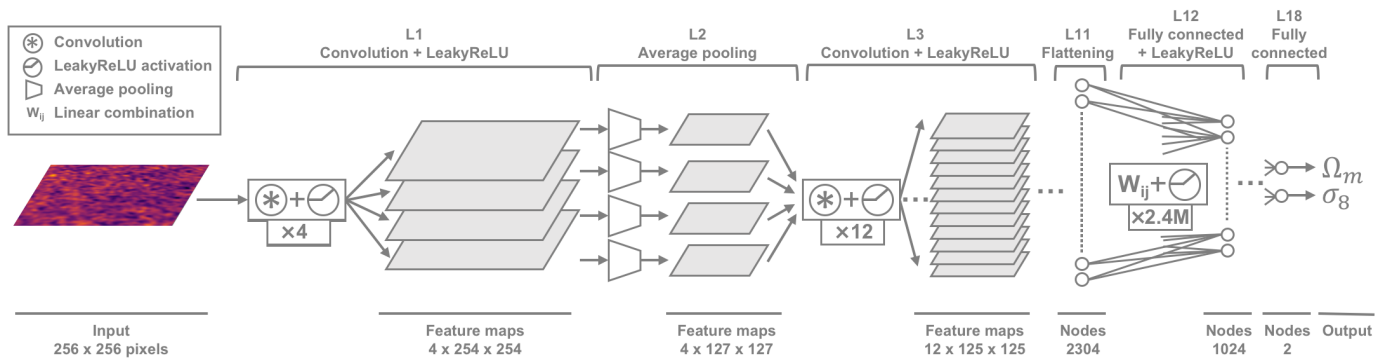
- ▶ According to Stone-Weierstrass theorem, two-layer architecture

$$f(x) = \alpha^\top \sigma(Wx)$$

for $\sigma = \exp$, $\alpha \in \mathbb{R}^p$, $W \in \mathbb{R}^{p \times d}$ is “sufficient” for almost any purpose (provided p is large enough)

- ▶ Other architectural designs typically based on application-specific considerations (e.g., convnets), optimization-based empirical observations (e.g., resnets), and/or manual experimentation (e.g., SwiGLU)
- ▶ Typical approach: Start from existing architecture that someone already used for a similar problem, and consider small modifications

11 / 34



12 / 34

Problem: How to fit neural network f_θ (with parameters θ) to training data?

► Typical training objective:

$$J(\theta) := \sum_{i=1}^n \text{loss}(f_\theta(x^{(i)}), y^{(i)})$$

```
loss = torch.nn.NLLLoss(reduction='sum')
J = loss(f(training_features), training_labels)
```

► Objective function can also be computed using a straight-line program

Amazing fact: Can compute gradient of $J(\theta)$ with respect to θ in (roughly) same amount of time as computing $J(\theta)$ itself (using [autodiff](#), discussed later)

Therefore can use gradient-based methods to (try to) minimize training objective

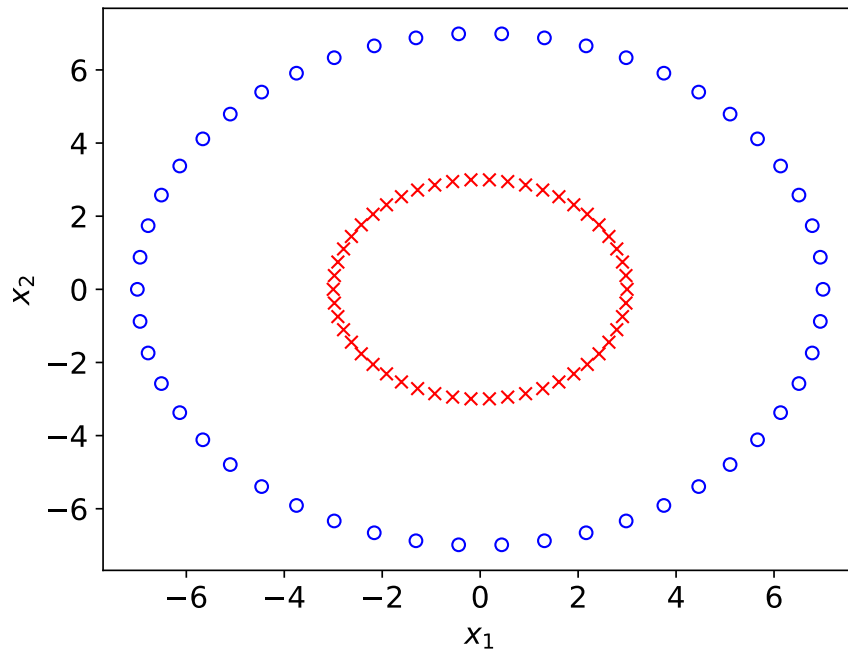
13 / 34

Main challenge: Objective function $J(\theta)$ might not be convex, so use of gradient-based optimization is more complicated (e.g., initialization, step sizes)

- ▶ Many tips and tricks (e.g., “Efficient BackProp”, LeCun et al, 1998)
- ▶ Ultimately need a lot of experimentation

Synthetic example

Data: classes are two concentric circles, 50 examples per class



15 / 34

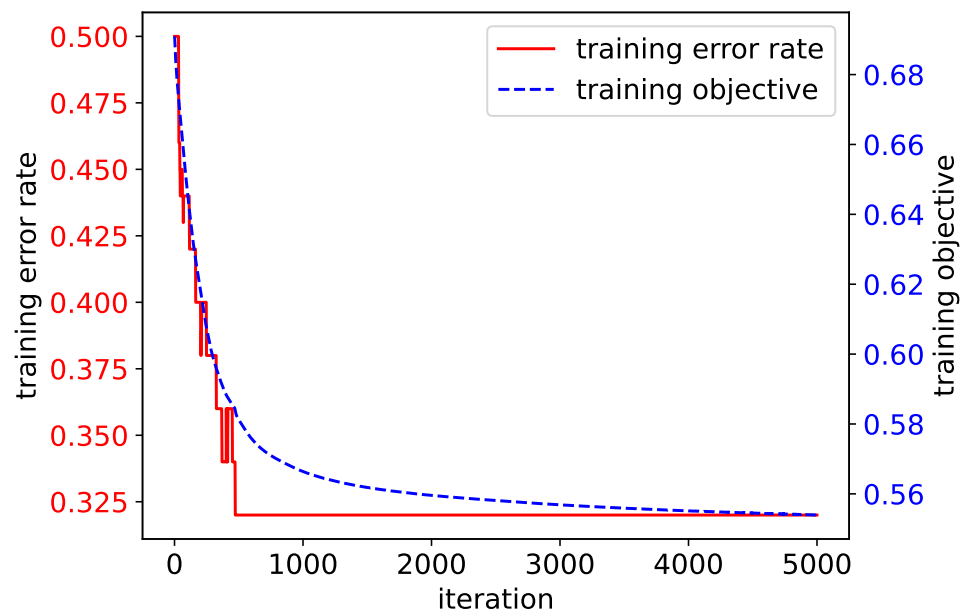
- ▶ Feature transformation: standardization
- ▶ Neural net: $f(x) = \text{softmax}(A \text{relu}(Wx + b) + c)$
 - ▶ Parameters: $W \in \mathbb{R}^{p \times 2}$, $b \in \mathbb{R}^p$, $A \in \mathbb{R}^{2 \times p}$, $c \in \mathbb{R}^2$
(We will vary the “width” p)
 - ▶ k -th output is prediction of $\Pr(Y = k \mid X = x)$
- ▶ Use gradient descent on average logarithmic loss on training data
 - ▶ Random initialization:

$$W_{i,j}, b_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{3}), \quad A_{i,j}, c_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{2}{p+1})$$

- ▶ Step size: $\eta_t = 0.1$

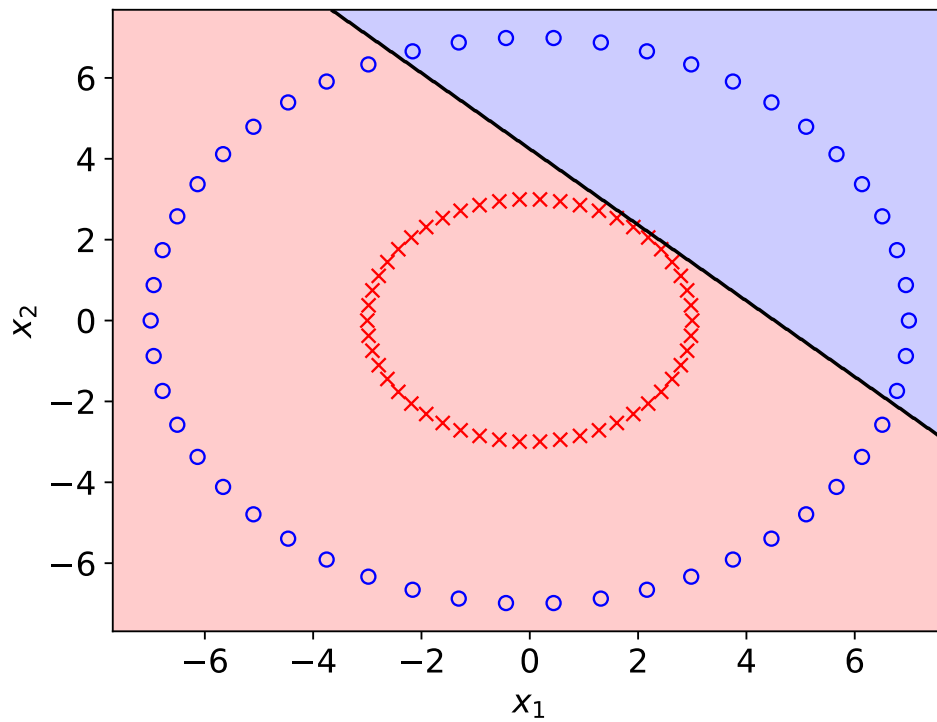
16 / 34

Results: $p = 2$



17 / 34

Results: $p = 2$



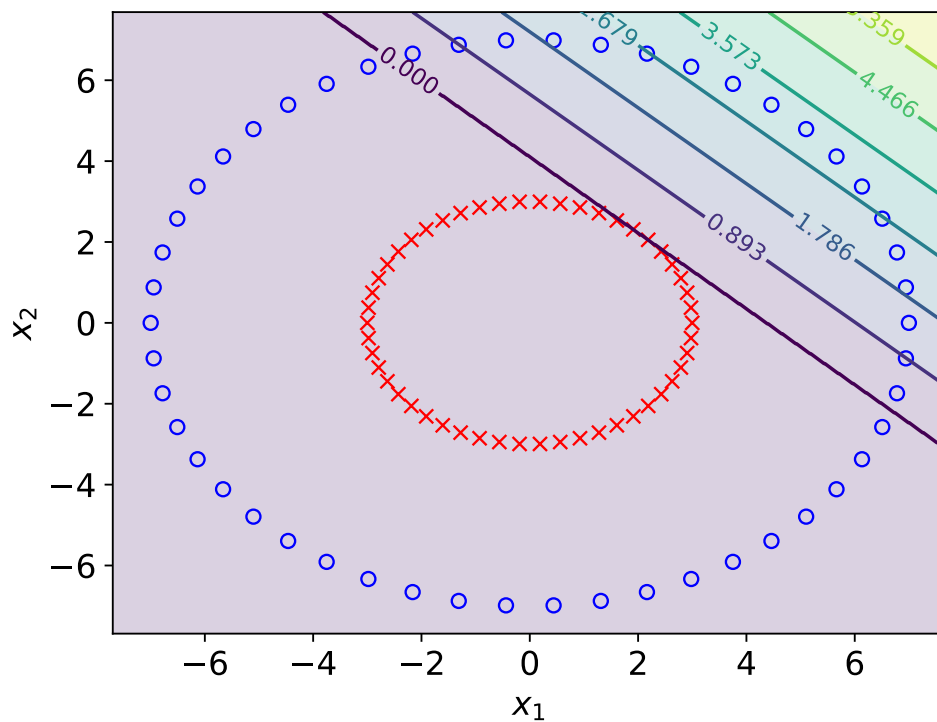
18 / 34

Results: $p = 2$

- ▶ Behaves like a linear classifier
- ▶ First component of $\text{relu}(Wx + b)$ is constant (0) over training data
- ▶ Only second component of $\text{relu}(Wx + b)$ varies over training data

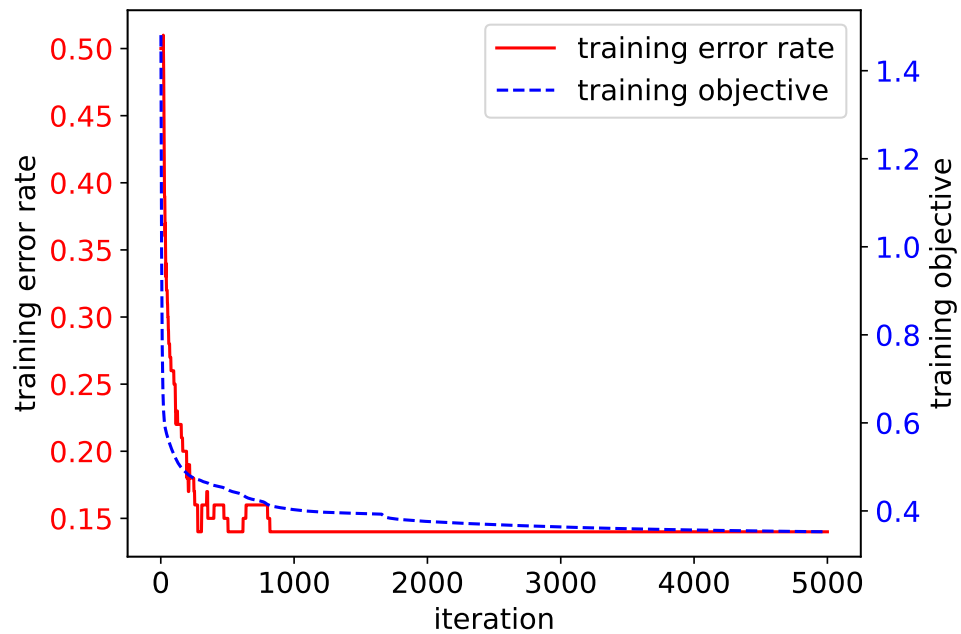
19 / 34

Results: $p = 2$ — second component of $\text{relu}(Wx + b)$



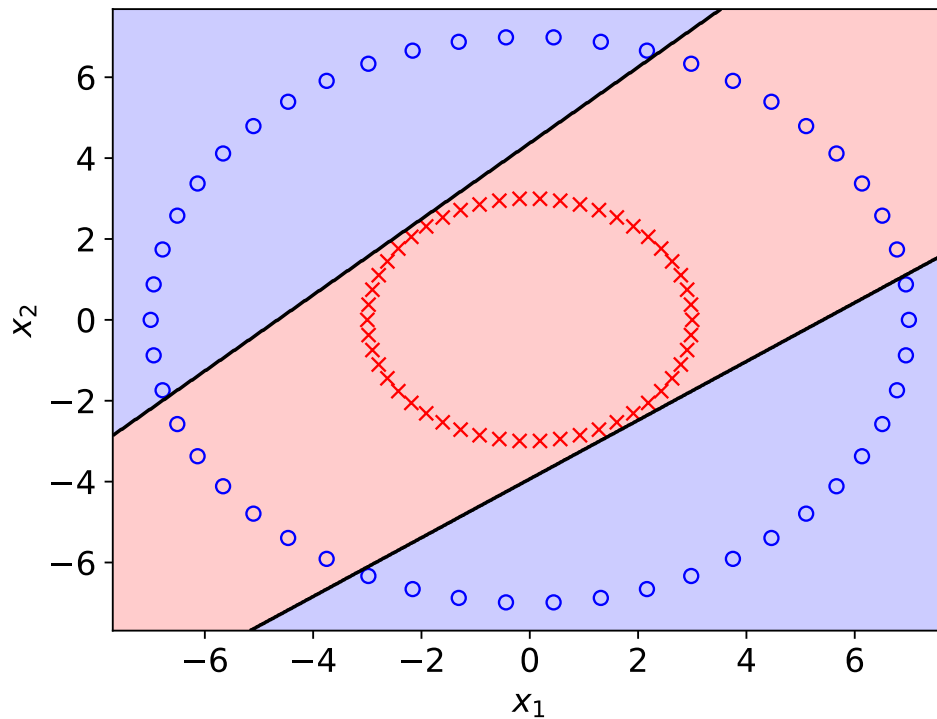
20 / 34

Results: $p = 2$ (different initialization)



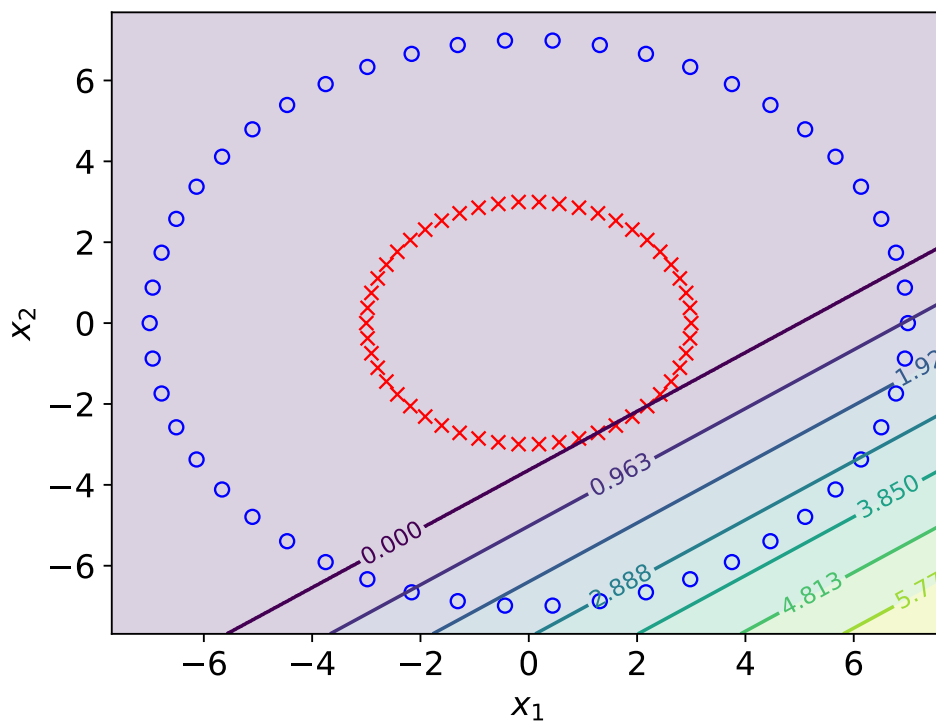
21 / 34

Results: $p = 2$ (different initialization)



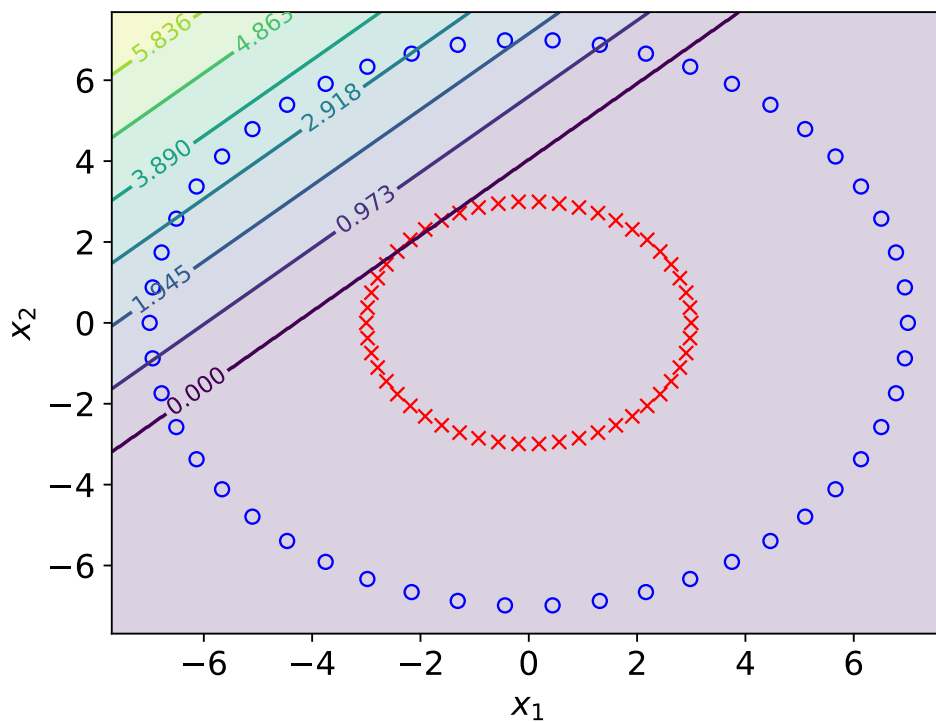
22 / 34

Results: $p = 2$ (different initialization) — first component of $\text{relu}(Wx + b)$



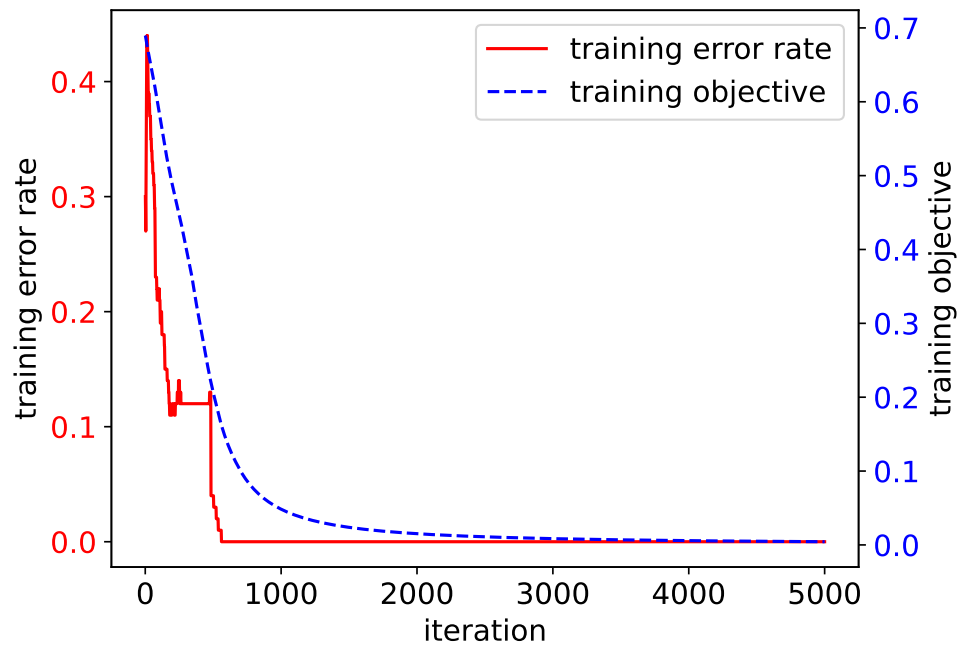
23 / 34

Results: $p = 2$ (different initialization) — second component of $\text{relu}(Wx + b)$



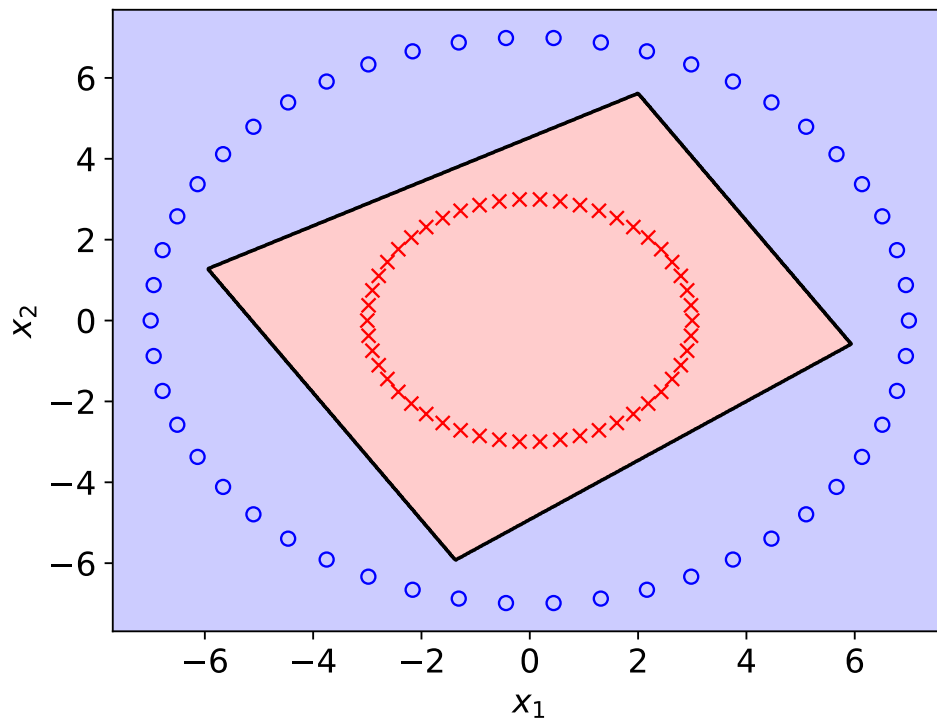
24 / 34

Results: $p = 3$



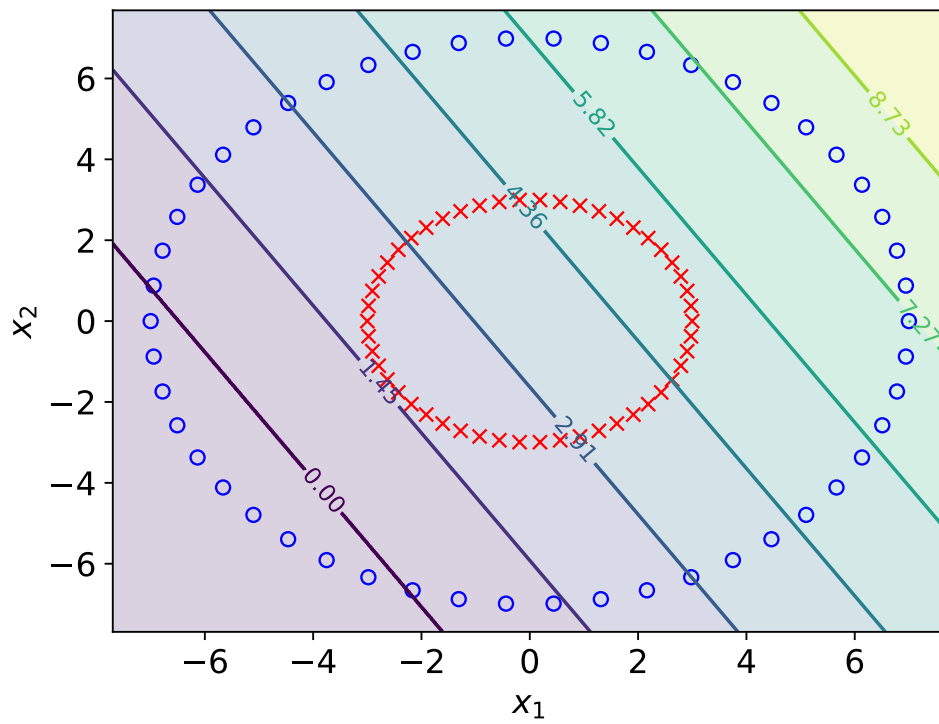
25 / 34

Results: $p = 3$



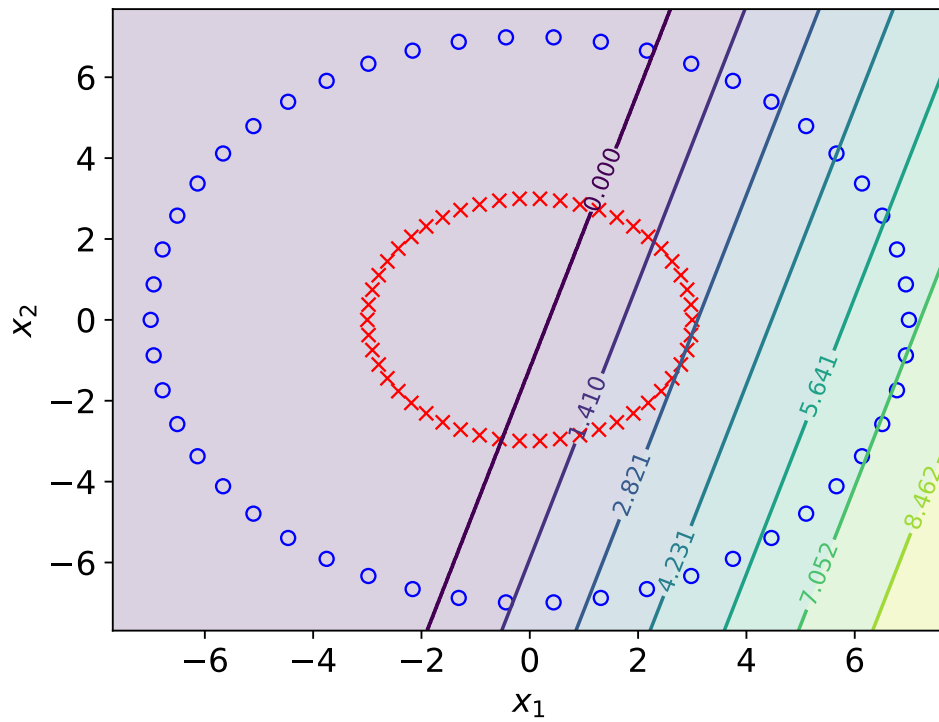
26 / 34

Results: $p = 3$ — first component of $\text{relu}(Wx + b)$



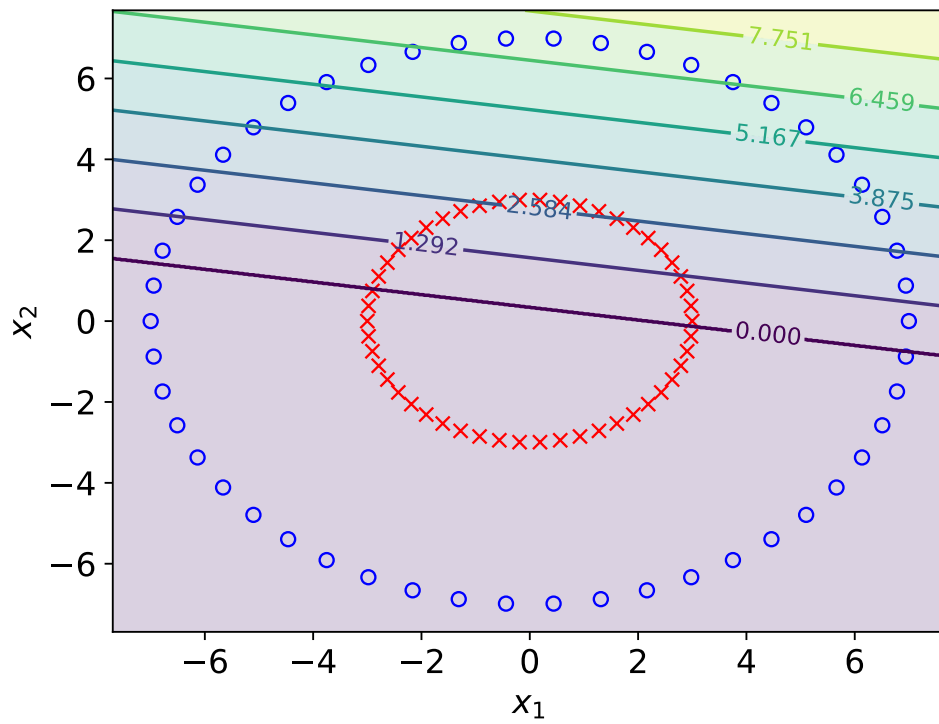
27 / 34

Results: $p = 3$ — second component of $\text{relu}(Wx + b)$



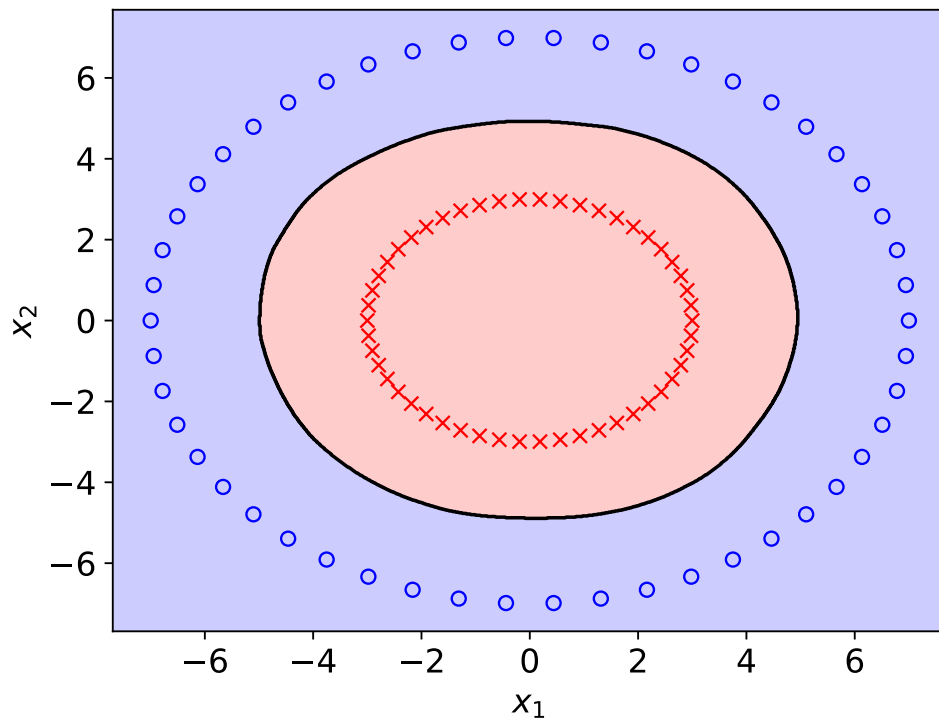
28 / 34

Results: $p = 3$ — third component of $\text{relu}(Wx + b)$



29 / 34

Results: $p = 1000$



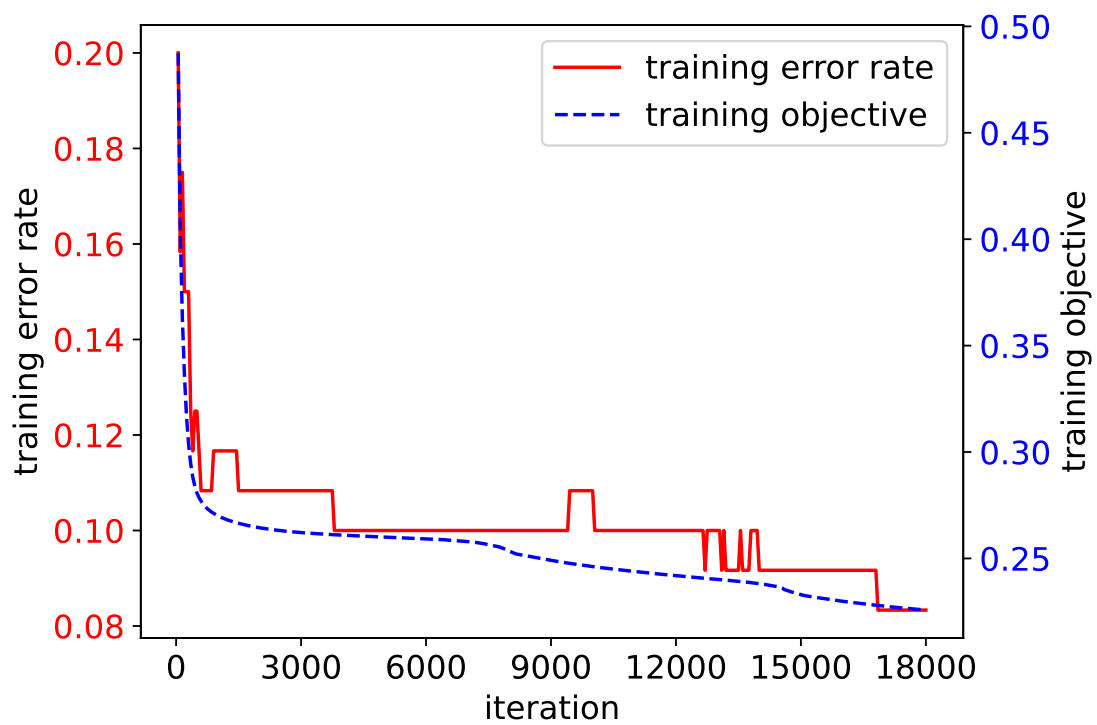
30 / 34

Iris data classifier

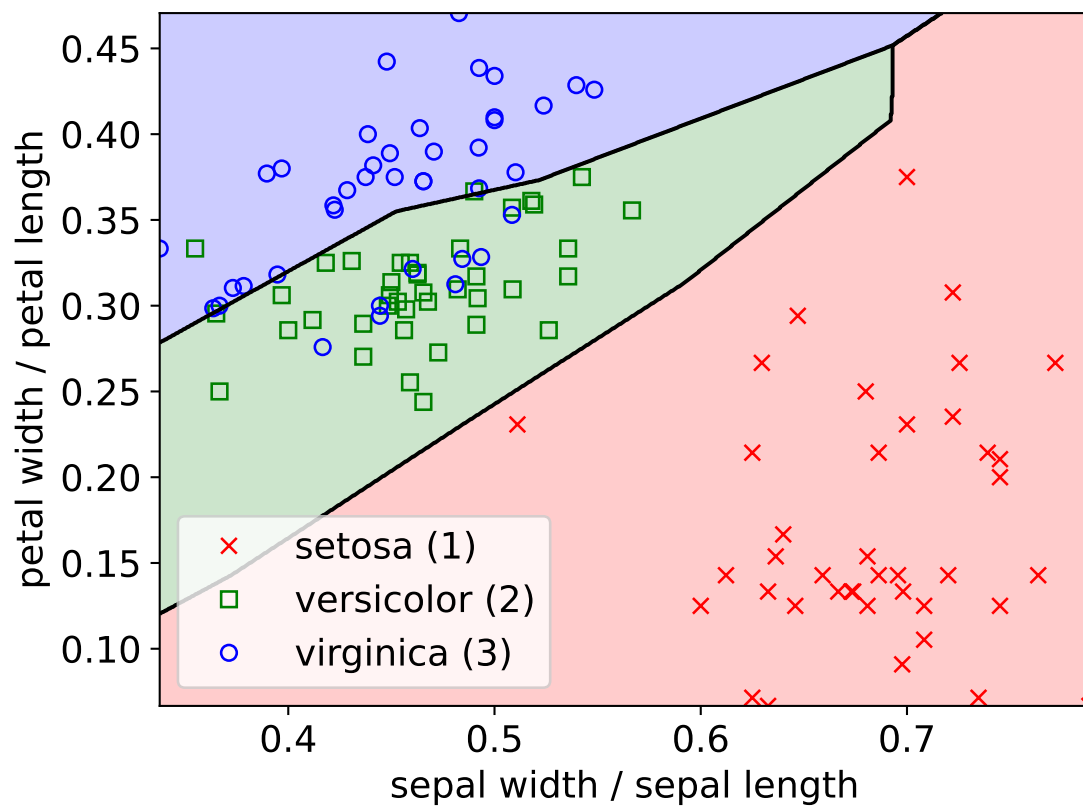
- Features:

$$x_1 = \text{sepal width} / \text{sepal length}, \quad x_2 = \text{petal width} / \text{petal length}$$

- Neural net: $f(x) = \text{softmax}(A \text{relu}(Wx + b) + c)$
 - Parameters: $W \in \mathbb{R}^{10 \times 2}$, $b \in \mathbb{R}^{10}$, $A \in \mathbb{R}^{2 \times 10}$, $c \in \mathbb{R}^2$
 - k -th output is prediction of $\Pr(Y = k \mid X = x)$
- Feature transformation and training procedure: same as in synthetic example
- Training error rate: 8.33%, test error rate: 10.0%



32 / 34



33 / 34

Deep learning lifestyle

- ▶ Since 2006–2012, use of neural networks has exploded in machine learning
- ▶ Called “deep learning” due to use of large and “deep” neural networks
 - ▶ Use of multiple layers is inspired by some biological systems (e.g., visual system) for processing natural signals
 - ▶ Conventional wisdom: very wide but shallow neural network can be replaced by moderately wide but deeper neural network
- ▶ Key factors in latest resurgence and success:
 - ▶ Graphics processing units (GPUs) to speed-up matrix operations
 - ▶ Easy-to-use numerical software with autodiff (e.g., pytorch)
 - ▶ Large benchmark datasets (e.g., ImageNet)