# Unit 8
# Neural Networks

**Prof. Phil Schniter**

THE OHIO STATE UNIVERSITY

**ECE 5307: Introduction to Machine Learning, Sp23**

# Learning objectives

- Understand 2-layer feedforward neural networks:
  - Motivation: learning feature transformations
  - Network architecture: linear and nonlinear layers
  - Choice of activation functions and training loss

- Understand mini-batch training and stochastic gradient descent

- Understand the back-propagation approach to gradient computation

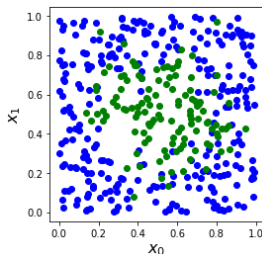- Know how to implement a neural network using PyTorch

# Outline

- Motivating Example: Learning a Feature Transformation

- Feed-Forward Neural Networks

- Training via Stochastic Gradient Descent

- Gradient Computation via Back-Propagation

- Implementing and Training Neural Nets with PyTorch

# Dealing with data that is not linearly separable

- Consider the data $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$ on the right:
  - Features $\boldsymbol{x}_i = [x_{i0}, x_{i1}]^{\mathsf{T}} \in \mathbb{R}^2$
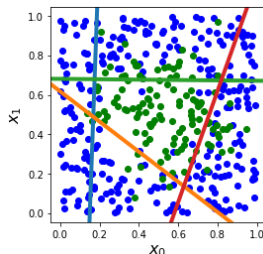  - Labels $y_i \in \{0, 1\}$
  - Not linearly separable!



- Could use a feature transformation to enhance linear separability, and then apply linear classification on the transformed features
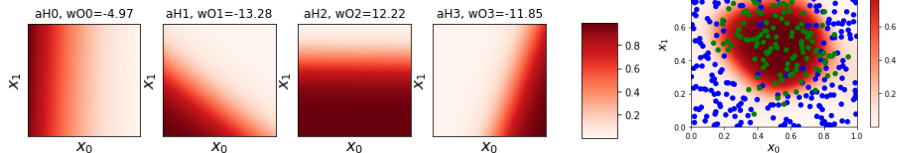  - But what if we don't know a good transformation?
  - Can we learn one?
  - Yes!

# A two-stage approach to classification

- A two-stage approach:
  1) Learn 4 transformed features, each the soft output of a linear classifier
     - "soft" output means in the interval $[0, 1]$
  2) Apply linear classification to the transformed features, giving a final soft output
     - high for one intersection of half-spaces



- Plot of transformed features and final soft output vs. $x$:



- The overall approach is successful in classifying the data!
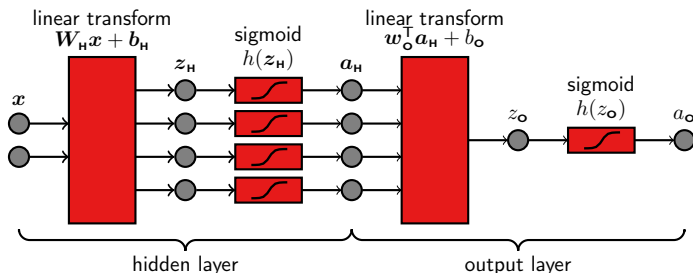
# Details of two-stage approach

- Stage 1: "Hidden layer"
  - scores: $z_{\mathsf{H}} = W_{\mathsf{H}} x + b_{\mathsf{H}} \in \mathbb{R}^{d_{\mathsf{H}}}$
  - soft outputs: $a_{\mathsf{H}} = h(z_{\mathsf{H}}) \in [0,1]^{d_{\mathsf{H}}}$
  - $a_{\mathsf{H}}$ are the learned features

- Stage 2: "Output layer"
  - score: $z_{\mathsf{O}} = w_{\mathsf{O}}^{\mathsf{T}} a_{\mathsf{H}} + b_{\mathsf{O}} \in \mathbb{R}$
  - soft output: $a_{\mathsf{O}} = h(z_{\mathsf{O}}) \in [0,1]$
  - $a_{\mathsf{O}}$ is the final $\Pr\{y = 1 \,|\, x\}$

- If we use $h(z) = \frac{1}{1+e^{-z}}$, a sigmoid, the output layer is logistic regression

- This is a multi-layer perceptron, or 2-stage feed-forward neural network!
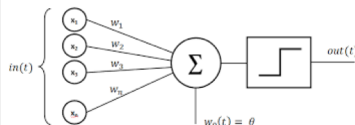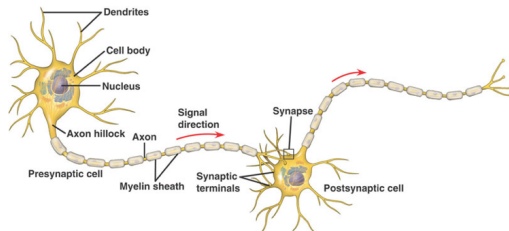
# Training the model

- Let's collect the model parameters into $\boldsymbol{\theta} \triangleq (\boldsymbol{W}_{\mathsf{H}}, \boldsymbol{b}_{\mathsf{H}}, \boldsymbol{w}_{\mathsf{o}}, b_{\mathsf{o}})$

- We fit these parameters using training data $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$

- Since we used logistic regression, the likelihood function would be

$$\Pr\{y_i = 1 \mid \boldsymbol{x}_i; \boldsymbol{\theta}\} = \frac{1}{1 + e^{-z_{\mathsf{o},i}}} \quad \text{for} \quad z_{\mathsf{o},i} = F(\boldsymbol{x}_i; \boldsymbol{\theta})$$

  - $F(\boldsymbol{x}; \boldsymbol{\theta})$ describes the network from input $\boldsymbol{x}$ to output score $z_{\mathsf{o}}$

- Then the maximum-likelihood model parameters are

$$\widehat{\boldsymbol{\theta}}_{\mathsf{ml}} = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left[ -\ln p(y_i \mid \boldsymbol{x}_i; \boldsymbol{\theta}) \right]$$
$$= \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left( \ln[1 + e^{z_{\mathsf{o},i}}] - y_i z_{\mathsf{o},i} \right) \quad \text{(binary cross-entropy loss)}$$

  - We will discuss the details of this optimization later

# Why is it called a "neural" network?



- Simple model of neurons:
  - Dendrites: Input currents from other neurons
  - Soma: Cell body, accumulation of charge
  - Axon: Outputs to other neurons

- Operation:
  - Output when the sum of input currents reaches a threshold
  - Similar to (artificial) neural network: $a_{\mathsf{H},j} = h(\boldsymbol{w}_{\mathsf{H},j}^{\mathsf{T}} \boldsymbol{x} + b_{\mathsf{H},j})$
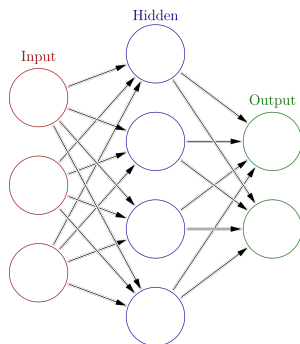
# History

- 1940s: Donald Hebb — Hebbian learning for neural plasticity
  - Hypothesized rule for updating synaptic weights in biological neurons

- 1950s: Frank Rosenblatt — Coined the term "perceptron"
  - Essentially a single-layer network, similar to logistic regression
  - Early computer implementations
  - But linear classifiers are limited, and so was compute power

- 1960s: Back-propagation — Efficient way to train multi-layer networks
  - We'll cover this later in this unit

- 1990s: Resurgence with greater computational power

- 2012-now: The deep-network revolution
  - Many more layers. Massive computational power and data
  - Breakthroughs in speech and image processing first, then many more fields . . .
  - We'll cover deep networks in the next unit

# Outline

- Motivating Example: Learning a Feature Transformation

- **Feed-Forward Neural Networks**

- Training via Stochastic Gradient Descent

- Gradient Computation via Back-Propagation

- Implementing and Training Neural Nets with PyTorch

# General structure of a feed-forward neural network

- Input: $\boldsymbol{x} = [x_1, \ldots, x_d]^{\mathsf{T}}$
  - $d =$ number of features (or inputs)

- Hidden layer:
  - $a_{\mathsf{H},l} = h_{\mathsf{H}}\big(b_{\mathsf{H},l} + \sum_{j=1}^{d} w_{\mathsf{H},lj} x_j\big),\ l = 1 \ldots d_{\mathsf{H}}$
  - $h_{\mathsf{H}}(\cdot)$ is a nonlinear "activation" function
  - $d_{\mathsf{H}} =$ number of hidden units

- Output layer:
  - $a_{\mathsf{O},k} = h_{\mathsf{O}}\big(b_{\mathsf{O},k} + \sum_{l=1}^{d_{\mathsf{H}}} w_{\mathsf{O},kl} a_{\mathsf{H},l}\big),\ k = 1 \ldots d_{\mathsf{O}}$
  - $h_{\mathsf{O}}(\cdot)$ is optional. If present, it must match loss
  - $d_{\mathsf{O}} =$ number of outputs

- Networks of this form also referred to as "multilayer perceptrons"

- Can use more than two layers ("deep network"), but for now we focus on two

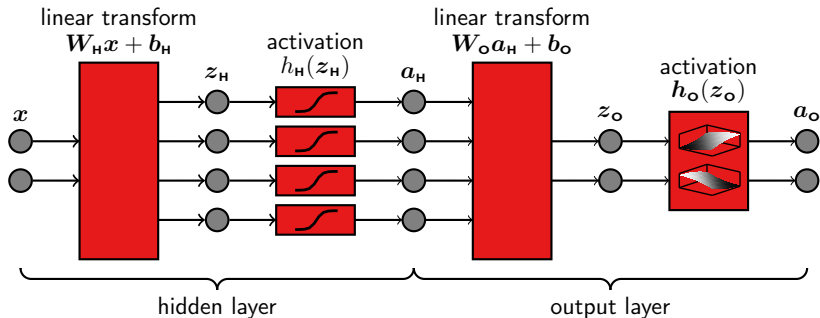# Vector/matrix representation of feed-forward neural network

- Hidden layer:
    - linear transform: $\boldsymbol{z_\mathsf{H}} = \boldsymbol{W_\mathsf{H}} \boldsymbol{x} + \boldsymbol{b_\mathsf{H}}$
    - nonlinear activation: $\boldsymbol{a_\mathsf{H}} = h_\mathsf{H}(\boldsymbol{z_\mathsf{H}}) \in \mathbb{R}^{d_\mathsf{H}}$, where $h_\mathsf{H}(\cdot)$ acts elementwise

- Output layer:
    - linear transform: $\boldsymbol{z_\mathsf{O}} = \boldsymbol{W_\mathsf{O}} \boldsymbol{a_\mathsf{H}} + \boldsymbol{b_\mathsf{O}}$
    - (optional) activation: $\boldsymbol{a_\mathsf{O}} = \boldsymbol{h_\mathsf{O}}(\boldsymbol{z_\mathsf{O}}) \in \mathbb{R}^{d_\mathsf{O}}$, where $\boldsymbol{h_\mathsf{O}}(\cdot)$ may act on full vector
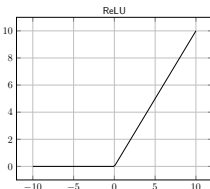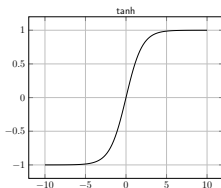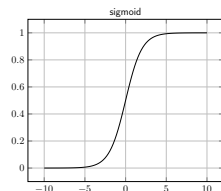
# Common choices for hidden-layer activation $h_{\text{H}}(\boldsymbol{z}_{\text{H}})$

- sigmoid: $h_{\text{H}}(z) = \frac{1}{1+e^{-z}} \triangleq \sigma(z)$
  - output is bounded, but not centered
  - sometimes used in shallow networks

- tanh: $h_{\text{H}}(z) = \tanh(z) = 2\sigma(z) - 1$
  - output is bounded and centered
  - often works better than sigmoid, because centered data is consistent with PyTorch weight initializations

- rectified linear unit (ReLU): $h_{\text{H}}(z) = \max\{0, z\}$
  - output is unbounded and not centered
  - most popular choice in deep networks
  - works well in shallow networks too

# Common choices for output-layer activation $h_\mathbf{o}(z_\mathbf{o})$

- Binary classification:
    - here $z_\mathbf{o} = z_\mathbf{o}$ is a scalar
    - option 1: use no output activation fxn (i.e., output $z_\mathbf{o}$)
    - option 2: $h_\mathbf{o}(z_\mathbf{o}) = \dfrac{1}{1 + e^{-z_\mathbf{o}}} = \Pr\{y=1 \,|\, \boldsymbol{x}\}$ (logistic)

- Multiclass classification with $K > 2$ classes:
    - here $z_\mathbf{o} = [z_{\mathbf{o},1}, \ldots, z_{\mathbf{o},K}]^\mathsf{T} \in \mathbb{R}^K$
    - option 1: use no output activation fxn (i.e., output $z_\mathbf{o}$)
    - option 2: $[h_\mathbf{o}(z_\mathbf{o})]_k = \dfrac{e^{z_{\mathbf{o},k}}}{\sum_{l=1}^{K} e^{z_{\mathbf{o},l}}} = \Pr\{y=k \,|\, \boldsymbol{x}\}$ for $k = 1...K$ (softmax)

- Regression with $K$-dimensional targets (i.e., $\boldsymbol{y}_i \in \mathbb{R}^K$):
    - here $z_\mathbf{o} = [z_{\mathbf{o},1}, \ldots, z_{\mathbf{o},K}]^\mathsf{T} \in \mathbb{R}^K$
    - use no output activation fxn (i.e., output $z_\mathbf{o}$)

# Loss functions for training

*The task determines the loss and the (optional) output activation $h_\mathbf{o}(\cdot)$!*

- Binary classification:
    - Likelihood: $\prod_{i=1}^{n} p(y_i | \boldsymbol{x}, \boldsymbol{\theta})$ with $\Pr\{y_i = 1 | \boldsymbol{x}\} = \dfrac{1}{1 + e^{-z_{\mathbf{o},i}}}$ (logistic)
    - Loss: $\mathcal{J}(\boldsymbol{\theta}) = \sum_{i=1}^{n} \left( \ln[1 + e^{z_{\mathbf{o},i}}] - y_i z_{\mathbf{o},i} \right)$ (binary cross entropy)
    - Or $\mathcal{J}(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left( y_i \ln a_{\mathbf{o},i} + [1 - y_i] \ln[1 - a_{\mathbf{o},i}] \right)$ if $a_{\mathbf{o},i} = h_\mathbf{o}(z_{\mathbf{o},i}) = \frac{1}{1 + e^{-z_{\mathbf{o},i}}}$

- Multiclass classification with $K > 2$ classes:
    - Likelihood: $\prod_{i=1}^{n} p(y_i | \boldsymbol{x}, \boldsymbol{\theta})$ with $\Pr\{y_i = k | \boldsymbol{x}\} = \dfrac{e^{z_{\mathbf{o},ik}}}{\sum_{l=1}^{K} e^{z_{\mathbf{o},il}}}$ (softmax)
    - Loss: $\mathcal{J}(\boldsymbol{\theta}) = \sum_{i=1}^{n} \left( \ln \sum_{k=1}^{K} e^{z_{\mathbf{o},ik}} - \sum_{k=1}^{K} y_{ik} z_{\mathbf{o},ik} \right)$ (cross entropy)
    - Or $\mathcal{J}(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \sum_{k=1}^{K} y_{ik} \ln a_{\mathbf{o},ik}$ if $a_{\mathbf{o},ik} = [\boldsymbol{h_\mathbf{o}}(\boldsymbol{z_{\mathbf{o},i}})]_k = \frac{e^{z_{\mathbf{o},ik}}}{\sum_{l=1}^{K} e^{z_{\mathbf{o},il}}}$

- Regression with $K$-dimensional targets:
    - Likelihood: $p(\boldsymbol{y} | \boldsymbol{x}) = \mathcal{N}(\boldsymbol{y}; \boldsymbol{z_o}, \sigma_\epsilon^2 \boldsymbol{I})$ (additive white Gaussian noise)
    - Loss: $\mathcal{J}(\boldsymbol{\theta}) = \sum_{i=1}^{n} \sum_{k=1}^{K} (y_{ik} - z_{\mathbf{o},ik})^2$ (quadratic)

# Outline

- Motivating Example: Learning a Feature Transformation

- Feed-Forward Neural Networks

- Training via Stochastic Gradient Descent

- Gradient Computation via Back-Propagation

- Implementing and Training Neural Nets with PyTorch

## Gradient descent

- Goal: Find model parameters $\boldsymbol{\theta}$ that minimize the loss function $\mathcal{J}(\boldsymbol{\theta})$:

$$\widehat{\boldsymbol{\theta}}_{\mathsf{ml}} = \arg\min_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) \quad \text{for} \quad \mathcal{J}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} J(\boldsymbol{\theta}, \boldsymbol{x}_i, \boldsymbol{y}_i),$$

where $J(\boldsymbol{\theta}, \boldsymbol{x}_i, \boldsymbol{y}_i)$ is the contribution from training sample $i$

- Can tackle this using the gradient descent (GD) algorithm:

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha_k \nabla \mathcal{J}(\boldsymbol{\theta}^k)$$
$$= \boldsymbol{\theta}^k - \frac{\alpha_k}{n} \sum_{i=1}^{n} \nabla J(\boldsymbol{\theta}^k, \boldsymbol{x}_i, \boldsymbol{y}_i)$$

  - Each iteration computes $n$ gradients
  - This is very expensive when $n$ (# training samples) is large!

# Stochastic gradient descent using mini-batches

- Main idea: In each step . . .
  - select a random mini-batch
  - evaluate gradient on mini-batch

- Details: At step $t = 1 \ldots T$:
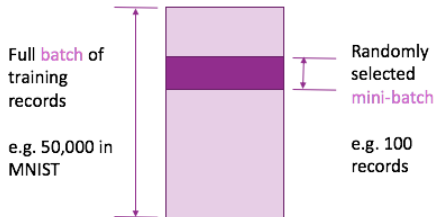  - randomly select subset of indices:
    $$I_t \subset \{1, \ldots, n\}$$
  - compute approximate gradient:
    $$g^t \triangleq \frac{1}{|I_t|} \sum_{i \in I_t} \nabla J(\theta^t, x_i, y_i)$$
  - update parameters:
    $$\theta^{t+1} = \theta^t - \alpha_t g^t$$

Full batch of training records

e.g. 50,000 in MNIST

Randomly selected mini-batch

e.g. 100 records

- Called stochastic gradient descent (SGD) because $g^t$ is a random approximation of the true gradient

# Justification for stochastic gradient descent

- The gradient approximation $\boldsymbol{g}^t$ is correct on average (i.e., unbiased):

$$\mathbb{E}\left\{\boldsymbol{g}^t \,\big|\, \boldsymbol{\theta}^t\right\} = \frac{1}{n}\sum_{i=1}^{n}\nabla J(\boldsymbol{\theta}^t, \boldsymbol{x}_i, \boldsymbol{y}_i) = \nabla\mathcal{J}(\boldsymbol{\theta}^t)$$

- Thus the SGD-updated parameters are also correct on average (i.e., unbiased):

$$\mathbb{E}\left\{\boldsymbol{\theta}^{t+1} \,\big|\, \boldsymbol{\theta}^t\right\} = \mathbb{E}\left\{\boldsymbol{\theta}^t - \alpha_t\boldsymbol{g}^t \,\big|\, \boldsymbol{\theta}^t\right\}$$
$$= \boldsymbol{\theta}^t - \alpha_t\nabla\mathcal{J}(\boldsymbol{\theta}^t)$$

- We can think of $\boldsymbol{g}^t$ as having "gradient noise" $\boldsymbol{\epsilon}^t$ that is zero mean:

$$\boldsymbol{g}^t = \nabla\mathcal{J}(\boldsymbol{\theta}^t) + \boldsymbol{\epsilon}^t \quad\text{with}\quad \mathbb{E}\{\boldsymbol{\epsilon}^t \,|\, \boldsymbol{\theta}^t\} = \boldsymbol{0}$$

  - This noise makes the SGD trajectory more noisy than the true GD trajectory
  - Can be mitigated by reducing the step-size $\alpha_t$

- For more details, see
  https://en.wikipedia.org/wiki/Stochastic_gradient_descent

# Implementing mini-batch

- Setup:
  - $n$ samples in the training dataset
  - $B$ samples in each (non-overlapping) mini-batch
  - $T = n/B$ mini-batches total
  - 1 SGD update step per mini-batch
  - $T$ updates in each "training epoch"
    - Each of the $n$ training samples is used once per training epoch

- Implementation: *In each epoch...*
  - Randomly shuffle all $n$ training indices $\{i\}$
    - This way, the mini-batches will change over the epochs!
  - Partition the shuffled indices $\{i\}$ into $T$ contiguous subsets $\{I_t\}_{t=1}^T$
  - Use subset $I_t$ to compute the gradient and cost in SGD update step $t$

# Outline

-

-

-

-

-

## Computing the gradient

- The network parameters $\boldsymbol{\theta} = (\boldsymbol{W_{\text{H}}}, \boldsymbol{b_{\text{H}}}, \boldsymbol{W_{\text{o}}}, \boldsymbol{b_{\text{o}}})$ must be trained

- To do this, we described the mini-batch SGD update

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \frac{\alpha_t}{|I_t|} \nabla \underbrace{\sum_{i \in I_t} J(\boldsymbol{\theta}^t; \boldsymbol{x}_i, \boldsymbol{y}_i)}_{\triangleq \, \mathcal{J}^t(\boldsymbol{\theta}^t)}$$

- But how do we compute the gradient?

  - Recall: the gradient is the partial derivative w.r.t. each parameter in $\boldsymbol{\theta}$
  - So we need to compute

  $$\frac{\partial \mathcal{J}^t(\boldsymbol{\theta}^t)}{\partial w_{\text{H},lj}} \; \forall l, j, \qquad \frac{\partial \mathcal{J}^t(\boldsymbol{\theta}^t)}{\partial b_{\text{H},l}} \; \forall l, \qquad \frac{\partial \mathcal{J}^t(\boldsymbol{\theta}^t)}{\partial w_{\text{o},kl}} \; \forall k, l, \qquad \frac{\partial \mathcal{J}^t(\boldsymbol{\theta}^t)}{\partial b_{\text{o},k}} \; \forall k$$

  - Going forward, we simplify our notation by dropping the batch index "$t$"

## The computation graph



- observed variable
- trainable variable
- derived variable

- The computation graph will help us to organize our computations

- The loss $\mathcal{J}$ can be computed using a forward pass through the graph:

$$z_{\mathsf{H},i} = W_{\mathsf{H}} x_i + b_{\mathsf{H}}, \quad \forall i = 1 \ldots B$$
$$a_{\mathsf{H},i} = h_{\mathsf{H}}(z_{\mathsf{H},i}), \quad \forall i = 1 \ldots B$$
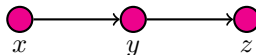$$z_{\mathsf{o},i} = W_{\mathsf{o}} a_{\mathsf{H},i} + b_{\mathsf{o}}, \quad \forall i = 1 \ldots B$$
$$\mathcal{J} = \sum_{i=1}^{B} J(z_{\mathsf{o},i}; y_i)$$

- The gradients can then be computed using a backward pass, as described next

- Note: we write the loss $J(\cdot, y_i)$ in terms of $z_{\mathsf{o},i}$, not $a_{\mathsf{o},i}$
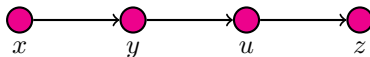  - PyTorch allows either, as discussed on page 44

# Review of the chain rule

- Suppose that $x, y, z$ are scalars. Then the chain rule says

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y}\frac{\partial y}{\partial x}$$

- Now suppose that another variable $u$ is inserted. Previous expression still holds!
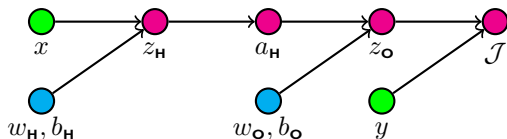
$$\frac{\partial z}{\partial x} = \underbrace{\frac{\partial z}{\partial u}\frac{\partial u}{\partial y}}_{\frac{\partial z}{\partial y}}\frac{\partial y}{\partial x}$$

- Now suppose that there are several intermediate variables $\{y_i\}_{i=1}^B$. The multivariable chain rule says

$$\frac{\partial z}{\partial x} = \sum_{i=1}^B \frac{\partial z}{\partial y_i}\frac{\partial y_i}{\partial x}$$

# Simple case 1: Scalar variables and batch size $B = 1$



$$z_\mathsf{H} = w_\mathsf{H} x + b_\mathsf{H}$$
$$a_\mathsf{H} = h_\mathsf{H}(z_\mathsf{H})$$
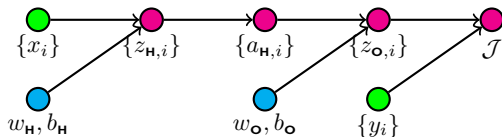$$z_\mathsf{O} = w_\mathsf{O} a_\mathsf{H} + b_\mathsf{O}$$
$$\mathcal{J} = J(z_\mathsf{O}; y)$$

Gradients follow from the chain rule. Previous computations can be reused!

- grad w.r.t. variables :
$$\frac{\partial \mathcal{J}}{\partial z_\mathsf{O}} = J'(z_\mathsf{O}), \qquad \frac{\partial \mathcal{J}}{\partial a_\mathsf{H}} = \frac{\partial \mathcal{J}}{\partial z_\mathsf{O}} \underbrace{\frac{\partial z_\mathsf{O}}{\partial a_\mathsf{H}}}_{=w_\mathsf{O}}, \qquad \frac{\partial \mathcal{J}}{\partial z_\mathsf{H}} = \frac{\partial \mathcal{J}}{\partial a_\mathsf{H}} \underbrace{\frac{\partial a_\mathsf{H}}{\partial z_\mathsf{H}}}_{=h'_\mathsf{H}(z_\mathsf{H})}$$

- grad w.r.t. parameters :
$$\frac{\partial \mathcal{J}}{\partial b_\mathsf{O}} = \frac{\partial \mathcal{J}}{\partial z_\mathsf{O}} \underbrace{\frac{\partial z_\mathsf{O}}{\partial b_\mathsf{O}}}_{=1}, \qquad \frac{\partial \mathcal{J}}{\partial w_\mathsf{O}} = \frac{\partial \mathcal{J}}{\partial z_\mathsf{O}} \underbrace{\frac{\partial z_\mathsf{O}}{\partial w_\mathsf{O}}}_{=a_\mathsf{H}} \qquad \text{\color{red}{can reuse}}\ \partial \mathcal{J}/\partial z_\mathsf{O}!$$

$$\frac{\partial \mathcal{J}}{\partial b_\mathsf{H}} = \frac{\partial \mathcal{J}}{\partial z_\mathsf{H}} \underbrace{\frac{\partial z_\mathsf{H}}{\partial b_\mathsf{H}}}_{=1}, \qquad \frac{\partial \mathcal{J}}{\partial w_\mathsf{H}} = \frac{\partial \mathcal{J}}{\partial z_\mathsf{H}} \underbrace{\frac{\partial z_\mathsf{H}}{\partial w_\mathsf{H}}}_{=x} \qquad \text{\color{red}{can reuse}}\ \partial \mathcal{J}/\partial z_\mathsf{H}!$$

# Simple case 2: Scalar variables and batch size $B > 1$



$$z_{\mathsf{H},i} = w_{\mathsf{H}} x_i + b_{\mathsf{H}}, \ \forall i$$
$$a_{\mathsf{H},i} = h_{\mathsf{H}}(z_{\mathsf{H},i}), \ \forall i$$
$$z_{\mathsf{o},i} = w_{\mathsf{o}} a_{\mathsf{H},i} + b_{\mathsf{o}}, \ \forall i$$
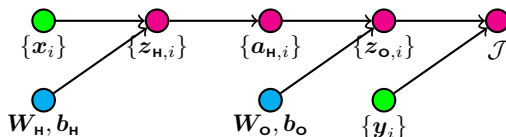$$\mathcal{J} = \sum_i J(z_{\mathsf{o},i}; y_i)$$

Now we need the multivariable chain rule for the parameter gradients:

- grad w.r.t. variables: $\dfrac{\partial \mathcal{J}}{\partial z_{\mathsf{o},i}} = J'(z_{\mathsf{o},i}), \quad \dfrac{\partial \mathcal{J}}{\partial a_{\mathsf{H},i}} = \dfrac{\partial \mathcal{J}}{\partial z_{\mathsf{o},i}} \underbrace{\dfrac{\partial z_{\mathsf{o},i}}{\partial a_{\mathsf{H},i}}}_{=w_{\mathsf{o}}}, \quad \dfrac{\partial \mathcal{J}}{\partial z_{\mathsf{H},i}} = \dfrac{\partial \mathcal{J}}{\partial a_{\mathsf{H},i}} \underbrace{\dfrac{\partial a_{\mathsf{H},i}}{\partial z_{\mathsf{H},i}}}_{=h_{\mathsf{H}}'(z_{\mathsf{H},i})}$

- grad w.r.t. parameters:

$$\frac{\partial \mathcal{J}}{\partial b_{\mathsf{o}}} = \sum_i \frac{\partial \mathcal{J}}{\partial z_{\mathsf{o},i}} \underbrace{\frac{\partial z_{\mathsf{o},i}}{\partial b_{\mathsf{o}}}}_{=1}, \qquad \frac{\partial \mathcal{J}}{\partial w_{\mathsf{o}}} = \sum_i \frac{\partial \mathcal{J}}{\partial z_{\mathsf{o},i}} \underbrace{\frac{\partial z_{\mathsf{o},i}}{\partial w_{\mathsf{o}}}}_{=a_{\mathsf{H},i}}$$

$$\frac{\partial \mathcal{J}}{\partial b_{\mathsf{H}}} = \sum_i \frac{\partial \mathcal{J}}{\partial z_{\mathsf{H},i}} \underbrace{\frac{\partial z_{\mathsf{H},i}}{\partial b_{\mathsf{H}}}}_{=1}, \qquad \frac{\partial \mathcal{J}}{\partial w_{\mathsf{H}}} = \sum_i \frac{\partial \mathcal{J}}{\partial z_{\mathsf{H},i}} \underbrace{\frac{\partial z_{\mathsf{H},i}}{\partial w_{\mathsf{H}}}}_{=x_i}$$

## Practical case: Vector variables and batch size $B > 1$



$$z_{\mathsf{H},i} = W_{\mathsf{H}} x_i + b_{\mathsf{H}}, \ \forall i$$
$$a_{\mathsf{H},i} = h_{\mathsf{H}}(z_{\mathsf{H},i}), \ \forall i$$
$$z_{\mathsf{O},i} = W_{\mathsf{O}} a_{\mathsf{H},i} + b_{\mathsf{O}}, \ \forall i$$
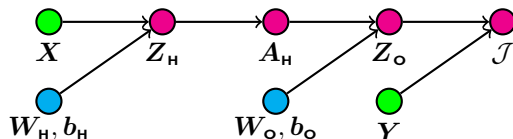$$\mathcal{J} = \sum_i J(z_{\mathsf{O},i}; y_i)$$

Now we also use the multivariable chain rule for gradient w.r.t. $a_{\mathsf{H},il}$:

- grad wrt variables :
$$\frac{\partial \mathcal{J}}{\partial z_{\mathsf{O},ik}} = \frac{\partial J(z_{\mathsf{O},i})}{\partial z_{\mathsf{O},ik}}, \quad \frac{\partial \mathcal{J}}{\partial a_{\mathsf{H},il}} = \sum_k \frac{\partial \mathcal{J}}{\partial z_{\mathsf{O},ik}} \underbrace{\frac{\partial z_{\mathsf{O},ik}}{\partial a_{\mathsf{H},il}}}_{=w_{\mathsf{O},kl}}, \quad \frac{\partial \mathcal{J}}{\partial z_{\mathsf{H},il}} = \frac{\partial \mathcal{J}}{\partial a_{\mathsf{H},il}} \underbrace{\frac{\partial a_{\mathsf{H},il}}{\partial z_{\mathsf{H},il}}}_{=h'_{\mathsf{H}}(z_{\mathsf{H},il})}$$

- grad wrt params :
$$\frac{\partial \mathcal{J}}{\partial b_{\mathsf{O},k}} = \sum_i \frac{\partial \mathcal{J}}{\partial z_{\mathsf{O},ik}} \underbrace{\frac{\partial z_{\mathsf{O},ik}}{\partial b_{\mathsf{O},k}}}_{=1}, \quad \frac{\partial \mathcal{J}}{\partial w_{\mathsf{O},kl}} = \sum_i \frac{\partial \mathcal{J}}{\partial z_{\mathsf{O},ik}} \underbrace{\frac{\partial z_{\mathsf{O},ik}}{\partial w_{\mathsf{O},kl}}}_{=a_{\mathsf{H},il}}$$

$$\frac{\partial \mathcal{J}}{\partial b_{\mathsf{H},l}} = \sum_i \frac{\partial \mathcal{J}}{\partial z_{\mathsf{H},il}} \underbrace{\frac{\partial z_{\mathsf{H},il}}{\partial b_{\mathsf{H},l}}}_{=1}, \quad \frac{\partial \mathcal{J}}{\partial w_{\mathsf{H},lj}} = \sum_i \frac{\partial \mathcal{J}}{\partial z_{\mathsf{H},il}} \underbrace{\frac{\partial z_{\mathsf{H},il}}{\partial w_{\mathsf{H},lj}}}_{=x_{ij}}$$

## Practical case: Matrix/vector formulation



$$Z_{\mathsf{H}}^{\mathsf{T}} = W_{\mathsf{H}} X^{\mathsf{T}} + b_{\mathsf{H}} \mathbf{1}_B^{\mathsf{T}}$$
$$A_{\mathsf{H}} = h_{\mathsf{H}}(Z_{\mathsf{H}})$$
$$Z_{\mathsf{O}}^{\mathsf{T}} = W_{\mathsf{O}} A_{\mathsf{H}}^{\mathsf{T}} + b_{\mathsf{O}} \mathbf{1}_B^{\mathsf{T}}$$
$$\mathcal{J} = \sum_{i=1}^{B} J(z_{\mathsf{O},i}; y_i)$$

- Define $X \triangleq [x_1, \ldots, x_B]^{\mathsf{T}}$, $\quad Z_{\mathsf{H}} \triangleq [z_{\mathsf{H},1}, \ldots, z_{\mathsf{H},B}]^{\mathsf{T}}$, $\quad A_{\mathsf{H}} \triangleq [a_{\mathsf{H},1}, \ldots, a_{\mathsf{H},B}]^{\mathsf{T}}$,
  $Y \triangleq [y_1, \ldots, y_B]^{\mathsf{T}}$, $\quad Z_{\mathsf{O}} \triangleq [z_{\mathsf{O},1}, \ldots, z_{\mathsf{O},B}]^{\mathsf{T}}$

- grad wrt variables : $\quad \left[\dfrac{\partial \mathcal{J}}{\partial Z_{\mathsf{O}}}\right]_{ik} = \dfrac{\partial J(z_{\mathsf{O},i})}{\partial z_{\mathsf{O},ik}}, \qquad \dfrac{\partial \mathcal{J}}{\partial A_{\mathsf{H}}} = \dfrac{\partial \mathcal{J}}{\partial Z_{\mathsf{O}}} W_{\mathsf{O}}, \qquad \dfrac{\partial \mathcal{J}}{\partial Z_{\mathsf{H}}} = \dfrac{\partial \mathcal{J}}{\partial A_{\mathsf{H}}} \odot h_{\mathsf{H}}'(Z_{\mathsf{H}})$

- grad wrt params : $\quad \dfrac{\partial \mathcal{J}}{\partial b_{\mathsf{O}}} = \left(\dfrac{\partial \mathcal{J}}{\partial Z_{\mathsf{O}}}\right)^{\mathsf{T}} \mathbf{1}_B \in \mathbb{R}^{d_{\mathsf{O}}}, \qquad \dfrac{\partial \mathcal{J}}{\partial W_{\mathsf{O}}} = \left(\dfrac{\partial \mathcal{J}}{\partial Z_{\mathsf{O}}}\right)^{\mathsf{T}} A_{\mathsf{H}} \in \mathbb{R}^{d_{\mathsf{O}} \times d_{\mathsf{H}}}$

  $\dfrac{\partial \mathcal{J}}{\partial b_{\mathsf{H}}} = \left(\dfrac{\partial \mathcal{J}}{\partial Z_{\mathsf{H}}}\right)^{\mathsf{T}} \mathbf{1}_B \in \mathbb{R}^{d_{\mathsf{H}}}, \qquad \dfrac{\partial \mathcal{J}}{\partial W_{\mathsf{H}}} = \left(\dfrac{\partial \mathcal{J}}{\partial Z_{\mathsf{H}}}\right)^{\mathsf{T}} X \in \mathbb{R}^{d_{\mathsf{H}} \times d}$

- Above, $\odot$ denotes the elementwise or "Hadamard" product

## Summary of forward and backward passes

So, to compute the gradients w.r.t. the parameters $\boldsymbol{\theta} = (\boldsymbol{W_{\mathsf{H}}}, \boldsymbol{b_{\mathsf{H}}}, \boldsymbol{W_{\mathsf{o}}}, \boldsymbol{b_{\mathsf{o}}})$, we . . .

- first perform the forward pass: $\quad \boldsymbol{Z_{\mathsf{H}}^{\mathsf{T}}} = \boldsymbol{W_{\mathsf{H}}} \boldsymbol{X^{\mathsf{T}}} + \boldsymbol{b_{\mathsf{H}}} \boldsymbol{1}_B^{\mathsf{T}}$

$$\boldsymbol{A_{\mathsf{H}}} = h_{\mathsf{H}}(\boldsymbol{Z_{\mathsf{H}}})$$
$$\boldsymbol{Z_{\mathsf{o}}^{\mathsf{T}}} = \boldsymbol{W_{\mathsf{o}}} \boldsymbol{A_{\mathsf{H}}^{\mathsf{T}}} + \boldsymbol{b_{\mathsf{o}}} \boldsymbol{1}_B^{\mathsf{T}}$$
$$\mathcal{J} = \sum_i J(\boldsymbol{z_{\mathsf{o}},i}; \boldsymbol{y}_i)$$

- then the backward pass:

$$\left[\frac{\partial \mathcal{J}}{\partial \boldsymbol{Z_{\mathsf{o}}}}\right]_{ik} = \frac{\partial J(\boldsymbol{z_{\mathsf{o}},i}; \boldsymbol{y}_i)}{\partial z_{\mathsf{o},ik}} \ \forall i = 1 \ldots B, \ k = 1 \ldots d_{\mathsf{o}}$$

$$\frac{\partial \mathcal{J}}{\partial \boldsymbol{b_{\mathsf{o}}}} = \left(\frac{\partial \mathcal{J}}{\partial \boldsymbol{Z_{\mathsf{o}}}}\right)^{\mathsf{T}} \boldsymbol{1}_B \ \text{ and } \ \frac{\partial \mathcal{J}}{\partial \boldsymbol{W_{\mathsf{o}}}} = \left(\frac{\partial \mathcal{J}}{\partial \boldsymbol{Z_{\mathsf{o}}}}\right)^{\mathsf{T}} \boldsymbol{A_{\mathsf{H}}}$$

$$\frac{\partial \mathcal{J}}{\partial \boldsymbol{Z_{\mathsf{H}}}} = \left(\frac{\partial \mathcal{J}}{\partial \boldsymbol{Z_{\mathsf{o}}}} \boldsymbol{W_{\mathsf{o}}}\right) \odot h_{\mathsf{H}}'(\boldsymbol{Z_{\mathsf{H}}})$$

$$\frac{\partial \mathcal{J}}{\partial \boldsymbol{b_{\mathsf{H}}}} = \left(\frac{\partial \mathcal{J}}{\partial \boldsymbol{Z_{\mathsf{H}}}}\right)^{\mathsf{T}} \boldsymbol{I}_B \ \text{ and } \ \frac{\partial \mathcal{J}}{\partial \boldsymbol{W_{\mathsf{H}}}} = \left(\frac{\partial \mathcal{J}}{\partial \boldsymbol{Z_{\mathsf{H}}}}\right)^{\mathsf{T}} \boldsymbol{X}$$

Called "back-propagation," since gradient computations work *backwards* from end!

## Examples of loss derivative

- With binary cross-entropy loss, we have $d_{\mathbf{o}} = 1$ and

$$\mathcal{J} = \sum_{i=1}^{B} J(z_{\mathbf{o},i}, y_i) \quad \text{with} \quad y_i \in \{0, 1\}$$
$$\text{where} \quad J(z_{\mathbf{o},i}, y_i) = \ln\left(1 + e^{z_{\mathbf{o},i}}\right) - y_i z_{\mathbf{o},i}$$
$$\Rightarrow \quad \frac{\partial J(z_{\mathbf{o},i}, y_i)}{\partial z_{\mathbf{o},i}} = \frac{e^{z_{\mathbf{o},i}}}{1 + e^{z_{\mathbf{o},i}}} - y_i$$

- With $K$-ary cross-entropy loss, we have $d_{\mathbf{o}} = K$ and

$$\mathcal{J} = \sum_{i=1}^{B} J(\boldsymbol{z}_{\mathbf{o},i}, \boldsymbol{y}_i) \quad \text{with} \quad y_{ik} = \begin{cases} 1 & \text{if } y_i = k \\ 0 & \text{if } y_i \neq k \end{cases}$$
$$\text{where} \quad J(\boldsymbol{z}_{\mathbf{o},i}, \boldsymbol{y}_i) = \ln\left(\sum_{l=1}^{K} e^{z_{\mathbf{o},il}}\right) - \sum_{l=1}^{K} y_{il} z_{\mathbf{o},il}$$
$$\Rightarrow \quad \frac{\partial J(\boldsymbol{z}_{\mathbf{o},i}, \boldsymbol{y}_i)}{\partial z_{\mathbf{o},ik}} = \frac{e^{z_{\mathbf{o},ik}}}{\sum_{l=1}^{K} e^{z_{\mathbf{o},il}}} - y_{ik}$$

# Outline

- Motivating Example: Learning a Feature Transformation

- Feed-Forward Neural Networks

- Training via Stochastic Gradient Descent

- Gradient Computation via Back-Propagation

- Implementing and Training Neural Nets with PyTorch

# Frameworks for implementing neural nets

The frameworks most popular today:

# Top 3 current frameworks

- TensorFlow (2015)
    - Widespread and mature for deployment
    - Original version was difficult to use. Newest version is better, but not great

- Keras (runs on top of TensorFlow or Theano or CNTK)
    - Very convenient interface for simple/standard tasks
    - Difficult to customize and debug for complex/nonstandard tasks

- PyTorch (2017)
    - Easier than TensorFlow to use and just as powerful
    - Most popular in research/academic community

We will use PyTorch! See the following for more about pros & cons:

- PyTorch vs. TensorFlow: Towards Data Science (6/20)
- PyTorch vs. TensorFlow vs. Keras: Simplilearn (9/20), The Startup (7/20)

# PyTorch recipe

1) Construct the dataset and dataloader objects

2) Construct the network
   - # hidden units, # output units, activations, etc . . .

3) Select the optimizer and loss criterion

4) Fit the network parameters

5) Apply the network

Great tutorials at https://pytorch.org/tutorials

# 1) Construct the dataset and dataloader objects

```python
import torch

import torch.utils.data

# Convert the numpy arrays to Tensor type
X_torch = torch.Tensor(X)
y_torch = torch.Tensor(y)

# Create a Dataset from the Tensors
dataset = torch.utils.data.TensorDataset(X_torch, y_torch)
# Create a DataLoader from the Dataset
loader = torch.utils.data.DataLoader(dataset, batch_size=100)
```

- `torch.Tensor`: the datatype used by PyTorch for multi-dimensional matrices

- `dataset`: the object used to hold the training data

- `DataLoader`: adds random sampling and multi-processor support to the Dataset

- Later we will draw mini-batches via:
```python
for batch, data in enumerate(loader):
    x_batch,y_batch = data
```

# 2) Construct the network

- The neural net is described by a `torch.nn.Module`
  - `__init__()` defines the network components
  - `forward()` defines one forward pass through the network

- We instantiate the neural net as "model"

```python
import torch.nn as nn

nh = 4

# nin: dimension of input data
# nh: number of hidden units
# nout: number of outputs = 1 since this is bi

class Net(nn.Module):
    def __init__(self,nin,nh,nout):
        super(Net,self).__init__()
        self.activation = nn.Sigmoid()
        self.Dense1 = nn.Linear(nin,nh)
        self.Dense2 = nn.Linear(nh,nout)

    def forward(self,x):
        x = self.activation(self.Dense1(x))
        out = self.activation(self.Dense2(x))
        return out

model = Net(nin=nx,nh=nh,nout=1)
```

# 3) Select the optimizer and loss criterion

```python
import torch.optim as optim

opt = optim.Adam(model.parameters(), lr=0.01)
criterion = nn.BCELoss()
```

- Learning algorithms are stored as "optimizer objects" in `torch.optim`
  - Many options, e.g., SGD, Rprop, RMSprop, AdaDelta, AdaGrad, Adam, etc.
  - Some descriptions can be found here, here, and here

- We instantiated the Adam optimizer as "opt" using . . .
  - the network parameters $\theta$ that we want to optimize
  - Adam's algorithmic parameters (i.e., learning rate)

- The `torch.nn` library includes many loss criteria
  - Examples: BCEloss, CrossEntropyLoss, MSELoss (i.e., RSS), L1Loss, etc.
  - These can be scaled and combined (e.g., "MSELoss $+ \lambda$ L1Loss" for Lasso)

# 4) Fit the network parameters

For each mini-batch $\{x_i\}$:

- compute outputs $\{a_{o,i}\}$

- compute loss $\mathcal{J}(\theta^t)$

- compute gradient $\nabla \mathcal{J}(\theta^t)$

- update parameters $\theta^{t+1}$

- record loss & accuracy

```python
num_epoch = 2000

print_intvl = 100
avg_loss = np.zeros([num_epoch])
avg_acc = np.zeros([num_epoch])

# Outer loop over epochs
for epoch in range(num_epoch):
    error = 0 # initialize error counter
    total = 0 # initialize total counter
    batch_loss = []
    # Inner loop over mini-batches
    for batch, data in enumerate(loader):
        x_batch,y_batch = data
        y_batch = y_batch.view(-1,1) # resizes y_batch to (batch_size,1)
        out = model(x_batch)
        # Compute loss
        loss = criterion(out,y_batch)
        batch_loss.append(loss.item())
        # Compute gradients using back-propagation
        opt.zero_grad()
        loss.backward()
        # Take an optimization 'step' (i.e., update parameters)
        opt.step()
        # Do hard decision
        guess = out.round()
        # Compute number of decision errors
        error += torch.sum(torch.abs(guess - y_batch))
        total += len(y_batch)

    acc = 100*(1-error/total) # Compute accuracy over epoch
    avg_loss[epoch] = np.mean(batch_loss) # Compute average loss over epoch
    avg_acc[epoch] = acc
```
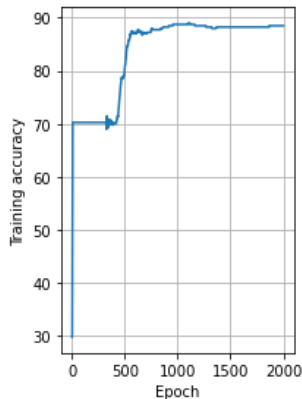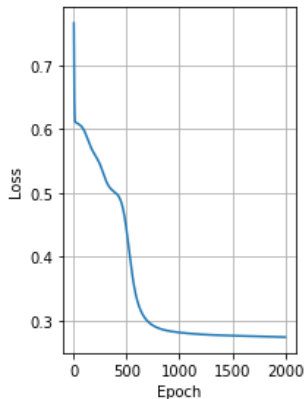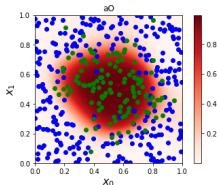
# Performance vs. epoch

Over the epochs . . .

- Loss gradually decreases

- Classification accuracy gradually increases

# 5) Applying the network: Visualizing the decision region

- `Xplot` created with rows densely sampling $x^{\mathsf{T}}$ over $[0,1]^2$

- `Xplot` recast as a `torch.Tensor`

- `model()` applies network

  - It calls `model.__call__()`, which eventually calls the `forward` method and does other book-keeping

- `.detach().numpy()` extracts the data and converts it to NumPy



```python
# Limits to plot the response.
xmin = [0,0]
xmax = [1,1]

# Use meshgrid to create the 2D input
nplot = 100
x0plot = np.linspace(xmin[0],xmax[1],nplot)
x1plot = np.linspace(xmin[0],xmax[1],nplot)
x0mat, x1mat = np.meshgrid(x0plot,x1plot)
Xplot = np.column_stack([x0mat.ravel(), x1mat.ravel()])
Xplot_tensor = torch.Tensor(Xplot)

# Compute the output and export to numpy
a0plot = model(Xplot_tensor).detach().numpy()
a0plot_mat = a0plot[:,0].reshape((nplot, nplot))

# Plot the recovered region
plt.imshow(np.flipud(a0plot_mat), extent=[xmin[0],xmax[0]
plt.colorbar()

# Overlay the samples
I0 = np.where(y==0)[0]
I1 = np.where(y==1)[0]
plt.plot(X[I0,0], X[I0,1], 'bo')
plt.plot(X[I1,0], X[I1,1], 'go')
plt.xlabel('$x_0$', fontsize=16)
plt.ylabel('$x_1$', fontsize=16)
plt.title('a0');
```
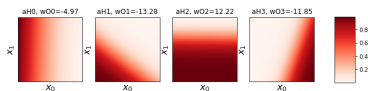
# 5) Applying the network: Visualizing the hidden layer

- `model.Dense1()` applies the first linear stage

- `model.activation()` applies the activation

- `.detach().numpy()` extracts the data and converts it to NumPy

- `model.state_dict()` exports the model parameters

```python
# Get the outputs of the hidden layer
ahid = model.activation(model.Dense1(Xplot_tensor))
ahid_plot = ahid.detach().numpy()
ahid_plot = ahid_plot.reshape((nplot,nplot,nh))

# Get the weights of the output layer
state_dict = model.state_dict()
Wo, bo = state_dict['Dense2.weight'], state_dict['Dense2.bias

fig = plt.figure(figsize=(10, 4))

for i in range(nh):
    plt.subplot(1,nh,i+1)
    ahid_ploti = np.flipud(ahid_plot[:,:,i])
    im = plt.imshow(ahid_ploti, extent=[xmin[0],xmax[0],xmin[
    plt.xticks([])
    plt.yticks([])
    plt.xlabel('$x_0$', fontsize=16)
    plt.ylabel('$x_1$', fontsize=16)
    plt.title('aH{0:d}, wO{0:d}={1:4.2f}'.format(i,Wo[0,i]))

fig.subplots_adjust(right=0.85)
cbar_ax = fig.add_axes([0.9, 0.30, 0.05, 0.4])
fig.colorbar(im, cax=cbar_ax);
```
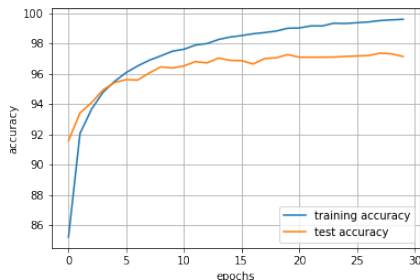
# Application to MNIST

- In a second demo, we use a 2-layer neural network for MNIST digit classification

- For this, we used
  - $d_{\mathbf{H}} = 100$ hidden nodes
  - $d_{\mathbf{O}} = 10$ output nodes (10 classes)
  - cross-entropy loss
  - sigmoidal activation fxn $h_{\mathbf{H}}(\cdot)$
  - the `Adam` optimizer with $\alpha = 10^{-3}$
  - 50,000 training samples
  - 10,000 test samples



- We achieved a test accuracy of $\approx 97\%$. Not bad!
  - Accuracy is similar to SVM. But neural net is much faster to apply!
  - And further fine-tuning could improve neural-net performance

# Fine-tuning the implementation

We have several decisions to make when implementing a 2-layer neural network:

- Network details
    - Number of hidden units $d_{\mathsf{H}}$
    - Type of hidden-layer activation functions $h_{\mathsf{H}}(\cdot)$

- Optimizer details
    - which optimizer (e.g., `Adam`)
    - learning rate (and how to schedule it over the epochs)
    - batch size
    - # epochs

- Regularization
    - L2 regularization can be implemented via optimizer's `weight_decay` option

Packages like Optuna can be used to tune all these hyperparameters

# PyTorch loss options: Be careful!

- Throughout these lecture notes, the loss $\mathcal{J}(\boldsymbol{\theta})$ has been defined using the output-layer linear scores $\boldsymbol{z_o}$, also called logits. For example,
  - Binary cross-entropy: $J(\boldsymbol{\theta}) = \sum_i \left( \ln[1 + e^{z_{\boldsymbol{o},i}}] - y_i z_{\boldsymbol{o},i} \right)$

- But the loss can also be written in terms of the output-layer activations $\boldsymbol{a_o}$:
  - Binary cross-entropy: $J(\boldsymbol{\theta}) = - \sum_i \left( y_i \ln a_{\boldsymbol{o},i} + [1 - y_i] \ln[1 - a_{\boldsymbol{o},i}] \right)$
    with sigmoid activation $a_{\boldsymbol{o},i} = \frac{1}{1 + e^{-z_{\boldsymbol{o},i}}}$

- PyTorch allows the loss to be defined in either way, so be careful!
  - `nn.BCEWithLogitsLoss`: takes logits $\{z_{\boldsymbol{o},i}\}$ as input
  - `nn.BCELoss`: takes sigmoid activations $\{a_{\boldsymbol{o},i}\}$ as input
  - `nn.CrossEntropyLoss`: takes logits $\{\boldsymbol{z_{o},i}\}$ as input
  - `nn.NLLLoss`: takes log-softmax activations $\{\boldsymbol{a_{o},i}\}$ as input

- For more discussion, see this blog post

# Learning objectives

- Understand 2-layer feedforward neural networks:
  - Motivation: learning feature transformations
  - Network architecture: linear and nonlinear layers
  - Choice of activation functions and training loss

- Understand mini-batch training and stochastic gradient descent

- Understand the back-propagation approach to gradient computation

- Know how to implement a neural network using PyTorch