

Gradient Descent and Training Neural Networks

Deep Learning
CS 435/635

Course Instructor: Chandresh
AI Lab Coordinator @IIT Indore

Course Content

Module I: History of Deep Learning, Sigmoid Neurons, Perceptrons, and learning algorithms. Multilayer Perceptrons (MLPs), Representation Power of MLPs,.

Module II: Feedforward Neural Networks. **Backpropagation**. first and second-order training methods. NN Training tricks, Regularization

Module III: Introduction to Autoencoders and their characteristics, relation to PCA, Regularization in autoencoders, and Types of autoencoders.

Module IV: Architecture of Convolutional Neural Networks (CNN), types of CNNs.

Module V: Architecture of Recurrent Neural Networks (RNN), Backpropagation through time. Encoder-Decoder Models, Attention Mechanism. Advanced Topics: Transformers and BERT.

Acknowledgement

- Neural Networks and Deep Learning course by Danna Gurari University of Colorado Boulder

Today's Topics

- Objective function: what to learn
- Gradient descent: how to learn
- Training a neural network: optimization
- Gradient descent for activation functions

Today's Topics

- Objective function: what to learn
- Gradient descent: how to learn
- Training a neural network: optimization
- Gradient descent for activation functions

Objective Function: Analogous to Learning...

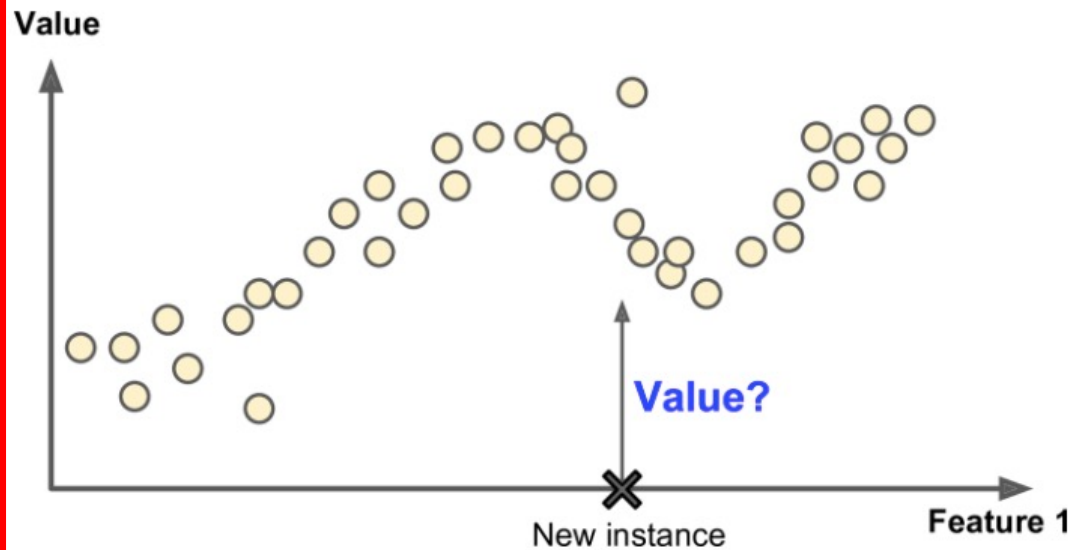
e.g., to walk



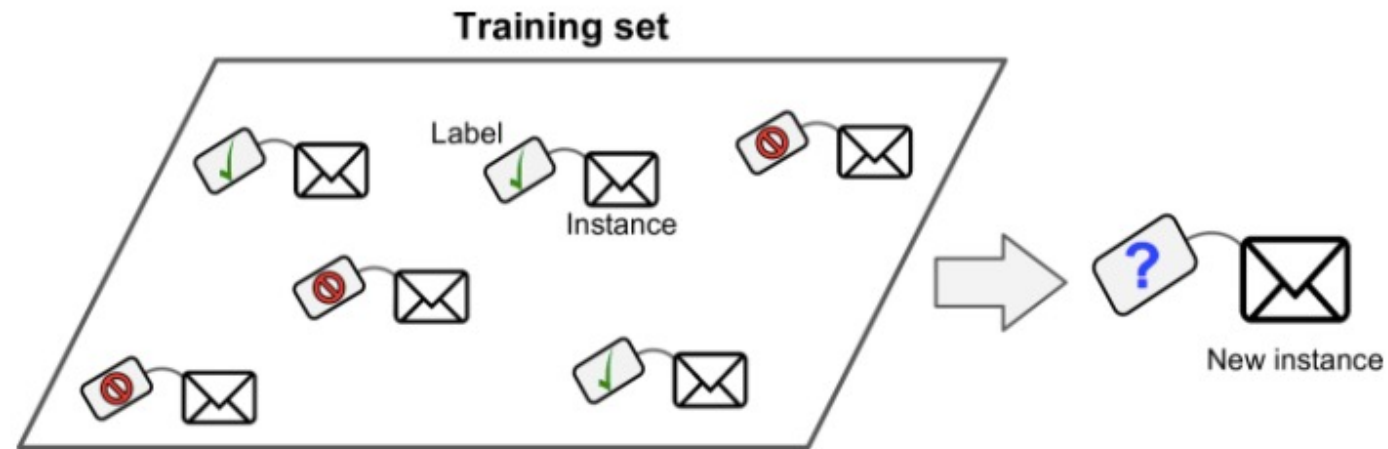
Key question: how do you measure/quantify task success?

Objective Function: Learn Model Parameters that Achieve a Specified (Measurable) Goal

Regression
(predict **continuous** value)

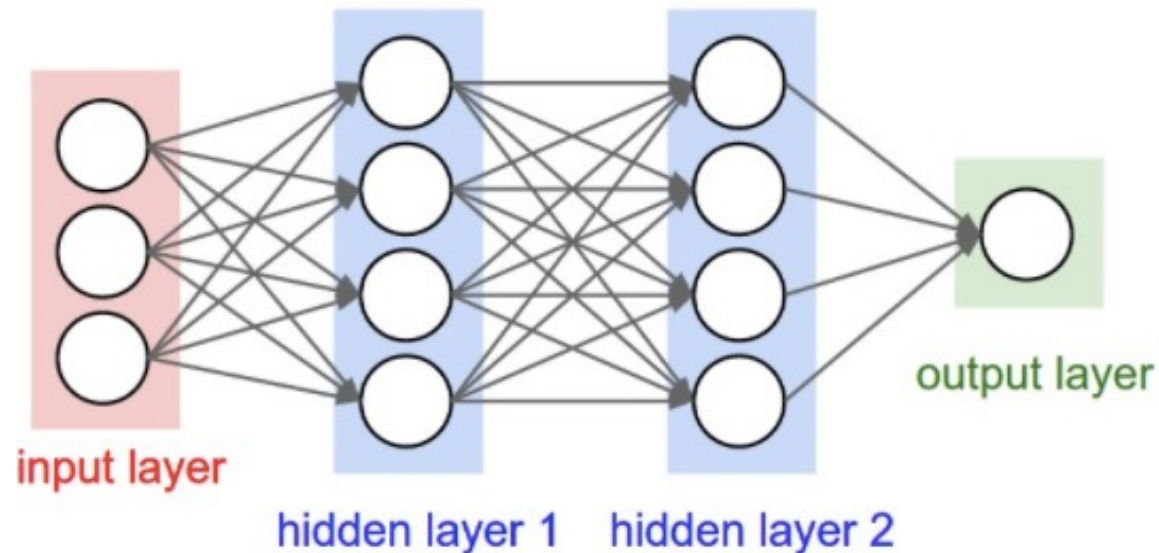


Classification
(predict **discrete** value)



Objective Function: Learn Model Parameters that Achieve a Specified (Measurable) Goal

e.g., make as small as possible the squared error (aka, L2 loss, quadratic loss)



Mean taken over n instances

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

True value

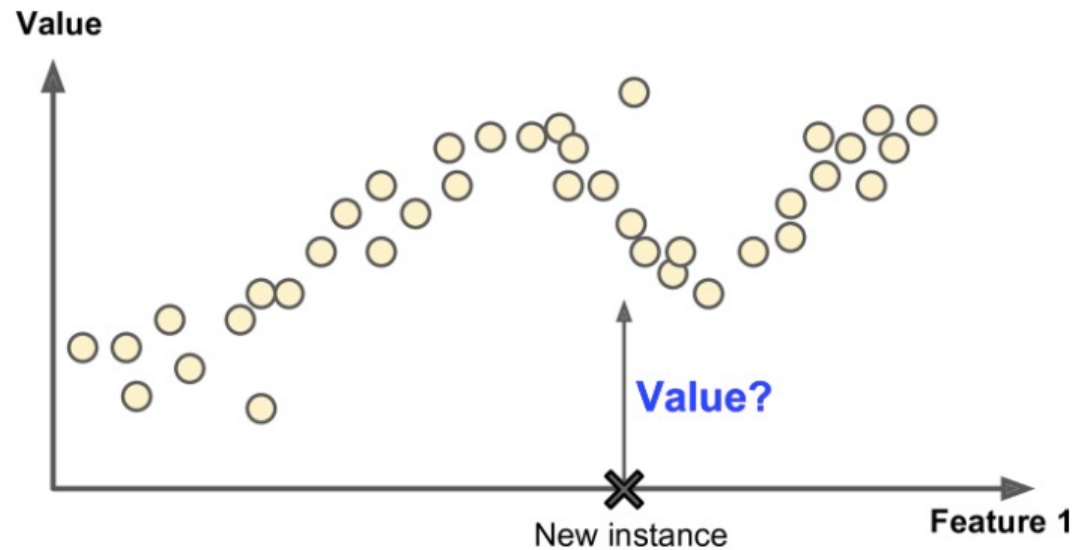
Predicted value

What is the range of possible values?

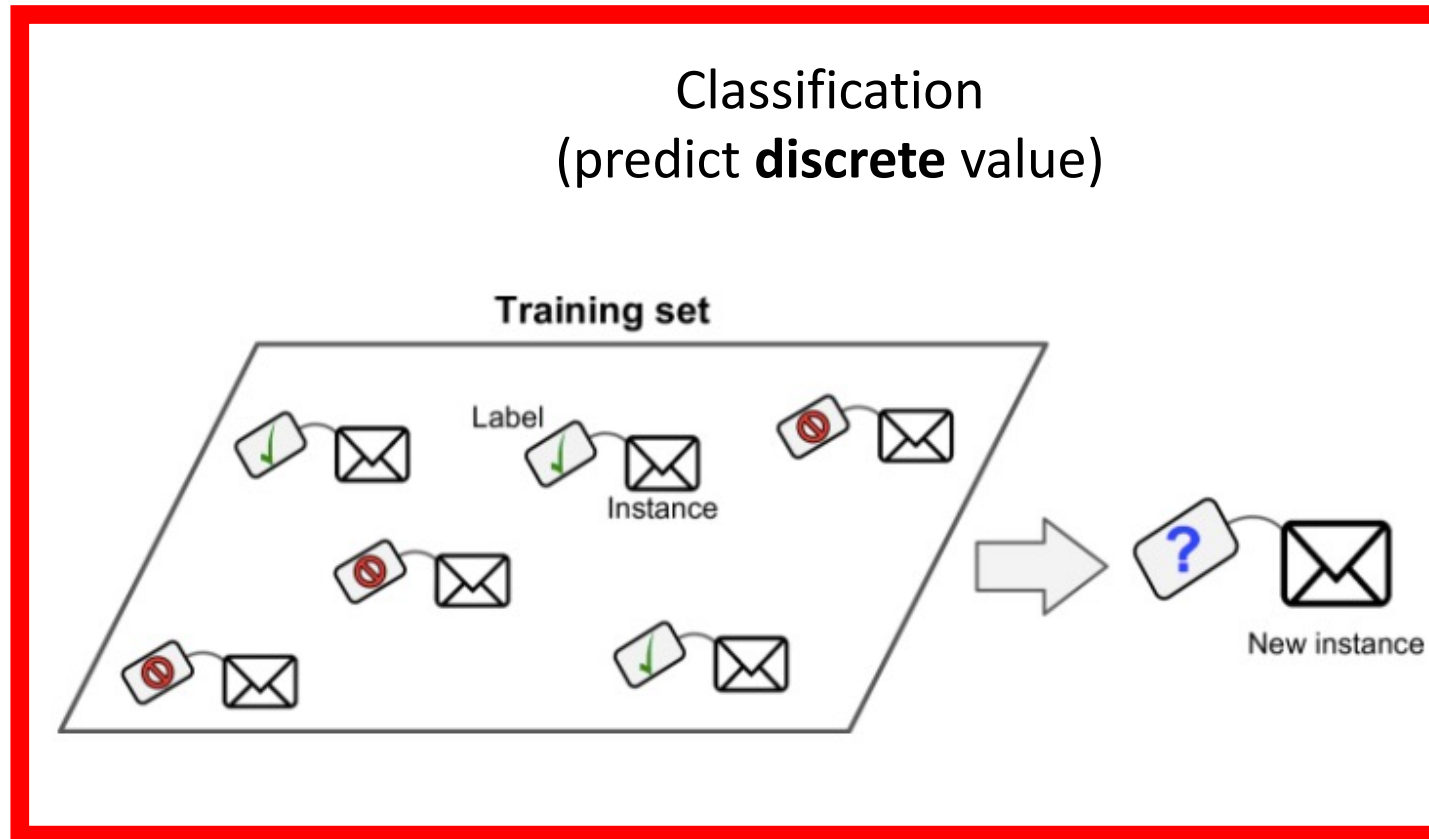
- Minimum: 0
 - i.e., all correct predictions
- Maximum: Infinity
 - i.e., incorrect predictions

Objective Function: Learn Model Parameters that Achieve a Specified (Measurable) Goal

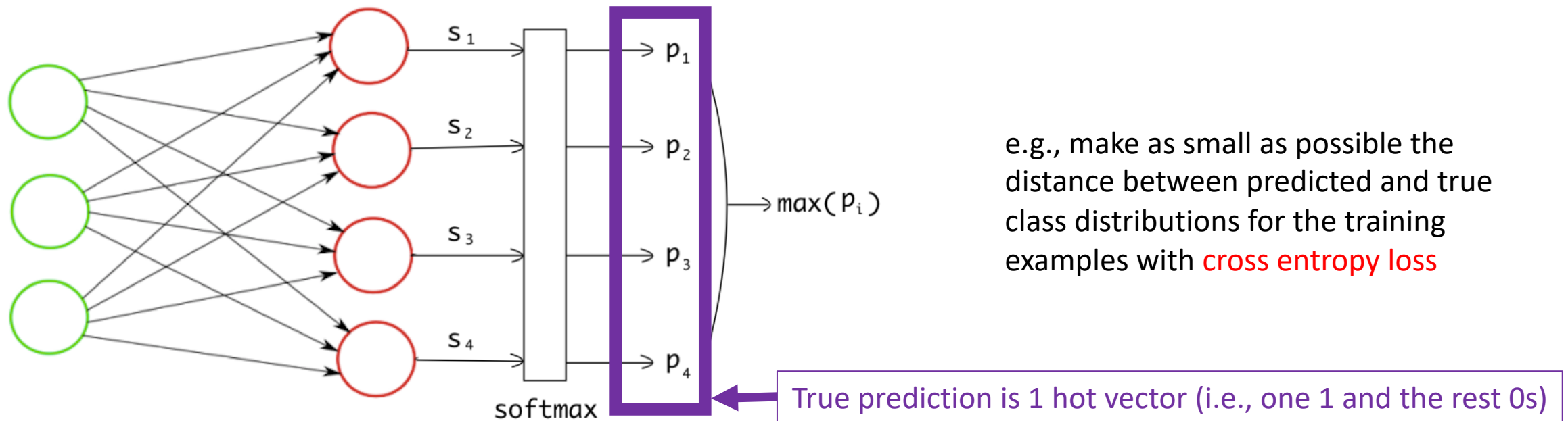
Regression
(predict **continuous** value)



Classification
(predict **discrete** value)



Objective Function: Learn Model Parameters that Achieve a Specified (Measurable) Goal



Objective Function: Learn Model Parameters that Achieve a Specified (Measurable) Goal

Probability distribution of true class

Probability distribution of predicted class

Number of classes

Recall, truth is set to 1 for one class and 0 otherwise

Observed features

Simplifies to the log of the predicted probability for the correct class (i.e., negative log likelihood loss)

$$\begin{aligned} L_{\text{CE}}(\hat{y}, y) &= - \sum_{k=1}^K y_k \log \hat{y}_k \\ &= - \sum_{k=1}^K y_k \log \hat{p}(y = k | x) \\ &= - \log \hat{y}_k, \quad (\text{where } k \text{ is the correct class}) \\ &= - \log \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)} \end{aligned}$$

Objective Function: Learn Model Parameters that Achieve a Specified (Measurable) Goal

Probability distribution of predicted class

Probability distribution of true class

Number of classes?

Recall, truth is set to 1 for one class and 0 otherwise

Observed features

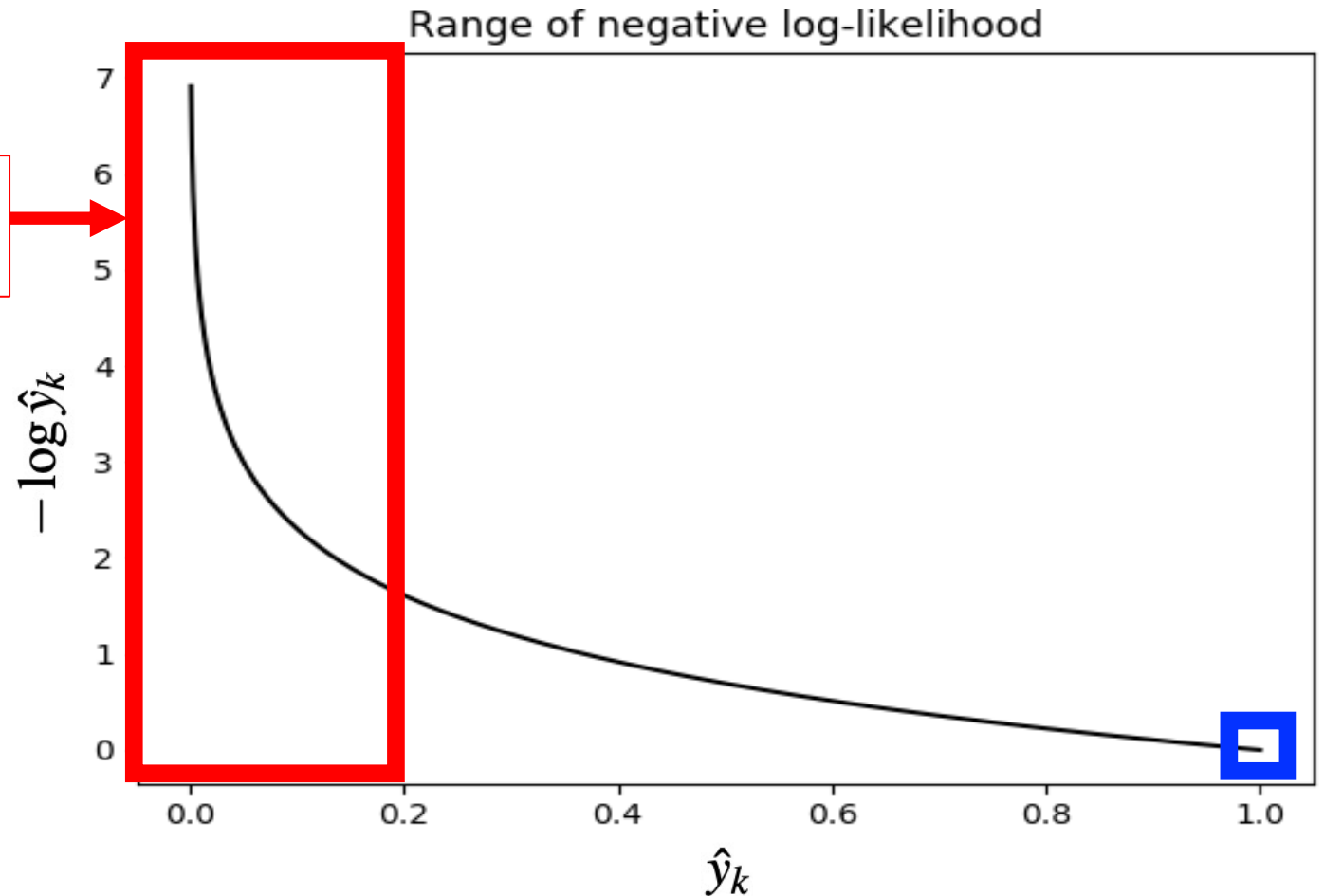
$$\begin{aligned} L_{\text{CE}}(\hat{y}, y) &= - \sum_{k=1}^K y_k \log \hat{y}_k \\ &= - \sum_{k=1}^K y_k \log \hat{p}(y = k | x) \\ &= - \log \hat{y}_k, \quad (\text{where } k \text{ is the correct class}) \\ &= - \log \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)} \end{aligned}$$

What is the range of possible values?

- Minimum: 0
 - i.e., correct prediction: negative log of 1
- Maximum: Infinity
 - i.e., incorrect prediction: negative log of 0

Objective Function: Learn Model Parameters that Achieve a Specified (Measurable) Goal

More confidently wrong predictions lead to greater error

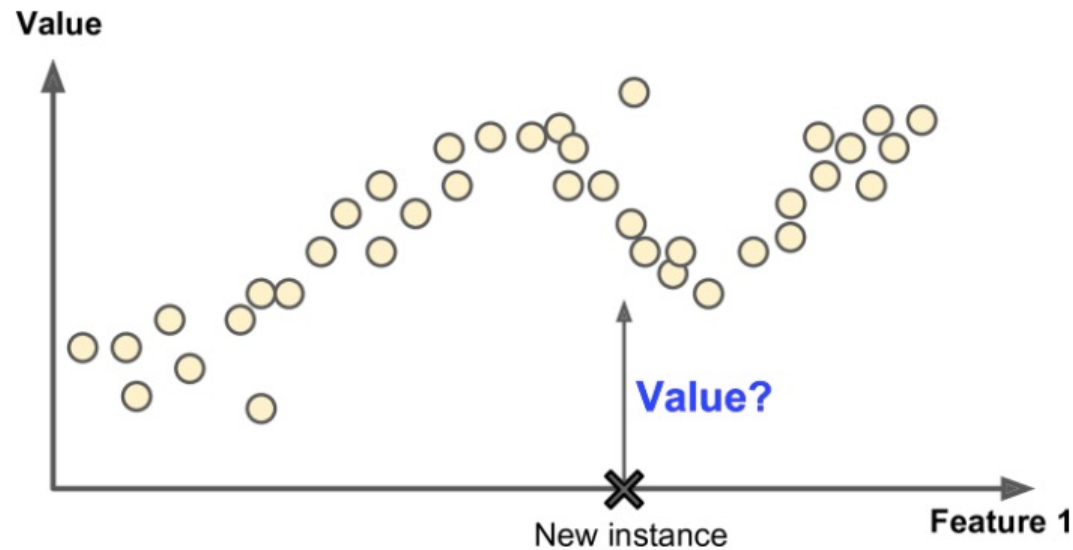


What is the range of possible values?

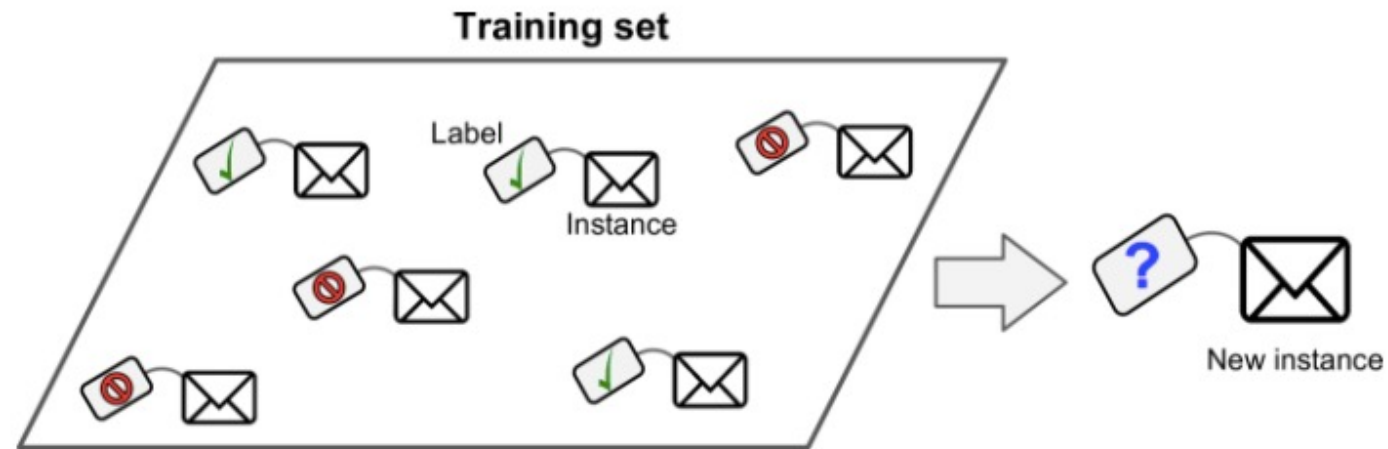
- Minimum: 0
 - i.e., correct prediction: negative log of 1
- Maximum: Infinity
 - i.e., incorrect prediction: negative log of 0

Objective Function: Learn Model Parameters that Achieve a Specified (Measurable) Goal

Regression
(predict **continuous** value)



Classification
(predict **discrete** value)

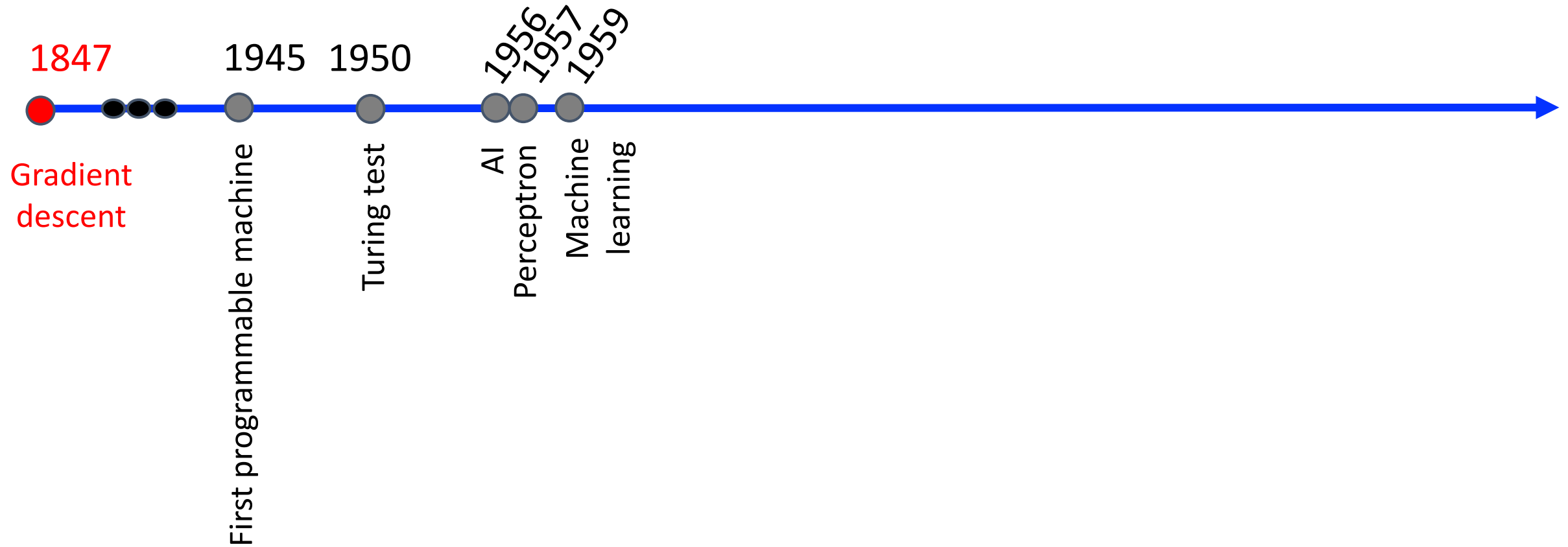


MANY objective functions exist, and we will examine popular ones in this course

Today's Topics

- Objective function: what to learn
- **Gradient descent: how to learn**
- Training a neural network: optimization
- Gradient descent for activation functions

Historical Context: Gradient Descent



Scalable way to train nonlinear models on “big data”

Gradient Descent: Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn linear model for converting kilometers to miles when only observing the input “miles” and output “kilometers”



Gradient Descent: Intuition

- Repeat:
 1. **Guess**
 2. Calculate error
- e.g., learn constant multiplier to convert US dollars to Israeli shekels

\$10 → Shekels = dollars x **constant**

Gradient Descent: Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert US dollars to Israeli shekels



Gradient Descent: Intuition

- Repeat:
 1. **Guess**
 2. Calculate error
- e.g., learn constant multiplier to convert US dollars to Israeli shekels

$\$10 \longrightarrow \text{Shekels} = \text{dollars} \times \text{constant}$

Gradient Descent: Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert US dollars to Israeli shekels



Gradient Descent: Intuition

- Repeat:
 1. **Guess**
 2. Calculate error
- e.g., learn constant multiplier to convert US dollars to Israeli shekels

$\$10 \longrightarrow \text{Shekels} = \text{dollars} \times \text{constant}$

Gradient Descent: Intuition

- Repeat:
 1. Guess
 2. Calculate error
- e.g., learn constant multiplier to convert US dollars to Israeli shekels



- Idea: iteratively adjust **constant (i.e., model parameter)** to try to reduce the error

Gradient Descent: Intuition

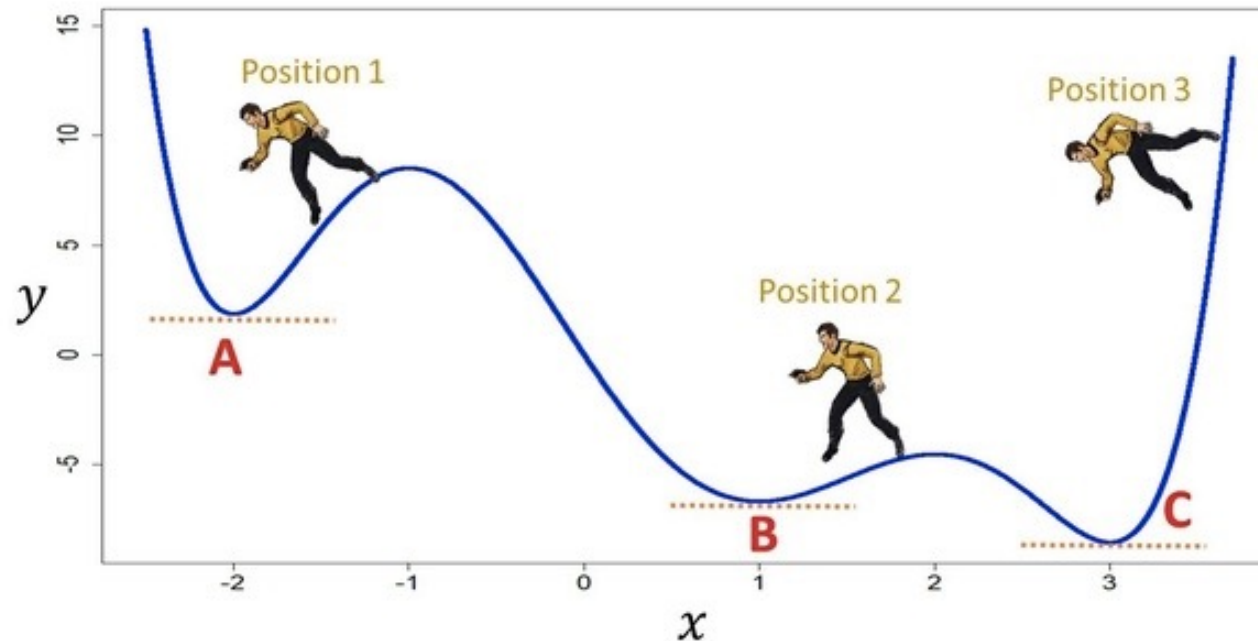
- Iteratively search for model parameters (i.e., weights and biases) that solve optimization problem (i.e., minimize or maximize an objective function)

Analogy: hiking to the bottom of a mountain range... blind or blindfolded!



Gradient Descent: Employs Calculus

- Idea: use derivatives!
 - Derivatives tells us how to change the input x to make a small change to the output $f(x)$
 - Gradient is a vector that indicates how $f(\mathbf{x})$ changes as each function variable changes (i.e., partial derivatives)
- Gradient descent:
 - Iteratively take steps in the opposite direction of the gradient to minimize the function

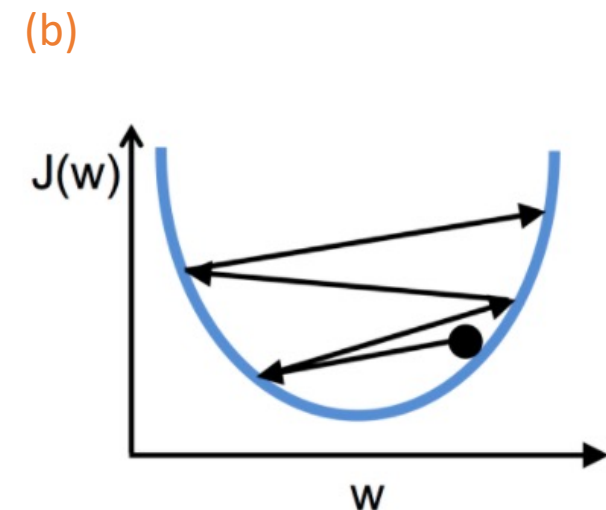
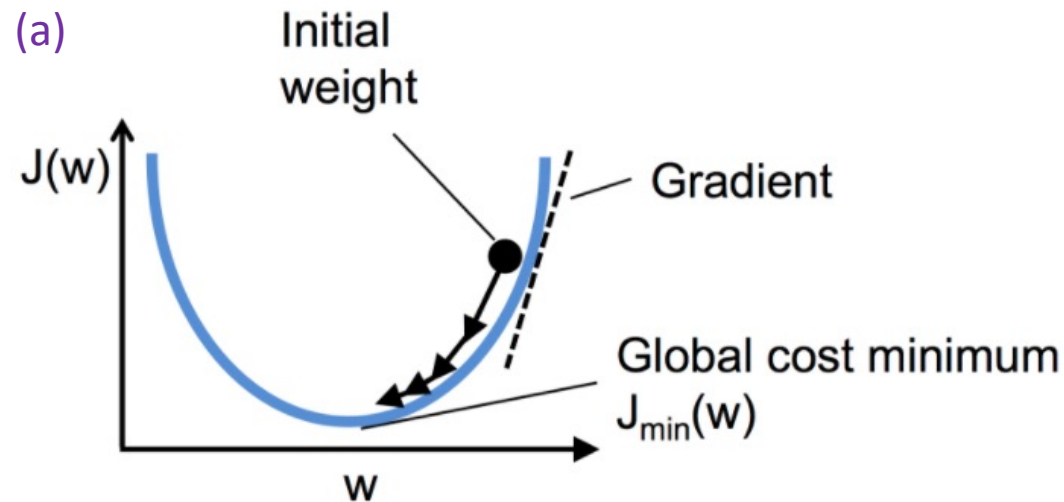


Which letter(s) are the global minima?

Which letter(s) are local minima?

Gradient Descent: How Much to Update?

- Step size = learning rate
 - (a) When learning rate is too small, convergence to good solution will be slow
 - (b) When learning rate is too large, convergence to a good solution is not possible



- Next lecture: examination of how to learn effectively using the gradients

Gradient Descent: How Often to Update?

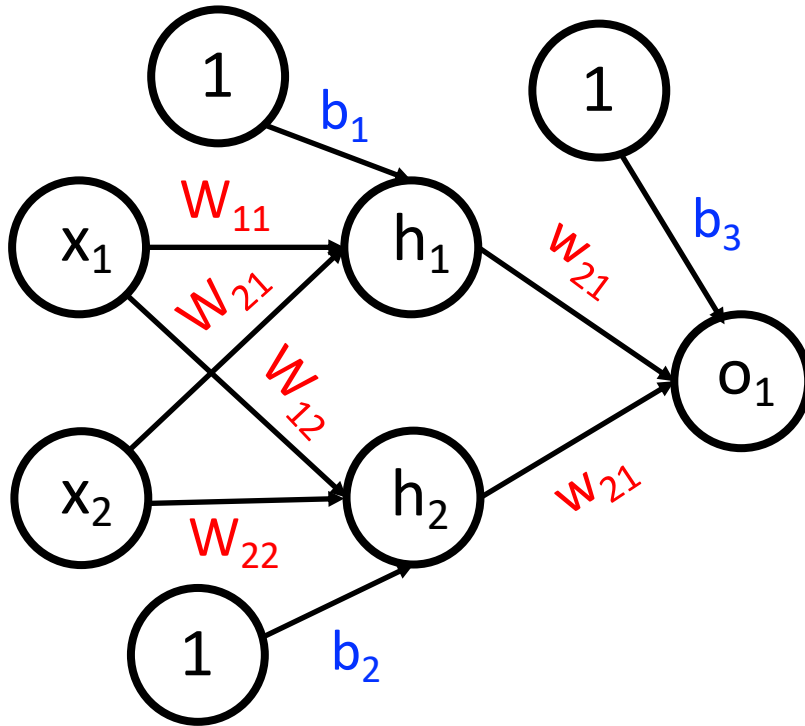
- Use calculations over ***all training examples*** (Batch gradient descent)
 - Less bouncing but can be slow or infeasible when dataset is large
- Use calculations from ***one training example*** (Stochastic gradient descent)
 - Fast to compute and can train using huge datasets (stores one instance in memory at each iteration) but updates are expected to bounce a lot
- Use calculations over ***subset of training examples*** (Mini-batch gradient descent)
 - Bounces less erratically than SGD and can train using huge datasets (store some instances in memory at each iteration) but can be slow or infeasible when dataset is large
- Often mini-batch gradient descent is used with maximum # of examples that fit in memory

Today's Topics

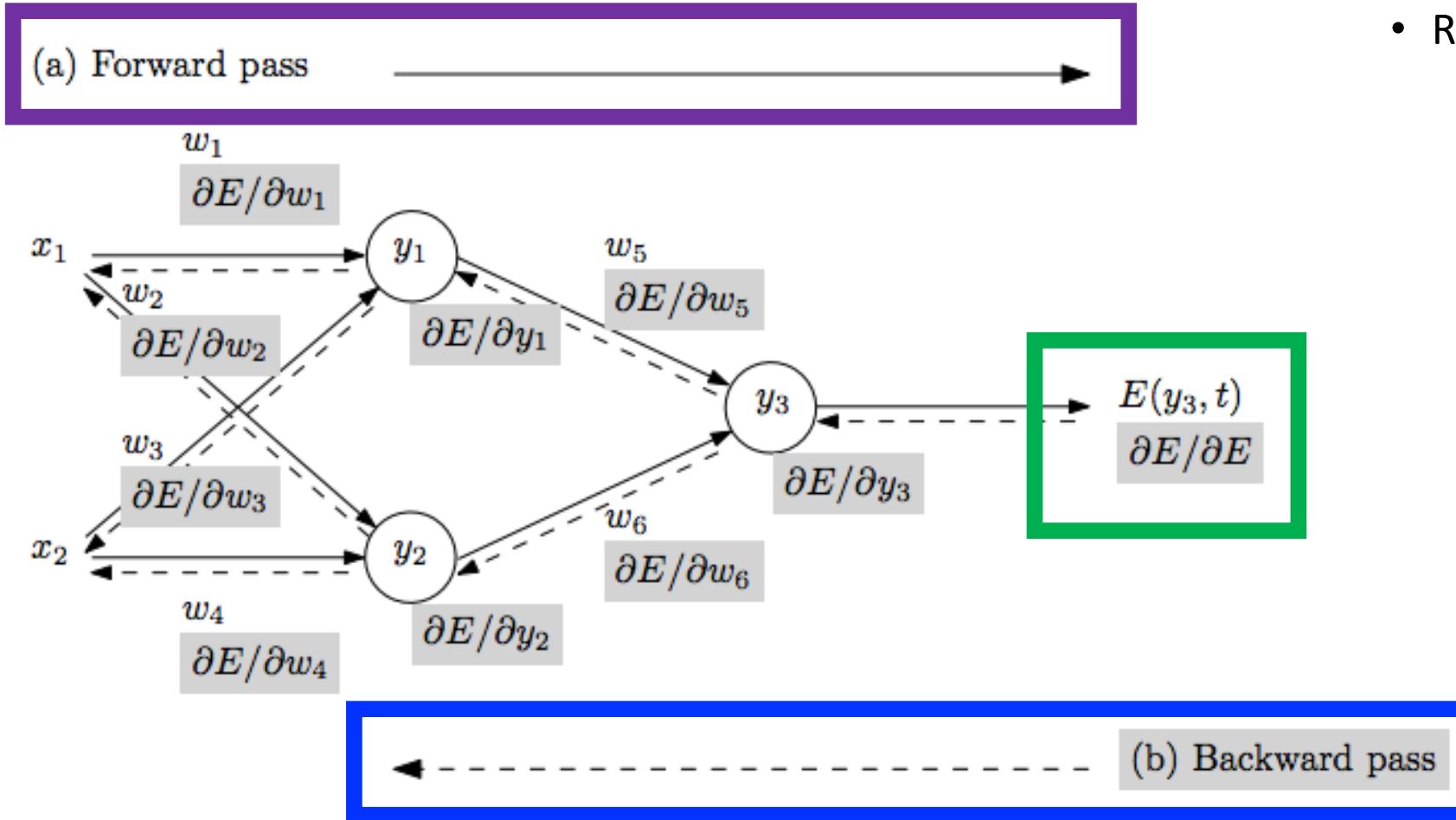
- Objective function: what to learn
- Gradient descent: how to learn
- Training a neural network: optimization
- Gradient descent for activation functions

Summary: Approach to Train Neural Network

- Learn model parameters (**weights**, **biases**) that minimize an objective function using gradient descent; e.g.,



Training: How Neural Networks Learn



- Repeat until stopping criterion met:

1. **Forward pass:** propagate training data through model to make prediction
2. Quantify the dissatisfaction with a model's results on the training data
3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
4. Update each parameter using calculated gradients

Training: How Neural Networks Learn

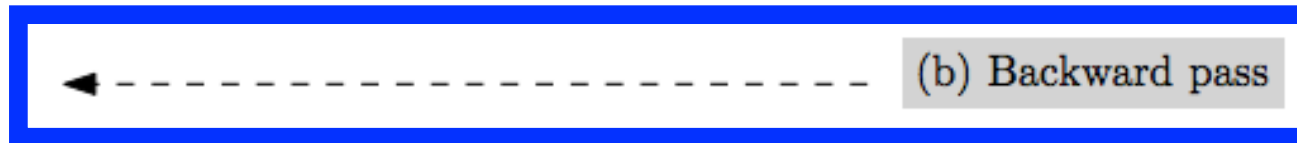
Key challenge: calculating gradients. Depends on:

1) Objective function

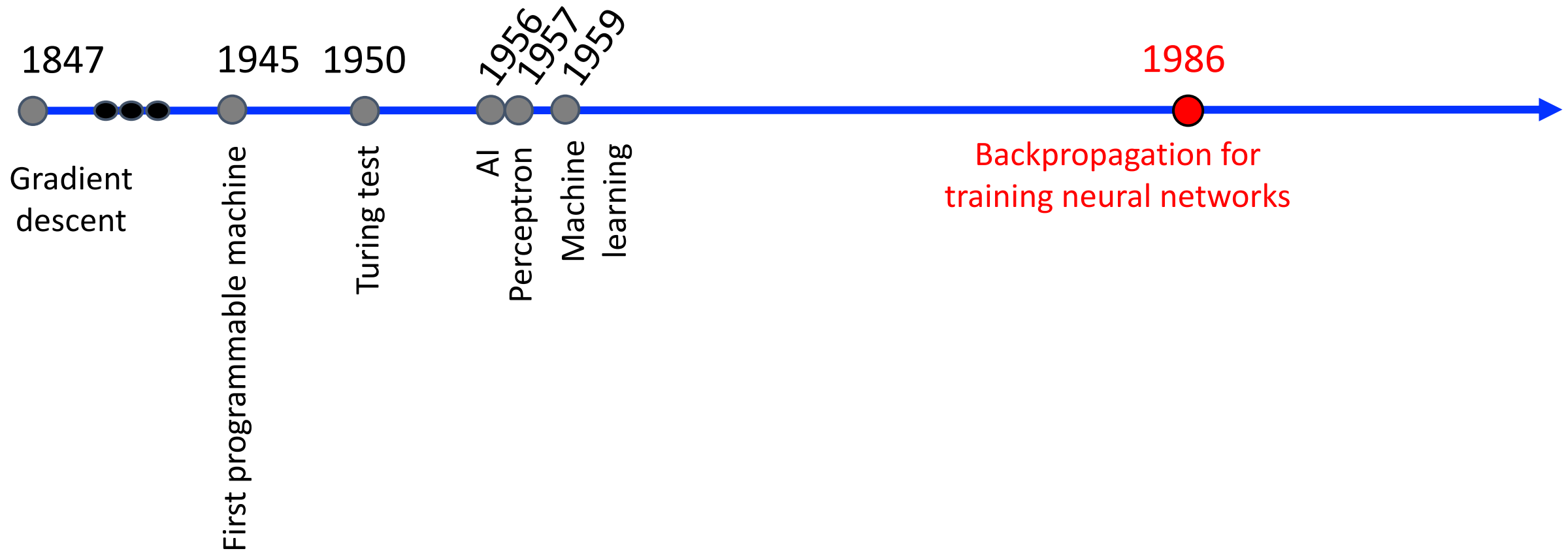
2) Activation functions

Repeat until stopping criterion met:

1. **Forward pass:** propagate training data through model to make prediction
2. Quantify the dissatisfaction with a model's results on the training data
3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
4. Update each parameter using calculated gradients



Solution: Backpropagation to Compute Gradients



D. Rulhart, G. Hinton, and R. Williams, Learning Internal Representations by Error Propagation, 1986.

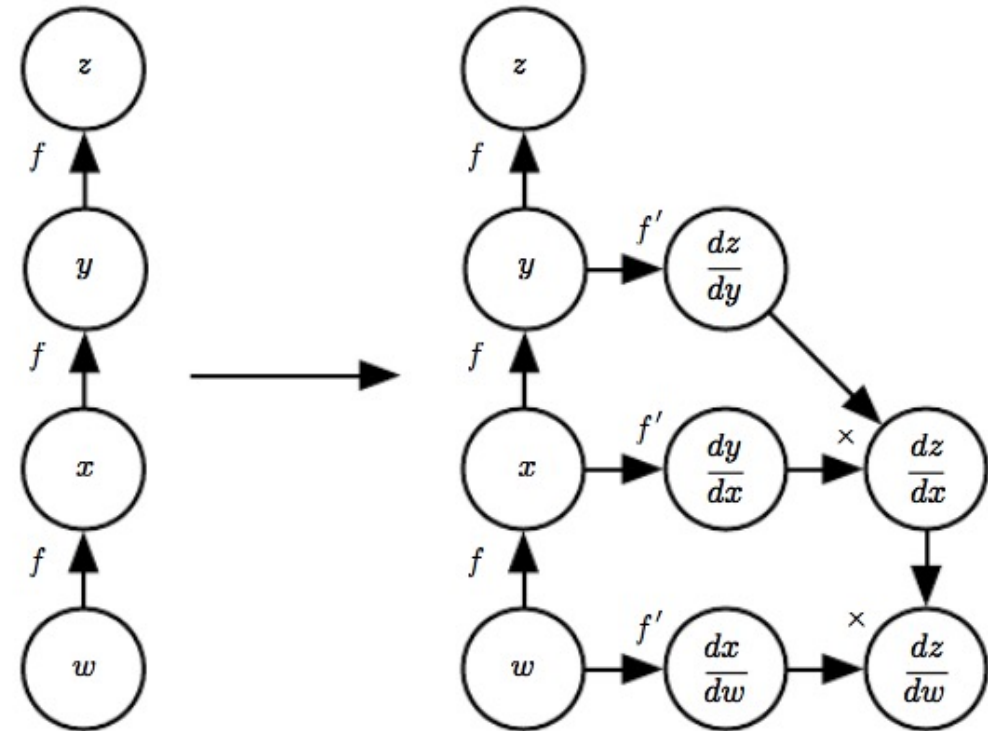
Solution: Backpropagation to Compute Gradients

- Idea: compute gradient on objective function to decide how to adjust each model parameter to get closer to solving the optimization problem
- Key observation: networks are functions connected in a chain

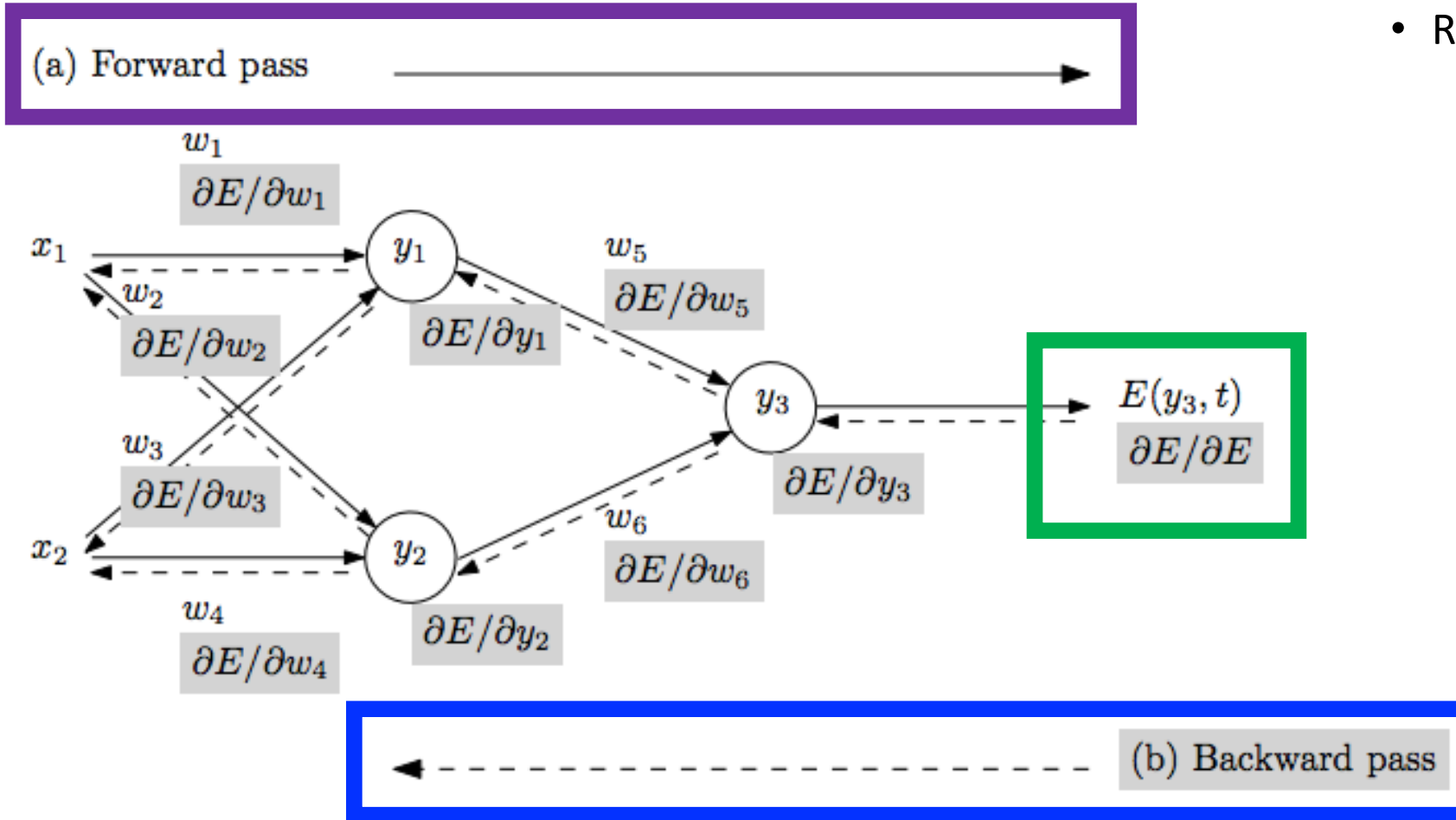
$$x = f(w), y = f(x), z = f(y)$$

Can use **chain rule of calculus** (and so compute from top to bottom where derivatives on the top are used to compute derivatives at the bottom);

$$\text{e.g., } \frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

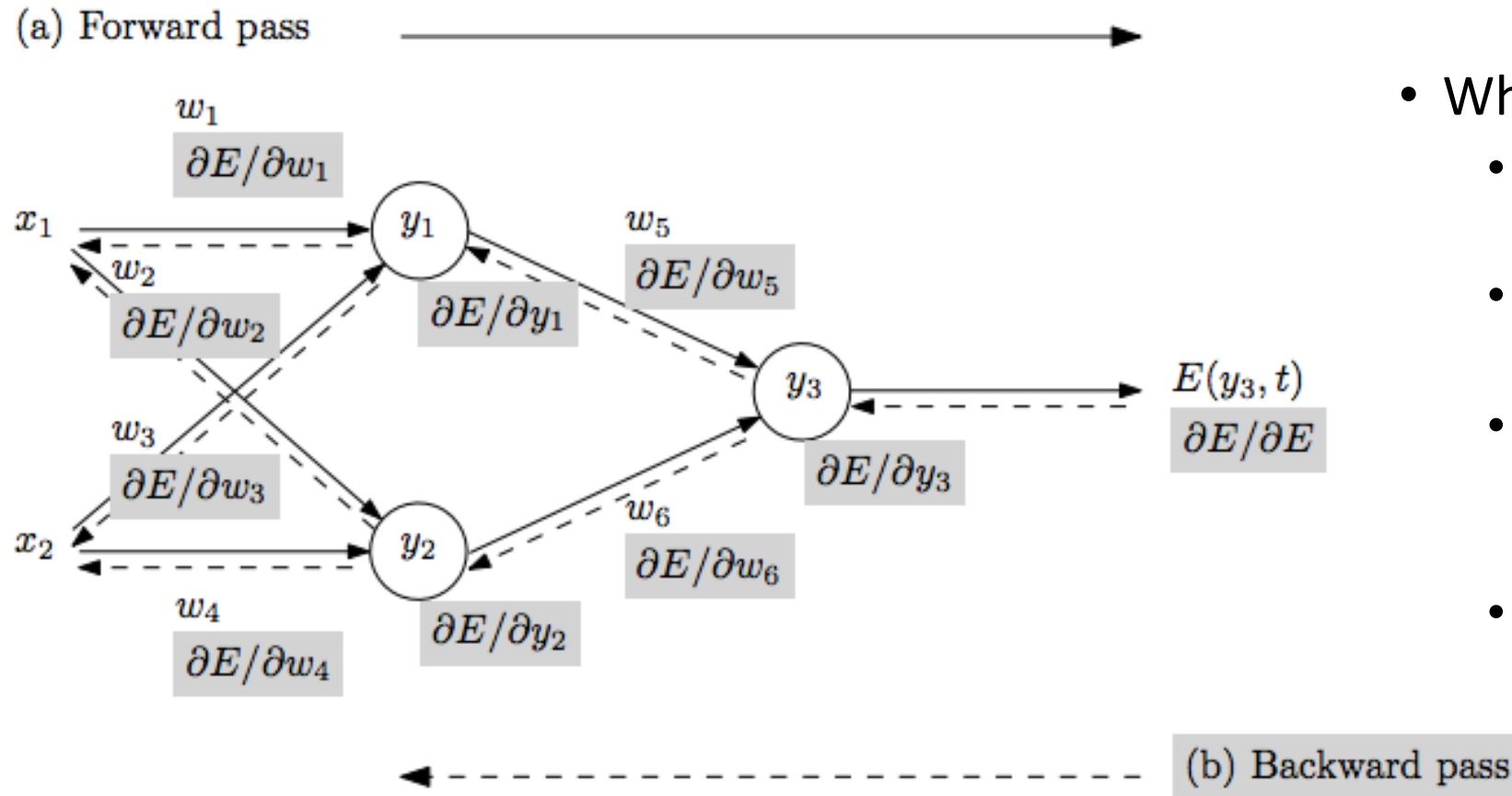


Training: How Neural Networks Learn



- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make prediction
 2. Quantify the dissatisfaction with a model's results on the training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

When to Stop Training Neural Networks?

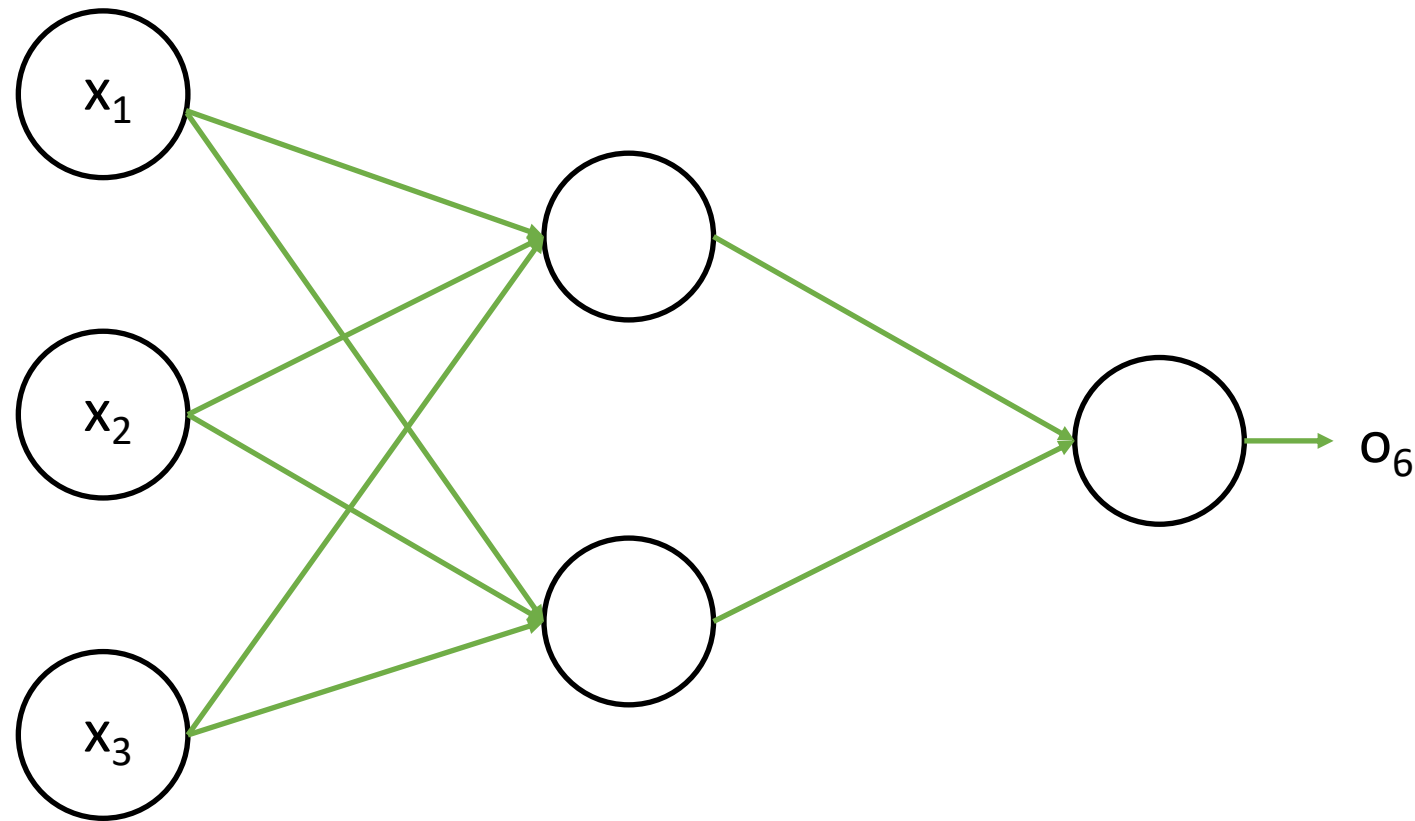


- What stopping criterion to use?
 - Weight changes are incredibly small
 - Finished a pre-specified number of epochs
 - Percentage of misclassified example is below some threshold
 - ...

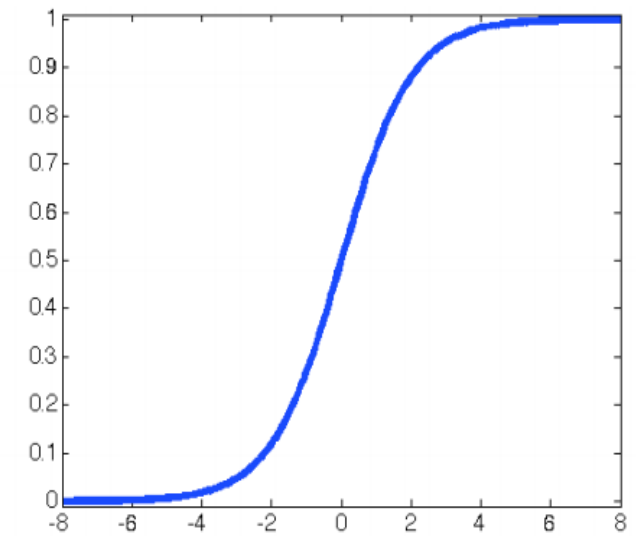
Example

- Binary classification: predict if a student will get a B- or better
- Inputs:
 - Percentage of assignments completed
 - Percentage of readings read
 - Percentage of lectures watched

Example: Choose Neural Network Architecture

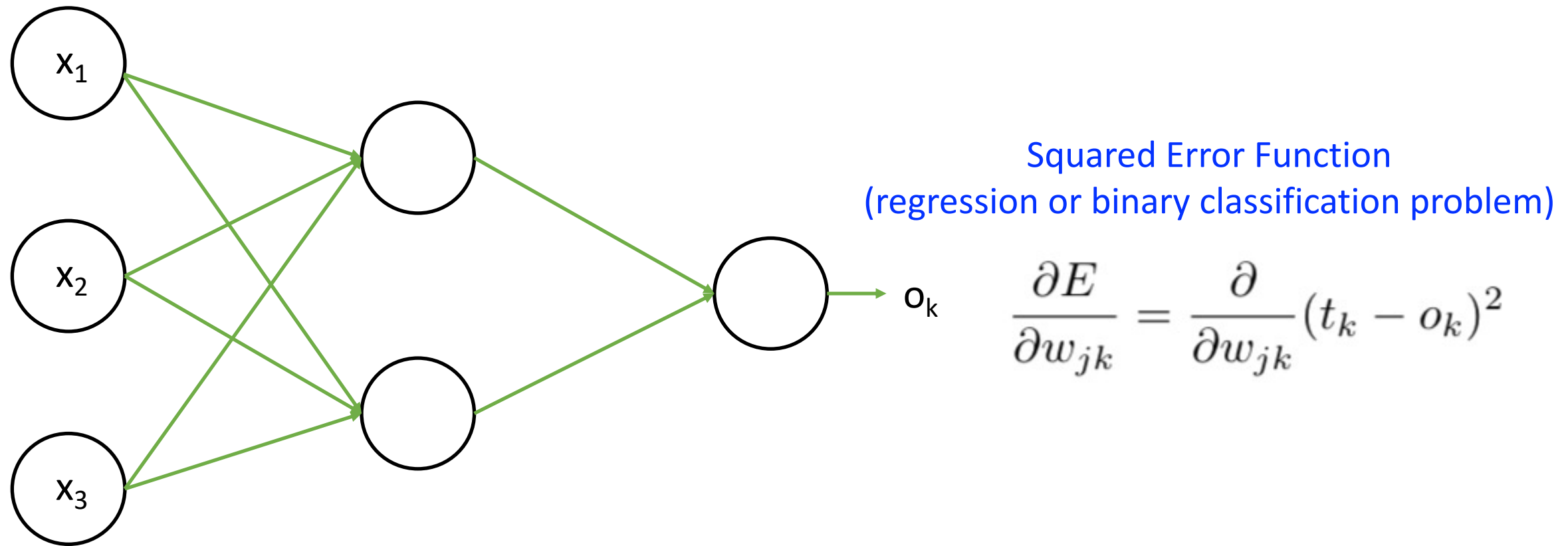


Sigmoid Activation
Function

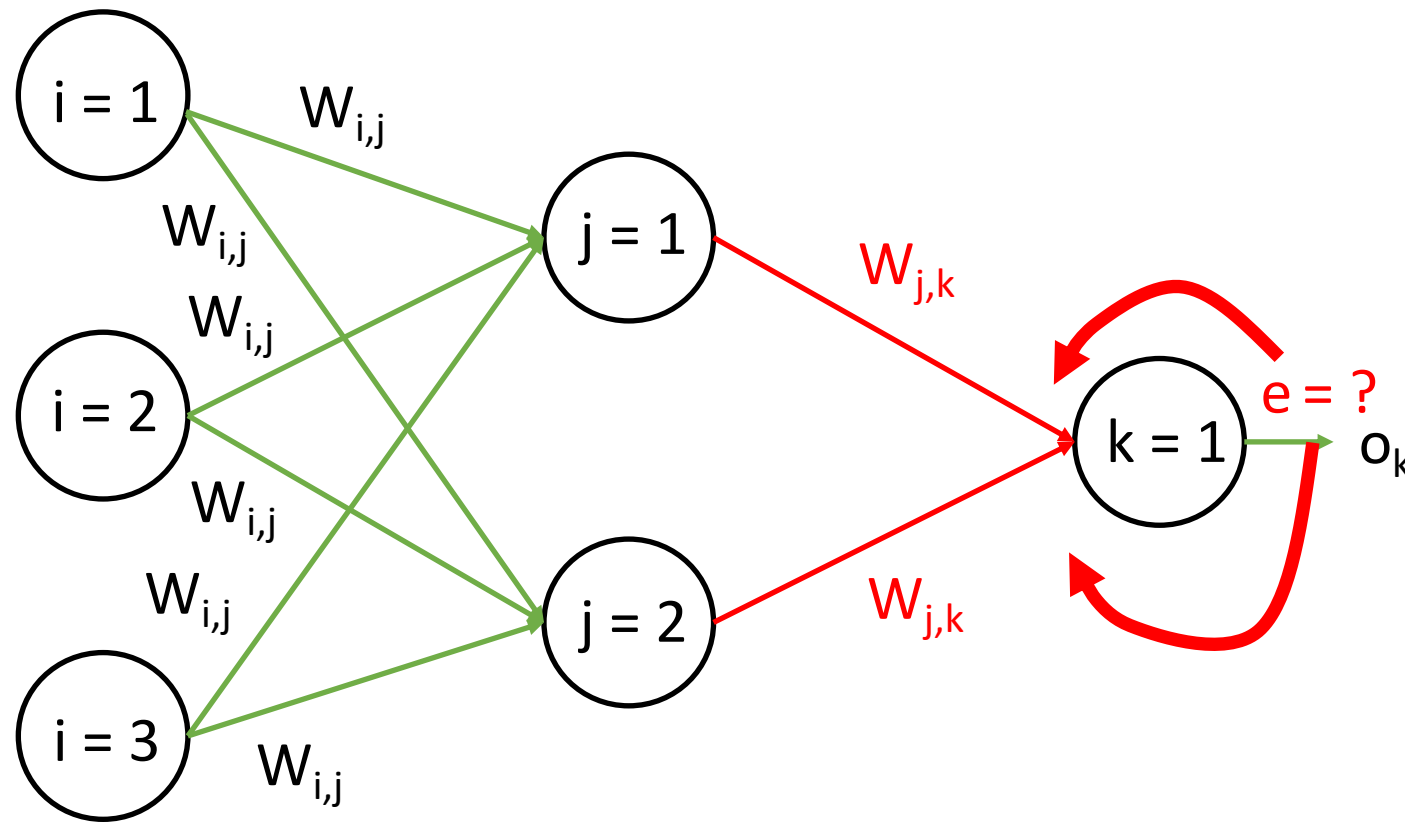


$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Example: Choose Loss Function for Training



Example: Resolve How to Compute Gradient? (Output Layer)



t is a constant value:

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} (t_k - o_k)^2$$

Example: Resolve How to Compute Gradient? (Output Layer)

t is a constant value:

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} (t_k - o_k)^2$$

Using the following chain rule :

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial o_k} * \frac{\partial o_k}{\partial w_{jk}}$$

$$\frac{\partial E}{\partial o_k} = -2(t_k - o_k)$$

Sigmoid activation function: $\sigma(x) = \frac{1}{1 + e^{-x}}$

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$\frac{d\sigma(x)}{dx} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right)$$

$$\frac{d\sigma(x)}{dx} = (1 - \sigma(x)) \sigma(x)$$

Example: Resolve How to Compute Gradient? (Output Layer)

t is a constant value:

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} (t_k - o_k)^2$$

Using the following chain rule :

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial o_k} * \frac{\partial o_k}{\partial w_{jk}}$$

$$\frac{\partial E}{\partial o_k} = -2(t_k - o_k)$$

Sigmoid activation function: $\sigma(x) = \frac{1}{1 + e^{-x}}$

$$\frac{d\sigma(x)}{dx} = (1 - \sigma(x)) \sigma(x)$$

We can rewrite our function as follows:

For efficiency, compute last

$$\frac{\partial E}{\partial w_{jk}} = -2(t_k - o_k) \text{sigmoid}\left(\sum_j w_{jk} * o_j\right) * (1 - \text{sigmoid}\left(\sum_j w_{jk} * o_j\right)) * o_j$$

Example: Resolve How to Compute Gradient? (Output Layer)

t is a constant value:

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} (t_k - o_k)^2$$

Using the following chain rule :

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial o_k} * \frac{\partial o_k}{\partial w_{jk}}$$

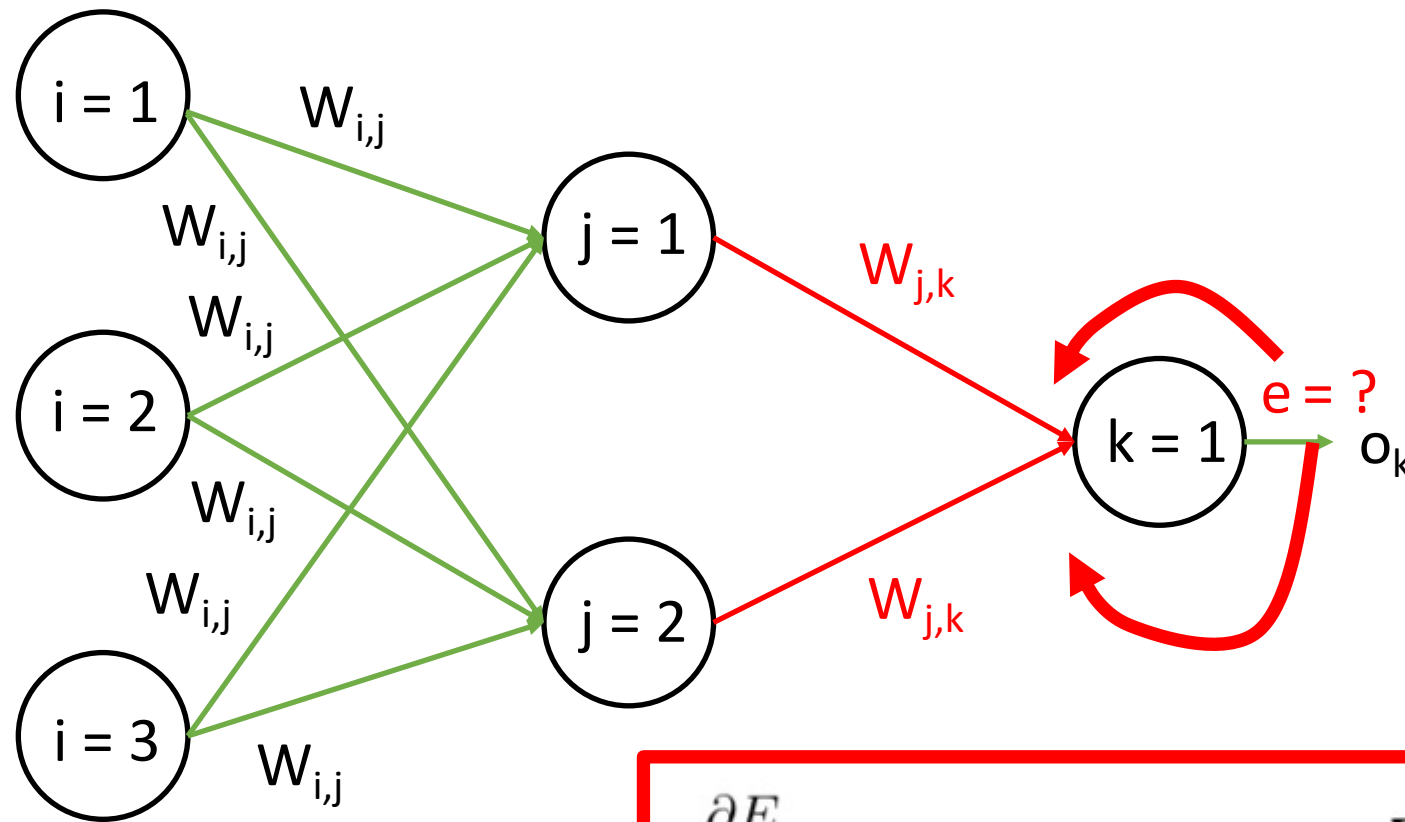
$$\frac{\partial E}{\partial o_k} = -2(t_k - o_k)$$

Sigmoid activation function: $\sigma(x) = \frac{1}{1 + e^{-x}}$

$$\frac{d\sigma(x)}{dx} = (1 - \sigma(x)) \sigma(x)$$

Key Observation: Possible because activation function
and loss function are **differentiable!!!**

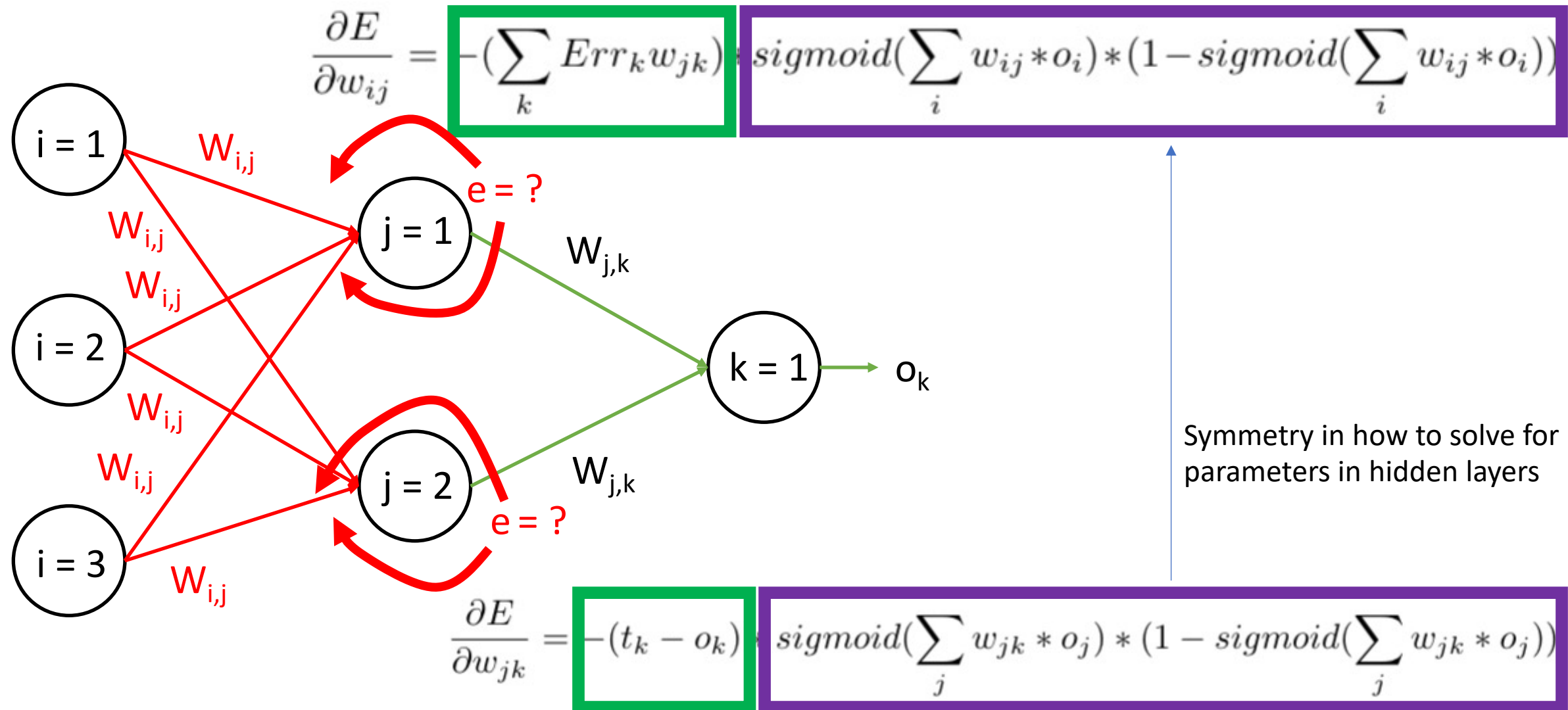
Example: Resolve How to Compute Gradient? (Output Layer)



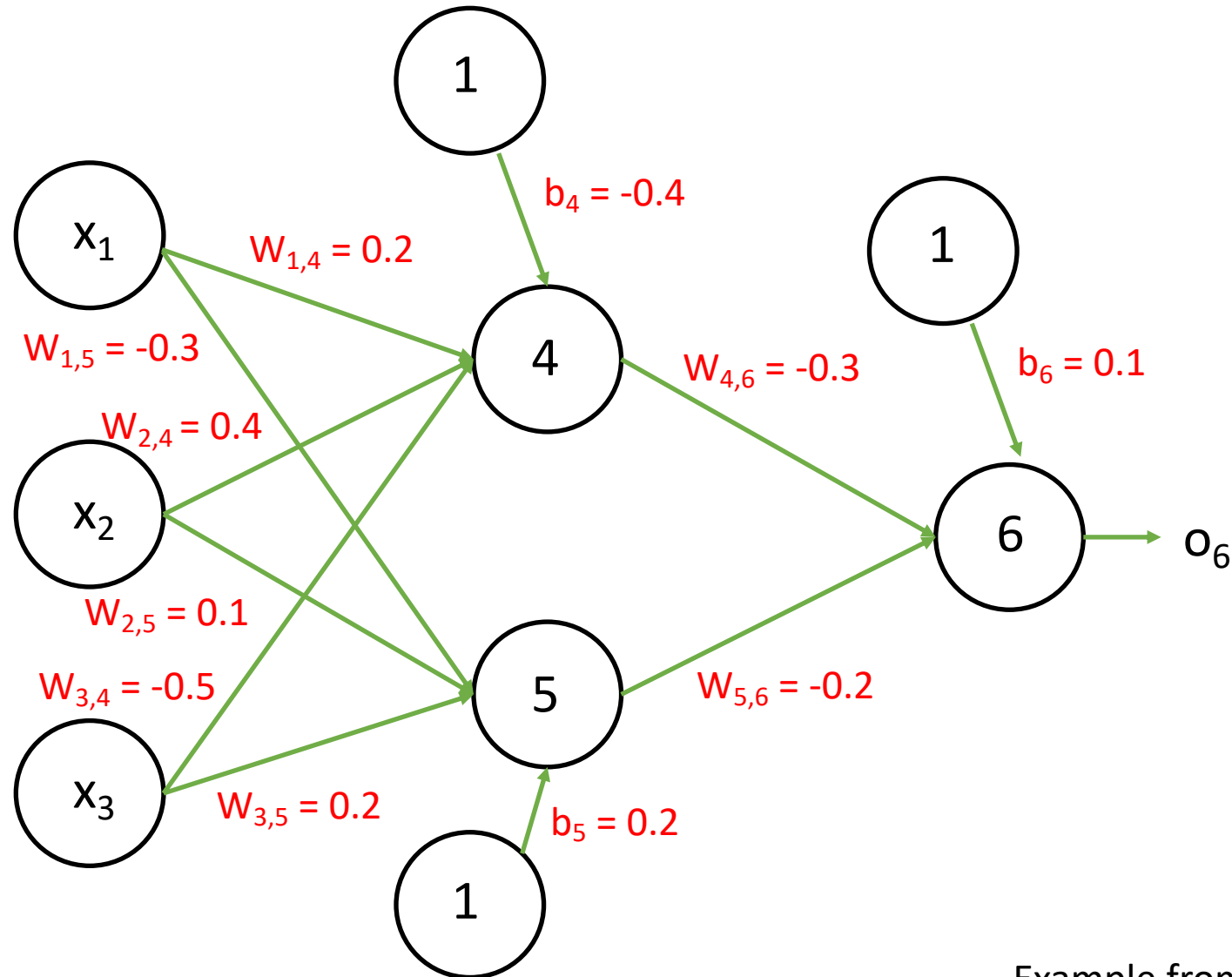
$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} (t_k - o_k)^2$$

$$\frac{\partial E}{\partial w_{jk}} = -(t_k - o_k) * \text{sigmoid}(\sum_j w_{jk} * o_j) * (1 - \text{sigmoid}(\sum_j w_{jk} * o_j))$$

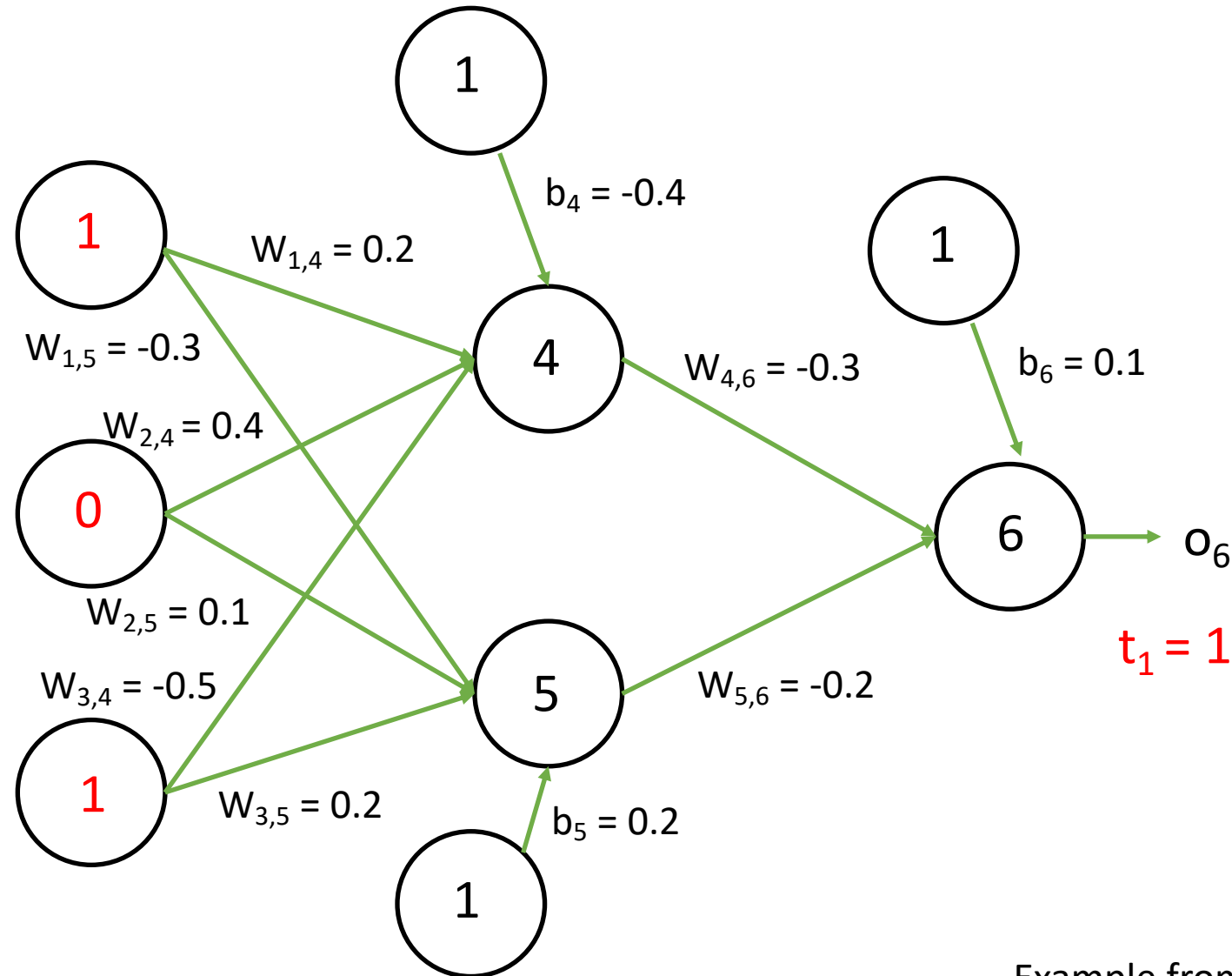
Example: How to Compute Gradient? (Hidden Layer)



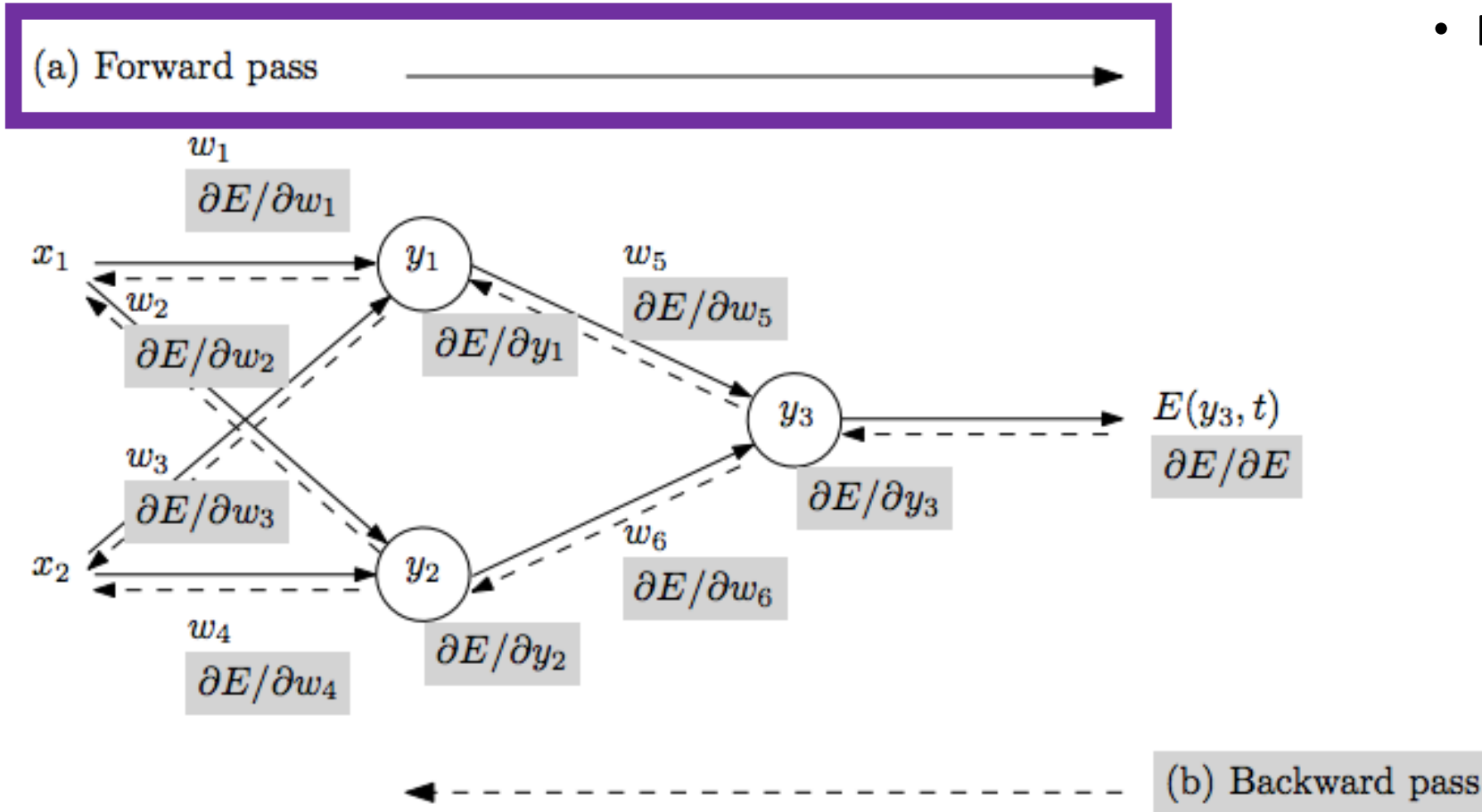
Example: Initialize Model Parameters



Example: Input Training Example

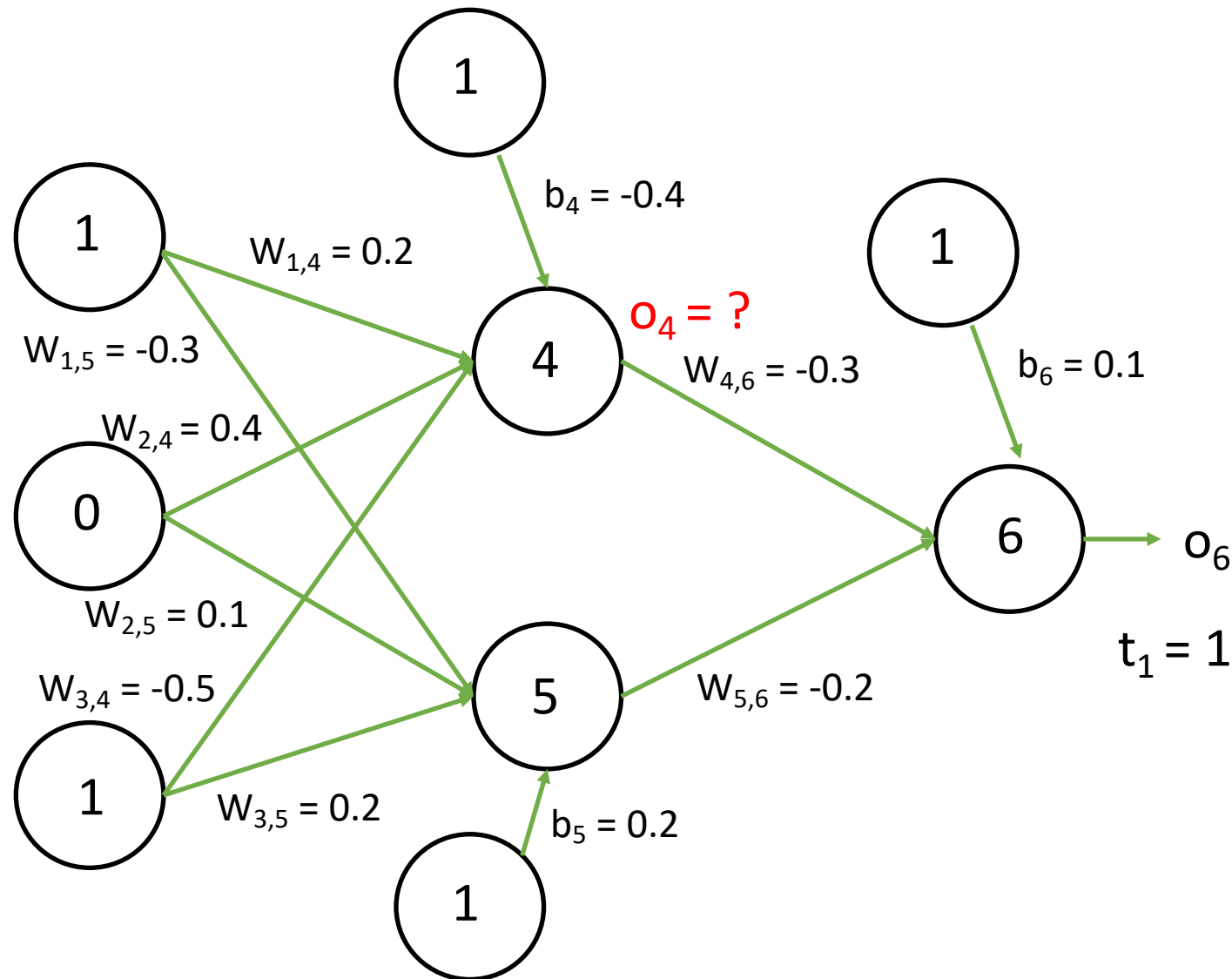


Training: How Neural Networks Learn



- Repeat until stopping criterion met:
 - Forward pass:** propagate training data through model to make prediction

Example: Step 1 – Forward Pass



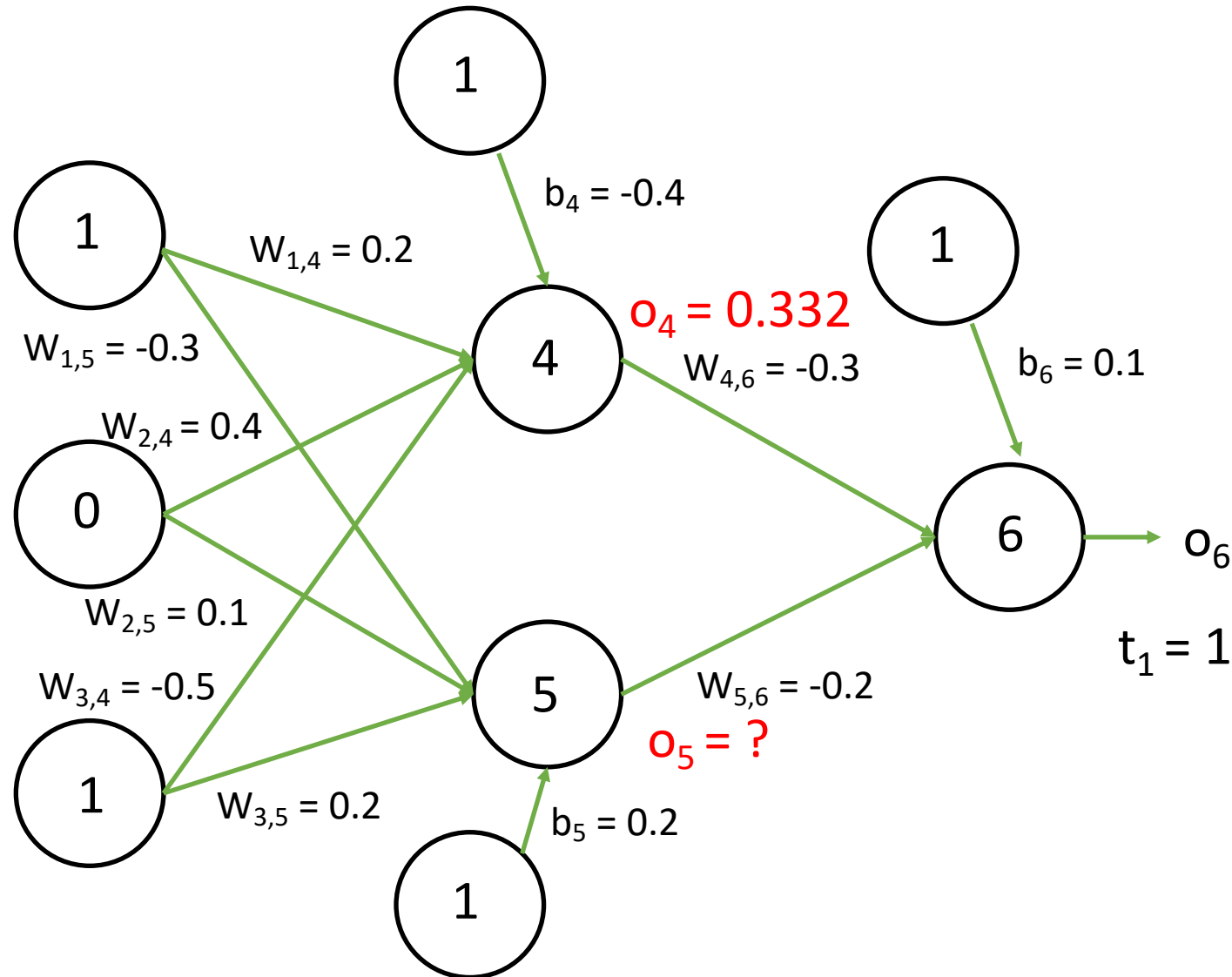
Input to node 4:

$$i_4 = (1 \times 0.2 + 0 \times 0.4 + 1 \times -0.5) - 0.4$$
$$i_4 = -0.7$$

Output of node 4 (sigmoid function):

$$o_4 = \text{sigmoid}(-0.7)$$
$$o_4 = 1/(1+e^{-(-0.7)})$$
$$o_4 = 0.332$$

Example: Step 1 – Forward Pass



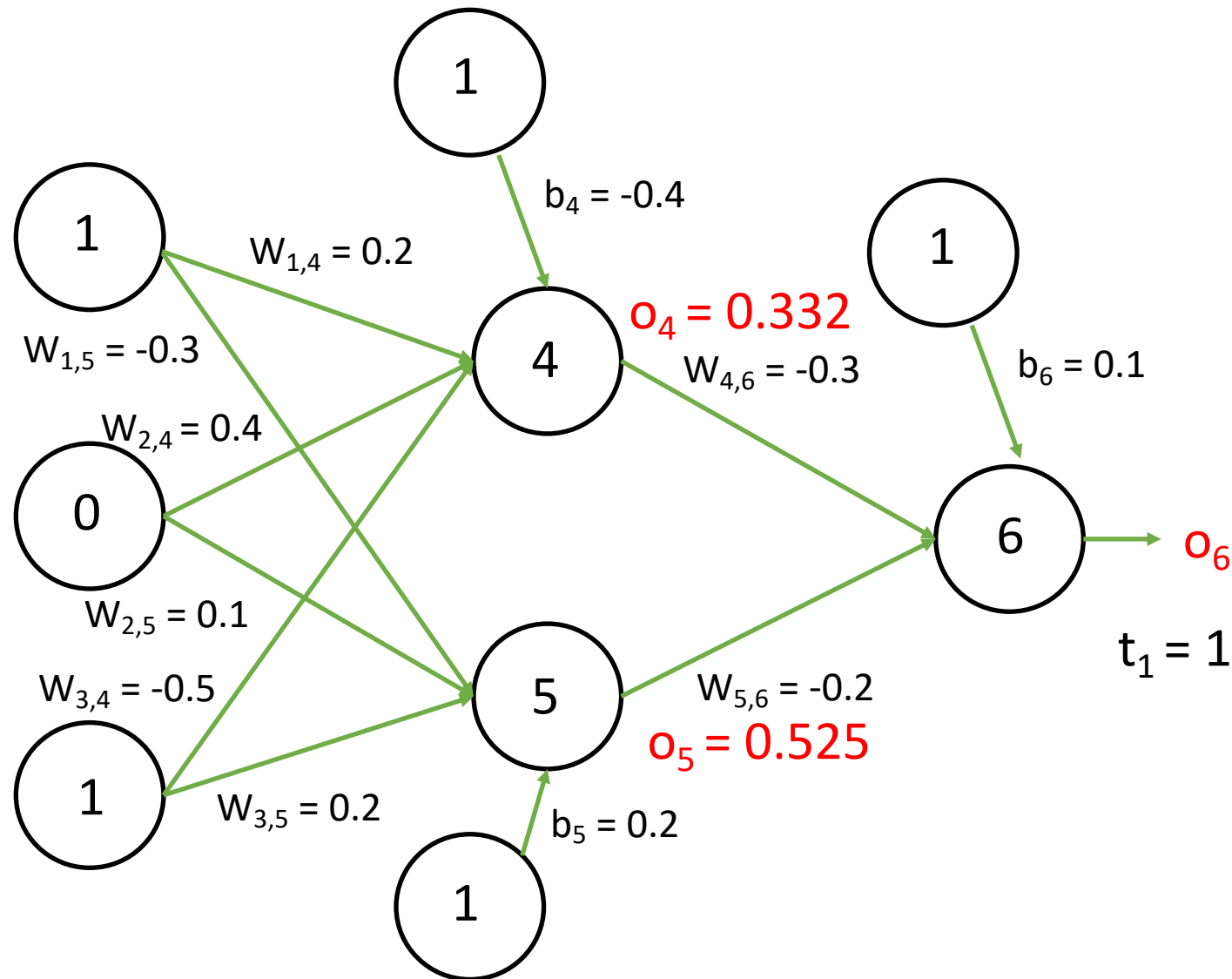
Input to node 5:

$$i_5 = (1 \times -0.3 + 0 \times 0.1 + 1 \times 0.2) + 0.2$$
$$i_5 = 0.1$$

Output of node 5 (sigmoid function):

$$o_5 = \text{sigmoid}(0.1)$$
$$o_5 = 1/(1+e^{-0.1})$$
$$o_5 = 0.525$$

Example: Step 1 – Forward Pass



Input to node 6:

$$i_6 = (0.332 \times -0.3 + 0.525 \times -0.2) + 0.1$$

$$i_6 = -0.105$$

Output of node 6 (sigmoid function):

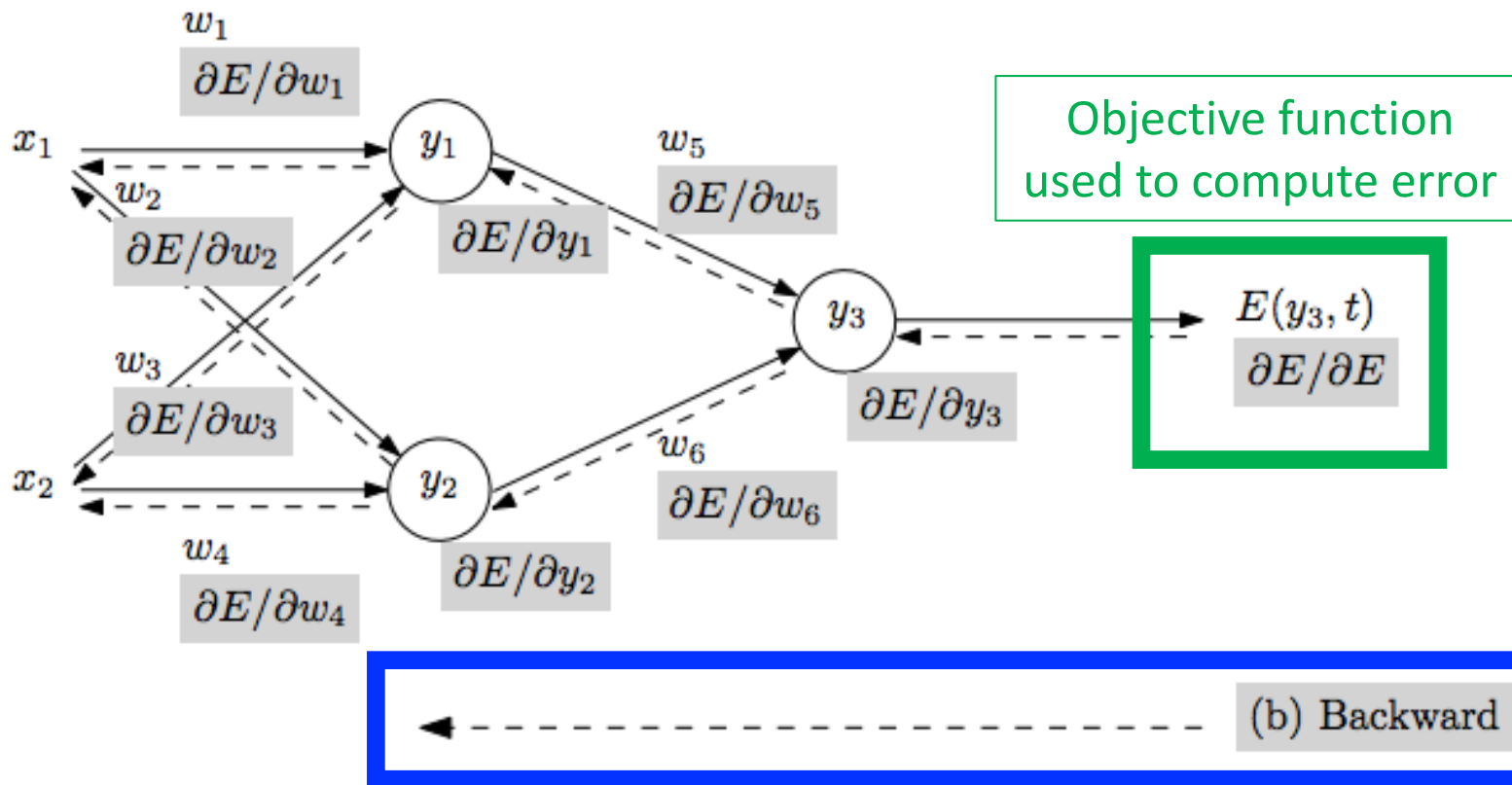
$$o_6 = \text{sigmoid}(-0.105)$$

$$o_6 = 1/(1+e^{-(-0.105)})$$

$$o_6 = 0.474$$

Training: How Neural Networks Learn

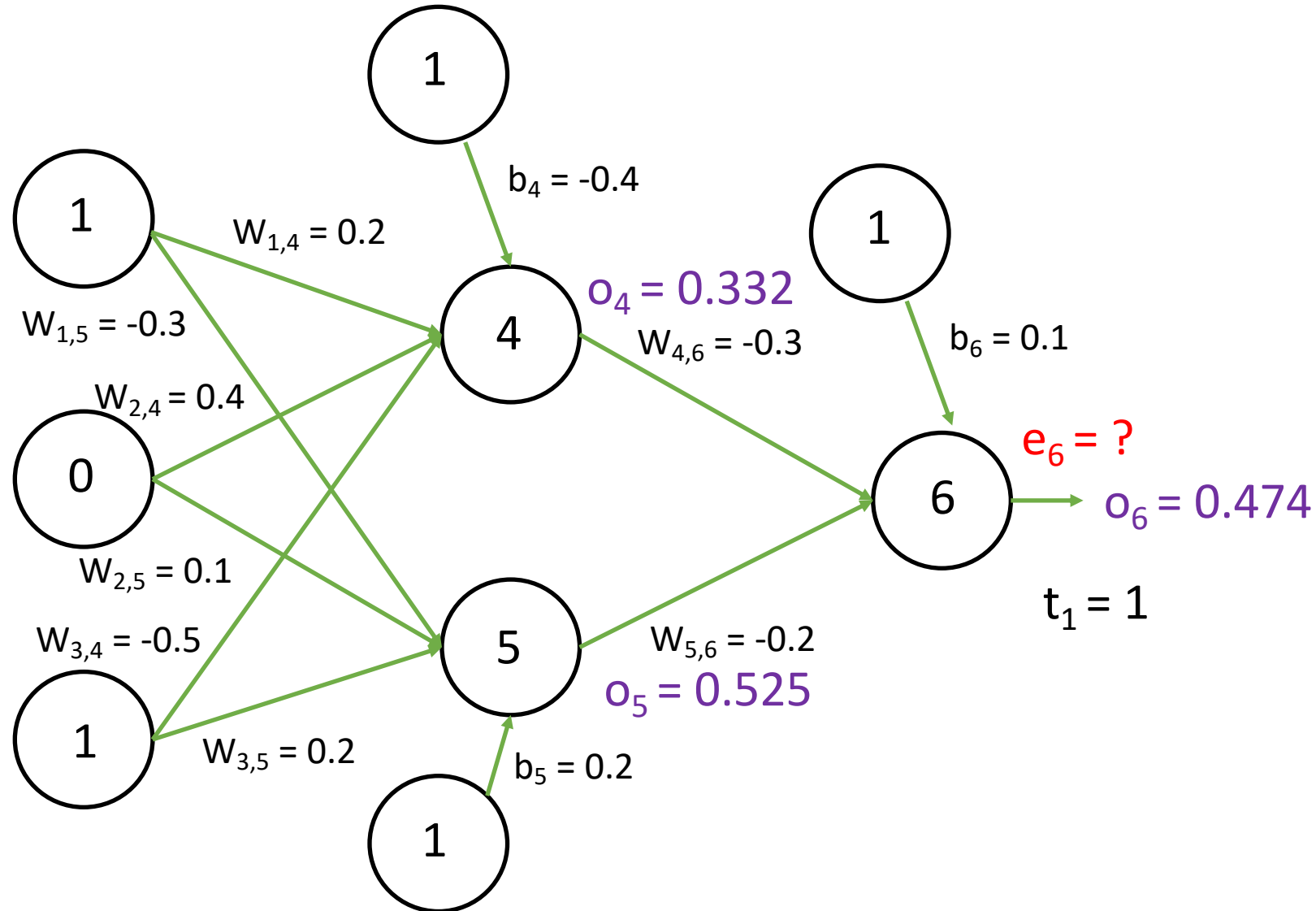
(a) Forward pass



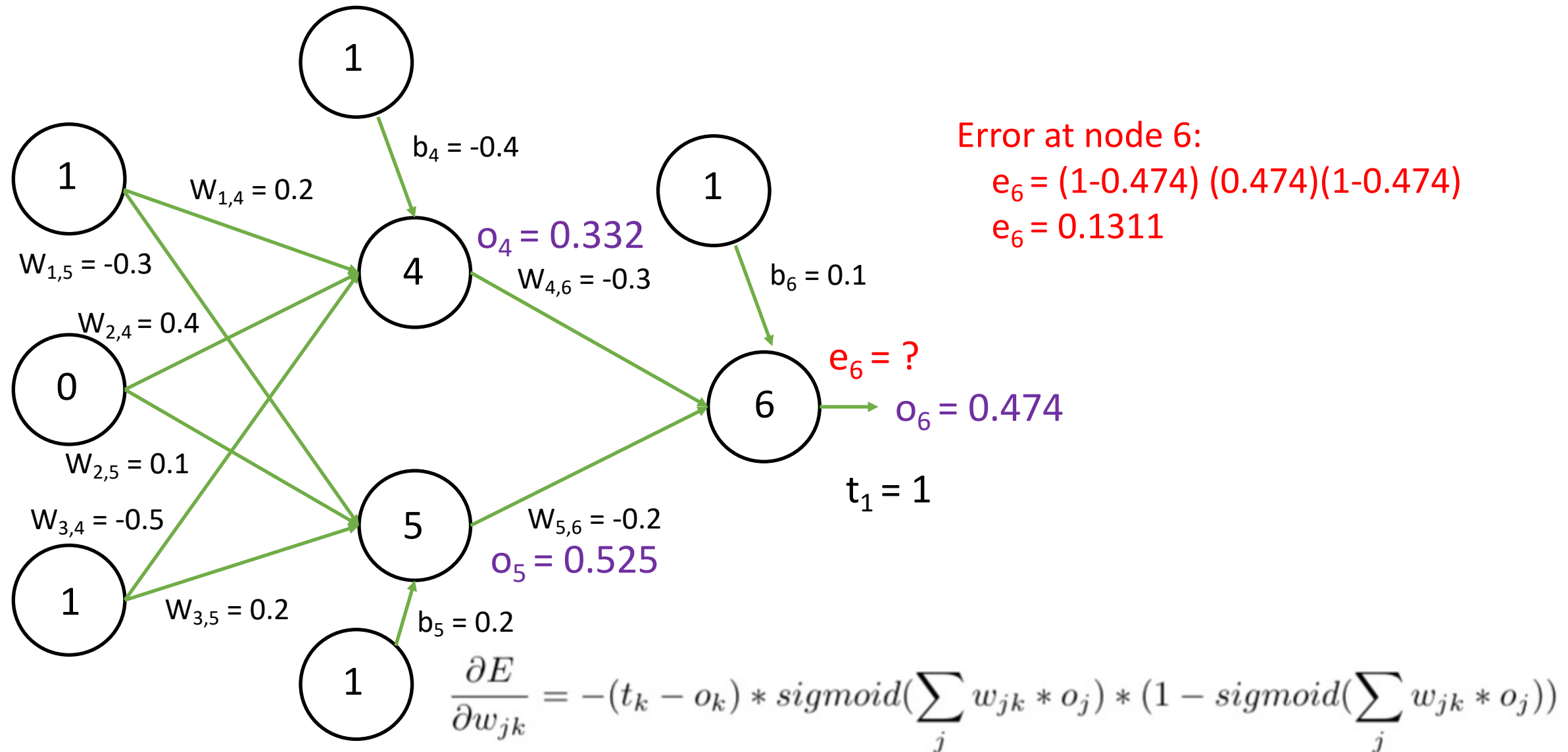
• Repeat until stopping criterion met:

1. **Forward pass:** propagate training data through model to make prediction
2. Quantify the dissatisfaction with a model's results on the training data
3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter

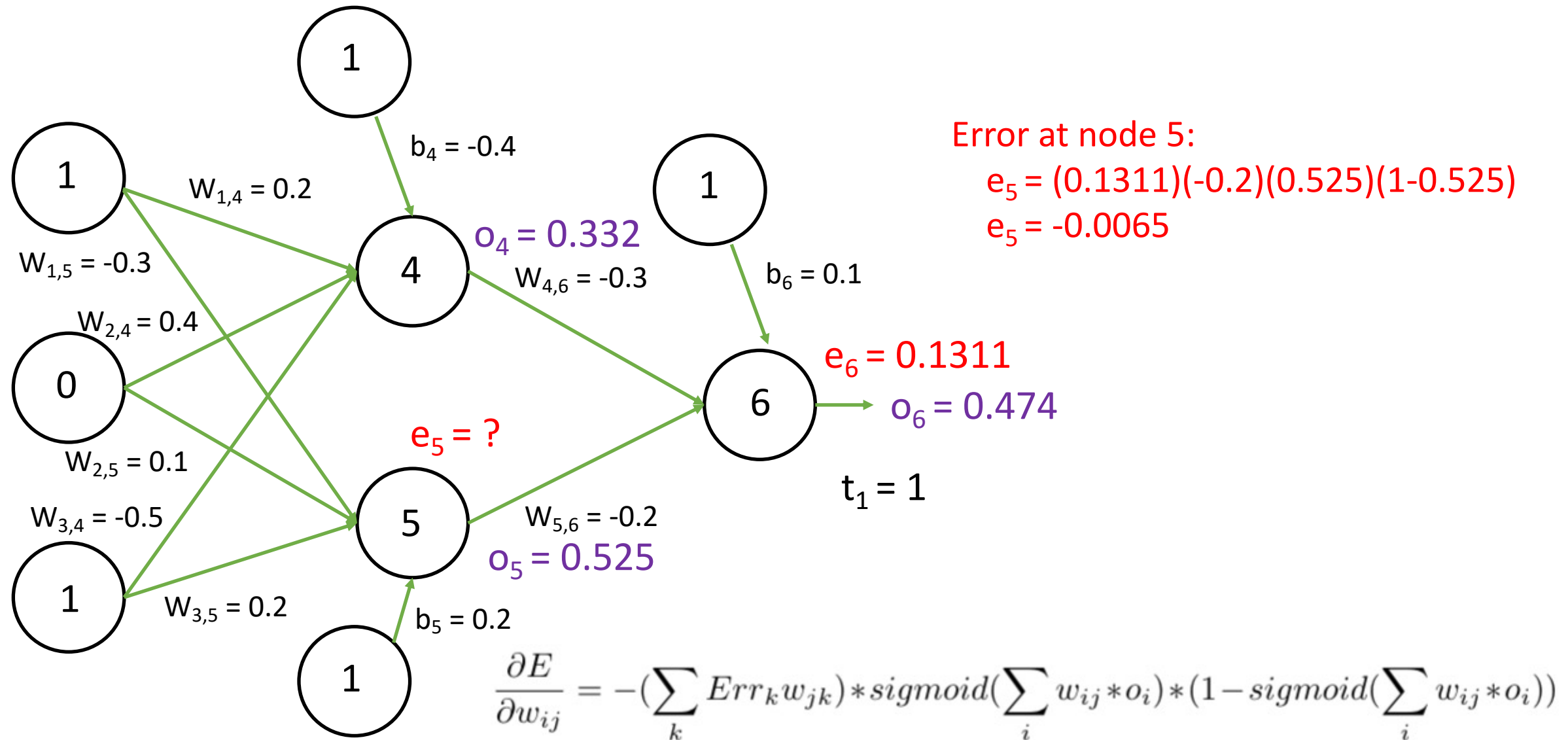
Example: Step 2 – Backward Pass



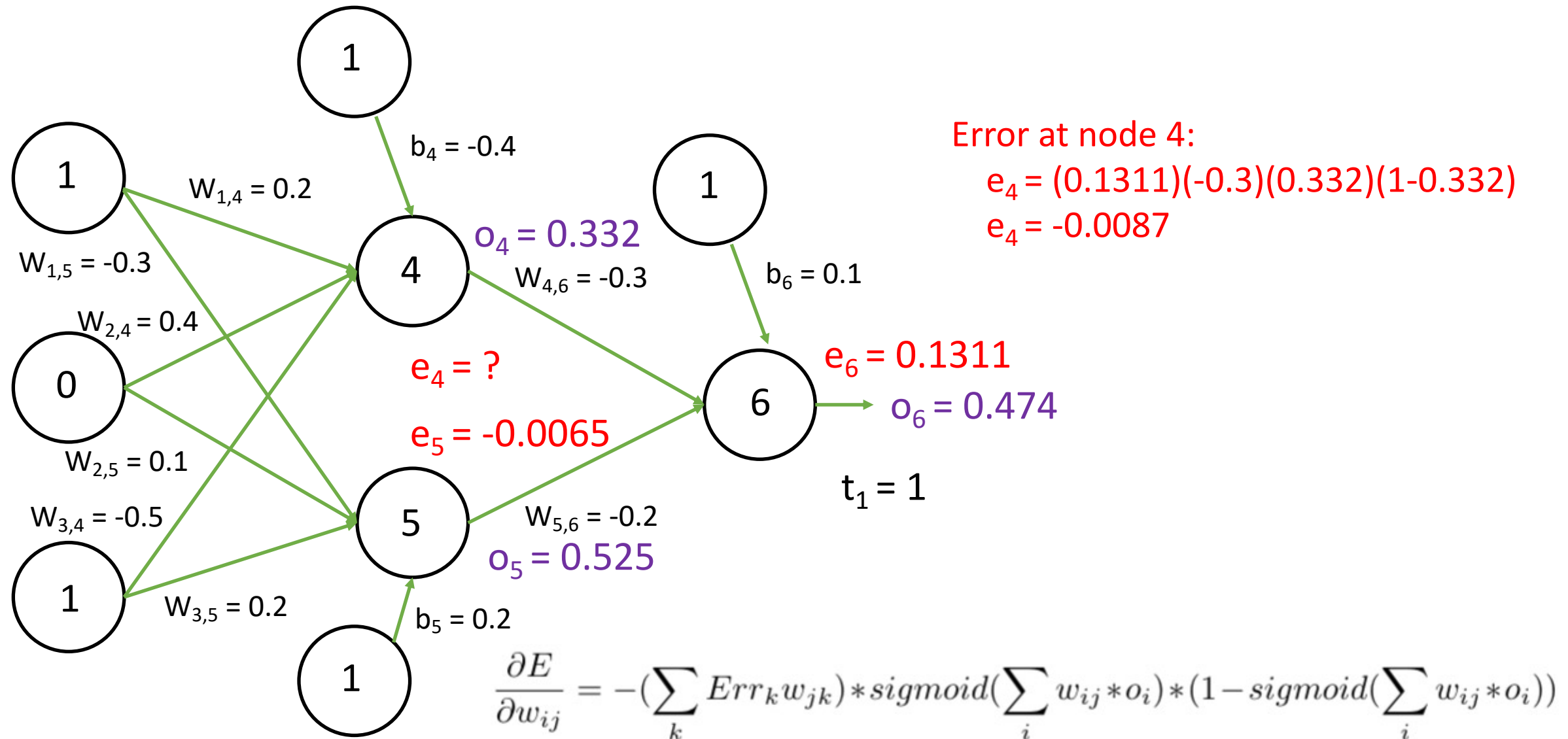
Example: Step 2 – Backward Pass



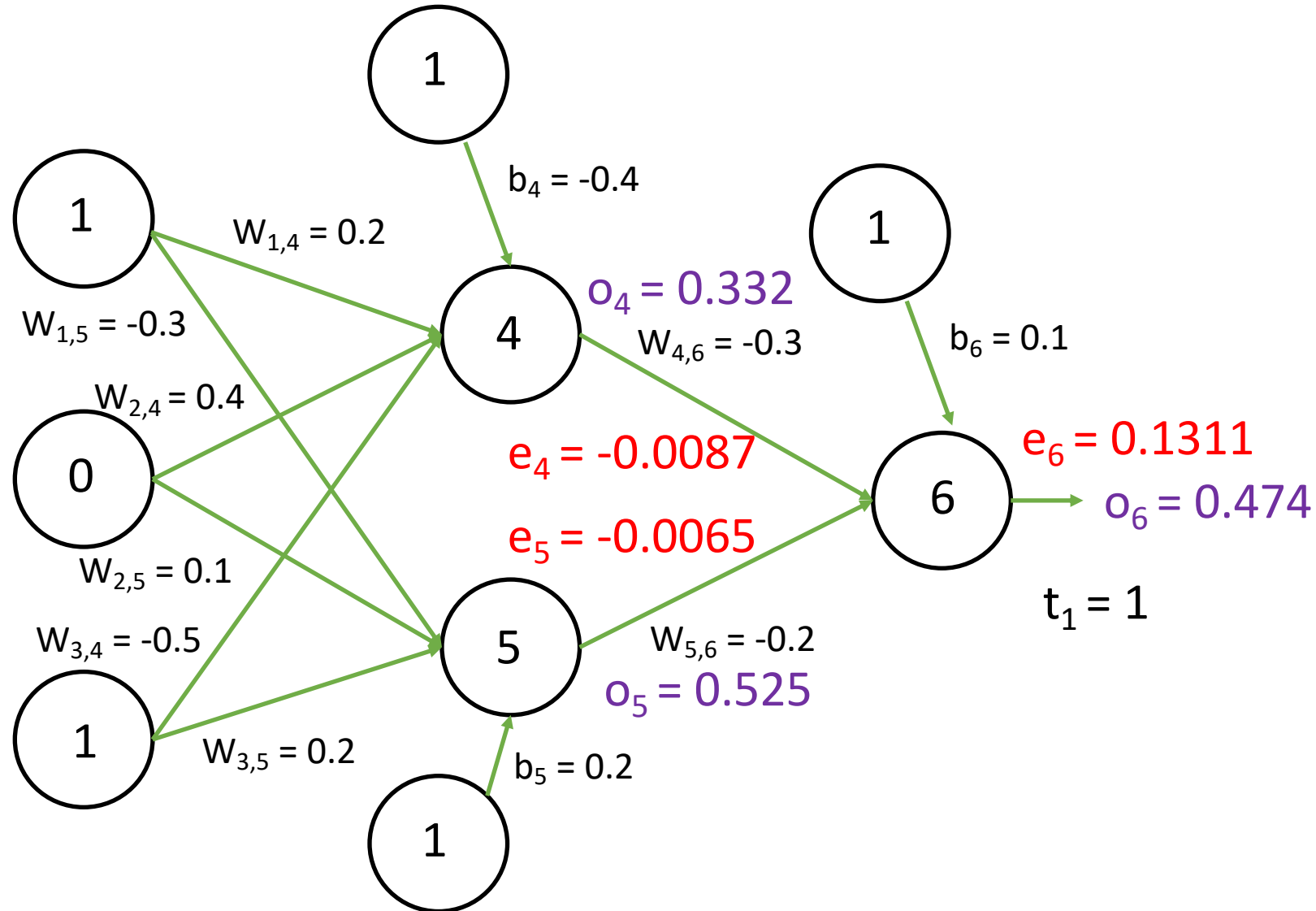
Example: Step 2 – Backward Pass



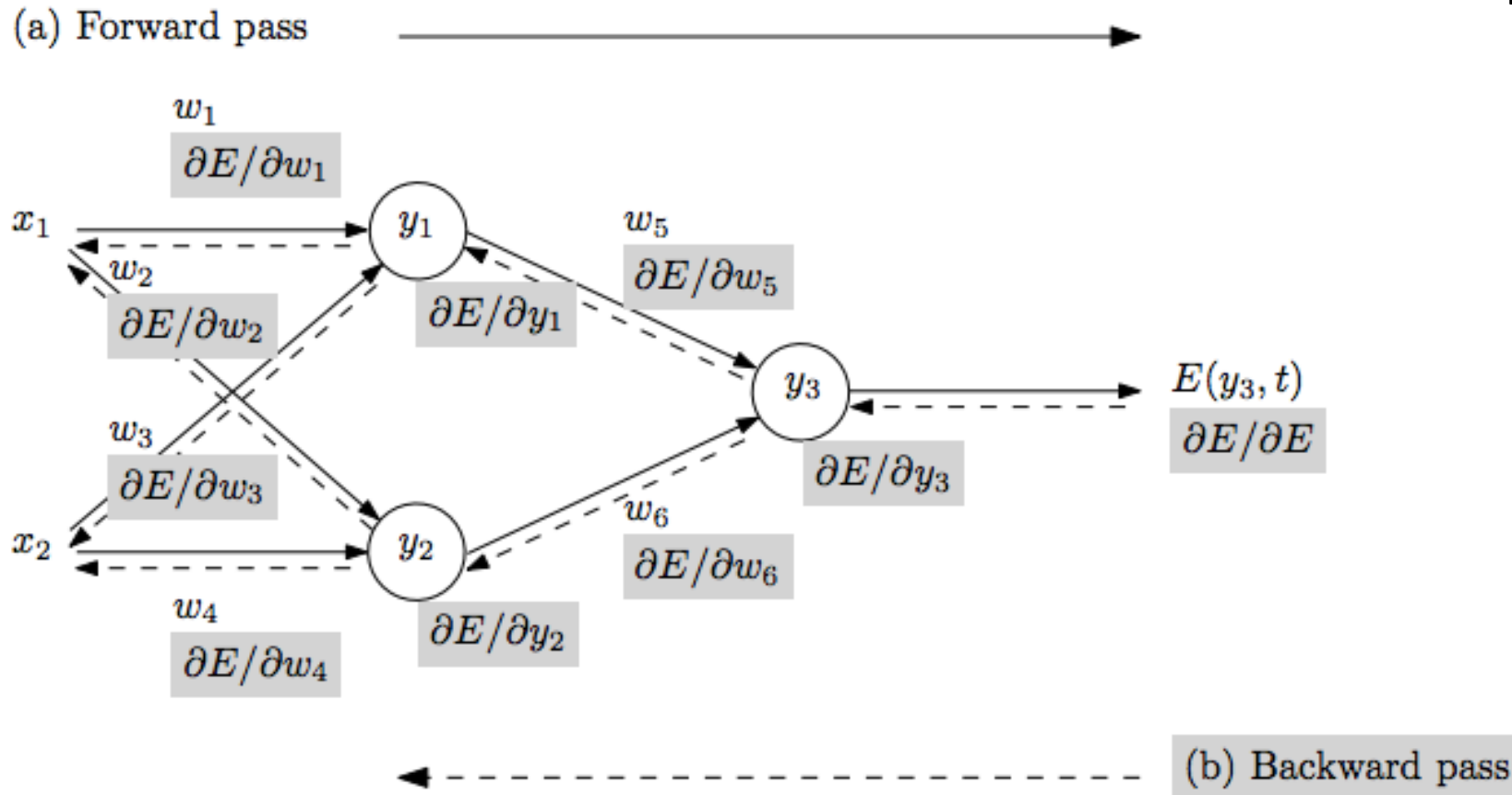
Example: Step 2 – Backward Pass



Example: Step 2 – Backward Pass

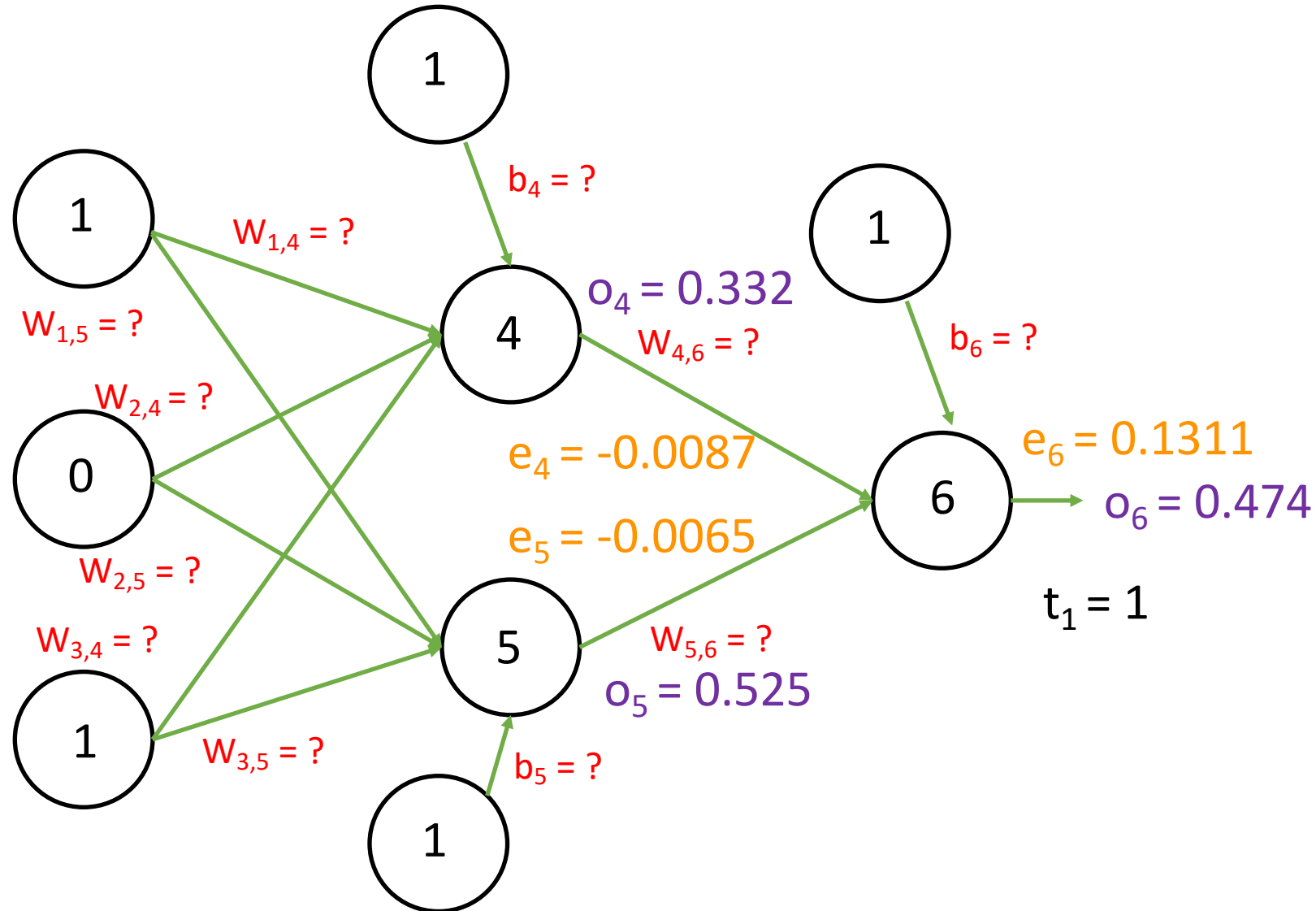


Training: How Neural Networks Learn

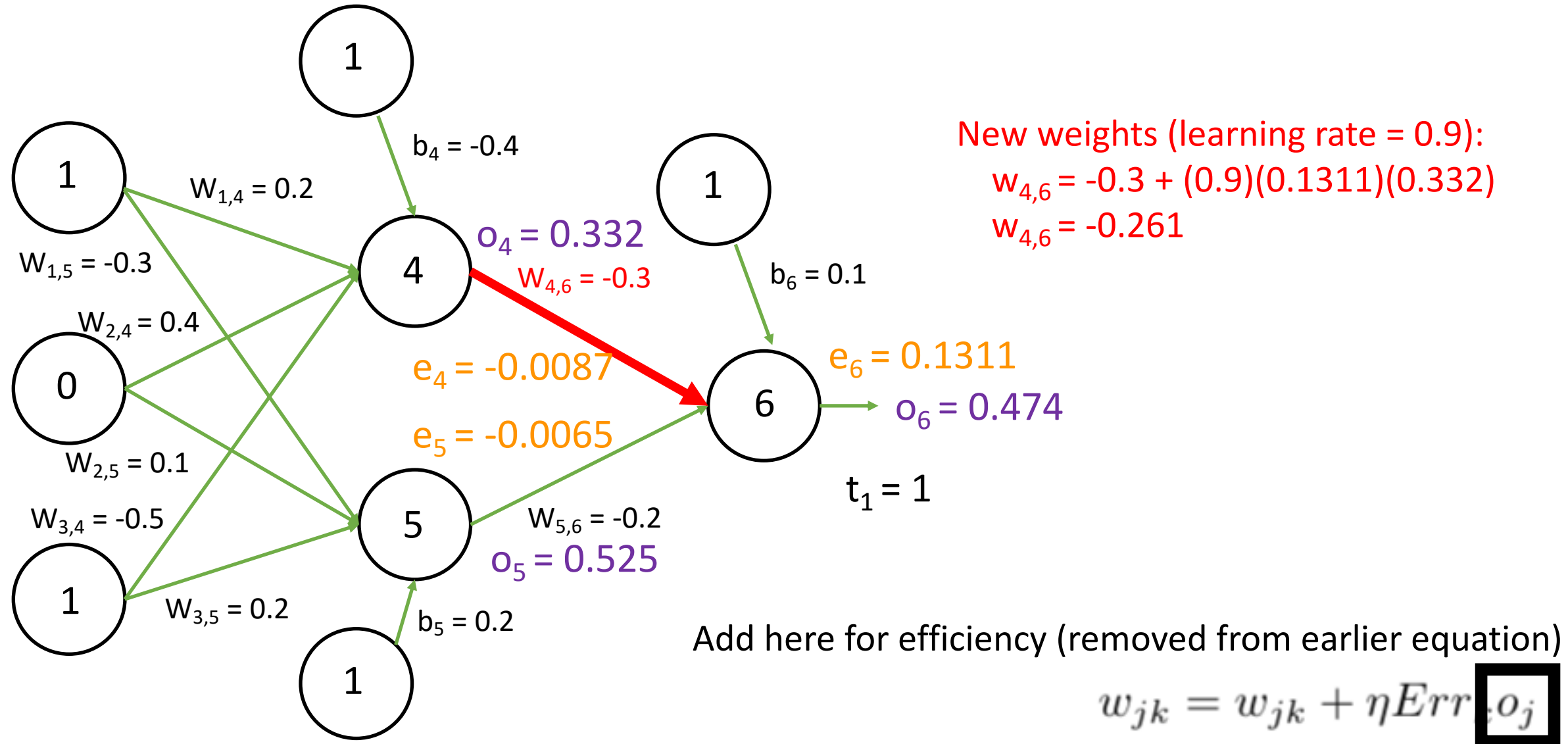


- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make prediction
 2. Quantify the dissatisfaction with a model's results on the training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. **Update each parameter using calculated gradients**

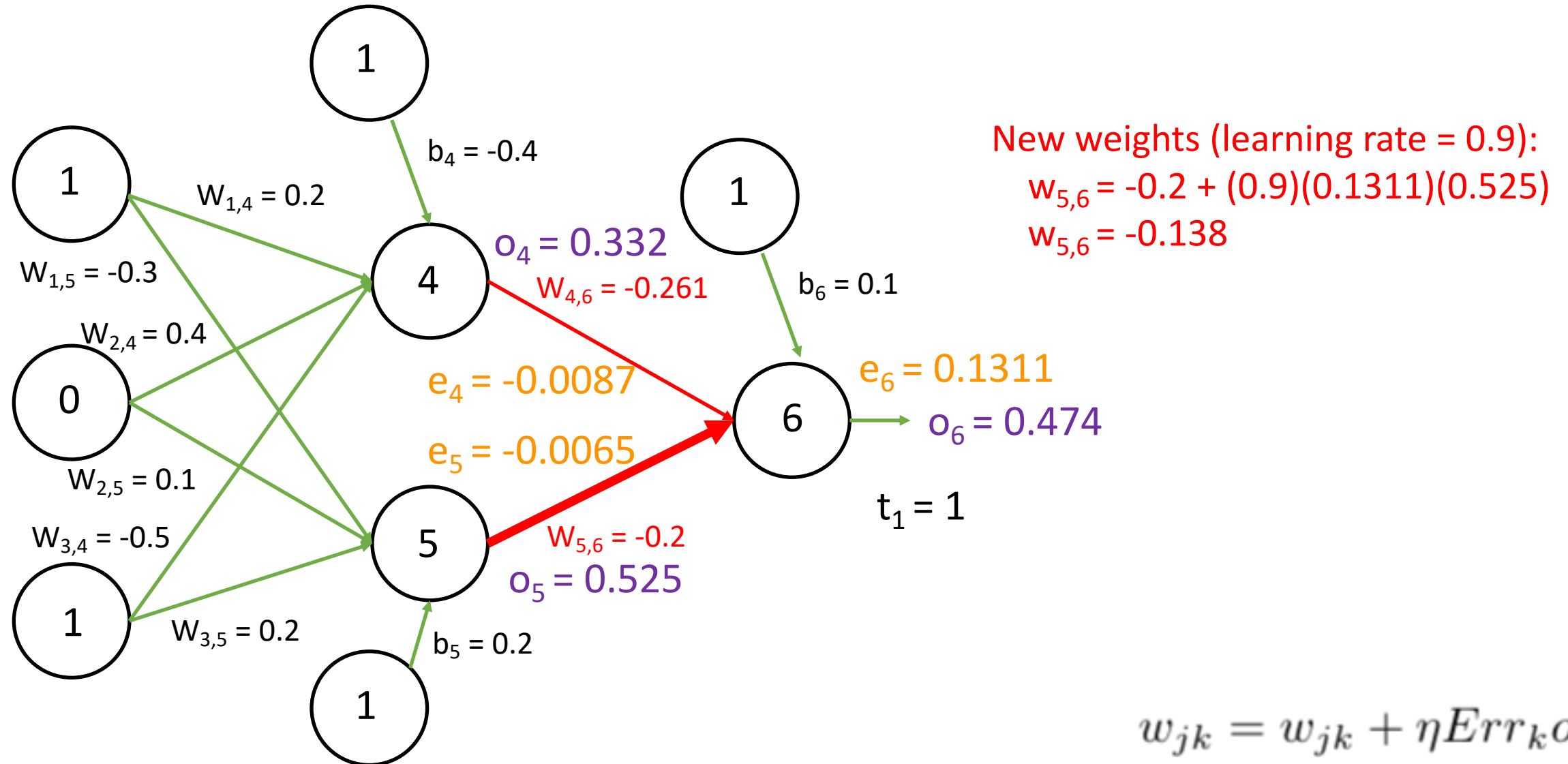
Example: Step 3 – Update Weights



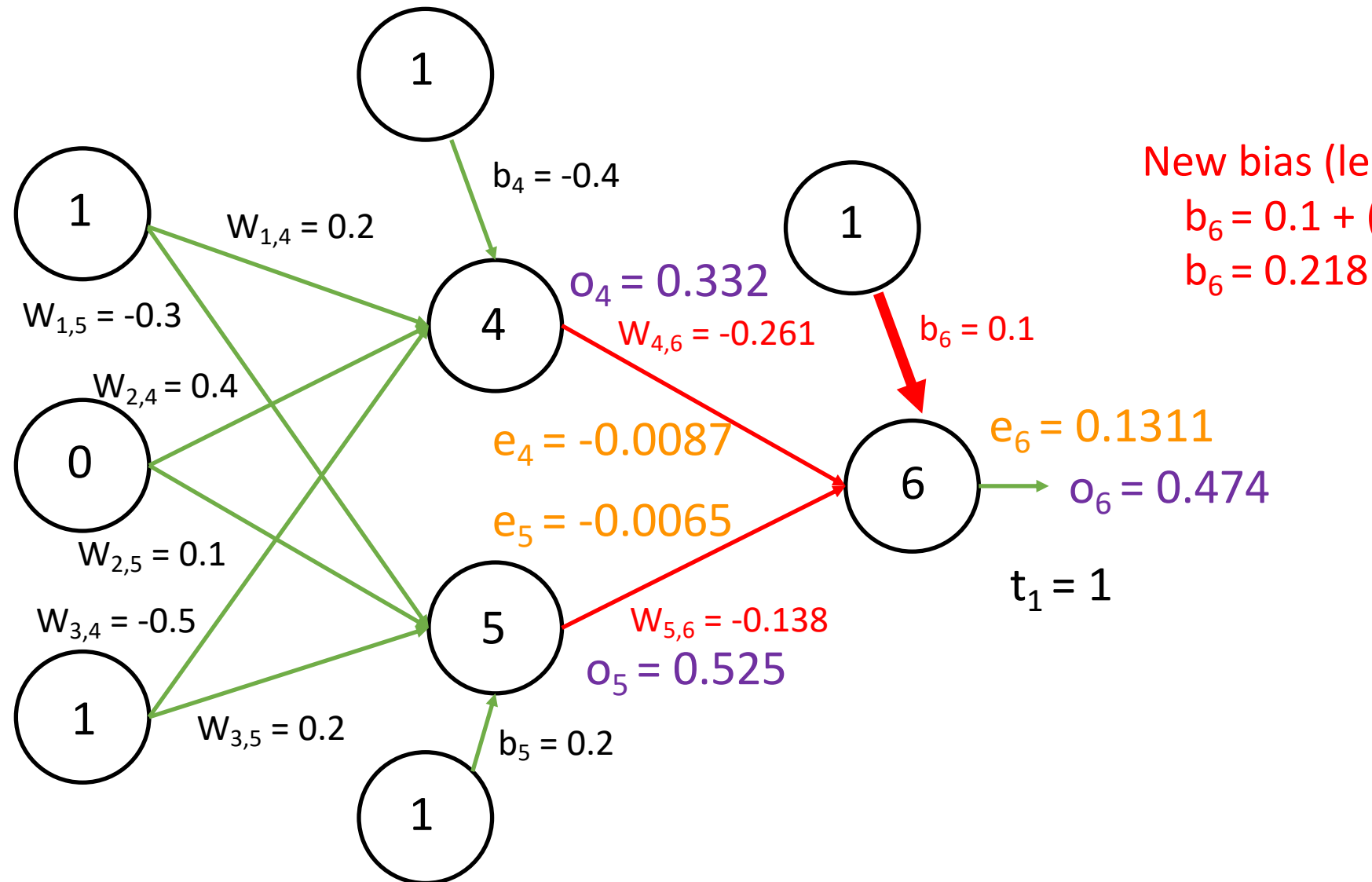
Example: Step 3 – Update Weights



Example: Step 3 – Update Weights

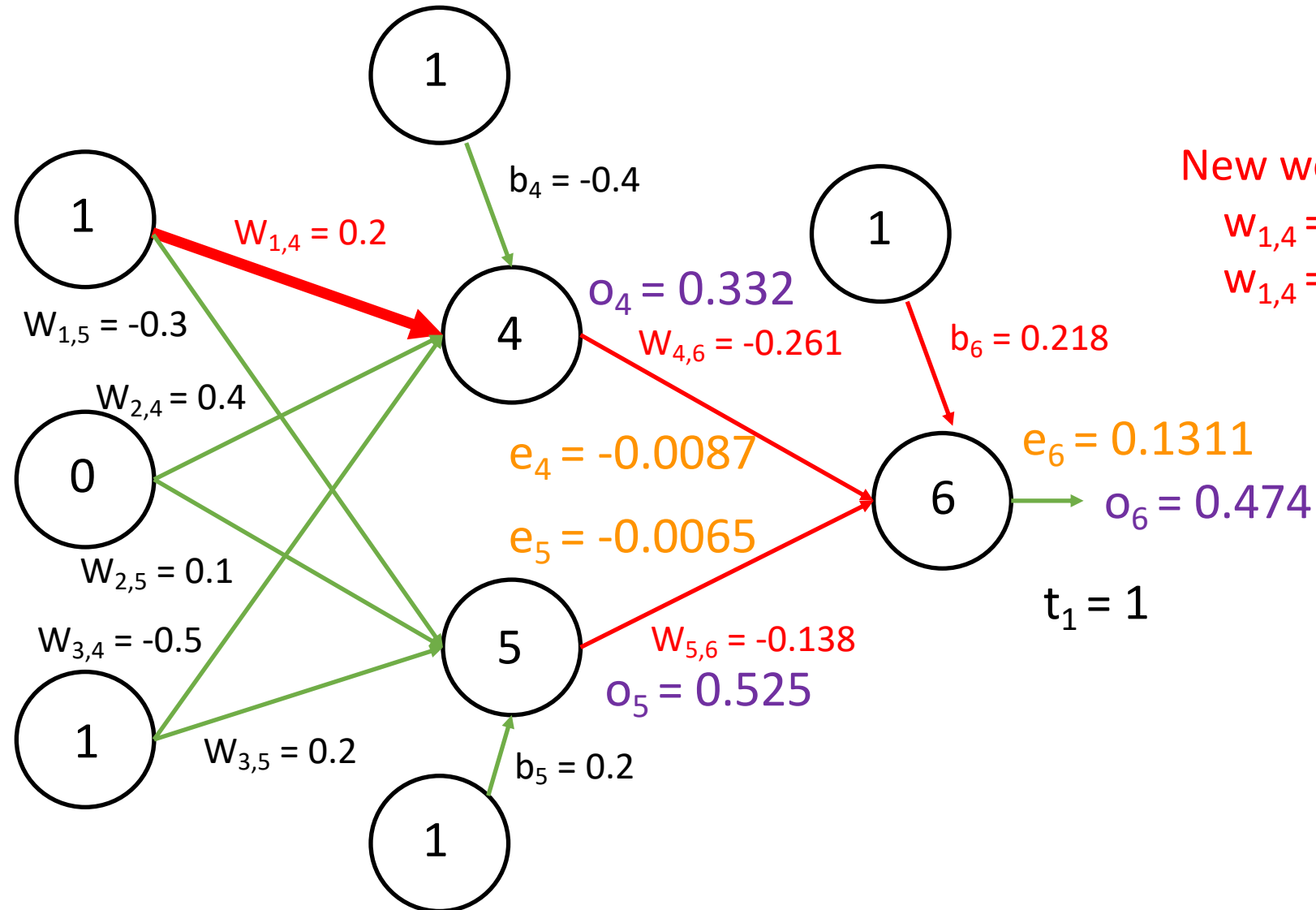


Example: Step 3 – Update Weights



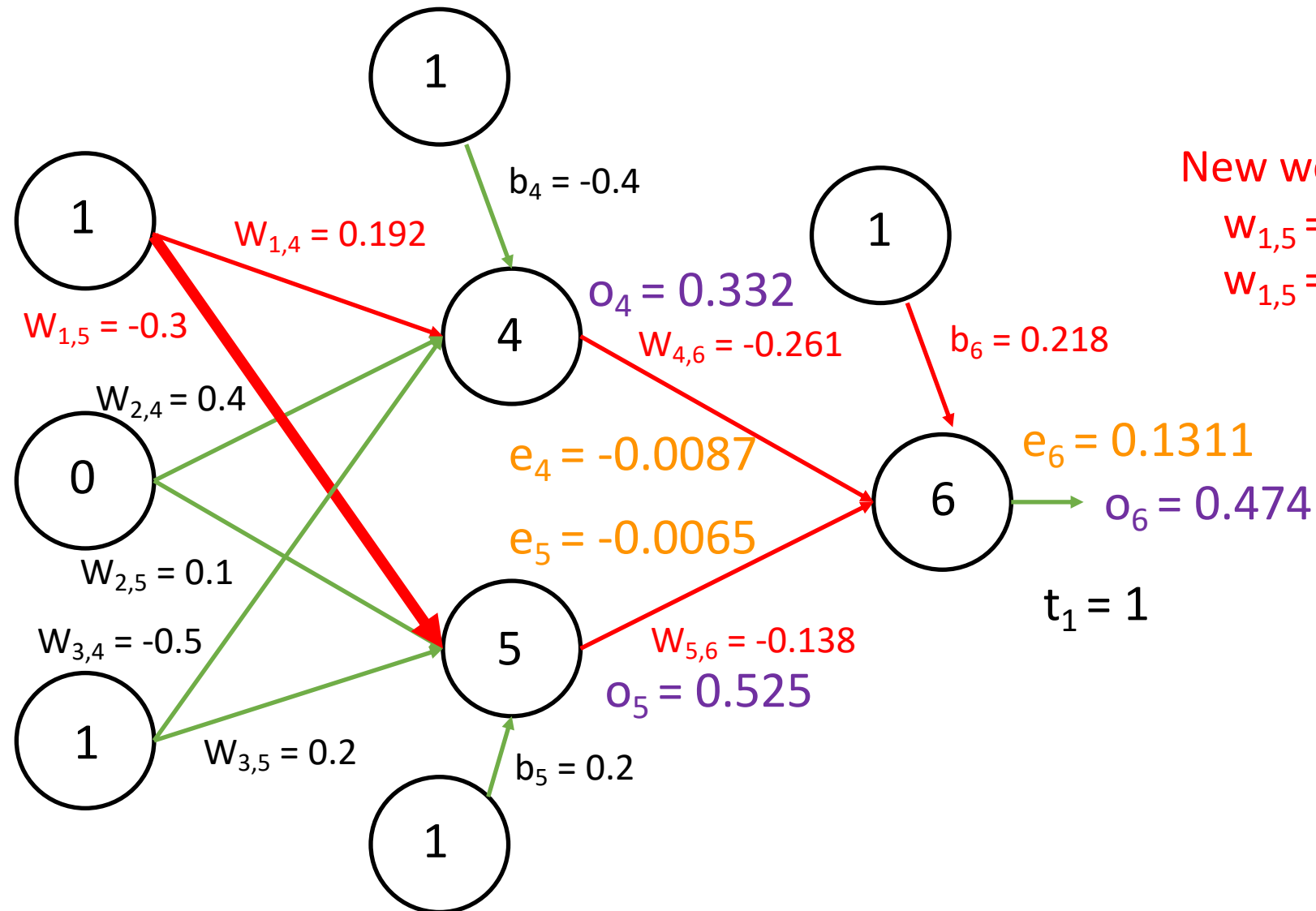
$$b_k = b_k + \eta Err_k$$

Example: Step 3 – Update Weights



$$w_{jk} = w_{jk} + \eta Err_k o_j$$

Example: Step 3 – Update Weights



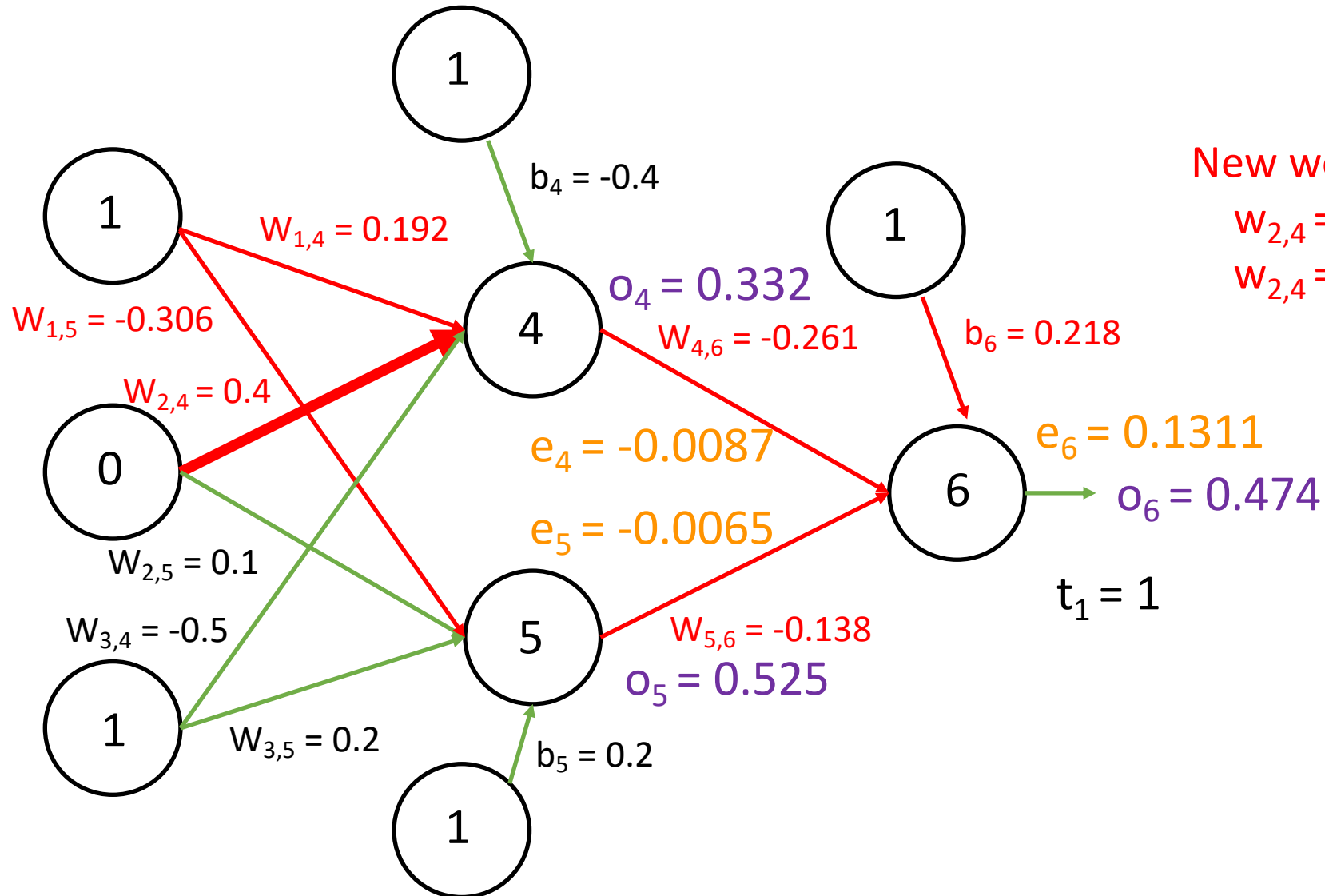
New weights (learning rate = 0.9):

$$w_{1,5} = -0.3 + (0.9)(-0.0065)(1)$$

$$w_{1,5} = -0.306$$

$$w_{jk} = w_{jk} + \eta Err_k o_j$$

Example: Step 3 – Update Weights



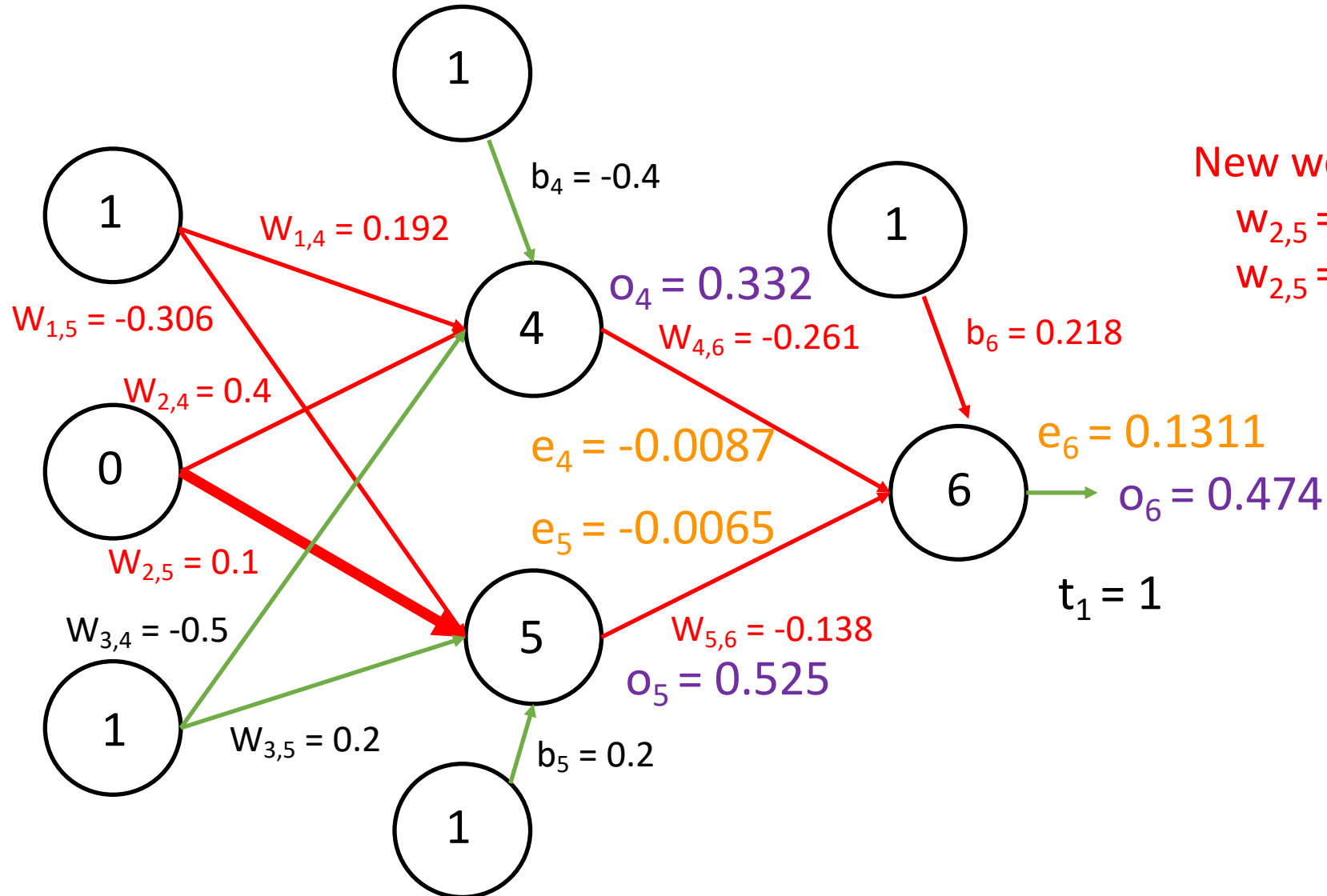
New weights (learning rate = 0.9):

$$w_{2,4} = 0.4 + (0.9)(-0.0087)(0)$$

$$w_{2,4} = 0.4$$

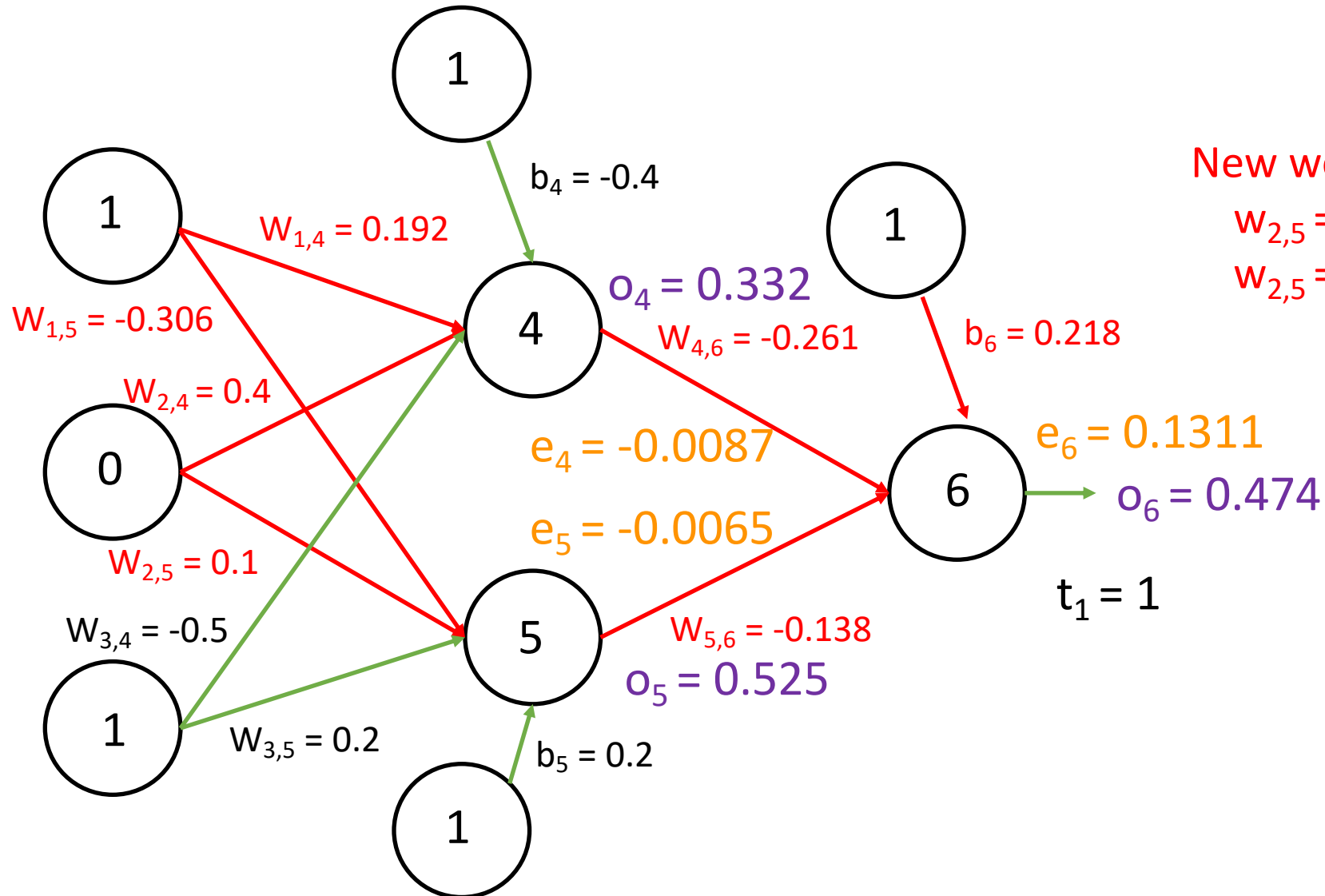
$$w_{jk} = w_{jk} + \eta Err_k o_j$$

Example: Step 3 – Update Weights



$$w_{jk} = w_{jk} + \eta Err_k o_j$$

Example: Step 3 – Update Weights



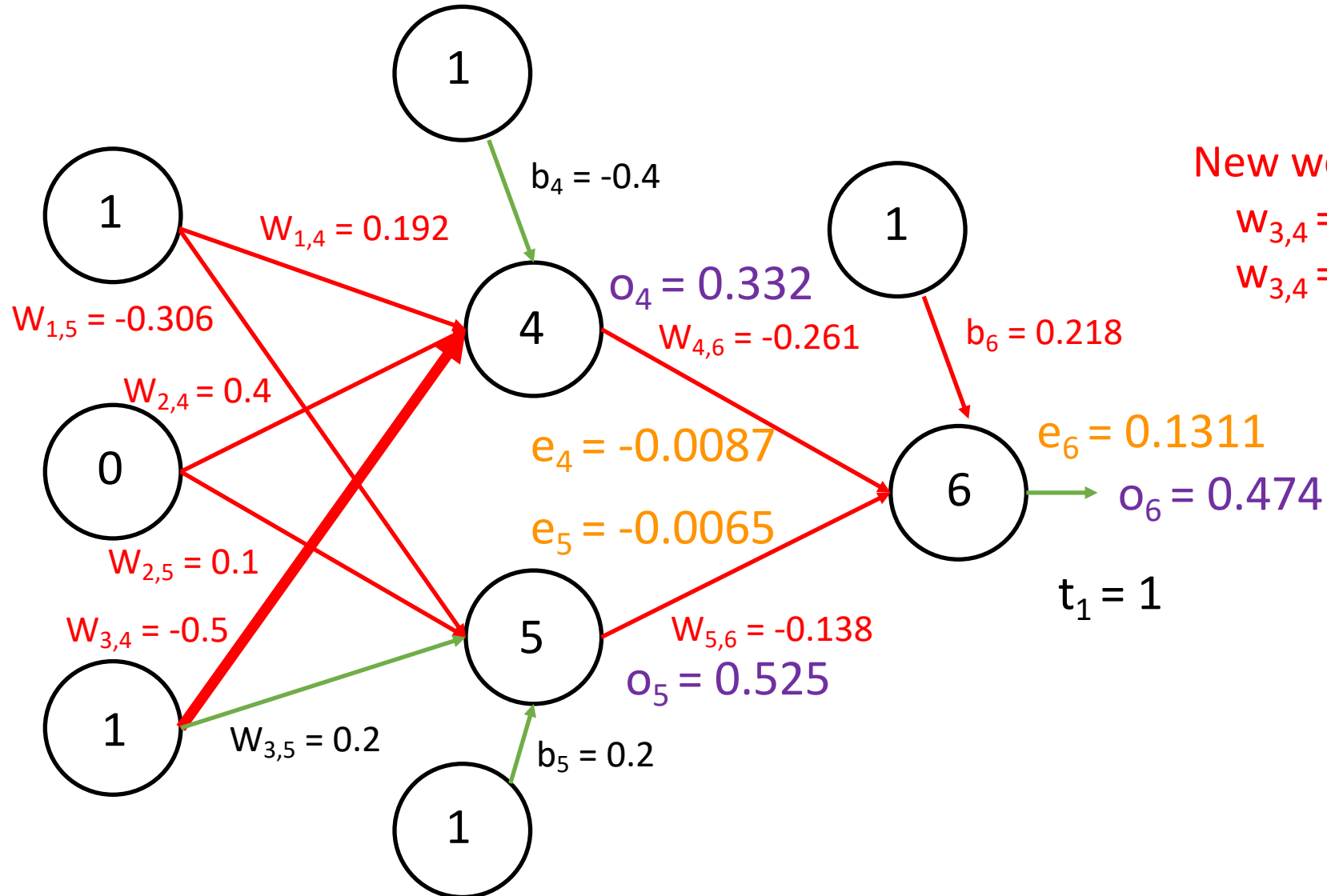
New weights (learning rate = 0.9):

$$w_{2,5} = 0.1 + (0.9)(-0.0065)(0)$$

$$w_{2,5} = 0.1$$

$$w_{jk} = w_{jk} + \eta Err_k o_j$$

Example: Step 3 – Update Weights



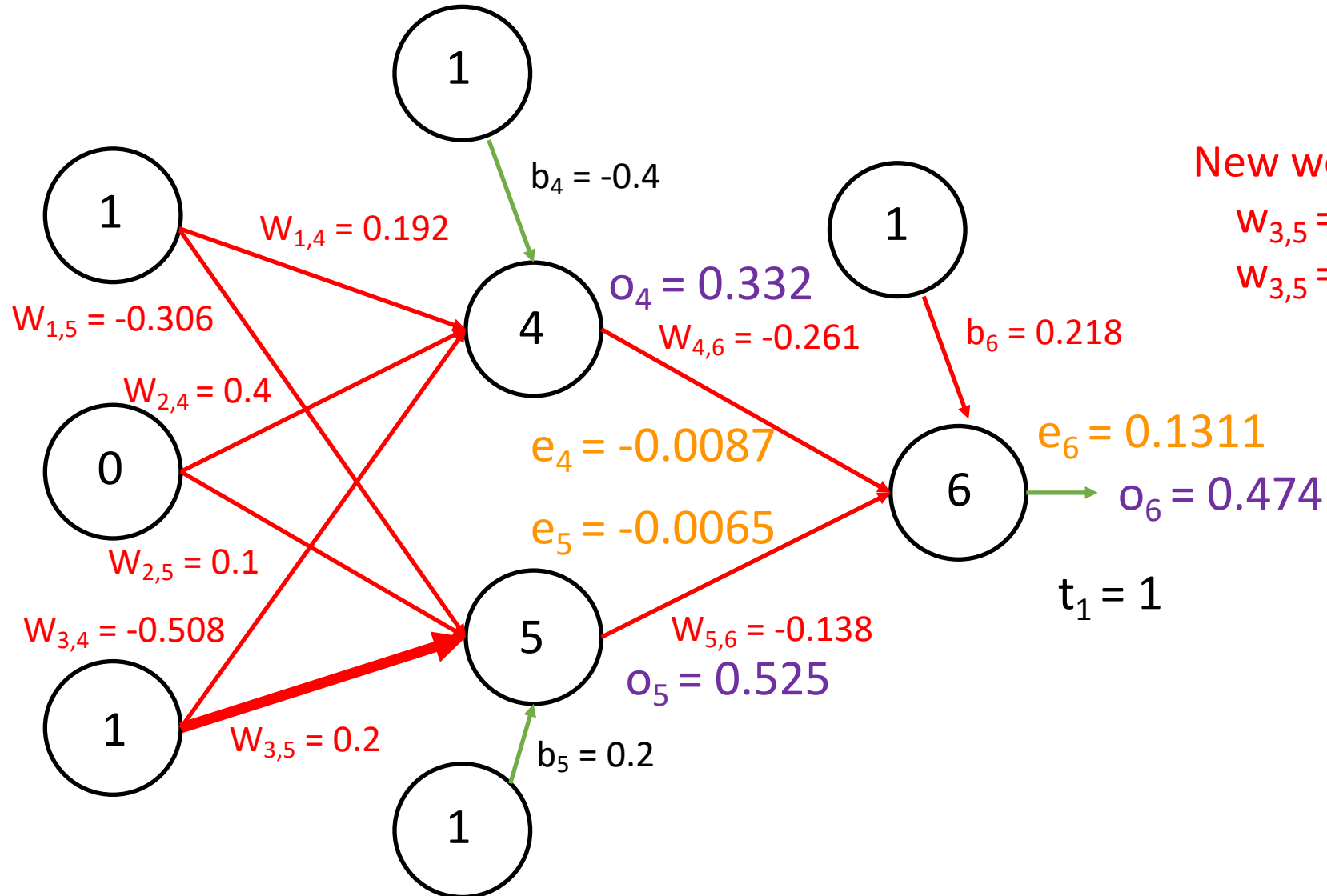
New weights (learning rate = 0.9):

$$w_{3,4} = -0.5 + (0.9)(-0.0087)(1)$$

$$w_{3,4} = -0.508$$

$$w_{jk} = w_{jk} + \eta Err_k o_j$$

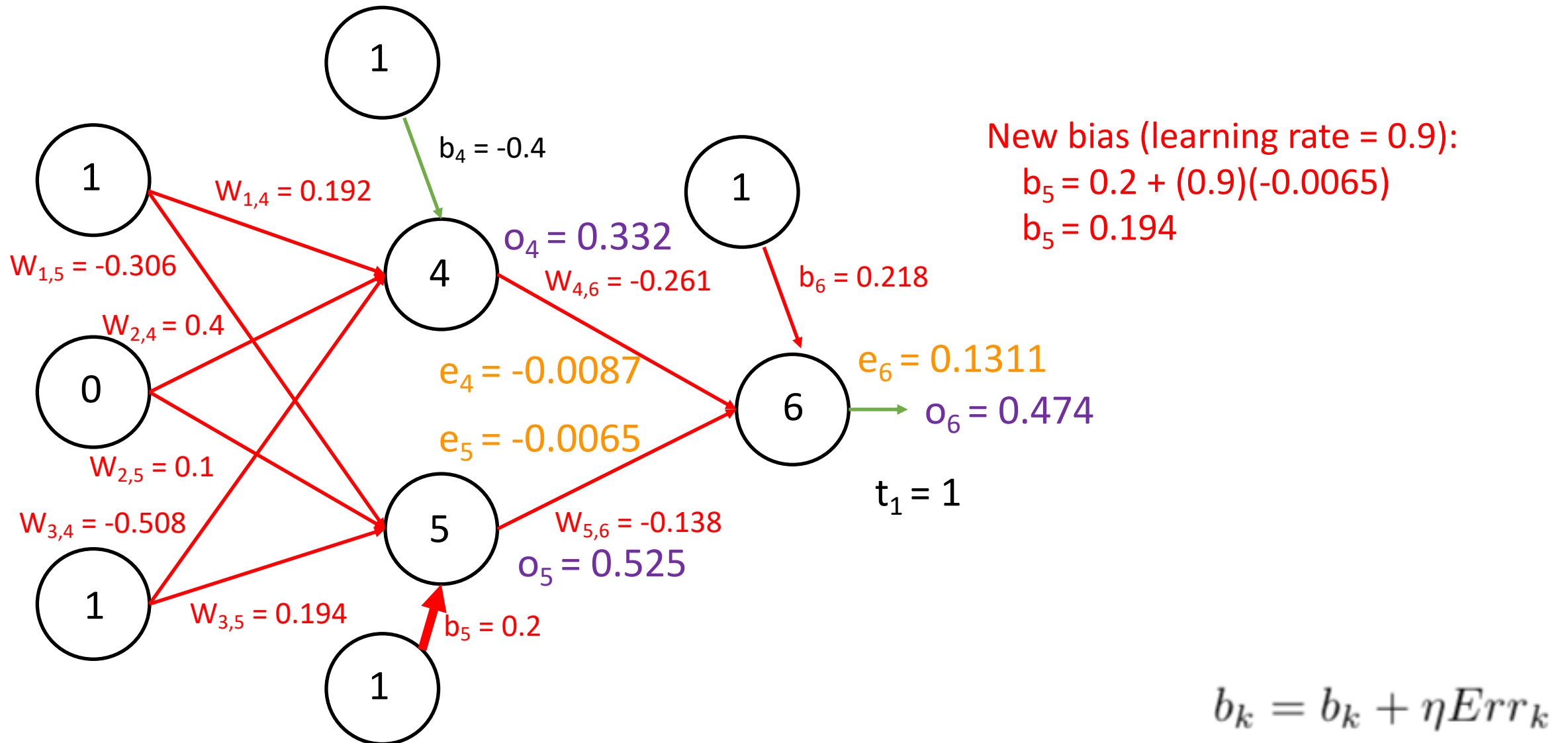
Example: Step 3 – Update Weights



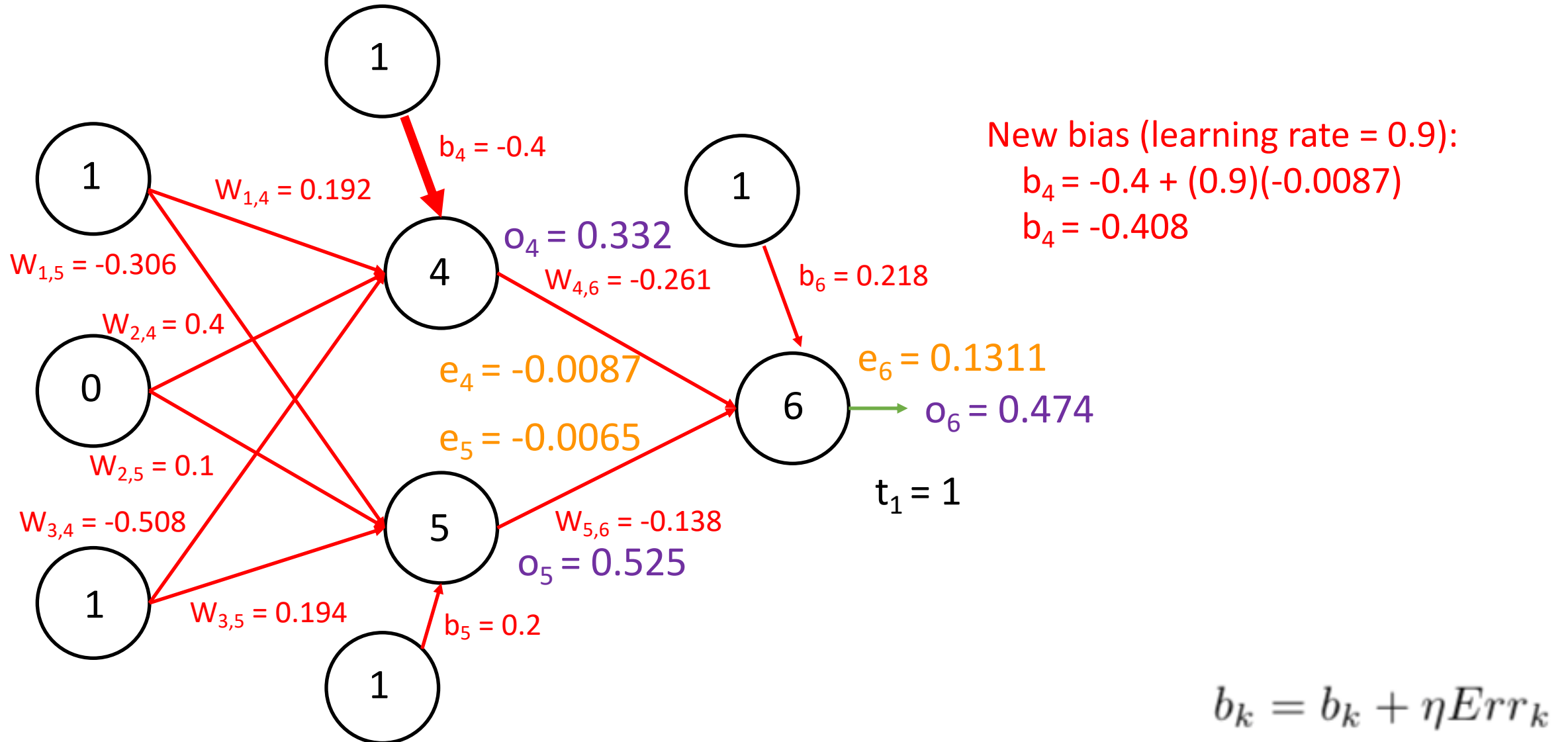
New weights (learning rate = 0.9):
 $w_{3,5} = 0.2 + (0.9)(-0.0065)(1)$
 $w_{3,5} = 0.194$

$$w_{jk} = w_{jk} + \eta Err_k o_j$$

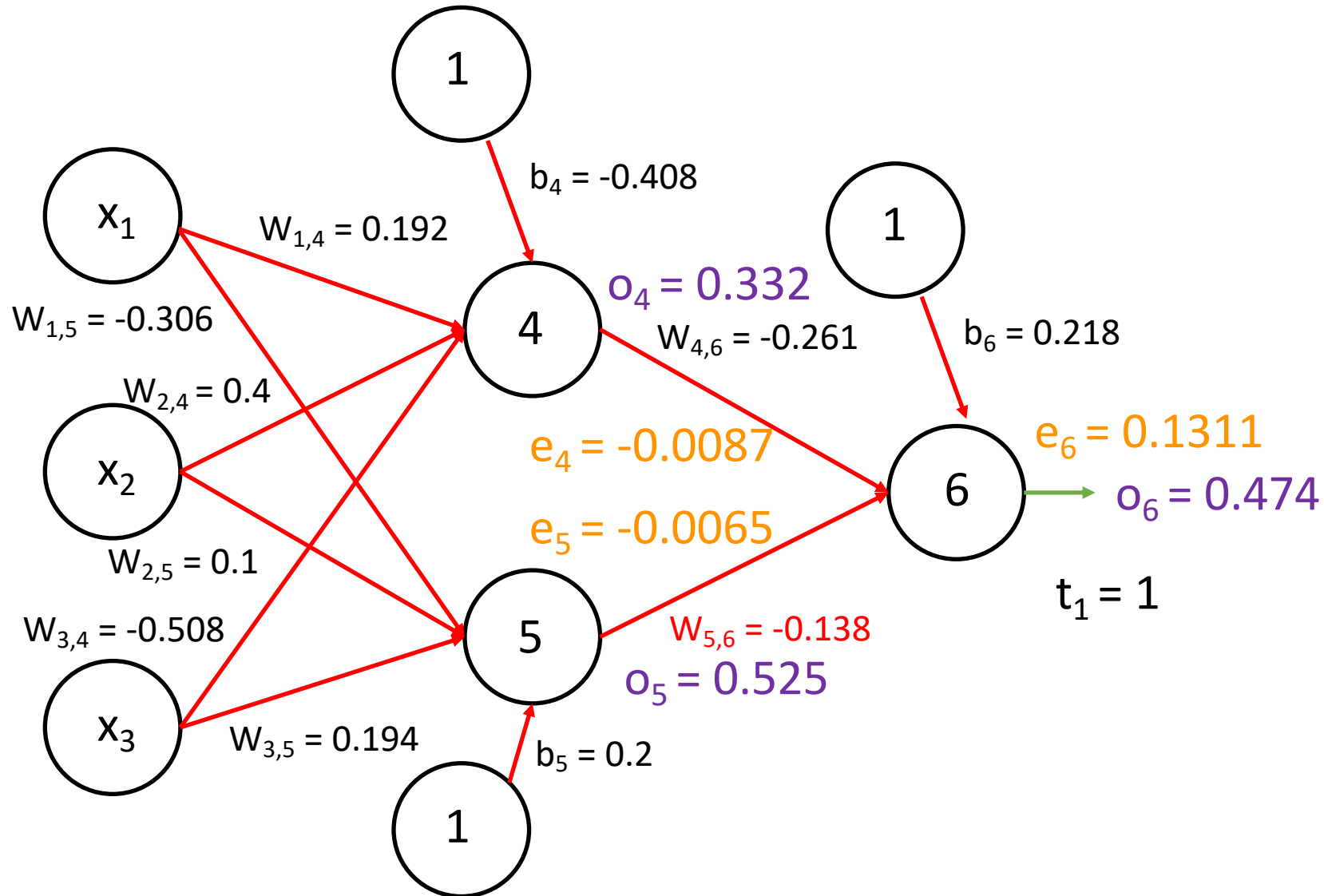
Example: Step 3 – Update Weights



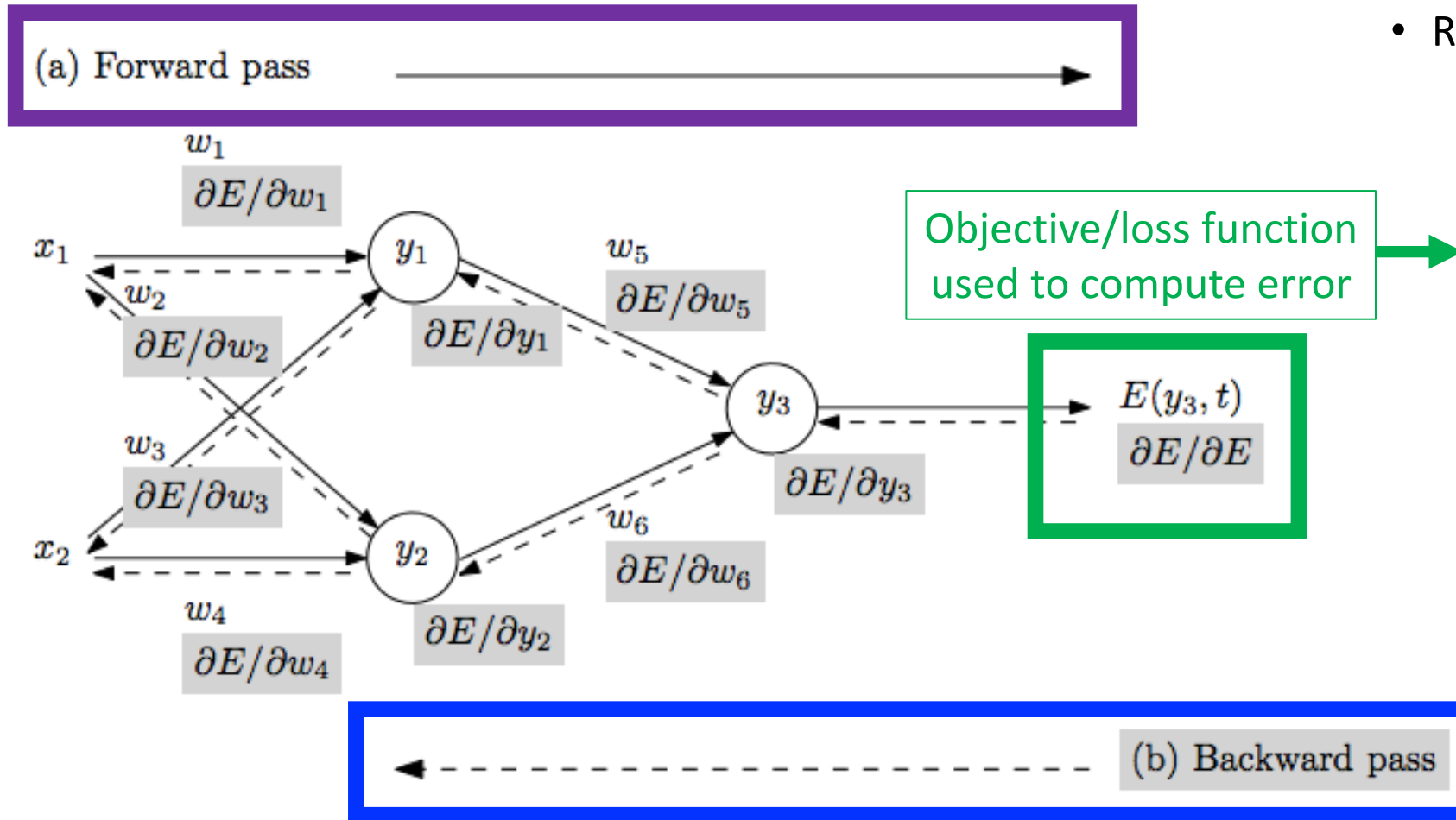
Example: Step 3 – Update Weights



Repeat Steps 1-3 With New Examples



Repeat Steps 1-4 With New Examples



- Repeat until stopping criterion met:

1. **Forward pass:** propagate training data through model to make prediction
2. Quantify the dissatisfaction with a model's results on the training data
3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
4. Update each parameter using calculated gradients

Repeat Steps 1-4 With New Examples

What type of gradient descent was used in the toy example?

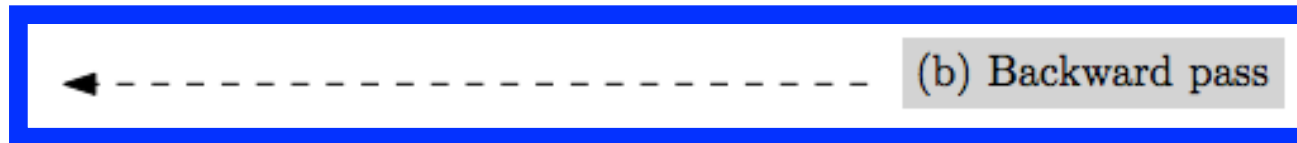
- a. Batch gradient descent
- b. Stochastic gradient descent
- c. Mini-batch gradient descent

- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make prediction
 2. Quantify the dissatisfaction with a model's results on the training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

Training: How Neural Networks Learn

The mean gradient is used for batch and mini-batch gradient descent

- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make prediction
 2. Quantify the dissatisfaction with a model's results on the training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

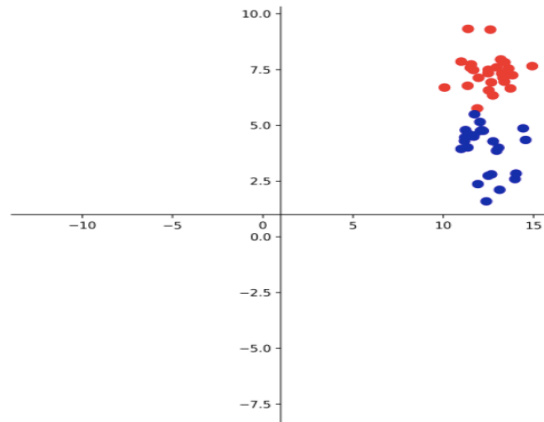


* Practical Detail (More in Future Lectures)

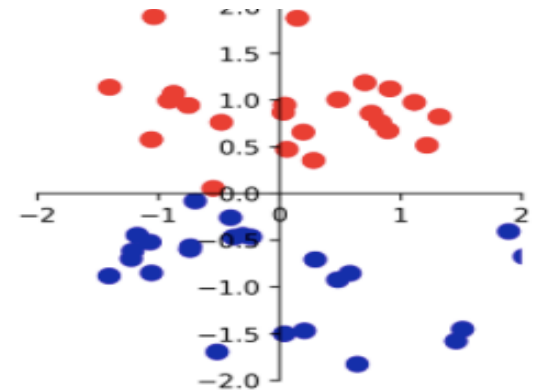
Basic data initialization approach:

- standardize so mean is 0 and standard deviation 1
- simplifies learning

Original data:



Standardized data:



Basic model parameter initialization:

- set weights to random values drawn from Gaussian or uniform distribution
- set biases to 0

* Practical Detail

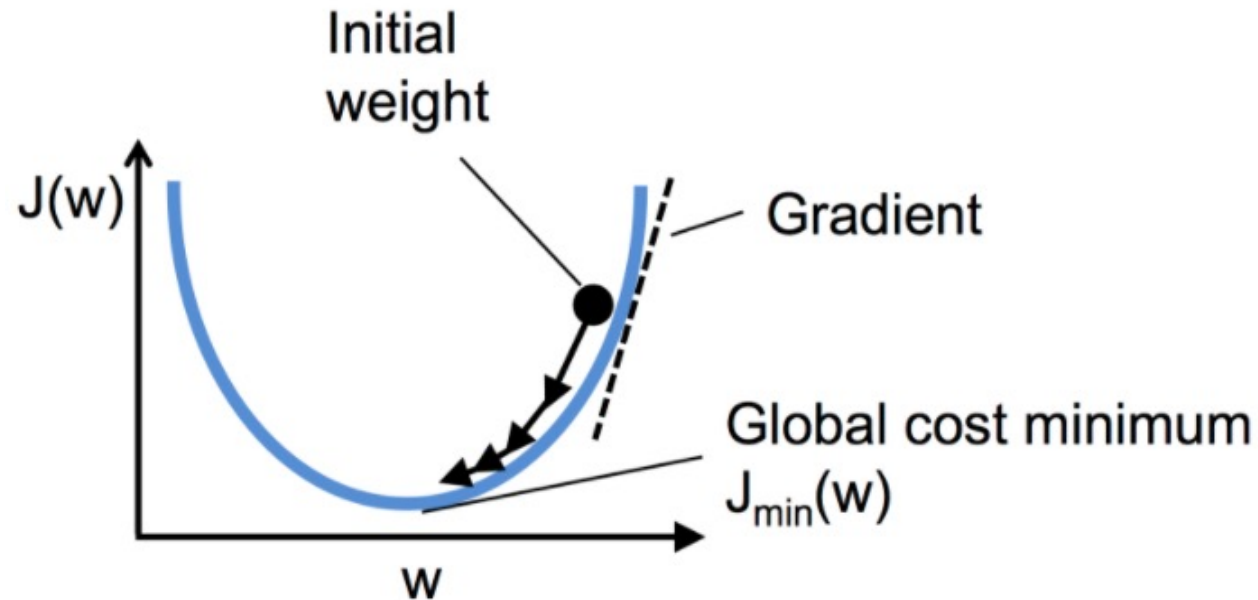
- When training neural networks, optimization function is often used interchangeably with loss function and cost function
 - Subtle nuances are discussed here: <https://www.baeldung.com/cs/cost-vs-loss-vs-objective-function>

Today's Topics

- Objective function: what to learn
- Gradient descent: how to learn
- Training a neural network: optimization
- Gradient descent for activation functions

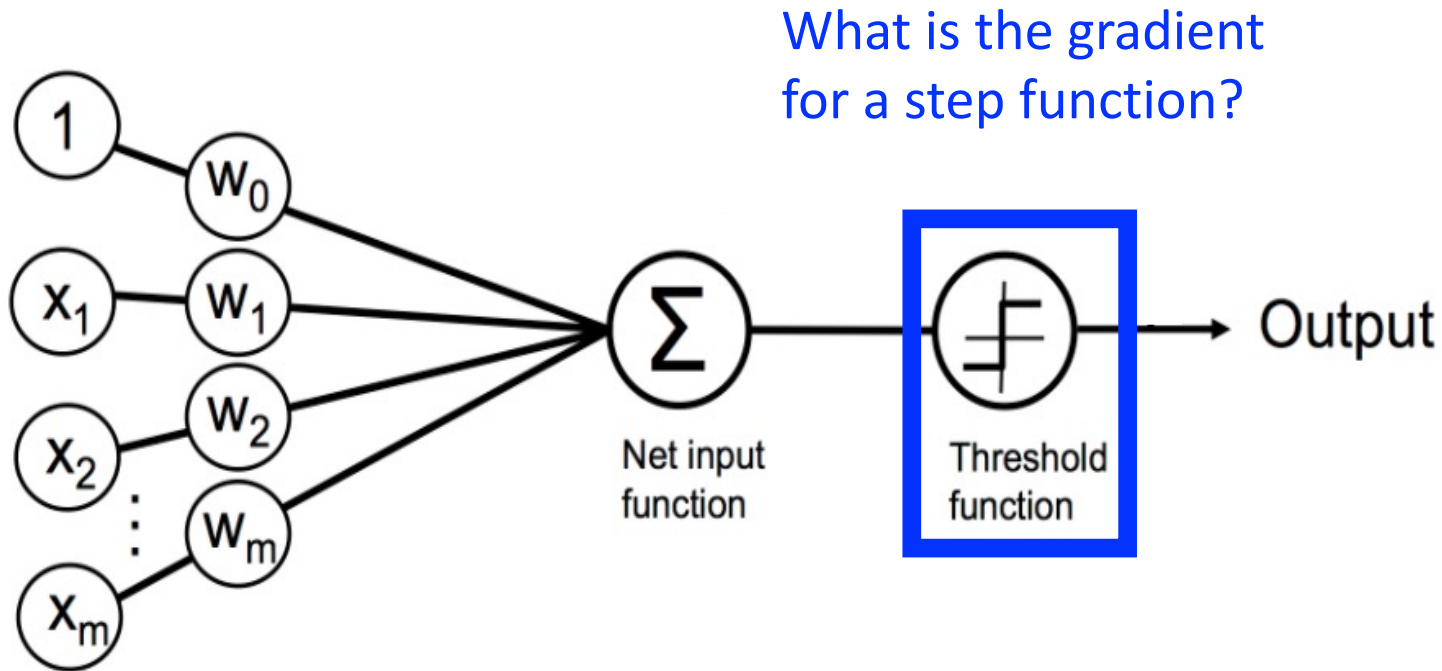
Activation Function Overview

- **Want:** function with a gradient large enough to support efficient learning

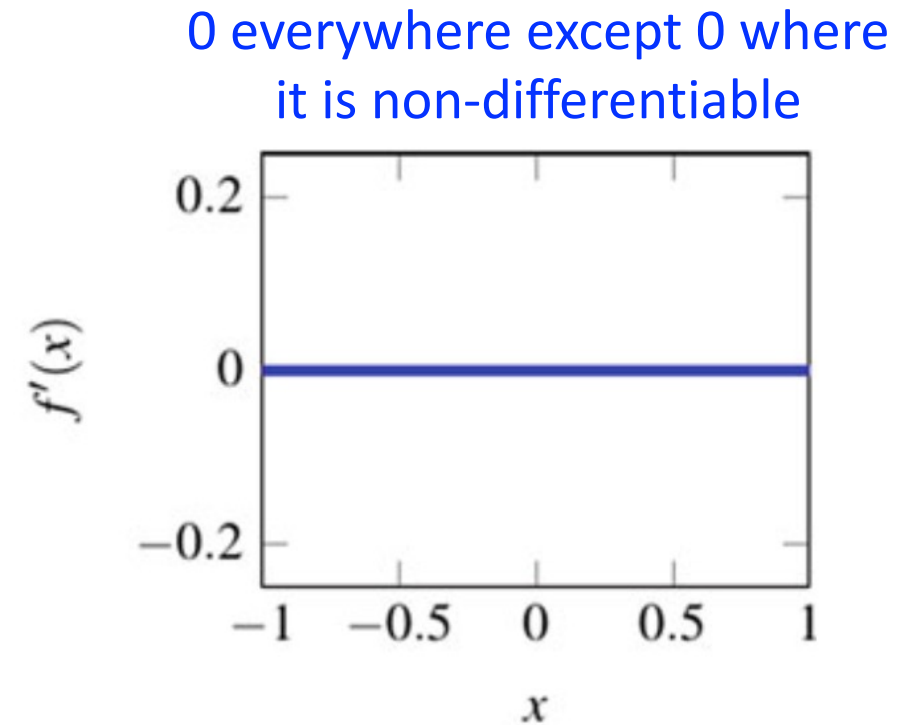


- **Implied requirement:** function should be differentiable

Activation Functions with Gradients: Revisiting Perceptron



Python Machine Learning; Raschka & Mirjalili

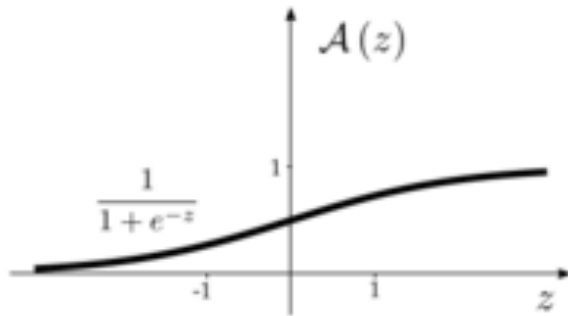


Deep Learning for NLP and Speech Recognition; Kamath, Liu, & Whitaker

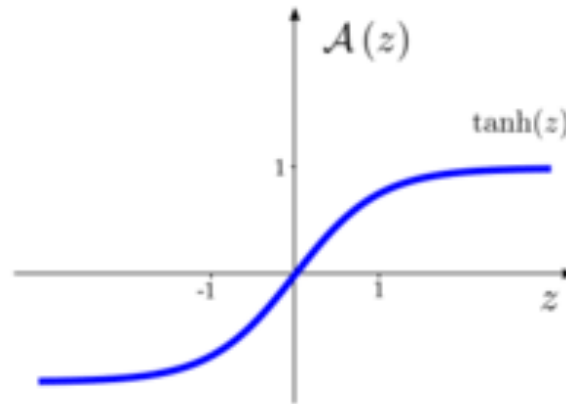
No gradient means model parameters wouldn't change with gradient descent!

Activation Functions with Gradients: Nonlinear Activation Functions

Sigmoid

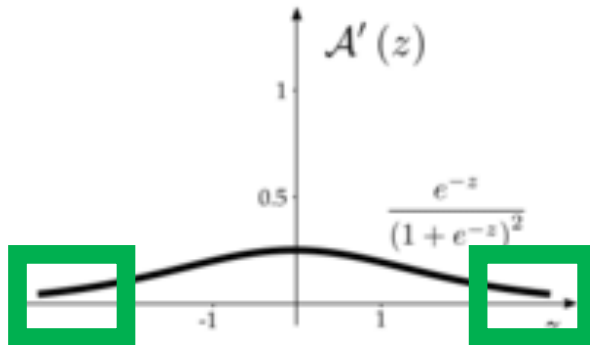


Tanh

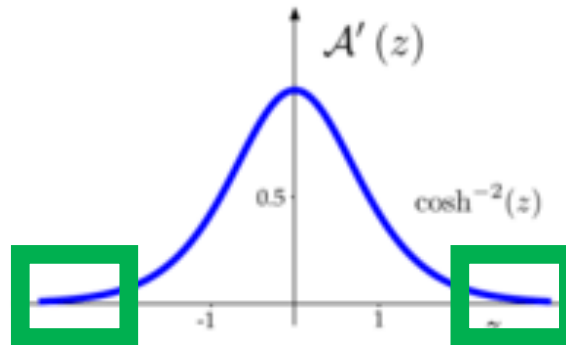


Problem: units with small or large “z” values lead to little/slow learning; why?

$\mathcal{A}'(z)$



$\mathcal{A}'(z)$



Reason: small gradients limit amount model parameters change with gradient descent

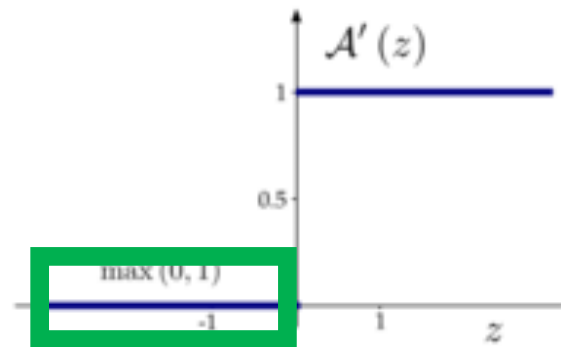
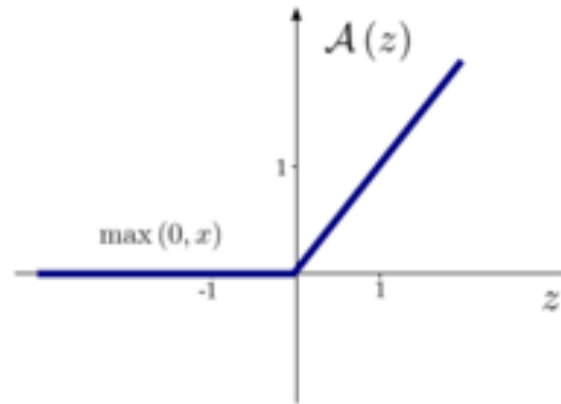
Activation Functions with Gradients: Nonlinear Activation Functions

Advantages:

- Fast to compute
- Large gradient when unit is “firing”

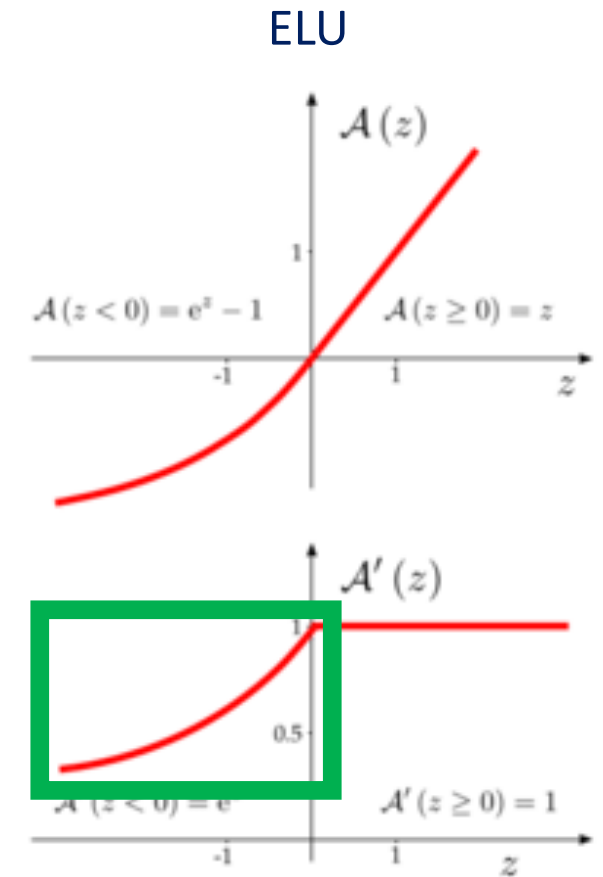
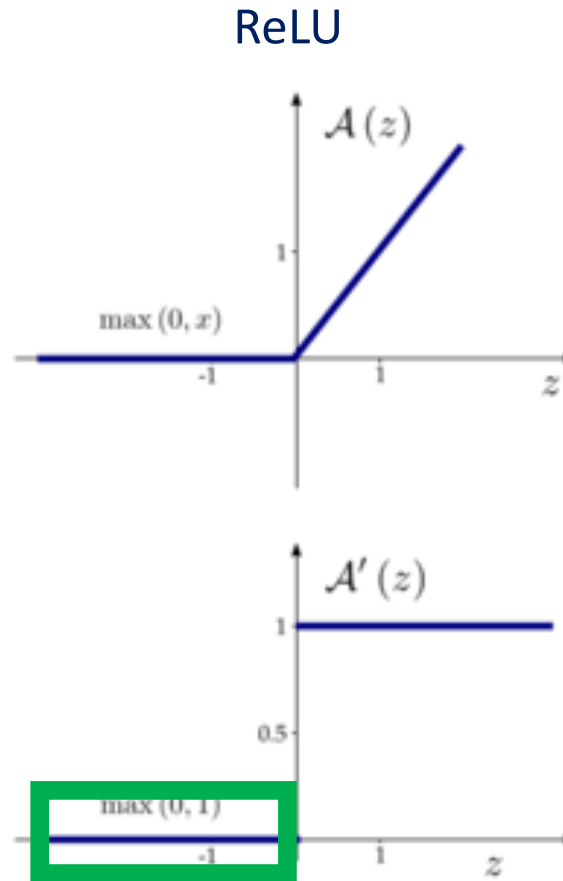
Problem: no gradient
means units can “die”

ReLU



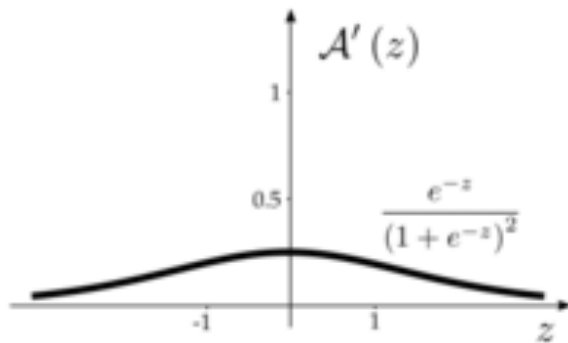
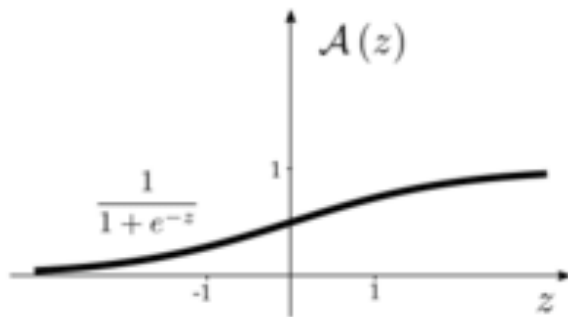
Activation Functions with Gradients: Nonlinear Activation Functions

Can avoid dying units by increasing computational complexity of ELU-based activation functions

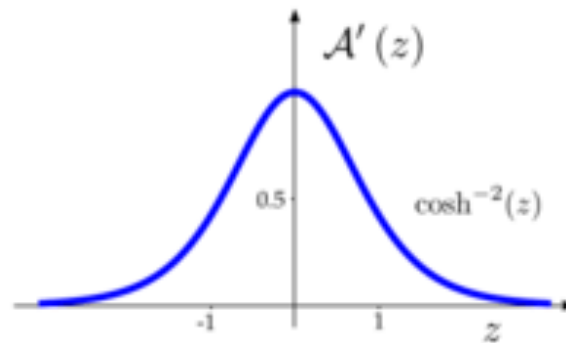
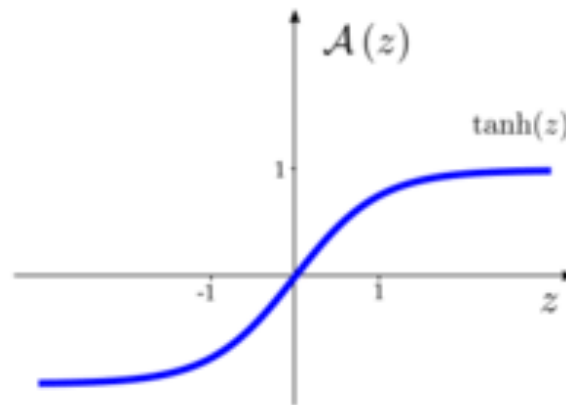


Activation Functions with Gradients: Nonlinear Activation Functions

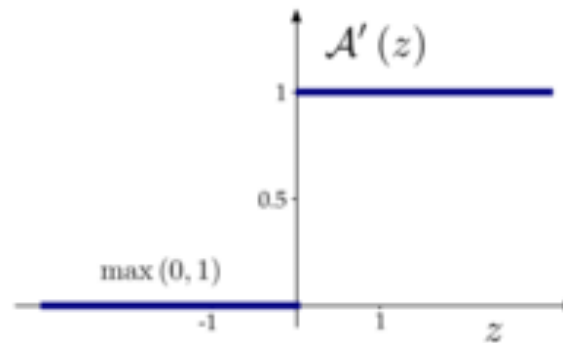
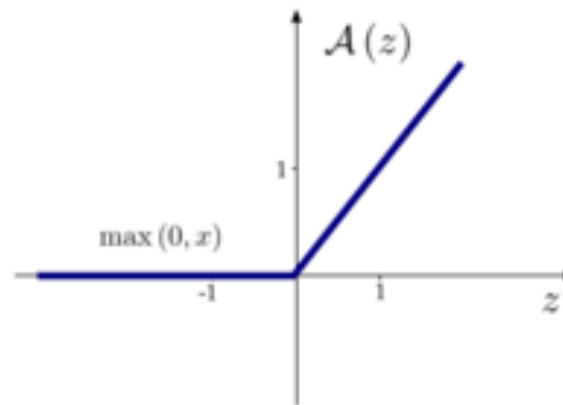
Sigmoid



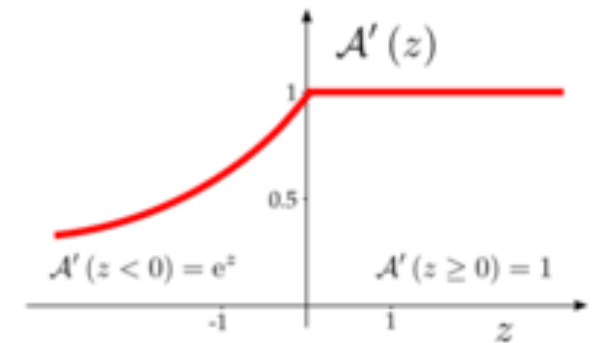
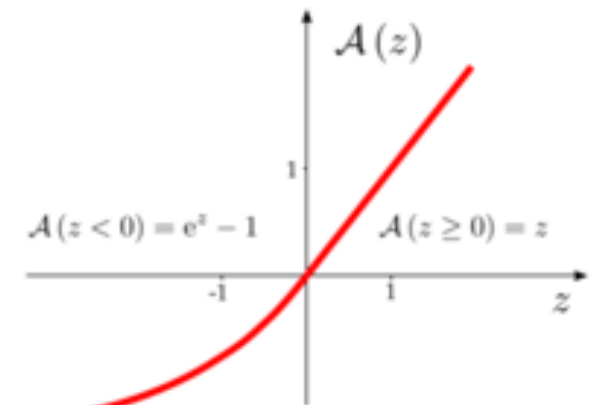
Tanh



ReLU

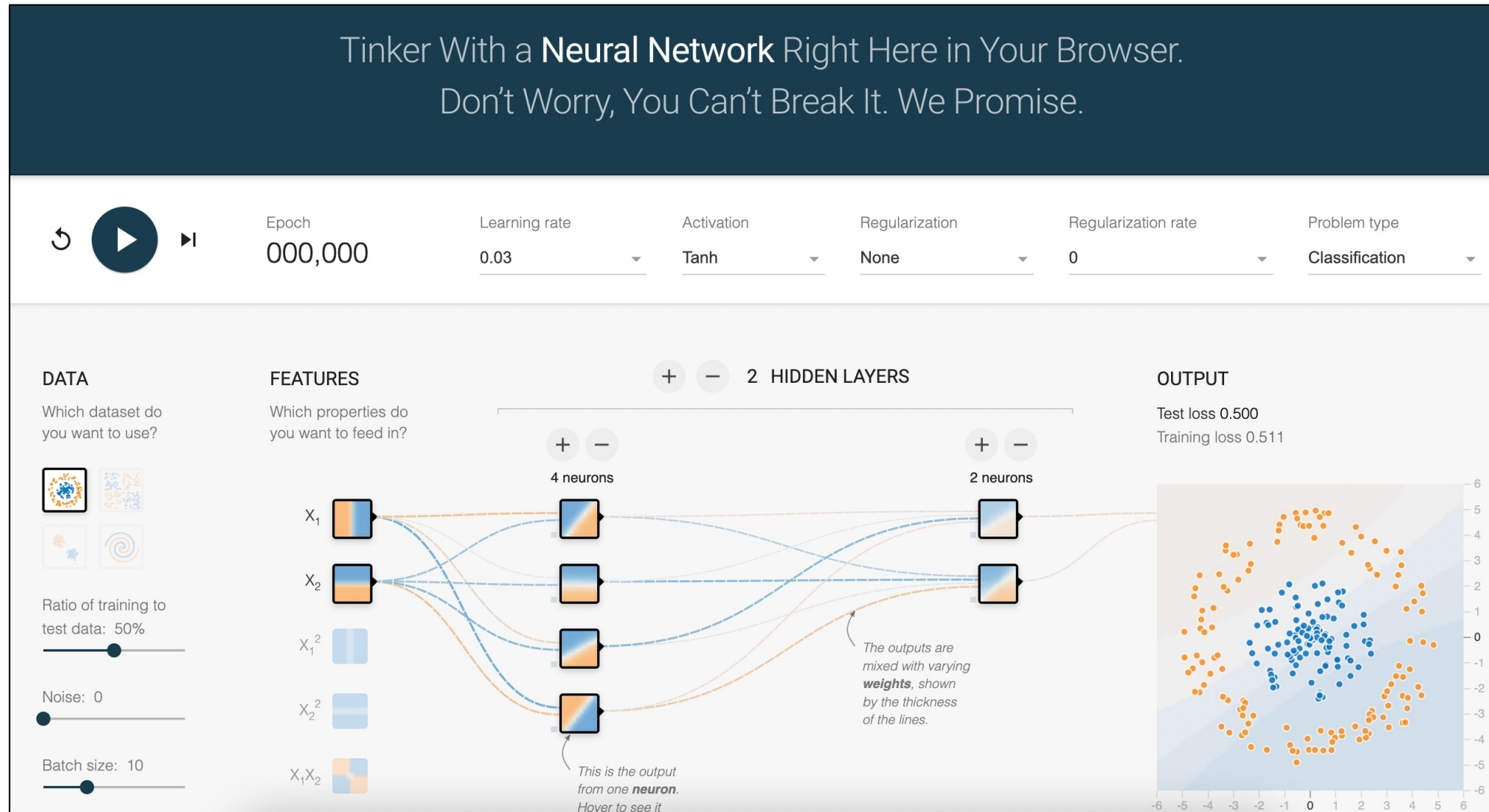


ELU



Goal: computationally-efficient functions with large gradients to support efficient learning

Demo: <https://playground.tensorflow.org/>



Today's Topics

- Objective function: what to learn
- Gradient descent: how to learn
- Training a neural network: optimization
- Gradient descent for activation functions



The End