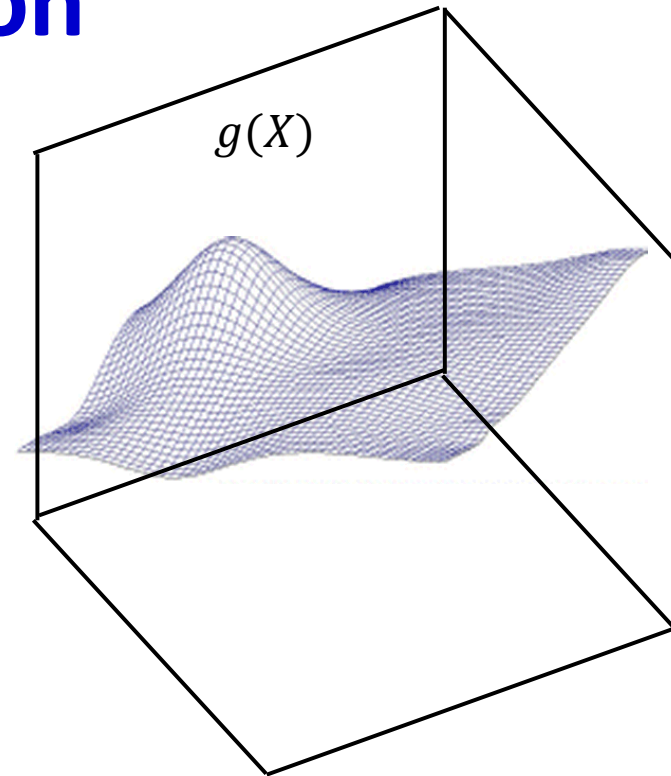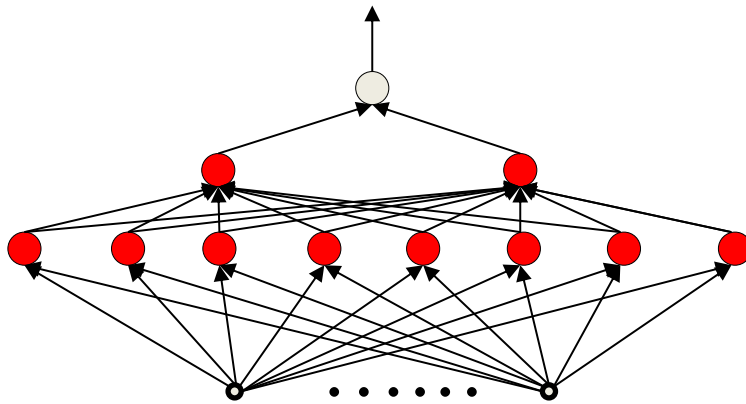# Neural Networks
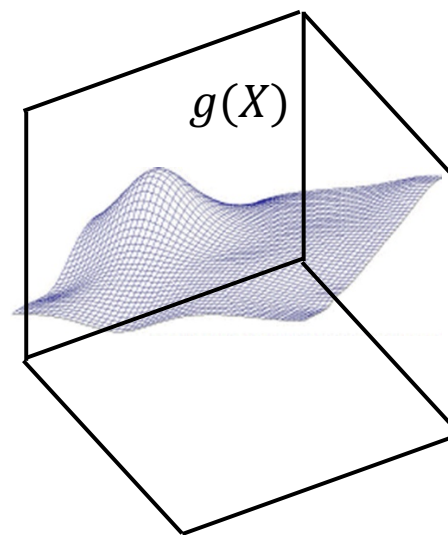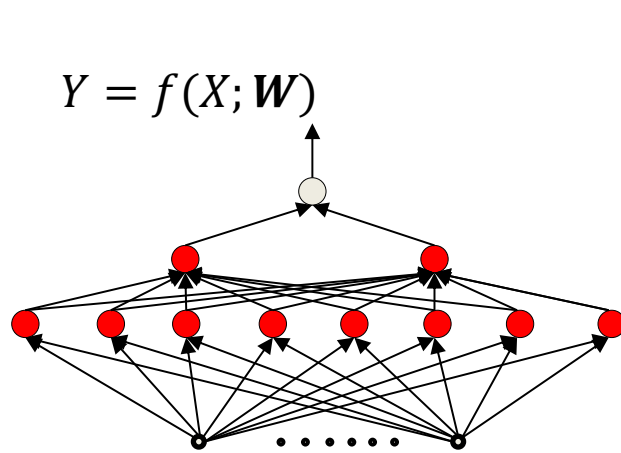# Learning the network: Backprop

11-785, Fall 2019

Lecture 4

# Recap: The MLP *can* represent any function



$g(X)$

- The MLP *can be constructed* to represent anything

- But *how* do we construct it?
  - *I.e.* how do we determine the weights (and biases) of the network to best represent a target function
    - *Assuming that the architecture of the network is given*

# Recap: How to learn the function
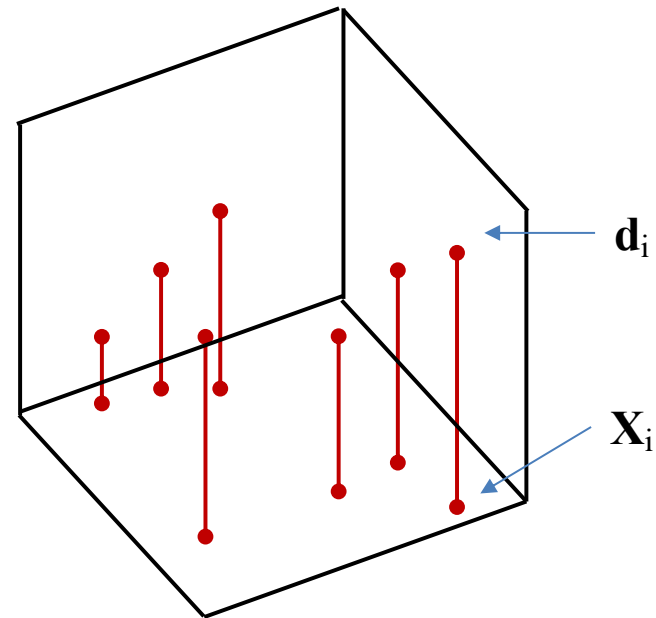
$Y = f(X; \boldsymbol{W})$

$g(X)$

- By minimizing expected error

$$\widehat{\boldsymbol{W}} = \operatorname*{argmin}_{W} \int_{X} div\big(f(X; W), g(X)\big)P(X)dX$$

$$= \operatorname*{argmin}_{W} E\big[div\big(f(X; W), g(X)\big)\big]$$

# Recap: Sampling the function



- $g(X)$ *is unknown, so sample it*

  – Basically, get input-output pairs for a number of samples of input $X_i$

  – Good sampling: the samples of $X$ will be drawn from $P(X)$

- Estimate function from the samples

# The *Empirical* risk



- The *empirical estimate* of the expected error is the *average* error over the samples

$$E\left[div(f(X;W), g(X))\right] \approx \frac{1}{T}\sum_{i=1}^{T} div(f(X_i;W), d_i)$$

- This approximation is an unbiased estimate of the *expected* divergence that we *actually* want to estimate
  - We can *hope* that minimizing the empirical loss will minimize the true loss
  - Caveat: This hope is generally not based on anything but, well, hope..

5

# Empirical Risk Minimization

$$Y = f(X; \boldsymbol{W})$$



- Given a training set of input-output pairs $(\boldsymbol{X}_1, \boldsymbol{d}_1), (\boldsymbol{X}_2, \boldsymbol{d}_2), \ldots, (\boldsymbol{X}_T, \boldsymbol{d}_T)$
  - Error on the i-th instance: $div(f(X_i; W), d_i)$
  - Empirical average error on all training data:

$$Loss(W) = \frac{1}{T} \sum_i div(f(X_i; W), d_i)$$

- Estimate the parameters to minimize the empirical estimate of expected error

$$\widehat{\boldsymbol{W}} = \underset{W}{\operatorname{argmin}} \, Loss(W)$$

  - I.e. minimize the *empirical error* over the drawn samples

6

# Empirical Risk Minimization

$$Y = f(X; \boldsymbol{W})$$

- Given a training set of input-output pairs $(\boldsymbol{X}_1, \boldsymbol{d}_1), (\boldsymbol{X}_2, \boldsymbol{d}_2), \ldots, (\boldsymbol{X}_T, \boldsymbol{d}_T)$
  - Error on the i-th instance: $div(f(X_i; W), d_i)$
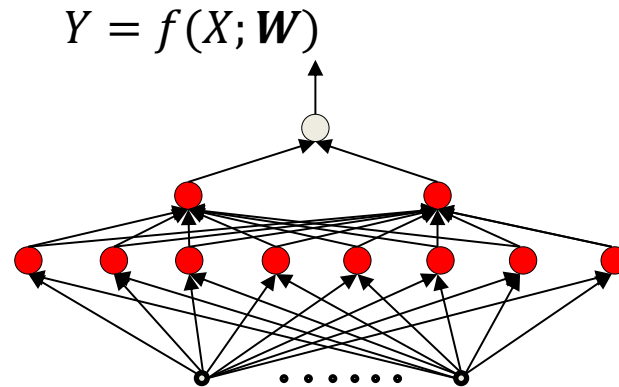  - Empirical average error on all training data:

$$Loss(W) = \frac{1}{T} \sum_i div(f(X_i; W), d_i)$$

- Estimate the parameters to minimize the empirical estimate of expected error

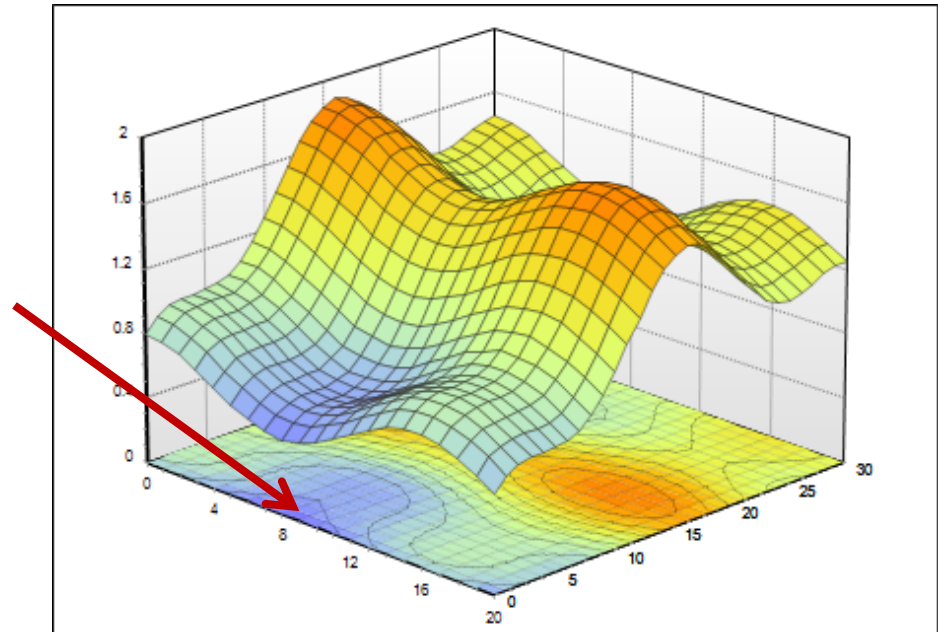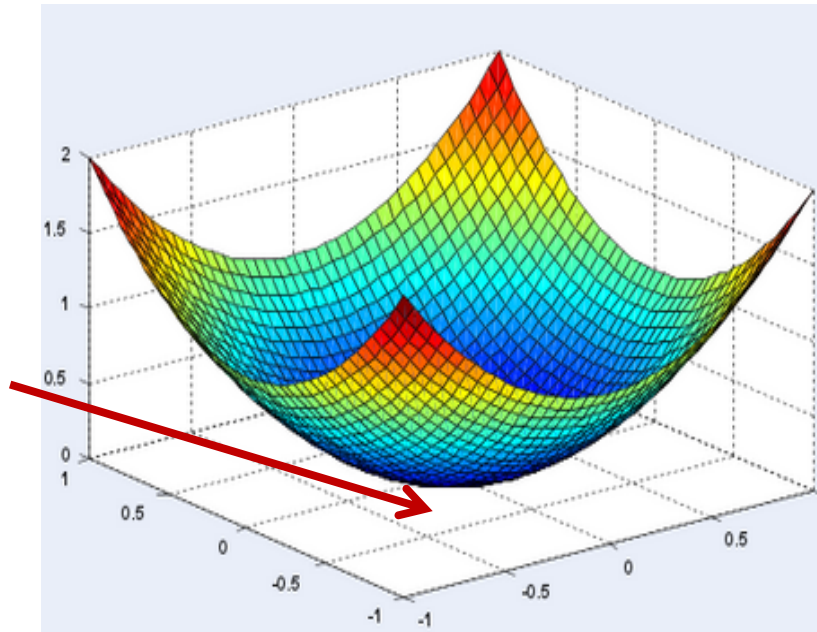$$\widehat{\boldsymbol{W}} = \underset{W}{\operatorname{argmin}} \; Loss(W)$$

  - I.e. minimize the *empirical error* over the drawn samples

7

- **A CRASH COURSE ON FUNCTION OPTIMIZATION**

# Finding the minimum of a scalar function of a multi-variate input



- The optimum point is a turning point – the gradient will be 0

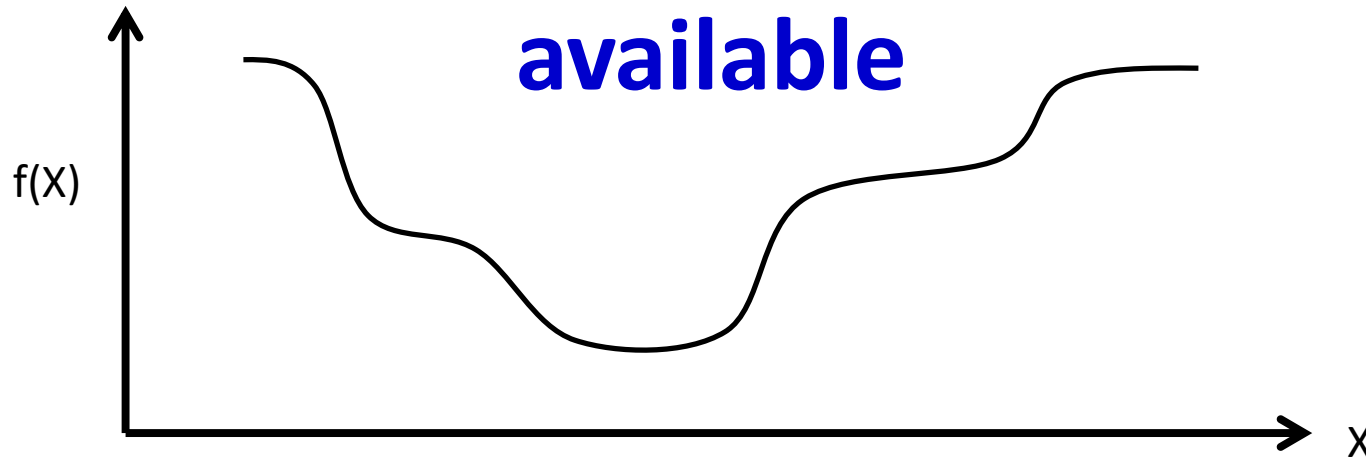# Unconstrained Minimization of function (Multivariate)

1. Solve for the $X$ where the gradient equation equals to zero

$$\nabla_X f(X) = 0$$

2. Compute the Hessian Matrix $\nabla^2 f(X)$ at the candidate solution and verify that

   – Hessian is positive definite (eigenvalues positive) -> to identify local minima

   – Hessian is negative definite (eigenvalues negative) -> to identify local maxima

# Closed Form Solutions are not always available



- Often it is not possible to simply solve $\nabla_X f(X) = 0$
  - The function to minimize/maximize may have an intractable form
- In these situations, iterative solutions are used
  - Begin with a "guess" for the optimal $X$ and refine it iteratively until the correct value is obtained

# Iterative solutions



- Iterative solutions
  - Start from an initial guess $x_0$ for the optimal $x$
  - Update the guess towards a (hopefully) "better" value of $f(x)$
  - Stop when $f(x)$ no longer decreases
- Problems:
  - Which direction to step in
  - How big must the steps be

# The Approach of Gradient Descent



- Iterative solution:  Trivial algorithm
  - Initialize $x^0$
  - While $\left\| \nabla_x f(x^k) \right\| > \varepsilon$ (or while $\left| f(x^{k+1}) - f(x^k) \right| > \varepsilon$)
    - $x^{k+1} = x^k - \eta^k \nabla_x f(x^k)^T$
      - $\eta^k$ is the "step size"

# Overall Gradient Descent Algorithm

- Initialize:
  - $x^0$
  - $k = 0$

- While $\left| f\left(x^{k+1}\right) - f\left(x^k\right) \right| > \varepsilon$
  - $x^{k+1} = x^k - \eta^k \nabla f(x^k)^T$
  - $k = k + 1$

# Convergence of Gradient Descent



- For appropriate step size, for convex (bowl-shaped) functions gradient descent will always find the minimum.

- For non-convex functions it will find a local minimum or an inflection point

- Returning to our problem..

# Problem Statement

- Given a training set of input-output pairs
$$(X_1, d_1), (X_2, d_2), \dots, (X_T, d_T)$$

- Minimize the following function
$$Loss(W) = \frac{1}{T} \sum_i div(f(X_i; W), d_i)$$

  w.r.t $W$

- This is problem of function minimization
  - An instance of optimization

# Problem Setup: Things to define

- Given a training set of input-output pairs
$$(X_1, d_1), (X_2, d_2), \ldots, (X_T, d_T)$$

- What are these input-output pairs?

$$Loss(W) = \frac{1}{T} \sum_i div(f(X_i; W), d_i)$$

# Problem Setup: Things to define

- Given a training set of input-output pairs
  $(X_1, d_1), (X_2, d_2), \ldots, (X_T, d_T)$

- What are these input-output pairs?

$$Loss(W) = \frac{1}{T} \sum_i div(f(X_i; W), d_i)$$

What is f() and what are its parameters W?

# Problem Setup: Things to define

- Given a training set of input-output pairs
$(X_1, d_1), (X_2, d_2), \ldots, (X_T, d_T)$

- What are these input-output pairs?

$$Loss(W) = \frac{1}{T} \sum_i div(f(X_i; W), d_i)$$

What is the divergence div()?

What is f() and what are its parameters W?

# Problem Setup: Things to define

- Given a training set of input-output pairs $(X_1, d_1), (X_2, d_2), \ldots, (X_T, d_T)$

- Minimize the following function

$$Loss(W) = \frac{1}{T} \sum_i div(f(X_i; W), d_i)$$

What is f() and what are its parameters W?

# What is f()? Typical network



- Multi-layer perceptron
- A *directed* network with a set of inputs and outputs
  - No loops
- Generic terminology
  - We will refer to the inputs as the *input units*
    - **No neurons here – the "input units" are just the inputs**
  - We refer to the outputs as the output units
  - Intermediate units are "hidden" units

23

# Typical network
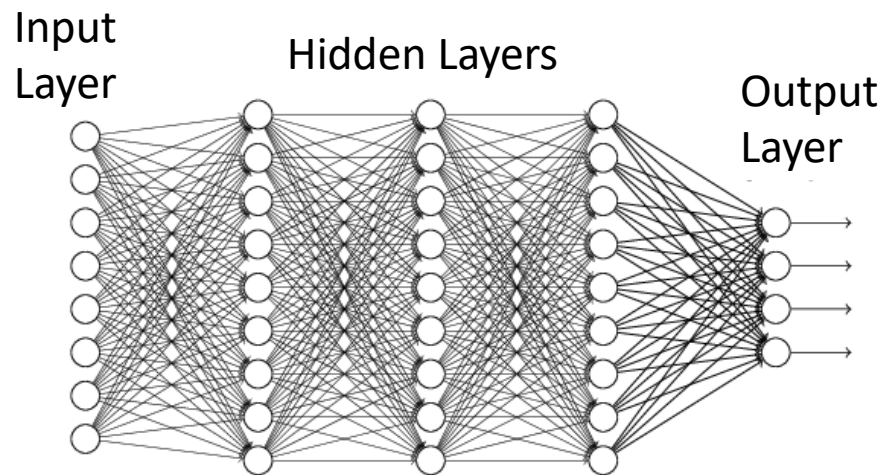


Input Layer

Hidden Layers

Output Layer

- We assume a "layered" network for simplicity
  - We will refer to the inputs as the *input layer*
    - No neurons here – the "layer" simply refers to inputs
  - We refer to the outputs as the output layer
  - Intermediate layers are "hidden" layers

# The individual neurons



- Individual neurons operate on a set of inputs and produce a single output

  - **Standard setup:** A differentiable activation function applied to an affine combination of the input

$$y = f\left(\sum_i w_i x_i + b\right)$$

  - More generally: *any* differentiable function
$$y = f(x_1, x_2, \ldots, x_N; W)$$

25

# The individual neurons



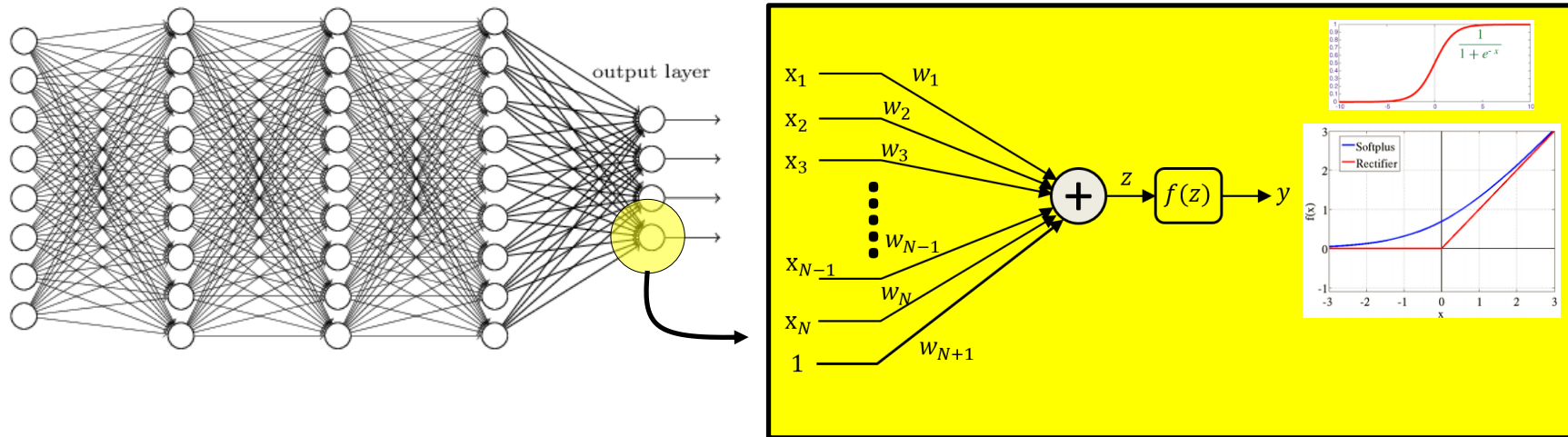- Individual neurons operate on a set of inputs and produce a single output

  - **Standard setup:** A differentiable activation function applied to an affine combination of the input

$$y = f\left(\sum_i w_i x_i + b\right)$$

  - More generally: *any* differentiable function

$$y = f(x_1, x_2, \ldots, x_N; W)$$

We will assume this unless otherwise specified

Parameters are weights $w_i$ and bias $b$

26

# Activations and their derivatives



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

$$f'(z) = f(z)(1 - f(z))$$



$$f(z) = \tanh(z)$$

$$f'(z) = (1 - f^2(z))$$



$$f(z) = \begin{cases} 0, & z < 0 \\ z, & z \geq 0 \end{cases}$$

[*] $\quad f'(z) = \begin{cases} 1, z \geq 0 \\ 0, z < 0 \end{cases}$

$$f(z) = \log(1 + \exp(z))$$

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

- Some popular activation functions and their derivatives

27

# Vector Activations



Input Layer    Hidden Layers    Output Layer

- We can also have neurons that have *multiple coupled* outputs

$$[y_1, y_2, \dots, y_l] = f(x_1, x_2, \dots, x_k; W)$$

  – Function $f()$ operates on set of inputs to produce set of outputs
  – Modifying a single parameter in $W$ will affect *all* outputs

# Vector activation example: Softmax



- Example: Softmax *vector* activation

$$z_i = \sum_j w_{ji} x_j + b_i$$

$$y = \frac{exp(z_i)}{\sum_j exp(z_j)}$$

Parameters are weights $w_{ji}$ and bias $b_i$

# Multiplicative combination: Can be viewed as a case of vector activations



$$z_i = \sum_j w_{ji} x_j + b_i$$

$$y_i = \prod_l (z_l)^{\alpha_{li}}$$

Parameters are weights $w_{ji}$ and bias $b_i$

- A layer of multiplicative combination is a special case of vector activation

# Typical network



Input Layer
Hidden Layers
Output Layer

- In a layered network, each layer of perceptrons can be viewed as a single vector activation

# Notation



- The input layer is the 0$^{\text{th}}$ layer

- We will represent the output of the i-th perceptron of the k$^{\text{th}}$ layer as $y_i^{(k)}$

  - **Input to network:** $y_i^{(0)} = x_i$

  - **Output of network:** $y_i = y_i^{(N)}$

- We will represent the weight of the connection between the i-th unit of the k-1th layer and the jth unit of the k-th layer as $w_{ij}^{(k)}$

  - The bias to the jth unit of the k-th layer is $b_j^{(k)}$

# Problem Setup: Things to define

- Given a training set of input-output pairs
$$(X_1, d_1), (X_2, d_2), \ldots, (X_T, d_T)$$

- What are these input-output pairs?

$$Loss(W) = \frac{1}{T} \sum_i div(f(X_i; W), d_i)$$

# Vector notation



- Given a training set of input-output pairs $(X_1, d_1), (X_2, d_2), \ldots, (X_T, d_T)$
- $X_n = [x_{n1}, x_{n2}, \ldots, x_{nD}]$ is the nth input vector
- $d_n = [d_{n1}, d_{n2}, \ldots, d_{nL}]$ is the nth desired output
- $Y_n = [y_{n1}, y_{n2}, \ldots, y_{nL}]$ is the nth vector of *actual* outputs of the network
- We will sometimes drop the first subscript when referring to a *specific* instance

# Representing the input



Input Layer / Hidden Layers / Output Layer

- Vectors of numbers
  - (or may even be just a scalar, if input layer is of size 1)
  - E.g. vector of pixel values
  - E.g. vector of speech features
  - E.g. real-valued vector representing text
    - We will see how this happens later in the course
  - Other real valued vectors

# Representing the output



Input Layer  Hidden Layers  Output Layer

- If the desired *output* is real-valued, no special tricks are necessary
  - Scalar Output : single output neuron
    - d = scalar (real value)
  - Vector Output : as many output neurons as the dimension of the desired output
    - d = [$d_1$ $d_2$ .. $d_L$] (vector of real values)

# Representing the output



- If the desired output is binary (is this a cat or not), use a simple 1/0 representation of the desired output
  - 1 = Yes it's a cat
  - 0 = No it's not a cat.

# Representing the output



- If the desired output is binary (is this a cat or not), use a simple 1/0 representation of the desired output

- Output activation: Typically a sigmoid

  - Viewed as the *probability* $P(Y = 1|X)$ of class value 1

    - Indicating the fact that for actual data, in general a feature value X may occur for both classes, but with different probabilities

    - Is differentiable

# Multi-class output: One-hot representations

- Consider a network that must distinguish if an input is a cat, a dog, a camel, a hat, or a flower

- We can represent this set as the following vector:

    $$[\text{cat} \; \text{dog} \; \text{camel} \; \text{hat} \; \text{flower}]^T$$

- For inputs of each of the five classes the desired output is:

    cat: $[1\,0\,0\,0\,0]^T$

    dog: $[0\,1\,0\,0\,0]^T$

    camel: $[0\,0\,1\,0\,0]^T$

    hat: $[0\,0\,0\,1\,0]^T$

    flower: $[0\,0\,0\,0\,1]^T$

- For an input of any class, we will have a five-dimensional vector output with four zeros and a single 1 at the position of that class

- This is a *one hot vector*

# Multi-class networks



Input Layer    Hidden Layers    Output Layer

- For a multi-class classifier with N classes, the one-hot representation will have N binary outputs
  - An N-dimensional binary vector
- The neural network's output too must ideally be binary (N-1 zeros and a single 1 in the right place)
- More realistically, it will be a probability vector
  - N probability values that sum to 1.

41

# Multi-class classification: Output



Input Layer   Hidden Layers   Output Layer

- Softmax *vector* activation is often used at the output of multi-class classifier nets

$$z_i = \sum_j w_{ji}^{(n)} y_j^{(n-1)}$$

$$y_i = \frac{exp(z_i)}{\sum_j exp(z_j)}$$

- This can be viewed as the probability $y_i = P(class = i | X)$

# Typical Problem Statement



- We are given a number of "training" data instances
- E.g. images of digits, along with information about which digit the image represents
- Tasks:
  - Binary recognition:  Is this a "2" or not
  - Multi-class recognition:  Which digit is this? Is this a digit in the first place?

# Typical Problem statement: binary classification

Training data



Input: vector of pixel values

Output: sigmoid

- Given, many positive and negative examples (training data),
  - learn all weights such that the network does the desired job

# Typical Problem statement: multiclass classification

Training data

$(5, 5)$  $(2, 2)$

$(2, 2)$  $(4, 4)$

$(0, 0)$  $(2, 2)$

Input Layer

Hidden Layers

Output Layer

softmax

Input: vector of pixel values

Output: Class prob

- Given, many positive and negative examples (training data),
  - learn all weights such that the network does the desired job

# Problem Setup: Things to define

- Given a training set of input-output pairs
  $$(X_1, d_1), (X_2, d_2), \ldots, (X_T, d_T)$$

- Minimize the following function

$$Loss(W) = \frac{1}{T}\sum_i div(f(X_i; W), d_i)$$

What is the divergence div()?

# Examples of divergence functions



- For real-valued output vectors, the (scaled) $L_2$ divergence is popular

$$Div(Y, d) = \frac{1}{2} \|Y - d\|^2 = \frac{1}{2} \sum_i (y_i - d_i)^2$$

  - Squared Euclidean distance between true and desired output
  - Note: this is differentiable

$$\frac{dDiv(Y, d)}{dy_i} = (y_i - d_i)$$

$$\nabla_Y Div(Y, d) = [y_1 - d_1, y_2 - d_2, \dots]$$

# For binary classifier



- For binary classifier with scalar output, $Y \in (0,1)$, $d$ is $0/1$, the cross entropy between the probability distribution $[Y, 1-Y]$ and the ideal output probability $[d, 1-d]$ is popular

$$Div(Y, d) = -d\log Y - (1-d)\log(1-Y)$$

- Minimum when $d = Y$

- Derivative

$$\frac{dDiv(Y,d)}{dY} = \begin{cases} -\dfrac{1}{Y} & if\ d = 1 \\ \dfrac{1}{1-Y} & if\ d = 0 \end{cases}$$

# For binary classifier



- For binary classifier with scalar output, $Y \in (0,1)$, $d$ is $0/1$, the cross entropy between the probability distribution $[Y, 1-Y]$ and the ideal output probability $[d, 1-d]$ is popular

$$Div(Y, d) = -d\log Y - (1-d)\log(1-Y)$$

  - Minimum when $d = Y$

- Derivative

$$\frac{dDiv(Y,d)}{dY} = \begin{cases} -\dfrac{1}{Y} & if \ d = 1 \\ \dfrac{1}{1-Y} & if \ d = 0 \end{cases}$$

Note: when $y = d$ the derivative is *not* 0

*Even though $div() = 0$* (minimum) *when y = d*

49

# For multi-class classification



- Desired output $d$ is a one hot vector $[0\ 0\ \dots 1\ \dots 0\ 0\ 0\ ]$ with the 1 in the $c$-th position (for class $c$)
- Actual output will be probability distribution $[y_1, y_2, \dots]$
- The cross-entropy between the desired one-hot output and actual output:

$$Div(Y, d) = -\sum_i d_i \log y_i = -\log y_c$$

- Derivative

$$\frac{dDiv(Y,d)}{dY_i} = \begin{cases} -\dfrac{1}{y_c} & for\ the\ c-th\ component \\ 0 & for\ remaining\ component \end{cases}$$

$$\nabla_Y Div(Y,d) = \left[0\ 0\ \dots \dfrac{-1}{y_c}\ \dots 0\ 0\right]$$

If $y_c < 1$, the slope is negative w.r.t. $y_c$

Indicates *increasing* $y_c$ will *reduce* divergence

50

# For multi-class classification



$d_1\, d_2\, d_3\, d_4$

KL Div() $\longrightarrow$ Div

- Desired output $d$ is a one hot vector $[0\ 0\ \dots 1\ \dots 0\ 0\ 0\,]$ with the 1 in the $c$-th position (for class $c$)
- Actual output will be probability distribution $[y_1, y_2, \dots]$
- The cross-entropy between the desired one-hot output and actual output:

$$Div(Y, d) = -\sum_i d_i \log y_i = -\log y_c$$

- Derivative

$$\frac{dDiv(Y,d)}{dY_i} = \begin{cases} -\dfrac{1}{y_c} & for\ the\ c-th\ component \\ 0 & for\ remaining\ component \end{cases}$$
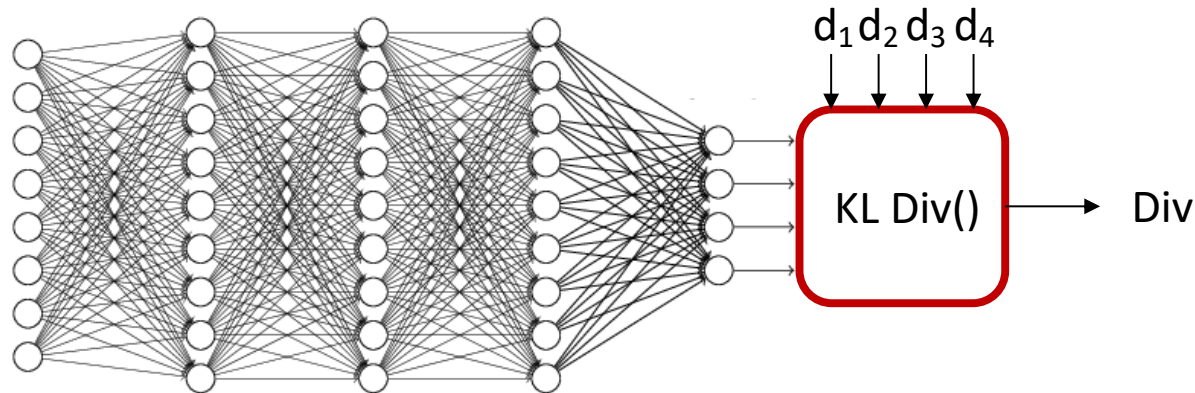
$$\nabla_Y Div(Y,d) = \left[0\ 0\ \dots \frac{-1}{y_c} \dots 0\ 0\right]$$

If $y_c < 1$, the slope is negative w.r.t. $y_c$

Indicates *increasing* $y_c$ will *reduce* divergence

Note: when $y = d$ the derivative is *not* 0

*Even though $div() = 0$ (minimum) when y = d*

# For multi-class classification



$d_1\, d_2\, d_3\, d_4$

KL Div() → Div

- It is sometimes useful to set the target output to $[\epsilon \quad \epsilon \,\ldots\, (1 - (K-1)\epsilon) \,\ldots\, \epsilon \quad \epsilon \quad \epsilon]$ with the value $1 - (\mathrm{K}-1)\epsilon$ in the $c$-th position (for class $c$) and $\epsilon$ elsewhere for some small $\epsilon$
  - "Label smoothing" -- aids gradient descent
- The cross-entropy remains:

$$Div(Y, d) = -\sum_i d_i \log y_i$$

- Derivative

$$\frac{dDiv(Y, d)}{dY_i} = \begin{cases} -\dfrac{1 - (K-1)\epsilon}{y_c} & for\ the\ c-th\ component \\[2mm] -\dfrac{\epsilon}{y_i} & for\ remaining\ components \end{cases}$$

52

# Problem Setup

- Given a training set of input-output pairs $(X_1, d_1), (X_2, d_2), \ldots, (X_T, d_T)$

- The error on the i$^{\text{th}}$ instance is $div(Y_i, d_i)$

- The loss

$$Loss = \frac{1}{T} \sum_i div(Y_i, d_i)$$

- Minimize $Loss$ w.r.t $\left\{ w_{ij}^{(k)}, b_j^{(k)} \right\}$

# Recap: Gradient Descent Algorithm

- In order to minimize any function $f(x)$ w.r.t. $x$
- Initialize:
  - $x^0$
  - $k = 0$

- While $\left| f\left(x^{k+1}\right) - f\left(x^k\right) \right| > \varepsilon$
  - $x^{k+1} = x^k - \eta^k \nabla f(x^k)^T$
  - $k = k + 1$

# Recap: Gradient Descent Algorithm

- In order to minimize any function $f(x)$ w.r.t. $x$
- Initialize:
  - $x^0$
  - $k = 0$

- While $\left| f\left(x^{k+1}\right) - f\left(x^k\right) \right| > \varepsilon$
  - For every component $i$
    - $x_i^{k+1} = x_i^k - \eta^k \dfrac{df}{dx_i}$  | Explicitly stating it by component
  - $k = k + 1$

# Training Neural Nets through Gradient Descent

**Total training Loss:**

$$Loss = \frac{1}{T} \sum_t Div(\boldsymbol{Y_t}, \boldsymbol{d_t})$$

- Gradient descent algorithm:

- Initialize all weights and biases $\left\{ w_{ij}^{(k)} \right\}$

  *Assuming the bias is also represented as a weight*

  – Using the extended notation: the bias is also a weight

- Do:

  – For every layer $k$ for all $i, j$, update:

    - $w_{i,j}^{(k)} = w_{i,j}^{(k)} - \eta \frac{dLoss}{dw_{i,j}^{(k)}}$

- Until *Loss* has converged

# Training Neural Nets through Gradient Descent

**Total training Loss:**

$$Loss = \frac{1}{T}\sum_t Div(\mathbf{Y_t}, \mathbf{d_t})$$

- Gradient descent algorithm:

- Initialize all weights $\left\{w_{ij}^{(k)}\right\}$

- Do:

  - For every layer $k$ for all $i, j$, update:

    - $w_{i,j}^{(k)} = w_{i,j}^{(k)} - \eta \dfrac{dLoss}{dw_{i,j}^{(k)}}$

- Until $Err$ has converged

# The derivative

Total training Loss:

$$Loss = \frac{1}{T} \sum_t Div(\mathbf{Y_t}, \mathbf{d_t})$$

- Computing the derivative

Total derivative:

$$\frac{dLoss}{dw_{i,j}^{(k)}} = \frac{1}{T} \sum_t \frac{dDiv(\mathbf{Y_t}, \mathbf{d_t})}{dw_{i,j}^{(k)}}$$

# The derivative

**Total training Loss:**

$$Loss = \frac{1}{T} \sum_t Div(\boldsymbol{Y_t}, \boldsymbol{d_t})$$

**Total derivative:**

$$\frac{d\boldsymbol{Loss}}{dw_{i,j}^{(k)}} = \frac{1}{T} \sum_t \frac{dDiv(\boldsymbol{Y_t}, \boldsymbol{d_t})}{dw_{i,j}^{(k)}}$$

- So we must first figure out how to compute the derivative of divergences of individual training inputs

# Calculus Refresher: Basic rules of calculus

For any differentiable function
$$y = f(x)$$
with derivative
$$\frac{dy}{dx}$$
the following must hold for sufficiently small $\Delta x$ $\Longrightarrow$ $\Delta y \approx \dfrac{dy}{dx}\Delta x$

For any differentiable function
$$y = f(x_1, x_2, \ldots, x_M)$$
with partial derivatives
$$\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \ldots, \frac{\partial y}{\partial x_M}$$
the following must hold for sufficiently small $\Delta x_1, \Delta x_2, \ldots, \Delta x_M$

$$\Delta y \approx \frac{\partial y}{\partial x_1}\Delta x_1 + \frac{\partial y}{\partial x_2}\Delta x_2 + \cdots + \frac{\partial y}{\partial x_M}\Delta x_M$$

61

# Calculus Refresher: Chain rule

For any nested function $y = f(g(x))$

$$\frac{dy}{dx} = \frac{\partial f}{\partial g(x)}\frac{dg(x)}{dx}$$

Check – we can confirm that :    $\Delta y = \dfrac{dy}{dx}\Delta x$

$z = g(x) \Longrightarrow \Delta z = \dfrac{dg(x)}{dx}\Delta x$

$y = f(z) \Longrightarrow \quad \Delta y = \dfrac{df}{dz}\Delta z = \dfrac{df}{dz}\dfrac{dg(x)}{dx}\Delta x$

# Calculus Refresher: Distributed Chain rule

$$y = f\big(g_1(x), g_1(x), \ldots, g_M(x)\big)$$

$$\frac{dy}{dx} = \frac{\partial f}{\partial g_1(x)}\frac{dg_1(x)}{dx} + \frac{\partial f}{\partial g_2(x)}\frac{dg_2(x)}{dx} + \cdots + \frac{\partial f}{\partial g_M(x)}\frac{dg_M(x)}{dx}$$
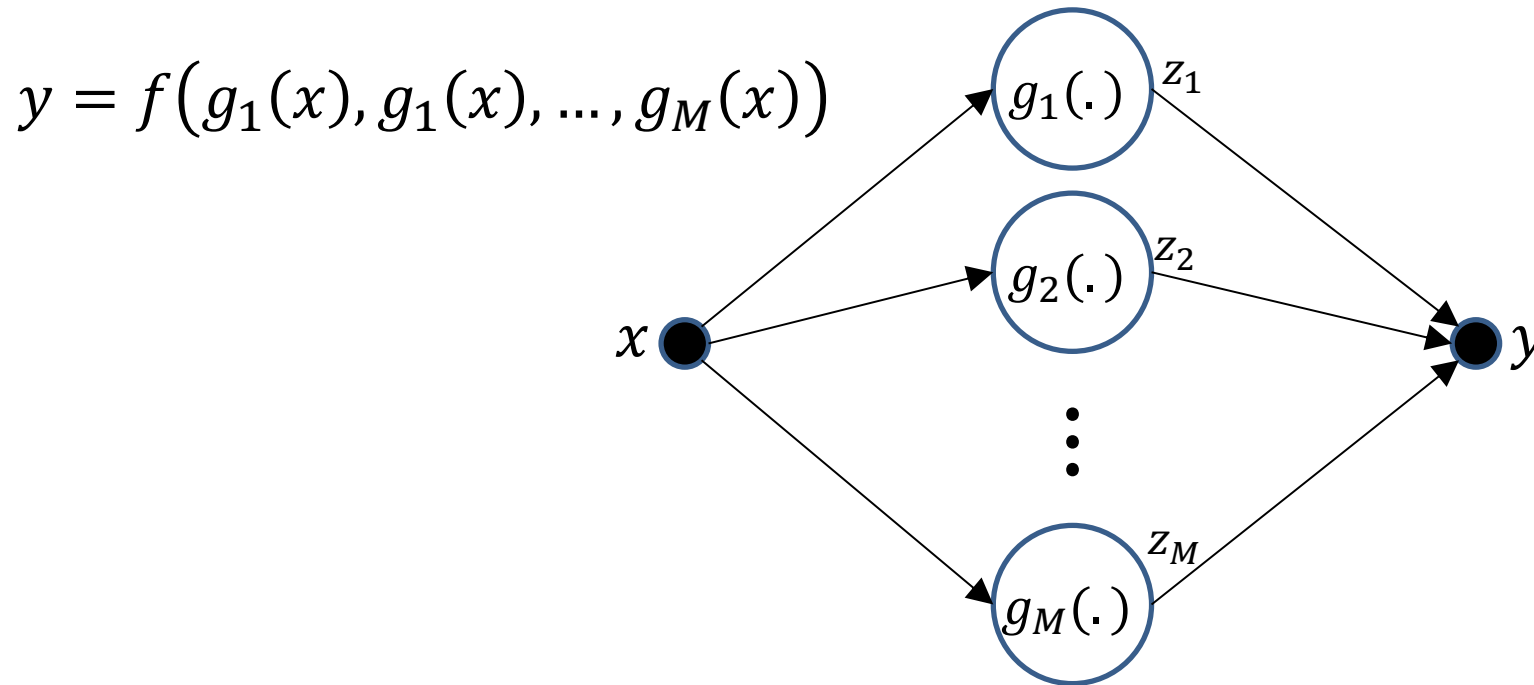
Check:  $\Delta y = \dfrac{dy}{dx}\Delta x$

$$\Delta y = \frac{\partial f}{\partial g_1(x)}\Delta g_1(x) + \frac{\partial f}{\partial g_2(x)}\Delta g_2(x) + \cdots + \frac{\partial f}{\partial g_M(x)}\Delta g_M(x)$$

$$\Delta y = \frac{\partial f}{\partial g_1(x)}\frac{dg_1(x)}{dx}\Delta x + \frac{\partial f}{\partial g_2(x)}\frac{dg_2(x)}{dx}\Delta x + \cdots + \frac{\partial f}{\partial g_M(x)}\frac{dg_M(x)}{dx}\Delta x$$

$$\Delta y = \left(\frac{\partial f}{\partial g_1(x)}\frac{dg_1(x)}{dx} + \frac{\partial f}{\partial g_2(x)}\frac{dg_2(x)}{dx} + \cdots + \frac{\partial f}{\partial g_M(x)}\frac{dg_M(x)}{dx}\right)\Delta x$$ ✓
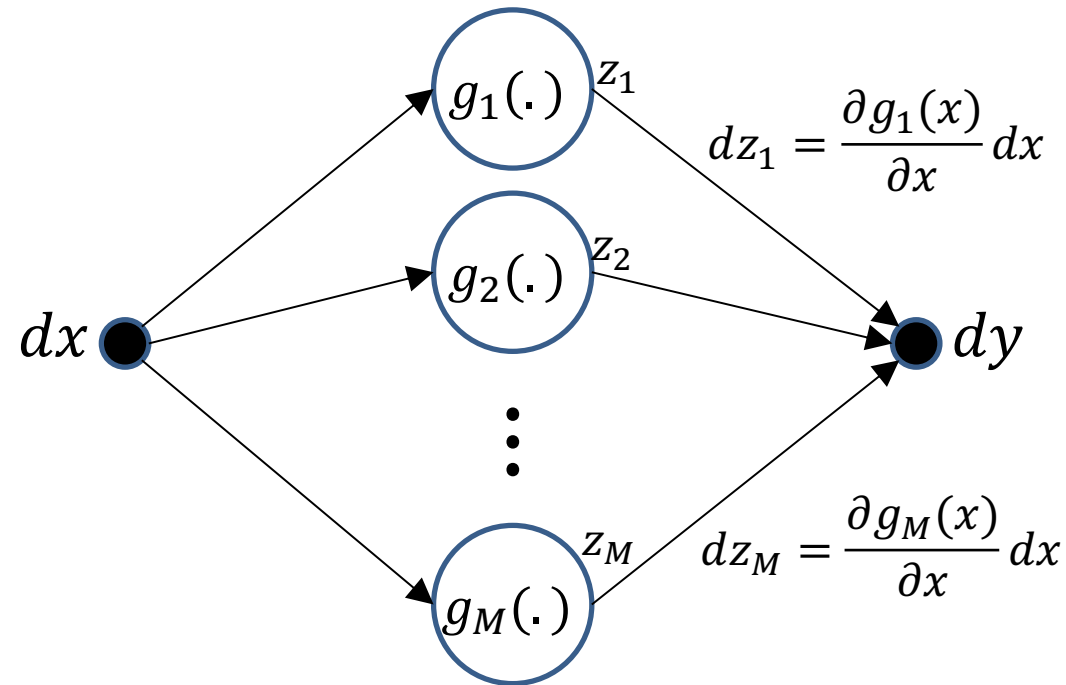
# Distributed Chain Rule: Influence Diagram

$$y = f(g_1(x), g_1(x), \ldots, g_M(x))$$



- $x$ affects $y$ through each of $g_1 \ldots g_M$

# Distributed Chain Rule: Influence Diagram



$$dz_1 = \frac{\partial g_1(x)}{\partial x} dx$$
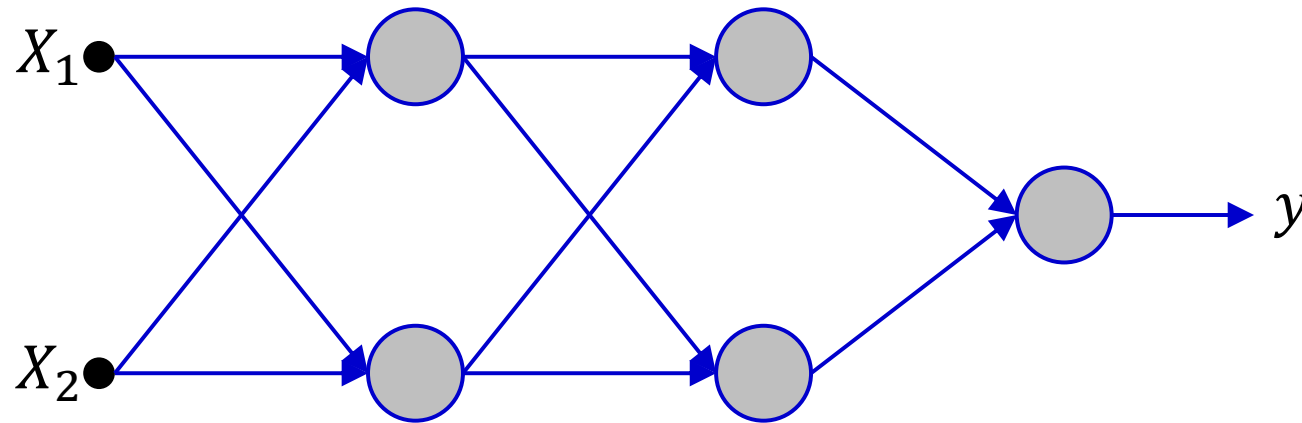
$$dz_M = \frac{\partial g_M(x)}{\partial x} dx$$

- Small perturbations in $x$ cause small perturbations in each of $g_1 \dots g_M$, each of which individually additively perturbs $y$

# Returning to our problem

- How to compute $\dfrac{d\boldsymbol{Div}(\boldsymbol{Y},\boldsymbol{d})}{dw_{i,j}^{(k)}}$

# A first closer look at the network



- Showing a tiny 2-input network for illustration
  - Actual network would have many more neurons and inputs

# A first closer look at the network



- Showing a tiny 2-input network for illustration
  - Actual network would have many more neurons and inputs
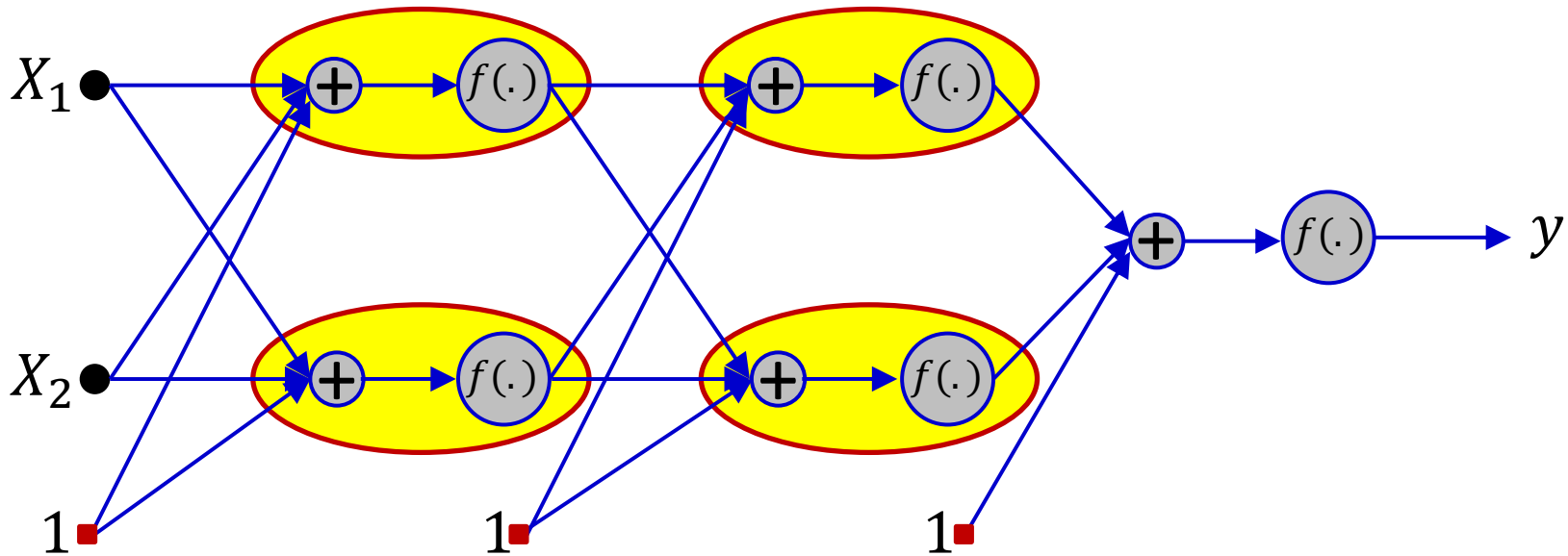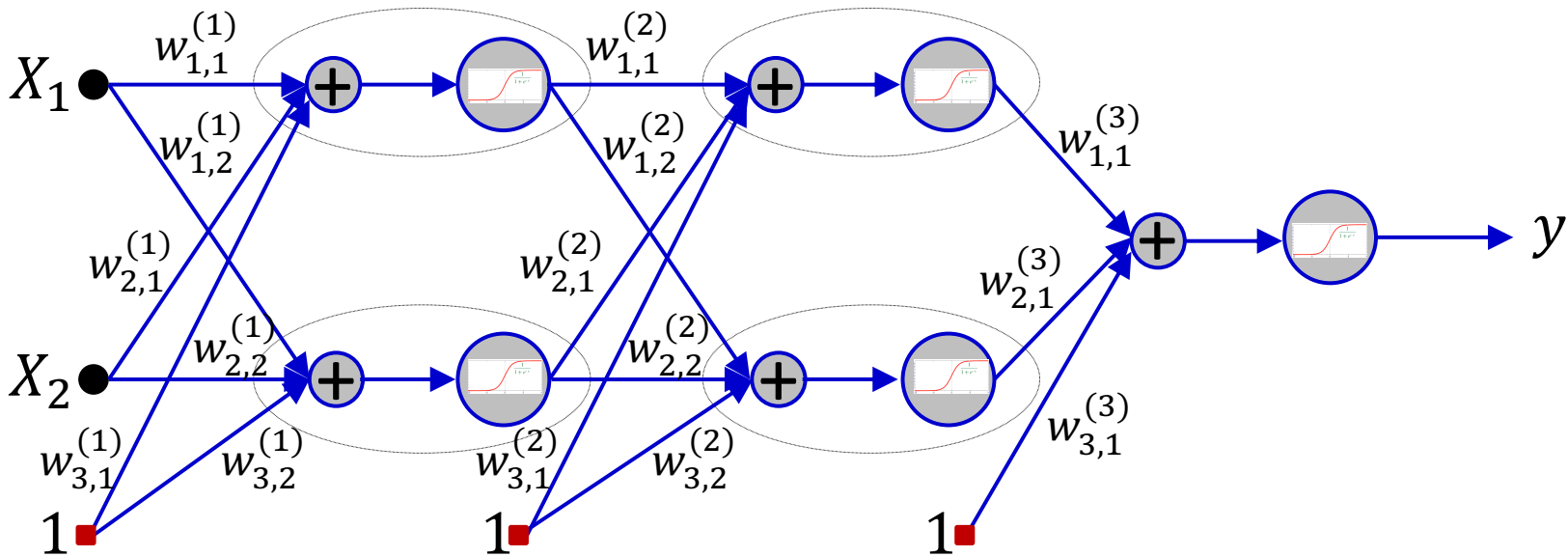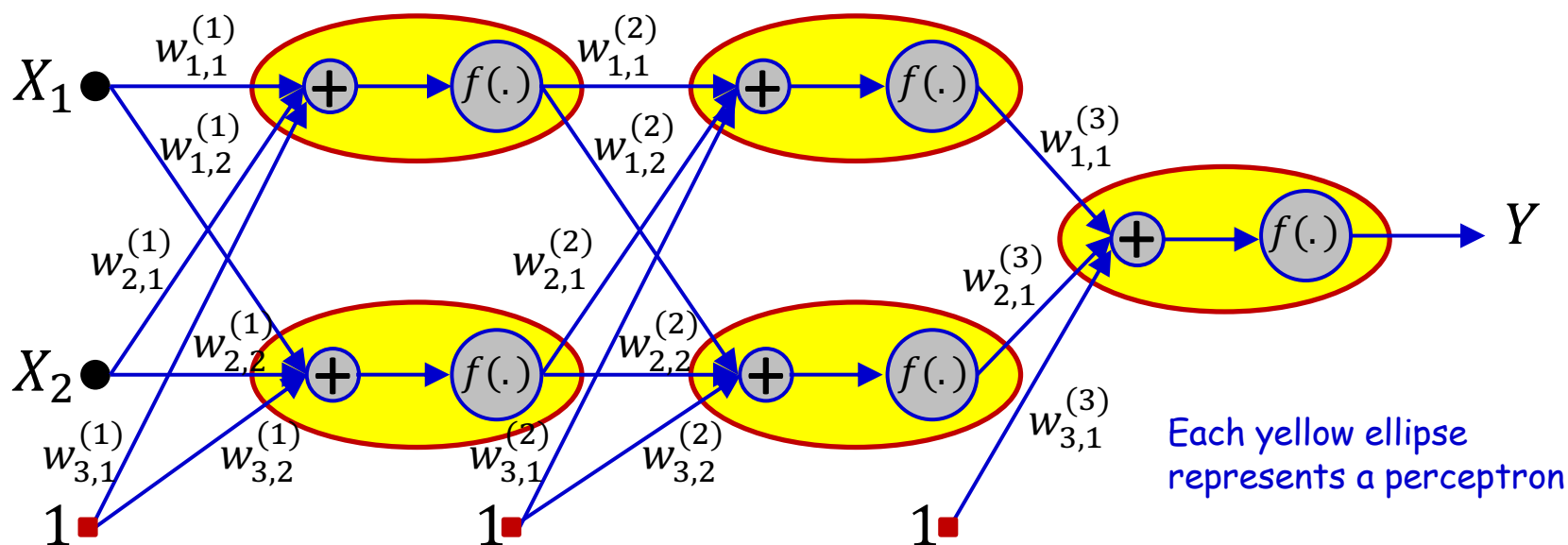- Explicitly separating the weighted sum of inputs from the activation

68

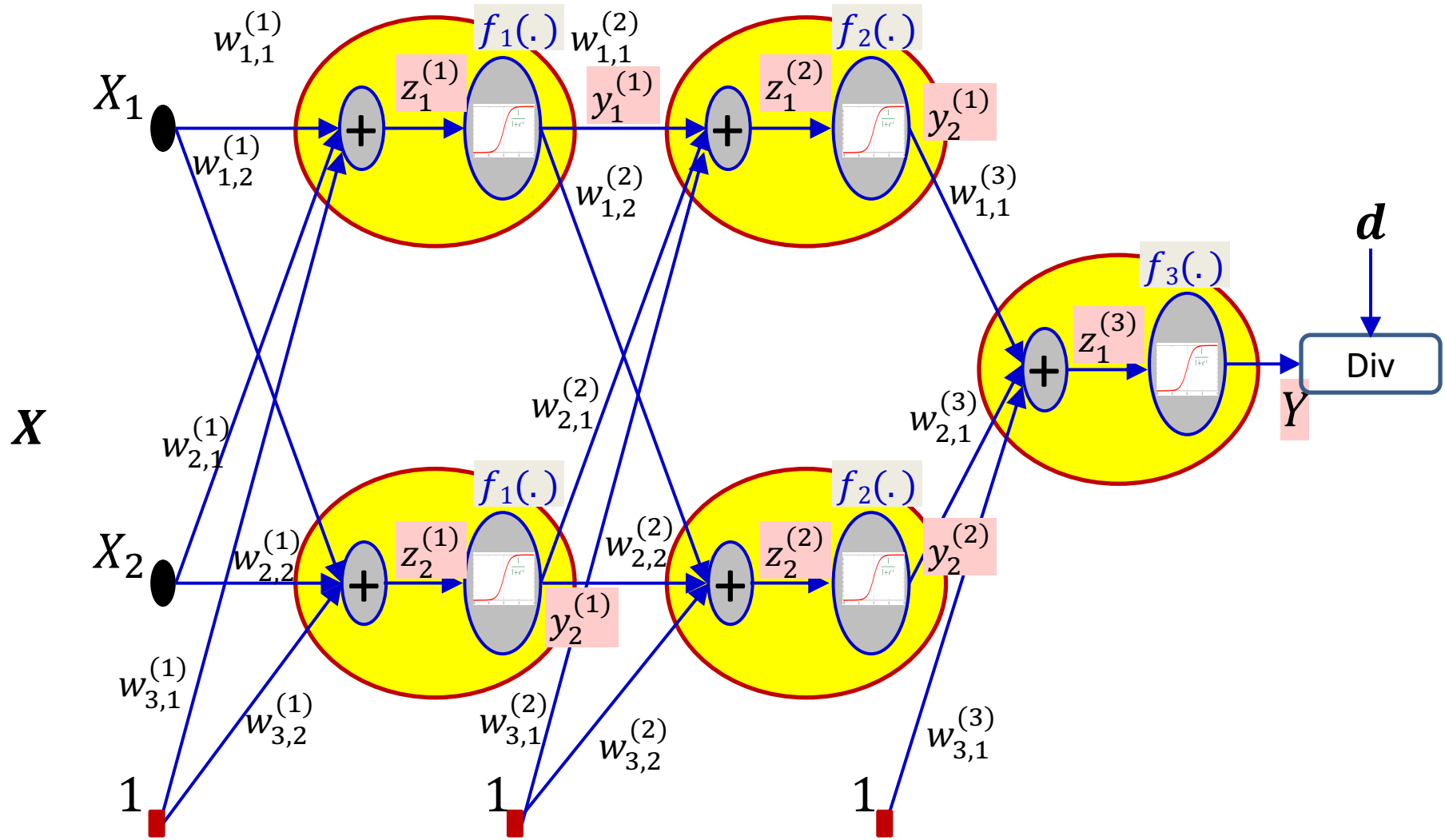# A first closer look at the network



- Showing a tiny 2-input network for illustration
  - Actual network would have many more neurons and inputs
- Expanded **with all weights and activations shown**
- The overall function is differentiable w.r.t every weight, bias and input

# Computing the derivative for a *single* input



Each yellow ellipse represents a perceptron

- Aim: compute derivative of $Div(Y, d)$ w.r.t. each of the weights

- But first, lets label *all* our variables and activation functions

70

# Computing the derivative for a *single* input

# Computing the gradient

- What is: $\dfrac{dDiv(\boldsymbol{Y},\boldsymbol{d})}{dw_{i,j}^{(k)}}$

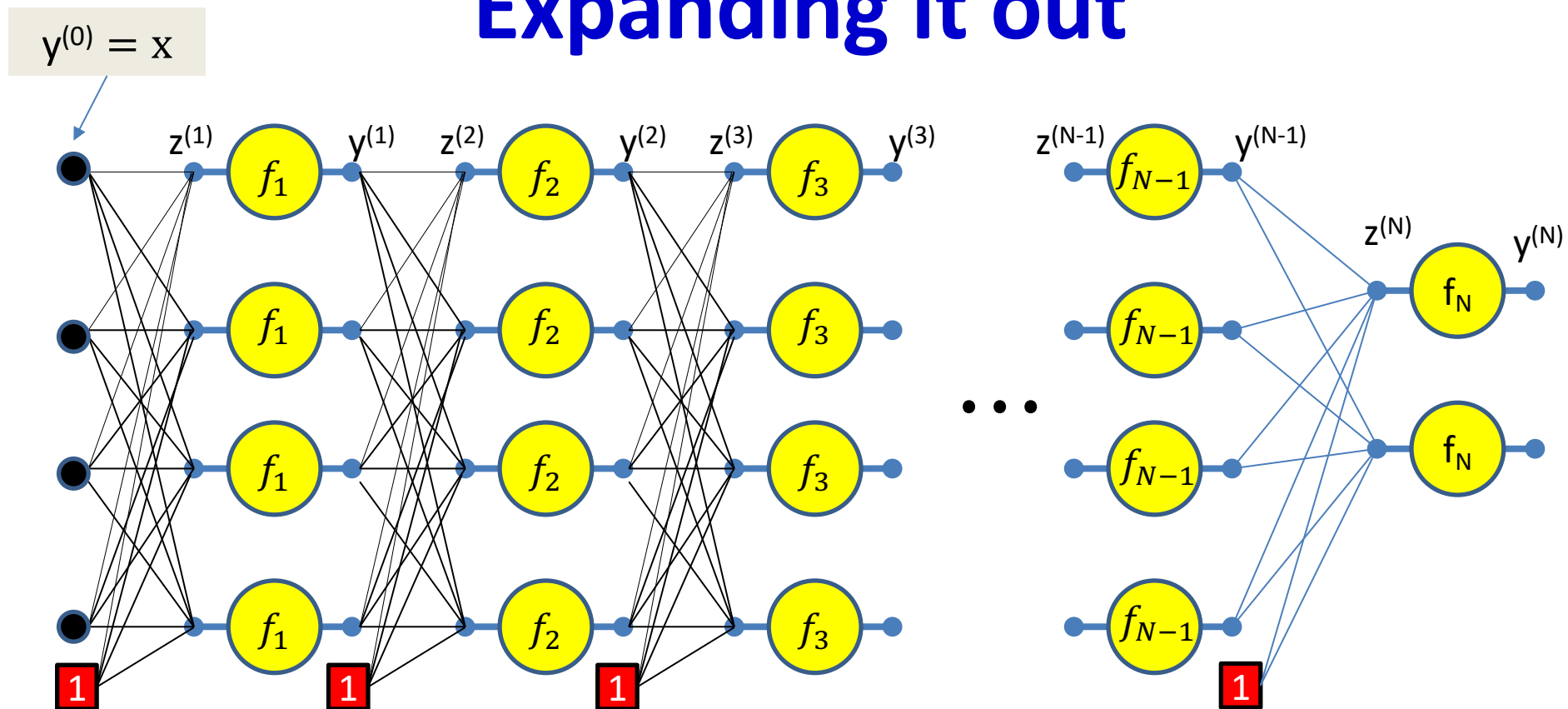  - Derive on board?

# Computing the gradient

- What is: $\dfrac{dDiv(\boldsymbol{Y},\boldsymbol{d})}{dw_{i,j}^{(k)}}$

- Derive on board?

- Note: computation of the derivative requires intermediate and final output values of the network in response to the input

# BP: Scalar Formulation



- The network again

# Expanding it out



Setting $y_i^{(0)} = x_i$ for notational convenience

Assuming $w_{0j}^{(k)} = b_j^{(k)}$ and $y_0^{(k)} = 1$ -- assuming the bias is a weight and extending the output of every layer by a constant 1, to account for the biases

# Expanding it out



$$z_1^{(1)} = \sum_i w_{i1}^{(1)} y_i^{(0)}$$

# Expanding it out



$y^{(0)} = x$

$z^{(1)}$

$y^{(1)}$ $z^{(2)}$ $y^{(2)}$ $z^{(3)}$ $y^{(3)}$

$z^{(N-1)}$ $y^{(N-1)}$ $z^{(N)}$ $y^{(N)}$

$f_1$ $f_2$ $f_3$ $f_{N-1}$ $f_N$

$$z_j^{(1)} = \sum_i w_{ij}^{(1)} y_i^{(0)}$$

$$y^{(0)} = x$$

$$z^{(1)}$$

$$y^{(1)}$$

$$z^{(2)} \quad y^{(2)} \quad z^{(3)} \qquad y^{(3)}$$

$$z^{(N-1)} \quad y^{(N-1)}$$

$$z^{(N)} \qquad y^{(N)}$$

$$z_j^{(1)} = \sum_i w_{ij}^{(1)} y_i^{(0)} \qquad y_j^{(1)} = f_1\left(z_j^{(1)}\right)$$

$$y^{(0)} = x$$

$$z^{(1)} \quad f_1 \quad y^{(1)} \quad z^{(2)}$$

$$z_j^{(1)} = \sum_i w_{ij}^{(1)} y_i^{(0)} \qquad y_j^{(1)} = f_1\left(z_j^{(1)}\right) \qquad z_j^{(2)} = \sum_i w_{ij}^{(2)} y_i^{(1)}$$

$$y^{(0)} = x$$

$$\mathbf{z^{(1)}} \quad \mathbf{y^{(1)}} \quad \mathbf{z^{(2)}} \quad \mathbf{y^{(2)}} \quad z^{(3)} \quad y^{(3)} \quad z^{(N-1)} \quad y^{(N-1)} \quad z^{(N)} \quad y^{(N)}$$

$$z_j^{(1)} = \sum_i w_{ij}^{(1)} y_i^{(0)} \qquad y_j^{(1)} = f_1\left(z_j^{(1)}\right) \qquad z_j^{(2)} = \sum_i w_{ij}^{(2)} y_i^{(1)} \qquad y_j^{(2)} = f_2\left(z_j^{(2)}\right)$$
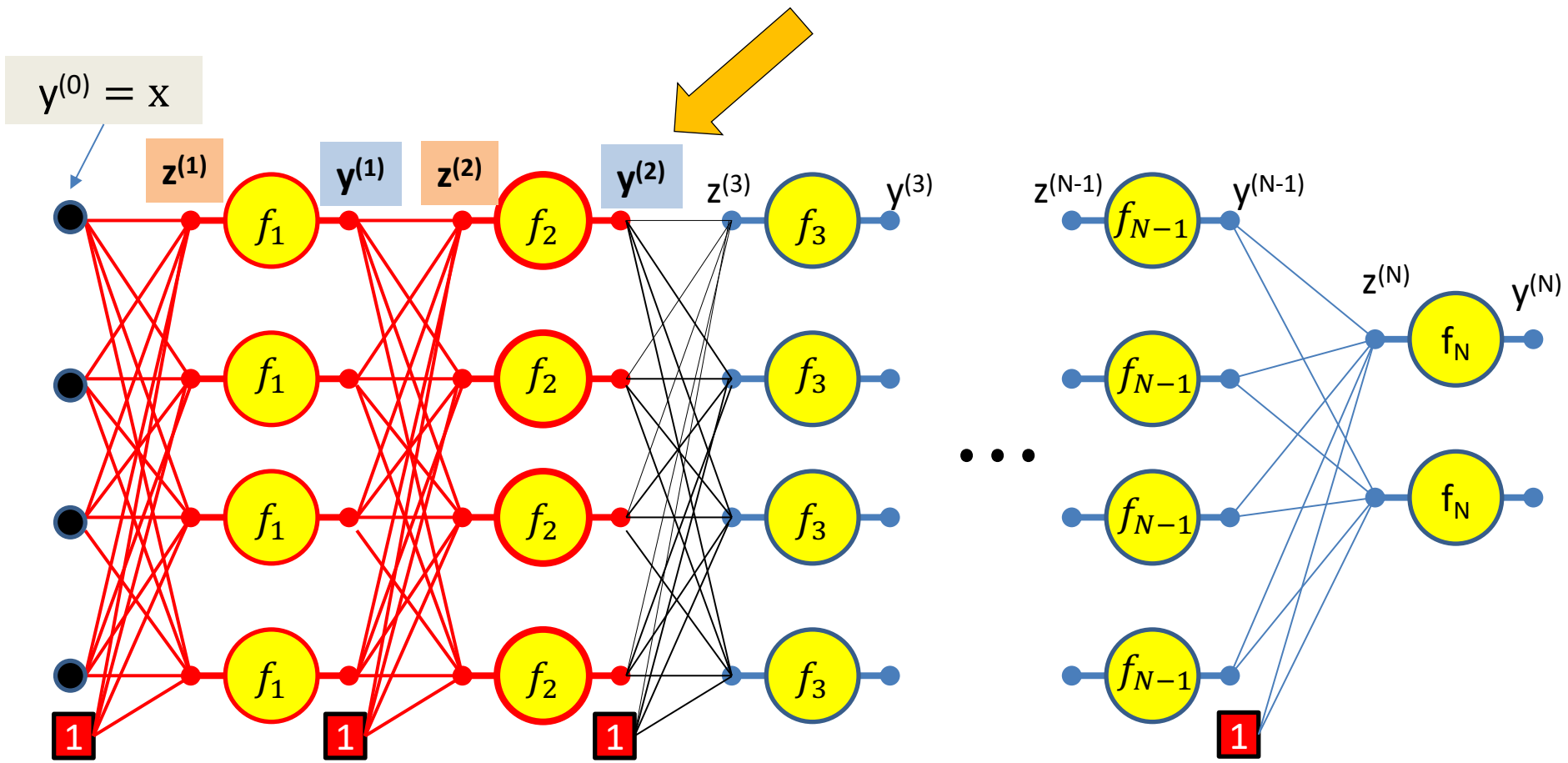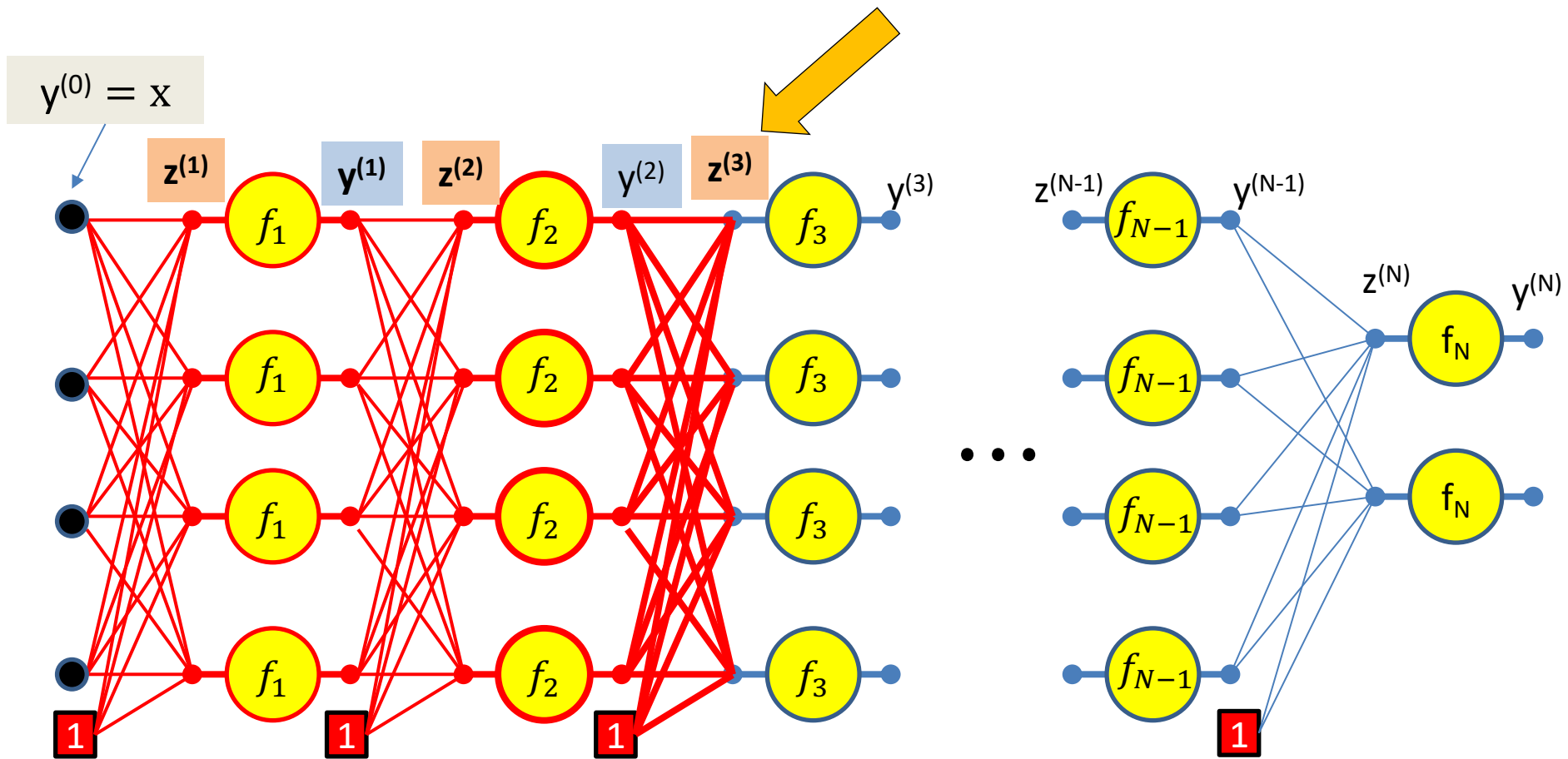
$$z_j^{(1)} = \sum_i w_{ij}^{(1)} y_i^{(0)}$$

$$y_j^{(1)} = f_1\left(z_j^{(1)}\right)$$

$$z_j^{(2)} = \sum_i w_{ij}^{(2)} y_i^{(1)}$$

$$y_j^{(2)} = f_2\left(z_j^{(2)}\right)$$

$$z_j^{(3)} = \sum_i w_{ij}^{(3)} y_i^{(2)}$$

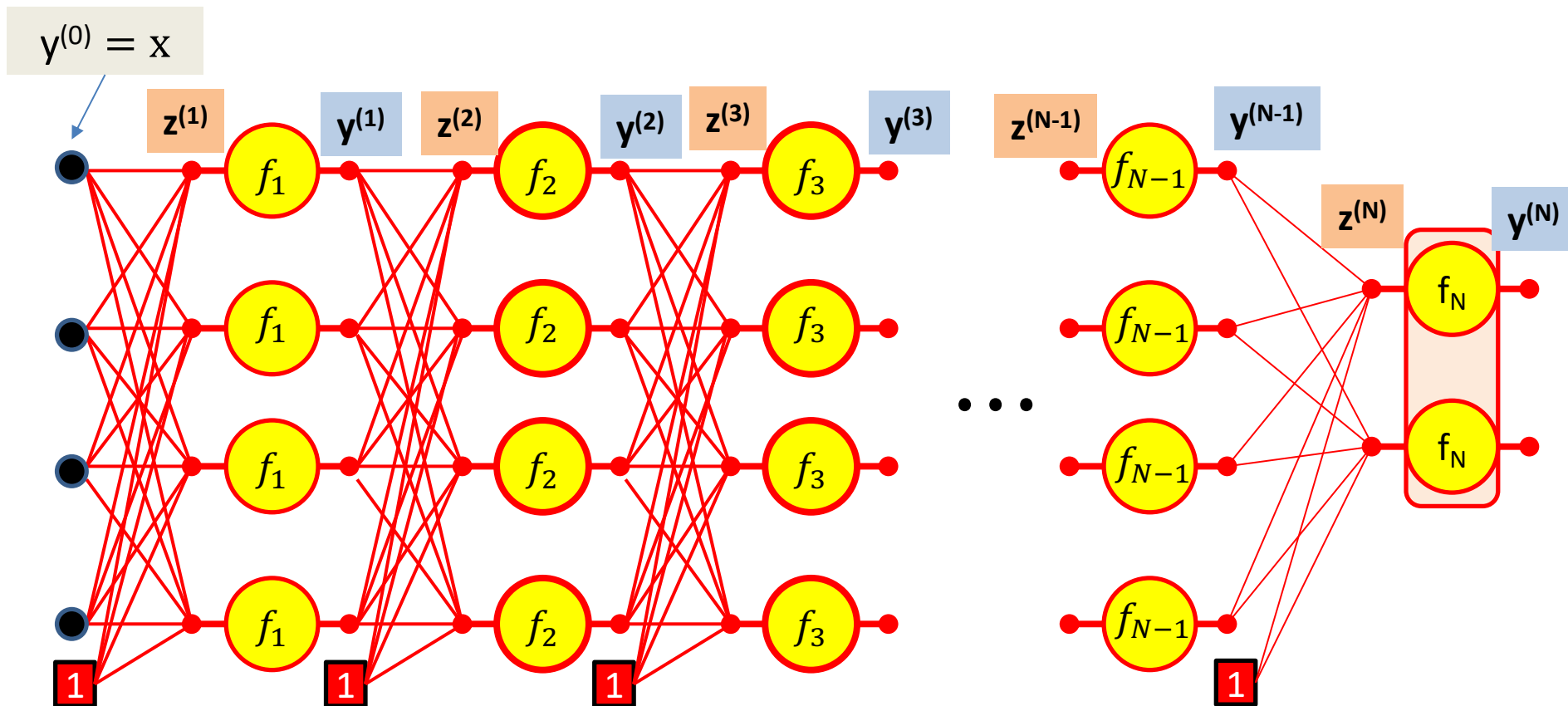$$z_j^{(1)} = \sum_i w_{ij}^{(1)} y_i^{(0)} \qquad y_j^{(1)} = f_1\left(z_j^{(1)}\right) \qquad z_j^{(2)} = \sum_i w_{ij}^{(2)} y_i^{(1)} \qquad y_j^{(2)} = f_2\left(z_j^{(2)}\right)$$

$$z_j^{(3)} = \sum_i w_{ij}^{(3)} y_i^{(2)} \qquad y_j^{(3)} = f_3\left(z_j^{(3)}\right) \qquad \cdots$$

$$y^{(0)} = x$$

$$z^{(1)} \quad y^{(1)} \quad z^{(2)} \quad y^{(2)} \quad z^{(3)} \quad y^{(3)} \quad z^{(N-1)} \quad y^{(N-1)}$$

$$z^{(N)} \quad y^{(N)}$$

$$y_j^{(N-1)} = f_{N-1}\left(z_j^{(N-1)}\right) \qquad z_j^{(N)} = \sum_i w_{ij}^{(N)} y_i^{(N-1)} \qquad \boldsymbol{y}^{(N)} = f_N\left(\boldsymbol{z}^{(N)}\right)$$

# Forward Computation



$y^{(0)} = x$

$$y_i^{(0)} = x_i$$

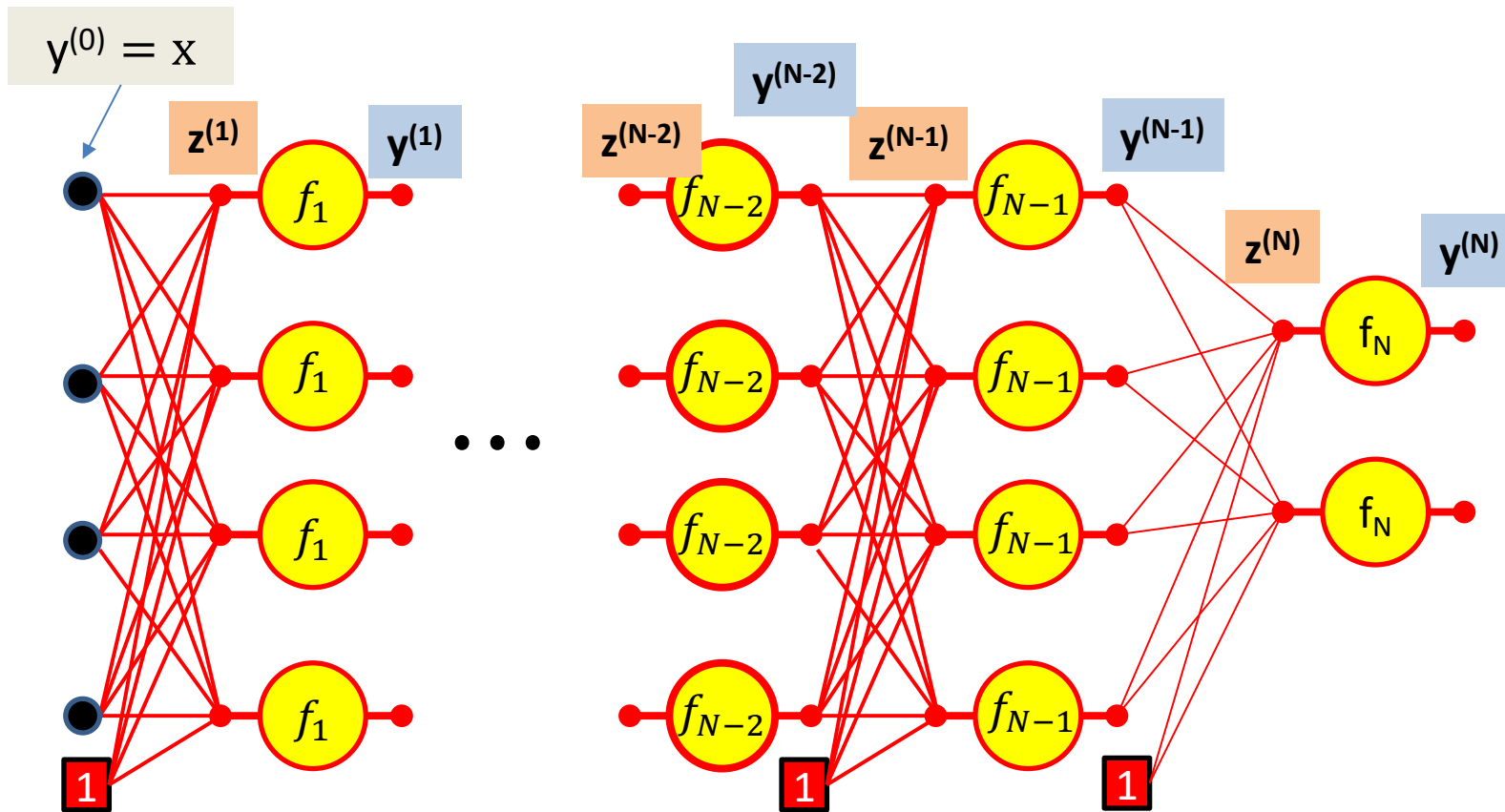ITERATE FOR k = 1:N

for j = 1:layer-width

$$z_j^{(k)} = \sum_i w_{ij}^{(k)} y_i^{(k-1)}$$

$$y_j^{(k)} = f_k\left(z_j^{(k)}\right)$$

# Forward "Pass"

- Input: $D$ dimensional vector $\mathbf{x} = [x_j, \ j = 1 \dots D]$
- Set:
  - $D_0 = D$, is the width of the $0^{\text{th}}$ (input) layer
  - $y_j^{(0)} = x_j, \ j = 1 \dots D; \qquad y_0^{(k=1\dots N)} = x_0 = 1$

- For layer $k = 1 \dots N$
  - For $j = 1 \dots D_k$    $D_k$ is the size of the kth layer
    - $z_j^{(k)} = \sum_{i=0}^{D_{k-1}} w_{i,j}^{(k)} y_i^{(k-1)}$
    - $y_j^{(k)} = f_k\left(z_j^{(k)}\right)$

- Output:
  - $Y = y_j^{(N)}, j = 1 .. D_N$

# Computing derivatives



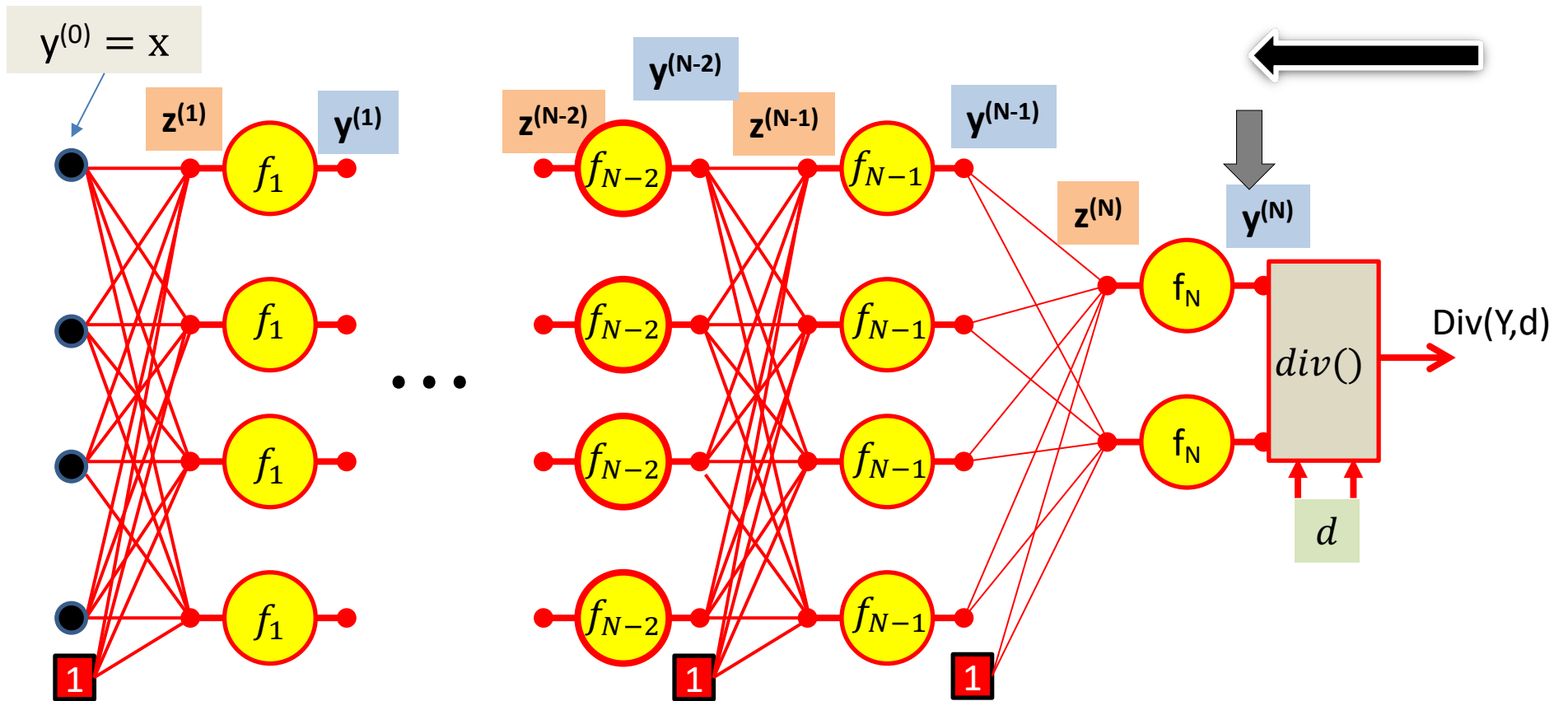We have computed all these intermediate values in the forward computation

We must remember them – we will need them to compute the derivatives

# Computing derivatives



First, we compute the divergence between the output of the net y = $y^{(N)}$ and the desired output $d$

# Computing derivatives



We then compute $\nabla_{Y^{(N)}} div(.)$ the derivative of the divergence w.r.t. the final output of the network $y^{(N)}$

# Computing derivatives



We then compute $\nabla_{Y^{(N)}} div(.)$ the derivative of the divergence w.r.t. the final output of the network $y^{(N)}$

We then compute $\nabla_{z^{(N)}} div(.)$ the derivative of the divergence w.r.t. the *pre-activation* affine combination $z^{(N)}$ using the chain rule

# Computing derivatives



Continuing on, we will compute $\nabla_{W^{(N)}} div(.)$ the derivative of the divergence with respect to the weights of the connections to the output layer

# Computing derivatives



Continuing on, we will compute $\nabla_{W^{(N)}} div(.)$ the derivative of the divergence with respect to the weights of the connections to the output layer

Then continue with the chain rule to compute $\nabla_{Y^{(N-1)}} div(.)$ the derivative of the divergence w.r.t. the output of the N-1th layer

# Computing derivatives



$y^{(0)} = x$

$z^{(1)}$    $y^{(1)}$

$z^{(N-2)}$    $y^{(N-2)}$    $z^{(N-1)}$    $y^{(N-1)}$

$z^{(N)}$    $y^{(N)}$

$f_1$    $f_1$    $f_1$    $f_1$

$f_{N-2}$    $f_{N-2}$    $f_{N-2}$    $f_{N-2}$

$f_{N-1}$    $f_{N-1}$    $f_{N-1}$    $f_{N-1}$

$f_N$    $f_N$

$div()$    Div(Y,d)

$d$

We continue our way backwards in the order shown

$\nabla_{z^{(N-1)}} div(.)$

$$y^{(0)} = x$$

$$z^{(1)} \quad y^{(1)}$$

$$y^{(N-2)}$$

$$z^{(N-2)} \quad z^{(N-1)} \quad y^{(N-1)}$$

$$z^{(N)} \quad y^{(N)}$$

$$f_1 \quad f_1 \quad f_1 \quad f_1$$

$$f_{N-2} \quad f_{N-1}$$

$$f_N$$

$$div()$$

$$Div(Y,d)$$

$$d$$

We continue our way backwards in the order shown

$$\nabla_{W^{(N-1)}} div(.)$$

$$y^{(0)} = x$$

$$z^{(1)}$$ $$y^{(1)}$$

$$f_1$$

$$z^{(N-2)}$$ $$y^{(N-2)}$$

$$f_{N-2}$$ $$f_{N-1}$$

$$z^{(N-1)}$$ $$y^{(N-1)}$$

$$z^{(N)}$$ $$y^{(N)}$$

$$f_N$$

$$div()$$

Div(Y,d)

$$d$$

1

We continue our way backwards in the order shown

$$\nabla_{Y^{(N-2)}} div(.)$$

$$y^{(0)} = x$$

$$z^{(1)}$$  $$y^{(1)}$$  $$f_1$$

$$y^{(N-2)}$$

$$z^{(N-2)}$$  $$f_{N-2}$$  $$z^{(N-1)}$$  $$f_{N-1}$$  $$y^{(N-1)}$$

$$z^{(N)}$$  $$y^{(N)}$$

$$f_N$$

$$div()$$

Div(Y,d)

$$d$$

We continue our way backwards in the order shown

$$\nabla_{z^{(N-2)}} div(.)$$

$y^{(0)} = x$

$z^{(1)}$  $y^{(1)}$  $y^{(N-2)}$  $z^{(N-2)}$  $z^{(N-1)}$  $y^{(N-1)}$

$z^{(N)}$  $y^{(N)}$

$f_1$  $f_{N-2}$  $f_{N-1}$  $f_N$

$div()$  Div(Y,d)

$d$

We continue our way backwards in the order shown

$\nabla_{Y^{(1)}} div(.)$

$y^{(0)} = x$

$z^{(1)}$ $y^{(1)}$ $y^{(N-2)}$ $z^{(N-2)}$ $z^{(N-1)}$ $y^{(N-1)}$ $z^{(N)}$ $y^{(N)}$

$f_1$ $f_1$ $f_1$ $f_1$

$f_{N-2}$ $f_{N-2}$ $f_{N-2}$ $f_{N-2}$

$f_{N-1}$ $f_{N-1}$ $f_{N-1}$ $f_{N-1}$

$f_N$ $f_N$

$div()$

Div(Y,d)

$d$

We continue our way backwards in the order shown

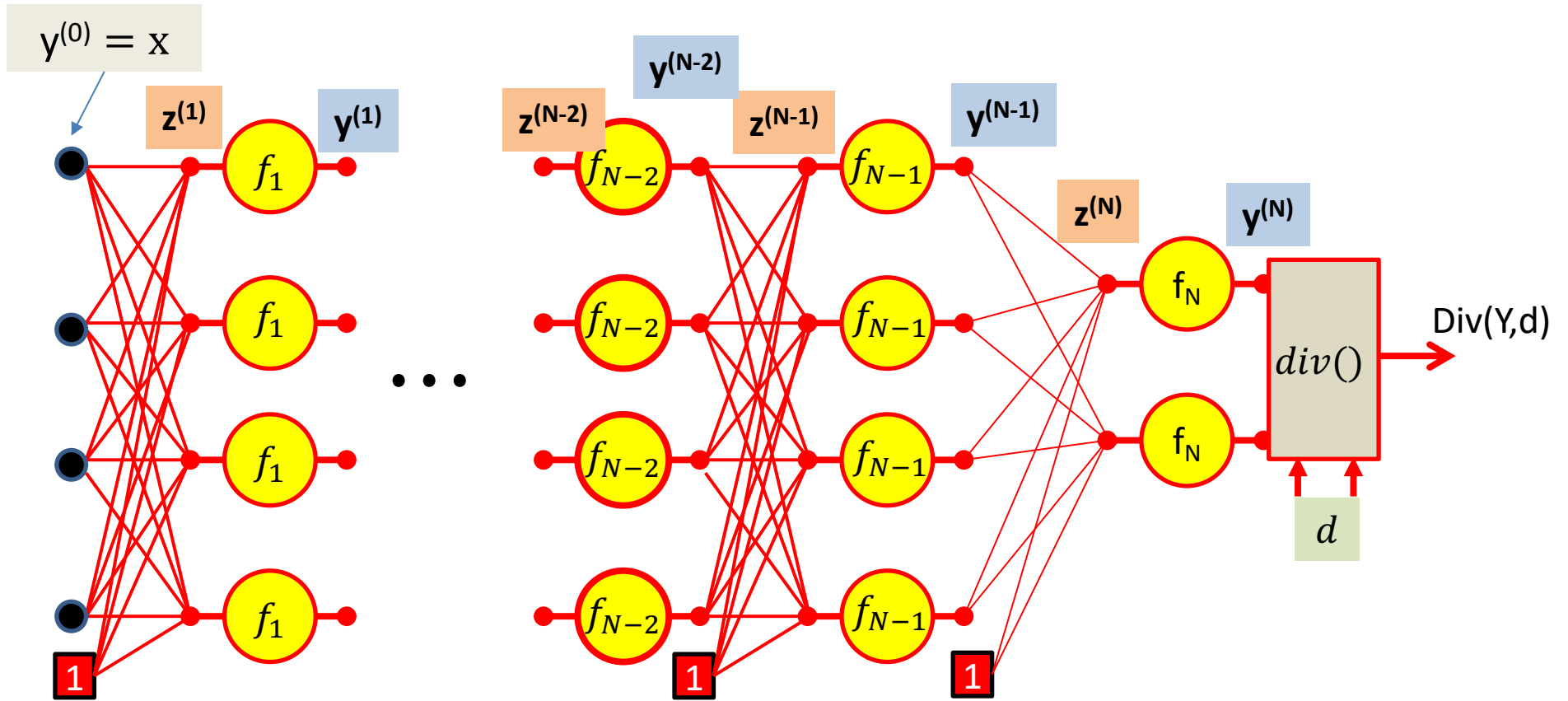$\nabla_{z^{(1)}} div(.)$

$$\nabla_{W^{(1)}} div(.)$$

We continue our way backwards in the order shown

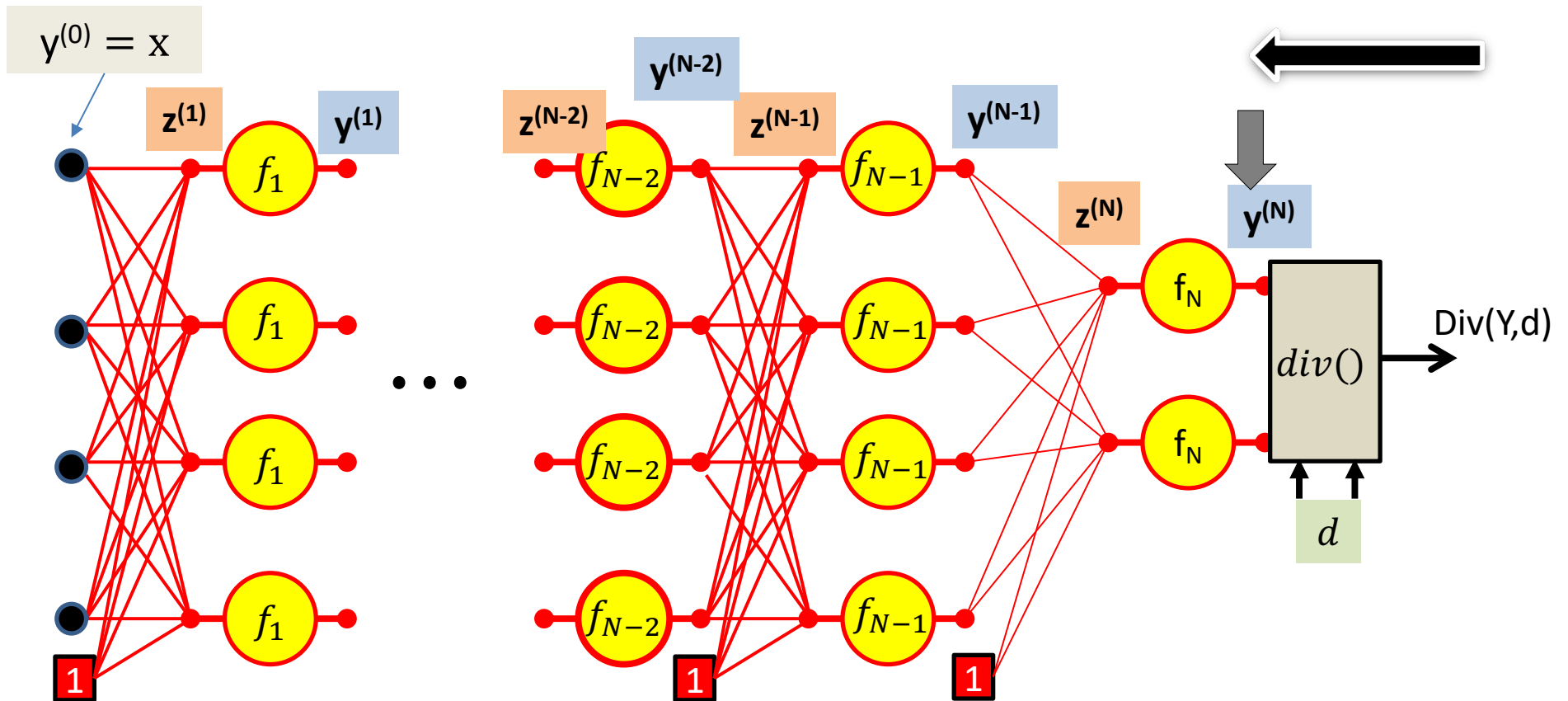# Backward Gradient Computation

- Lets actually see the math..

# Computing derivatives

# Computing derivatives
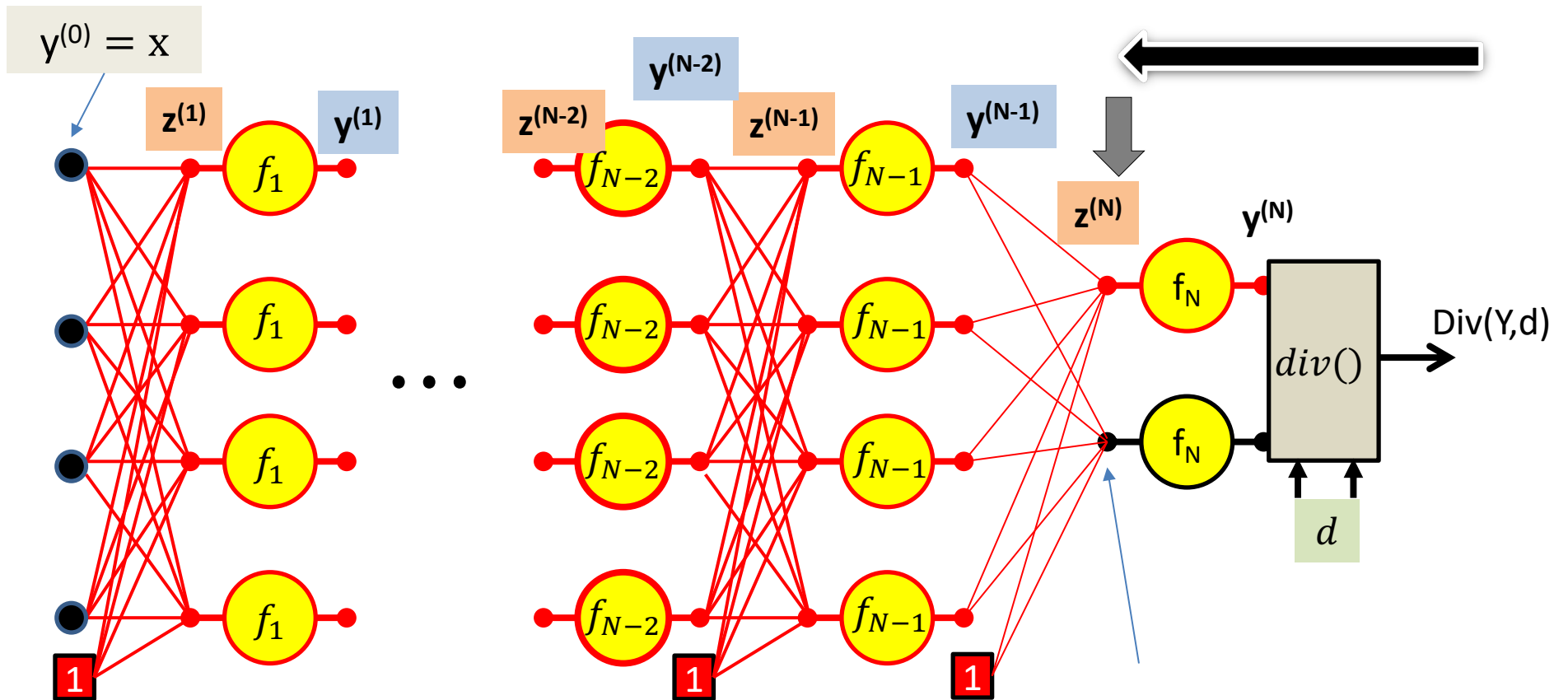


The derivative w.r.t the actual output of the network is simply the derivative w.r.t to the output of the final layer of the network

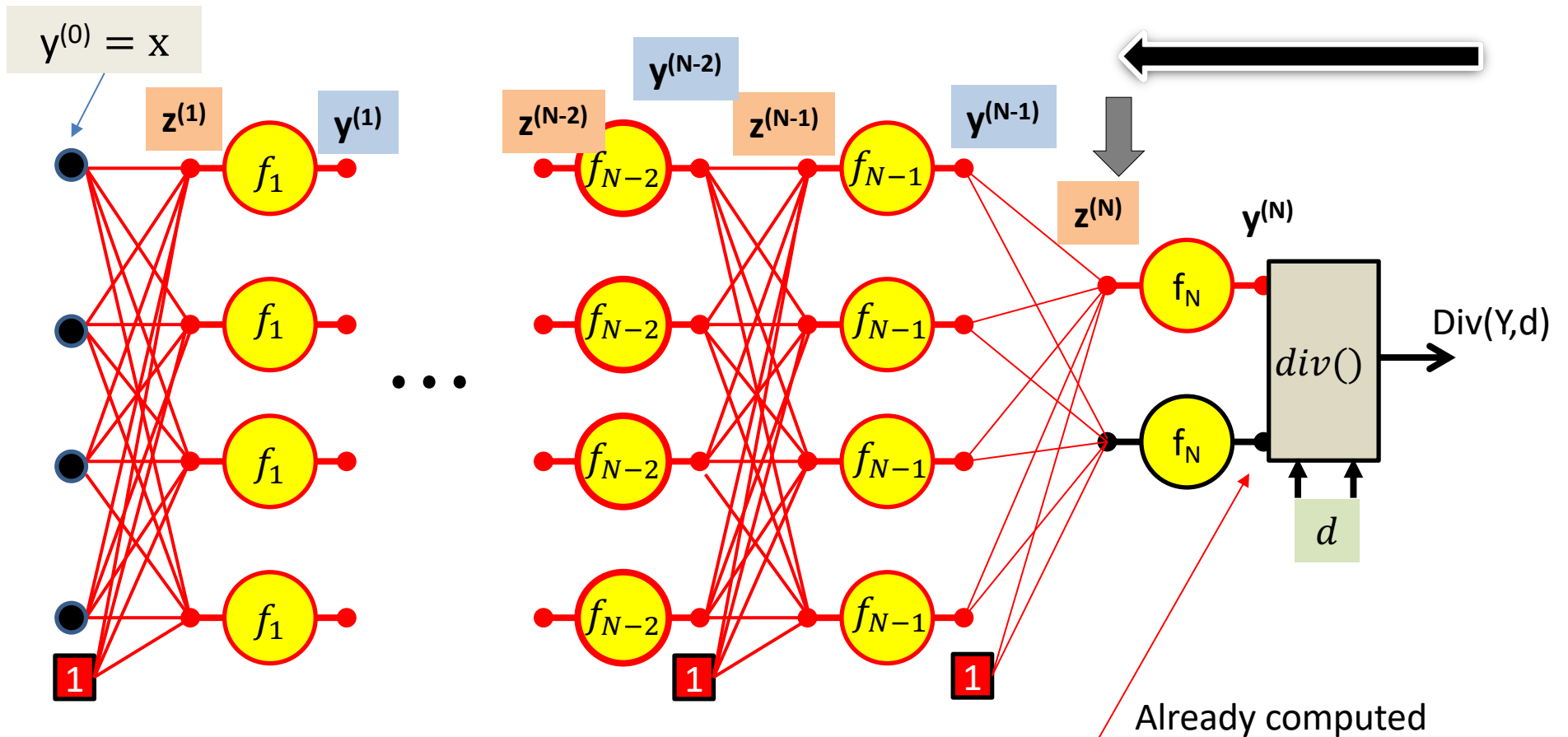$$\frac{\partial Div(Y,d)}{\partial y_i} = \frac{\partial Div(Y,d)}{\partial y_i^{(N)}}$$

# Computing derivatives



$$\frac{\partial Div}{\partial z_1^{(N)}} = \frac{\partial y_1^{(N)}}{\partial z_1^{(N)}} \frac{\partial Div}{\partial y_1^{(N)}}$$

# Computing derivatives



$$\frac{\partial Div}{\partial z_1^{(N)}} = \frac{\partial y_1^{(N)}}{\partial z_1^{(N)}} \frac{\partial Div}{\partial y_1^{(N)}}$$

# Computing derivatives



$$\frac{\partial Div}{\partial z_1^{(N)}} = \frac{\partial y_1^{(N)}}{\partial z_1^{(N)}} \frac{\partial Div}{\partial y_1^{(N)}}$$

$$f_N'\left(z_1^{(N)}\right)$$

Derivative of activation function

# Computing derivatives



$y^{(0)} = x$

$z^{(1)}$  $y^{(1)}$

$y^{(N-2)}$

$z^{(N-2)}$  $z^{(N-1)}$  $y^{(N-1)}$

$z^{(N)}$  $y^{(N)}$

$f_1$  $f_1$  $f_1$  $f_1$

$f_{N-2}$  $f_{N-1}$

$f_N$

$div()$

Div(Y,d)

$d$

$f_N' \left( z_1^{(N)} \right)$

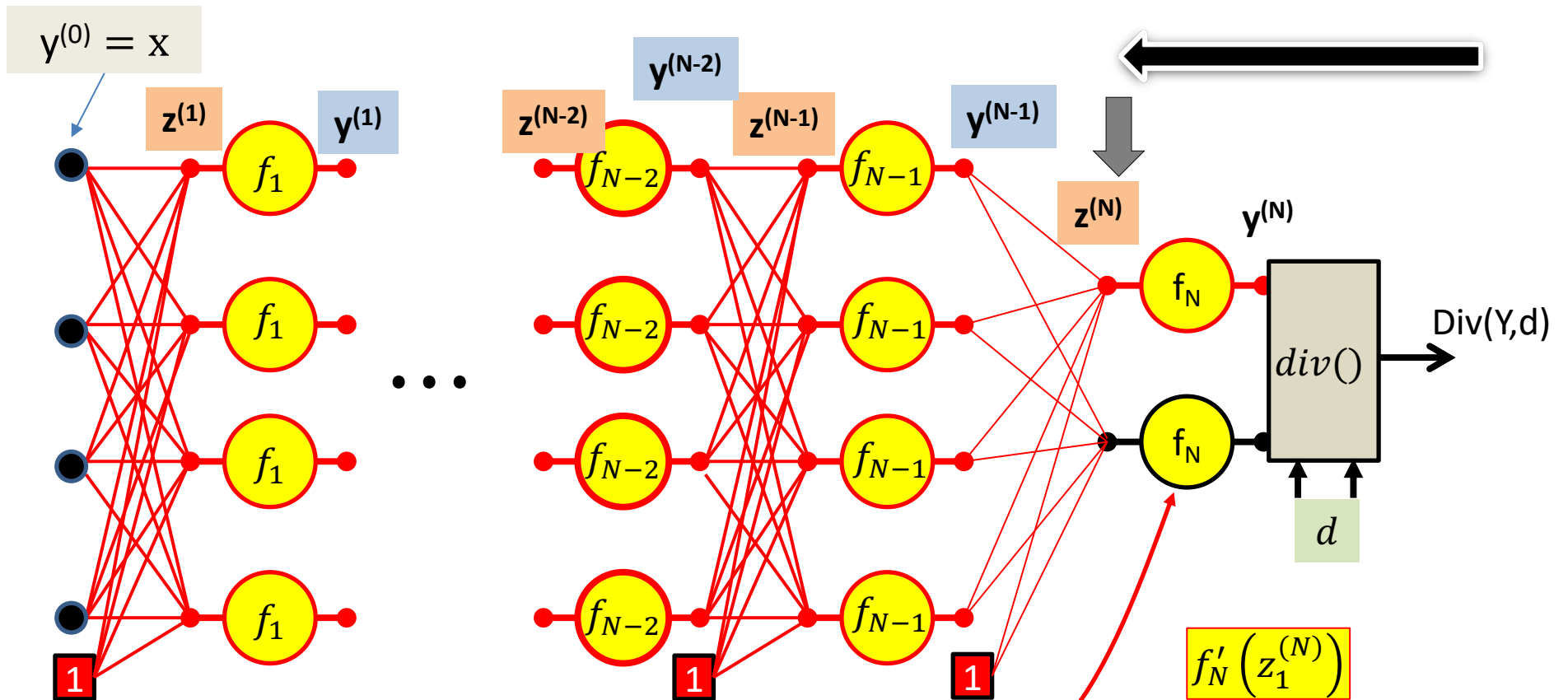Derivative of activation function

Computed in forward pass

$$\frac{\partial Div}{\partial z_1^{(N)}} = \frac{\partial y_1^{(N)}}{\partial z_1^{(N)}} \frac{\partial Div}{\partial y_1^{(N)}}$$

# Computing derivatives



$$\frac{\partial Div}{\partial z_1^{(N)}} = f_N'\left(z_1^{(N)}\right)\frac{\partial Div}{\partial y_1^{(N)}}$$

# Computing derivatives



$$\frac{\partial Div}{\partial z_i^{(N)}} = f_N'\left(z_i^{(N)}\right)\frac{\partial Div}{\partial y_i^{(N)}}$$

# Computing derivatives



$$\frac{\partial Div}{\partial w_{11}^{(N)}} = \frac{\partial z_1^{(N)}}{\partial w_{11}^{(N)}} \frac{\partial Div}{\partial z_1^{(N)}}$$
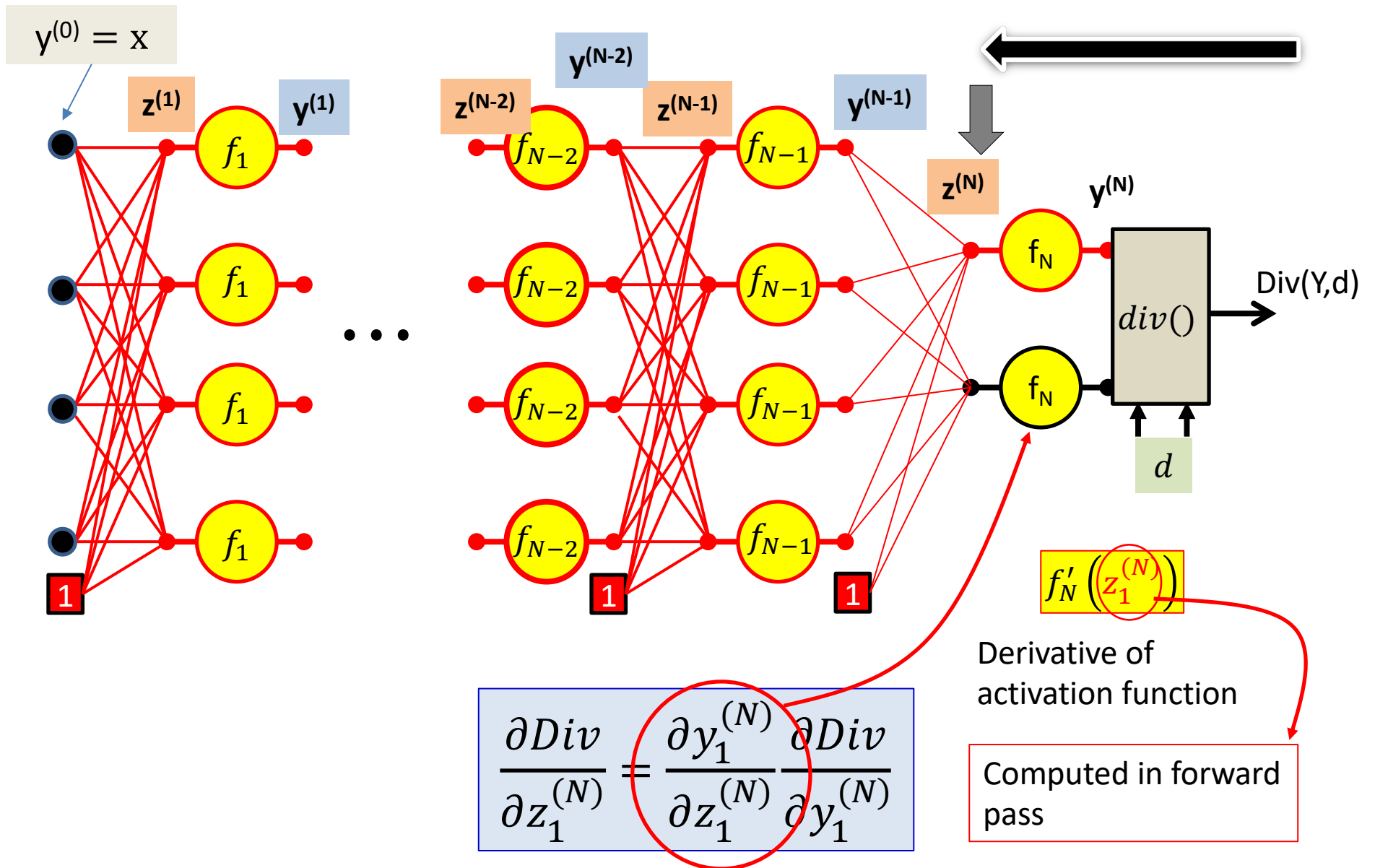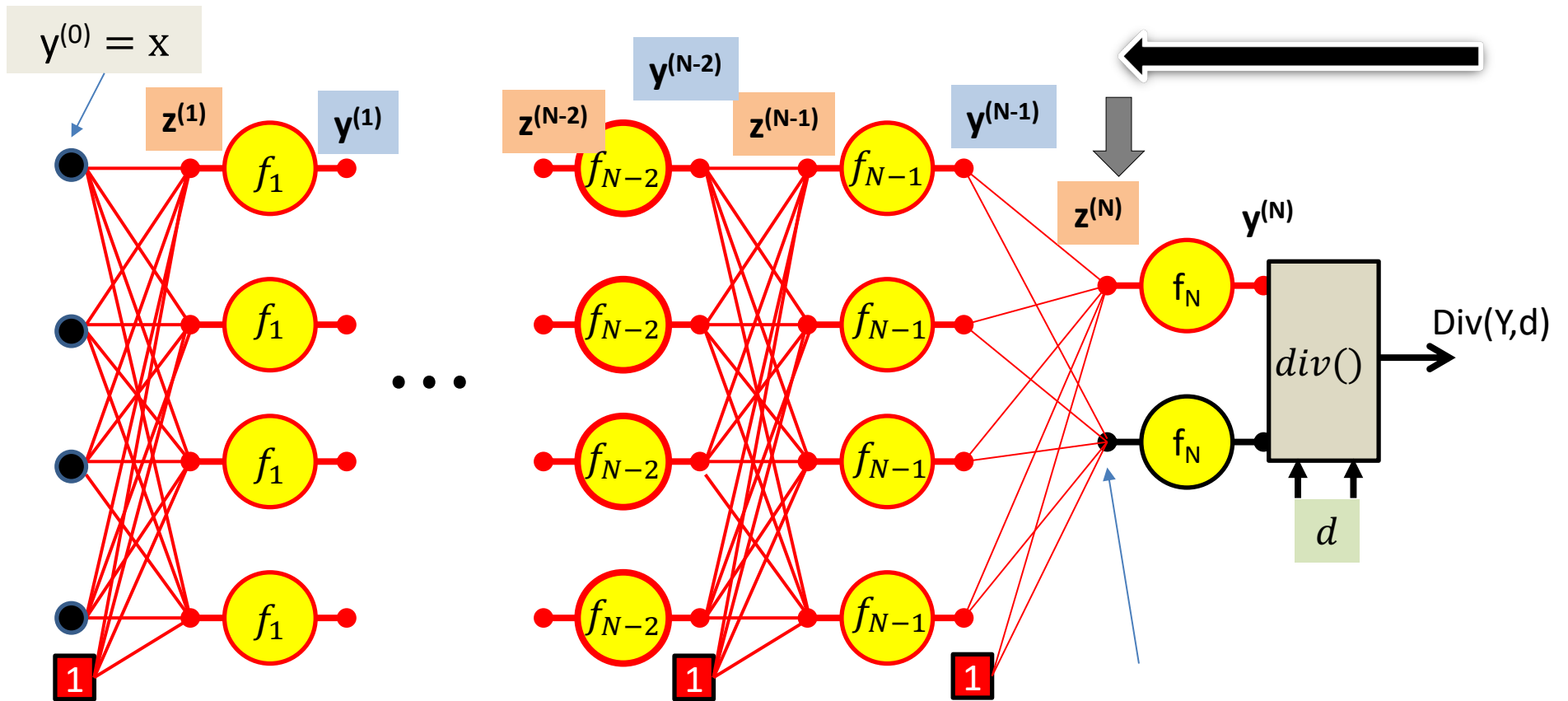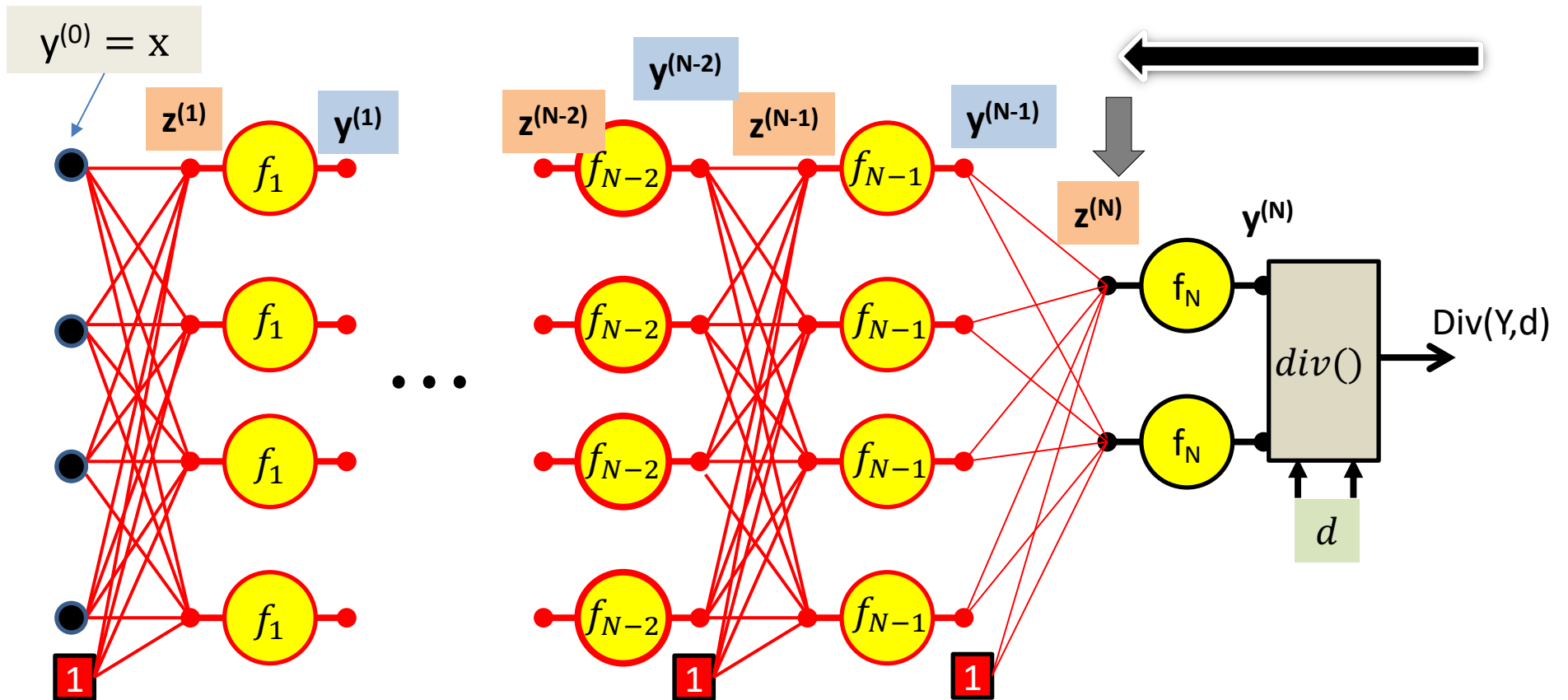
# Computing derivatives



$$\frac{\partial Div}{\partial w_{11}^{(N)}} = \frac{\partial z_1^{(N)}}{\partial w_{11}^{(N)}} \cdot \frac{\partial Div}{\partial z_1^{(N)}}$$

Just computed

# Computing derivatives



$$\frac{\partial Div}{\partial w_{11}^{(N)}} = \frac{\partial z_1^{(N)}}{\partial w_{11}^{(N)}} \frac{\partial Div}{\partial z_1^{(N)}}$$

$y_1^{(N-1)}$

Because
$$z_1^{(N)} = w_{11}^{(N)} y_1^{(N-1)} + \text{other terms}$$

# Computing derivatives



$$\frac{\partial Div}{\partial w_{11}^{(N)}} = \frac{\partial z_1^{(N)}}{\partial w_{11}^{(N)}} \frac{\partial Div}{\partial z_1^{(N)}}$$
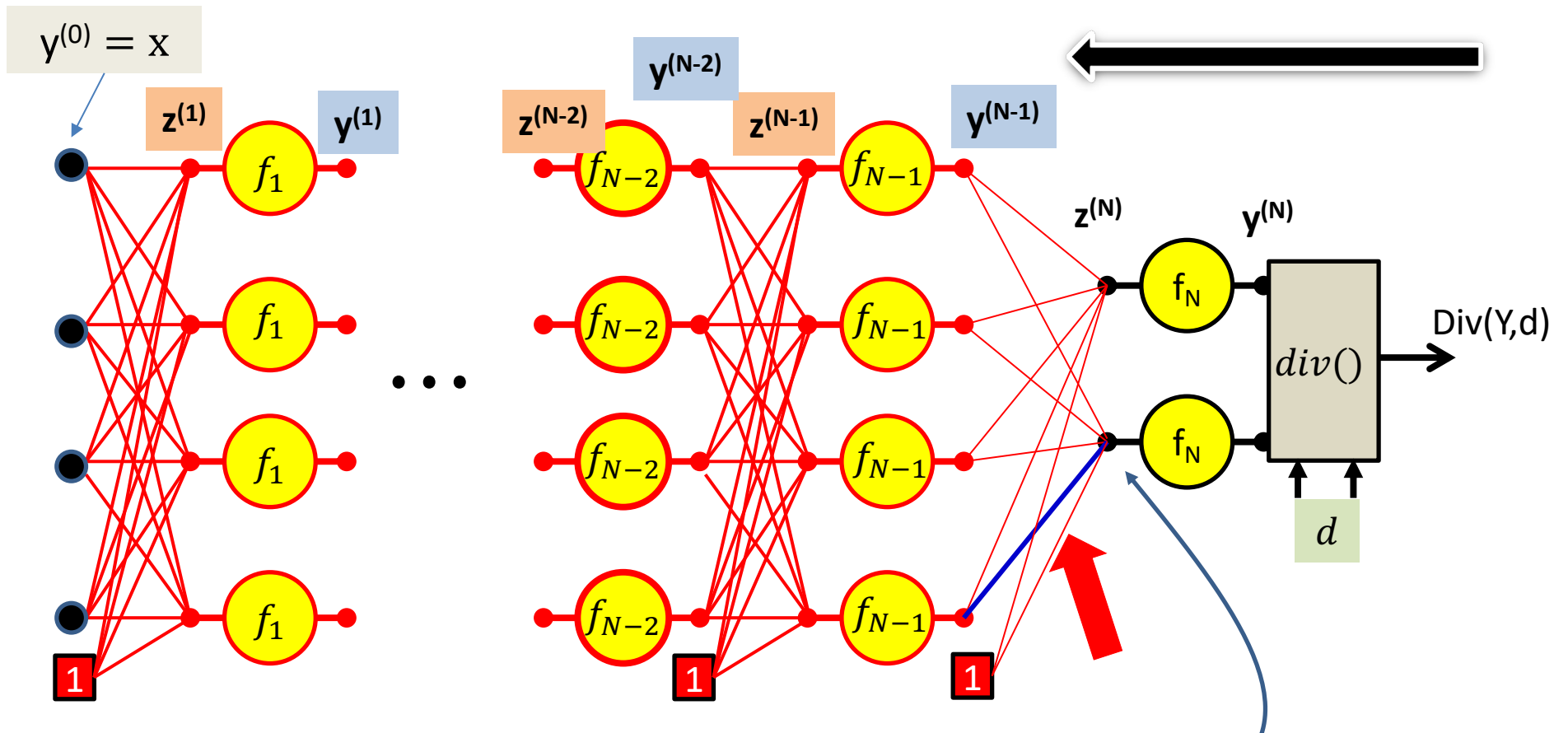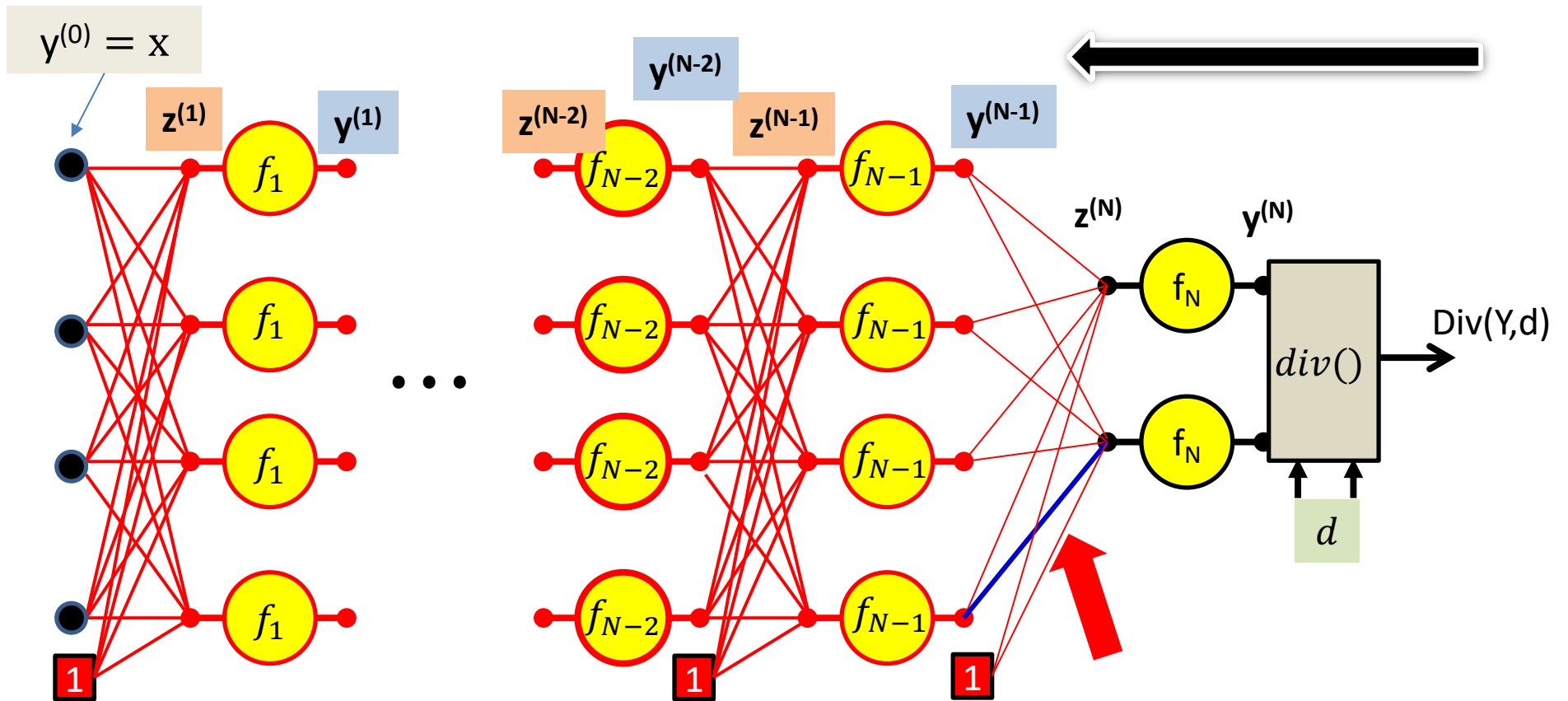
$y_1^{(N-1)}$

Computed in forward pass

Because
$$z_1^{(N)} = w_{11}^{(N)} y_1^{(N-1)} + \text{other terms}$$

# Computing derivatives



$y^{(0)} = x$

$z^{(1)}$ $y^{(1)}$ $y^{(N-2)}$ $z^{(N-2)}$ $z^{(N-1)}$ $y^{(N-1)}$

$z^{(N)}$ $y^{(N)}$

$f_1$ ... $f_{N-2}$ $f_{N-1}$ $f_N$ $div()$ Div(Y,d)

$d$

$$\frac{\partial Div}{\partial w_{11}^{(N)}} = y_1^{(N-1)} \frac{\partial Div}{\partial z_1^{(N)}}$$

# Computing derivatives



$$\frac{\partial Div}{\partial w_{ij}^{(N)}} = y_i^{(N-1)} \frac{\partial Div}{\partial z_j^{(N)}}$$

For the bias term $y_0^{(N-1)} = 1$

# Computing derivatives



$$\frac{\partial Div}{\partial y_1^{(N-1)}} = \sum_j \frac{\partial z_j^{(N)}}{\partial y_1^{(N-1)}} \frac{\partial Div}{\partial z_j^{(N)}}$$
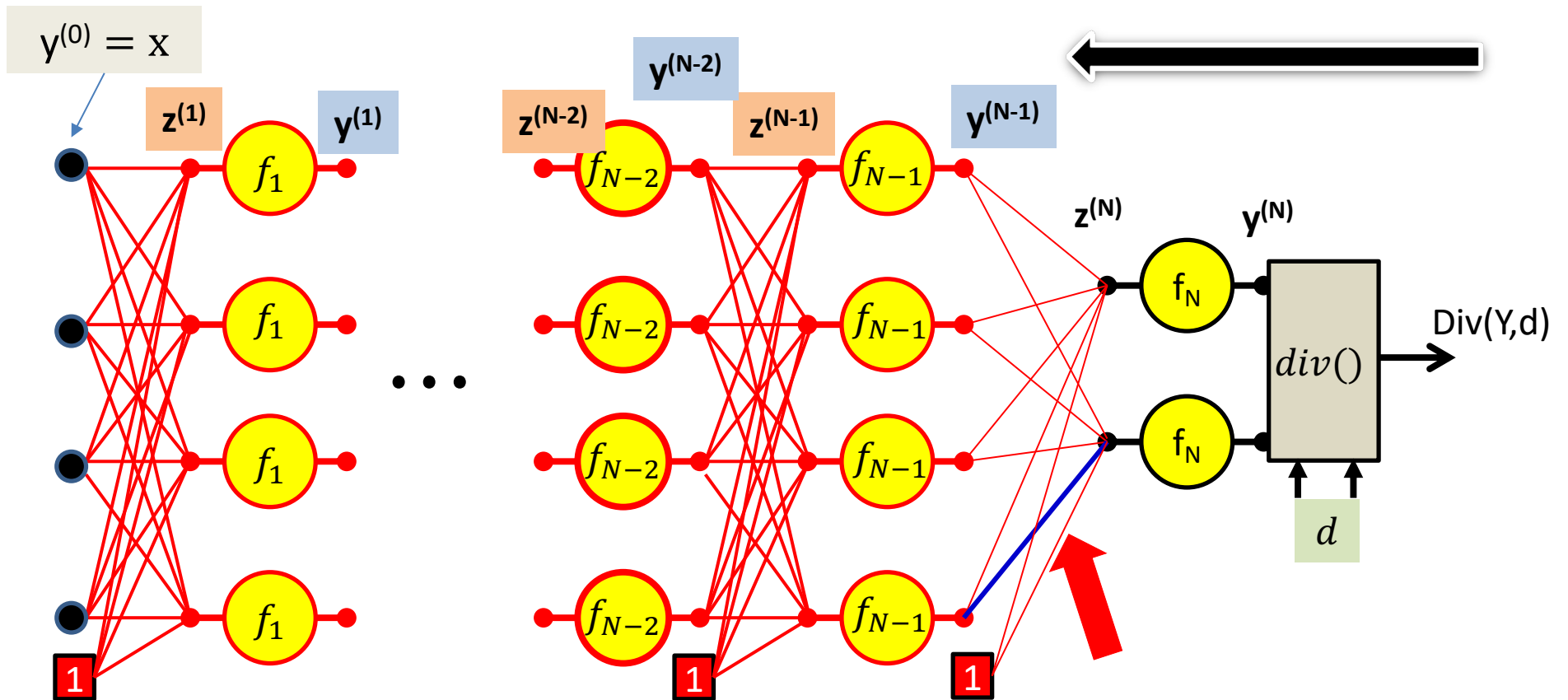
# Computing derivatives



$$\frac{\partial Div}{\partial y_1^{(N-1)}} = \sum_j \frac{\partial z_j^{(N)}}{\partial y_1^{(N-1)}} \frac{\partial Div}{\partial z_j^{(N)}}$$

Already computed

# Computing derivatives



$y^{(0)} = x$

$z^{(1)}$    $y^{(1)}$

$y^{(N-2)}$

$z^{(N-2)}$    $z^{(N-1)}$    $y^{(N-1)}$

$f_1$   $f_{N-2}$   $f_{N-1}$

$z^{(N)}$   $y^{(N)}$

$f_N$

$div()$   Div(Y,d)

$d$

$$\frac{\partial Div}{\partial y_1^{(N-1)}} = \sum_j \frac{\partial z_j^{(N)}}{\partial y_1^{(N-1)}} \frac{\partial Div}{\partial z_j^{(N)}}$$

$w_{1j}^{(N-1)}$

Because
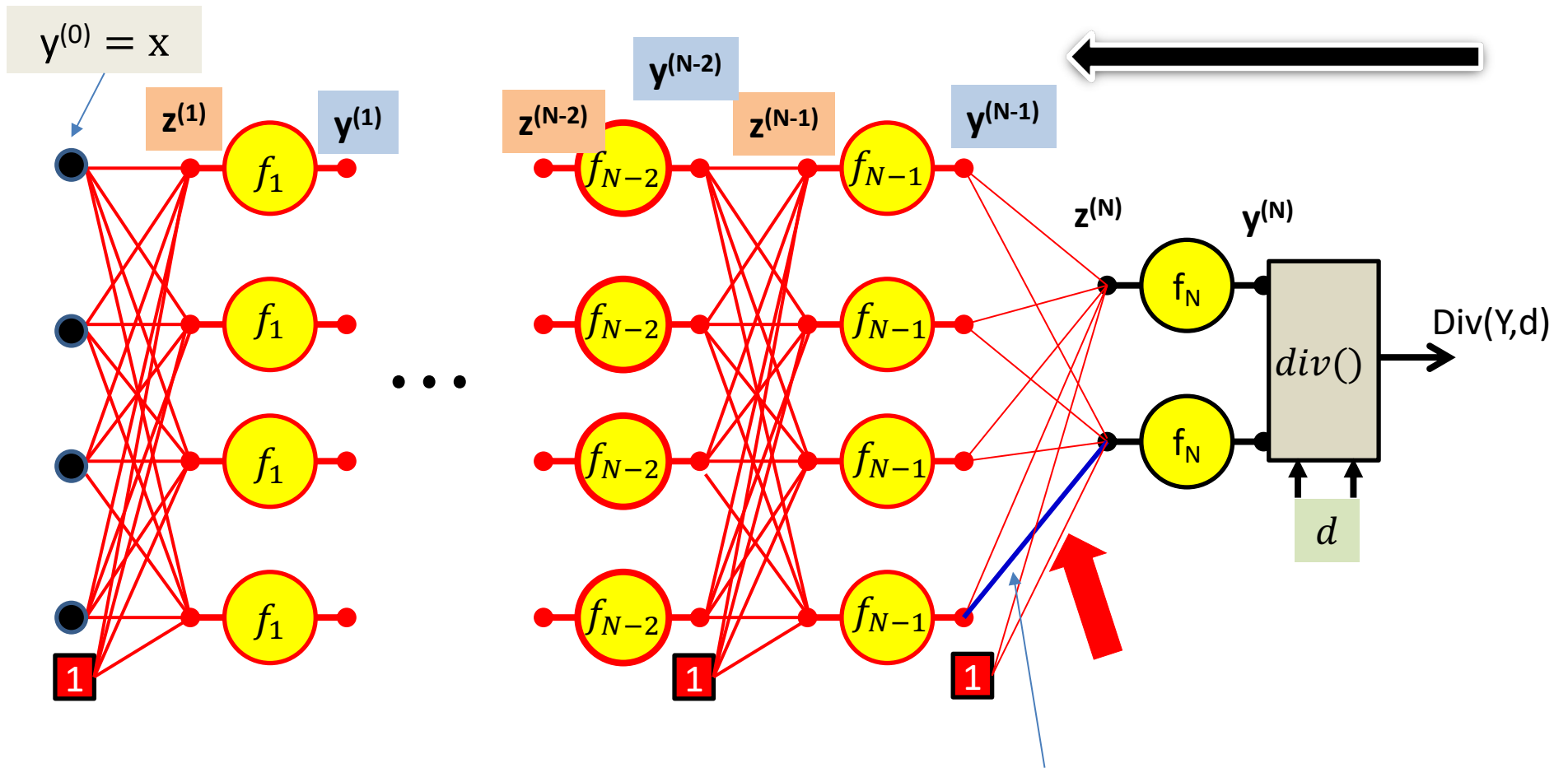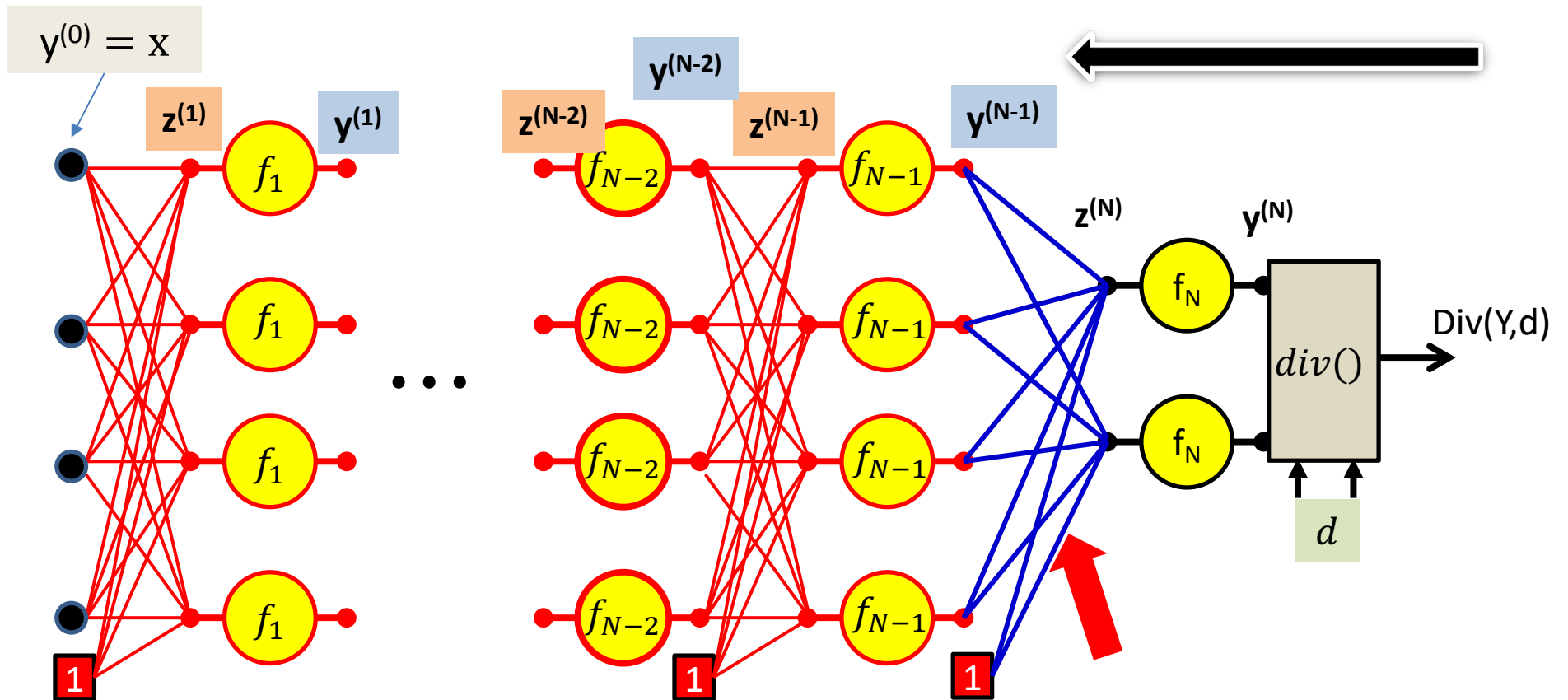$$z_j^{(N)} = w_{1j}^{(N)} y_1^{(N-1)} + \text{other terms}$$

# Computing derivatives



$$\frac{\partial Div}{\partial y_1^{(N-1)}} = \sum_j w_{1j}^{(N)} \frac{\partial Div}{\partial z_j^{(N)}}$$

# Computing derivatives



$y^{(0)} = x$

$$\frac{\partial Div}{\partial y_i^{(N-1)}} = \sum_j w_{ij}^{(N)} \frac{\partial Div}{\partial z_j^{(N)}}$$

# Computing derivatives



We continue our way backwards in the order shown

$$\frac{\partial Div}{\partial z_i^{(N-1)}} = f'_{N-1}\left(z_i^{(N-1)}\right)\frac{\partial Div}{\partial y_i^{(N-1)}}$$

$y^{(0)} = x$

$z^{(1)}$

$y^{(1)}$

$y^{(N-2)}$

$z^{(N-2)}$

$z^{(N-1)}$

$y^{(N-1)}$

$z^{(N)}$

$y^{(N)}$

$f_1$

$f_{N-2}$

$f_{N-1}$

$f_N$

$div()$

$Div(Y,d)$

$d$

We continue our way backwards in the order shown

$$\frac{\partial Div}{\partial w_{ij}^{(N-1)}} = y_i^{(N-2)} \frac{\partial Div}{\partial z_j^{(N-1)}}$$

For the bias term $y_0^{(N-2)} = 1$

We continue our way backwards in the order shown

$$\frac{\partial Div}{\partial y_i^{(N-2)}} = \sum_j w_{ij}^{(N-1)} \frac{\partial Div}{\partial z_j^{(N-1)}}$$

We continue our way backwards in the order shown

$$\frac{\partial Div}{\partial z_i^{(N-2)}} = f'_{N-2}\left(z_i^{(N-2)}\right)\frac{\partial Div}{\partial y_i^{(N-2)}}$$

We continue our way backwards in the order shown

$$\frac{\partial Div}{\partial y_1^{(1)}} = \sum_j w_{ij}^{(2)} \frac{\partial Div}{\partial z_j^{(2)}}$$

We continue our way backwards in the order shown

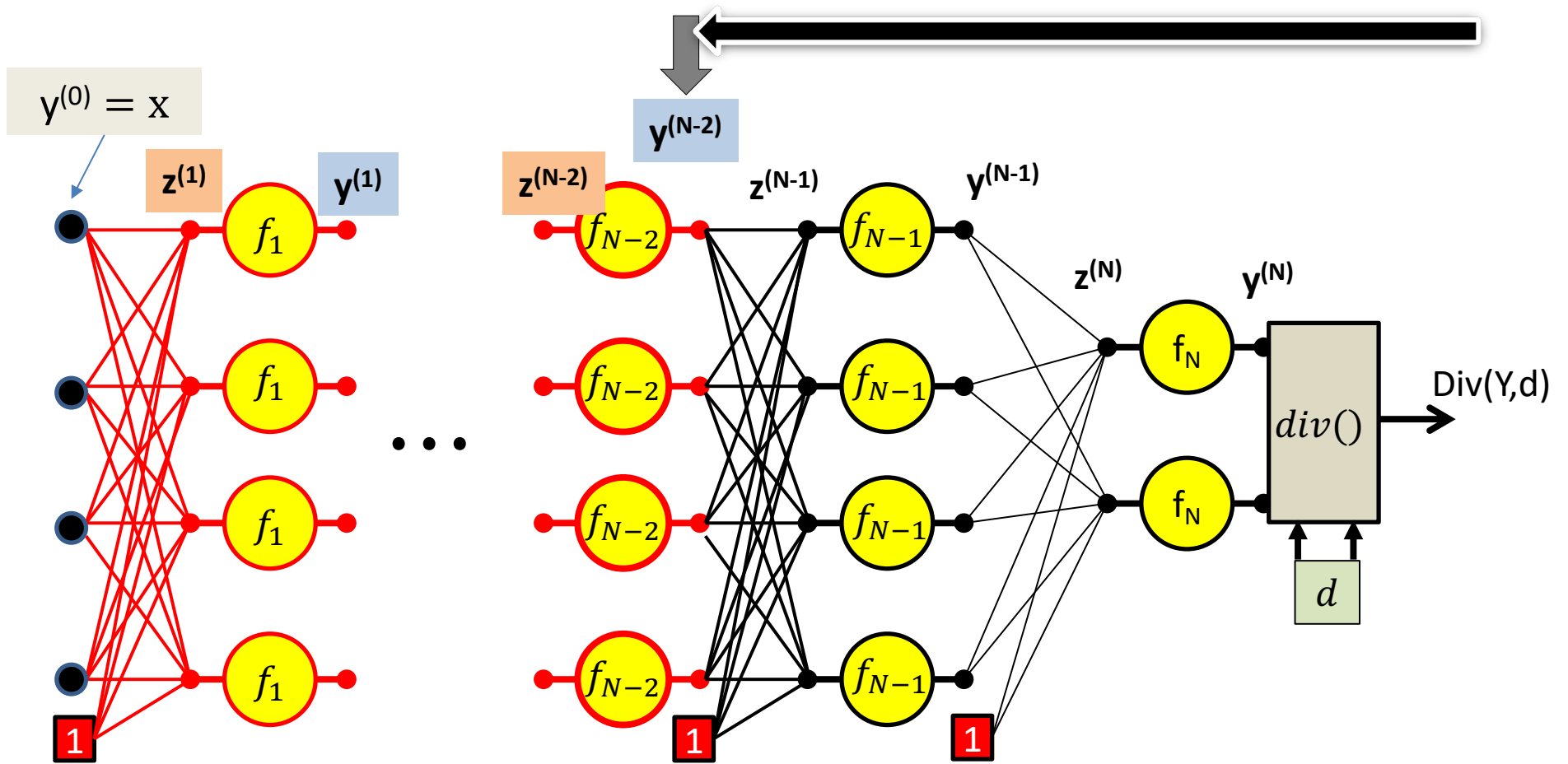$$\frac{\partial Div}{\partial z_i^{(1)}} = f_1'\left(z_i^{(1)}\right)\frac{\partial Div}{\partial y_i^{(1)}}$$
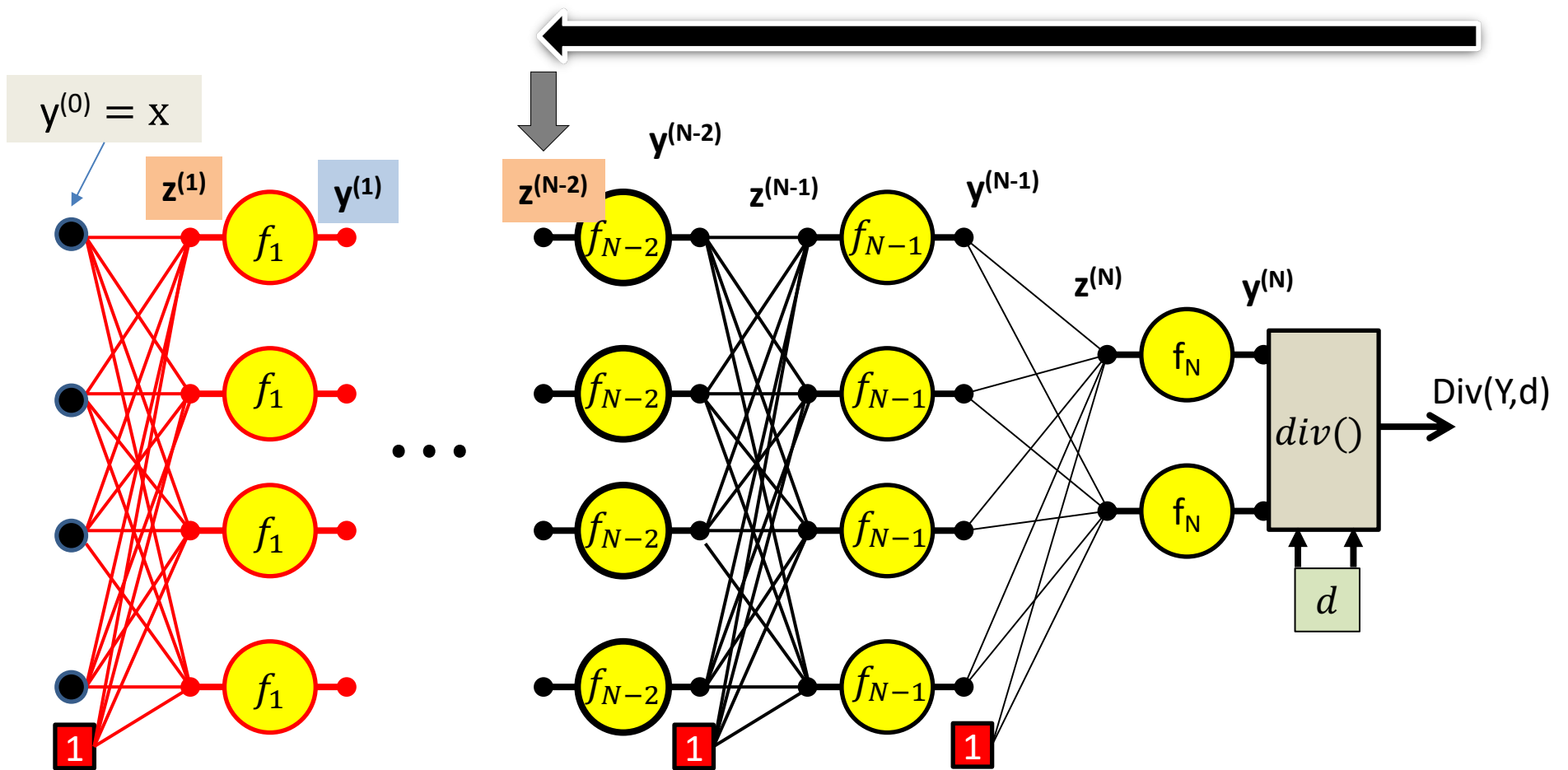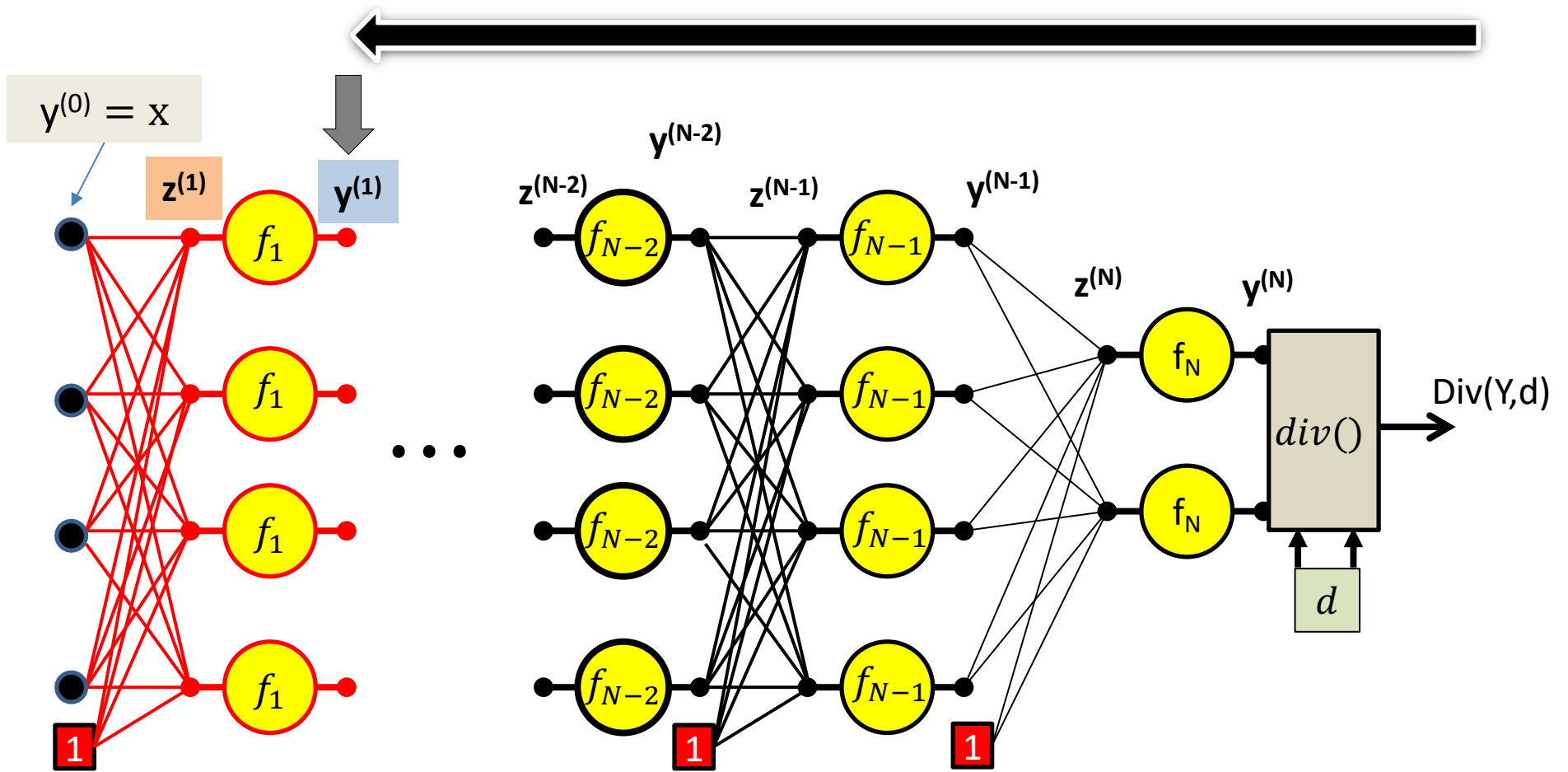
We continue our way backwards in the order shown

$$\frac{\partial Div}{\partial w_{ij}^{(1)}} = y_i^{(1)} \frac{\partial Div}{\partial z_j^{(1)}}$$

# Gradients: Backward Computation



$z^{(k-1)}$  $y^{(k-1)}$  $z^{(k)}$  $y^{(k)}$  $z^{(N-1)}$  $y^{(N-1)}$

$z^{(N)}$  $y^{(N)}$  $f_N$

Div(Y,d)

Div(Y,d)

$f_N$

Figure assumes, but does not show the "1" bias nodes

Initialize: Gradient w.r.t network output

$$\frac{\partial Div}{\partial y_i} = \frac{\partial Div(Y,d)}{\partial y_i^{(N)}}$$

$$\frac{\partial Div}{\partial z_i^{(N)}} = f_k'\left(z_i^{(N)}\right)\frac{\partial Div}{\partial y_i^{(N)}}$$

$For\ k\ =\ N-1..0$
  $For\ i\ =\ 1: layer\ width$

$$\frac{\partial Div}{\partial y_i^{(k)}} = \sum_j w_{ij}^{(k+1)} \frac{\partial Div}{\partial z_j^{(k+1)}} \qquad \frac{\partial Div}{\partial z_i^{(k)}} = f_k'\left(z_i^{(k)}\right)\frac{\partial Div}{\partial y_i^{(k)}}$$

$$\forall j\ \frac{\partial Div}{\partial w_{ij}^{(k+1)}} = y_i^{(k)}\frac{\partial Div}{\partial z_j^{(k+1)}}$$

# Backward Pass

- Output layer (N) :
  - For $i = 1 \dots D_N$

    - $\dfrac{\partial Div}{\partial y_i} = \dfrac{\partial Div(Y,d)}{\partial y_i^{(N)}}$

    - $\dfrac{\partial Div}{\partial z_i^{(N)}} = \dfrac{\partial Div}{\partial y_i^{(N)}} \dfrac{\partial y_i^{(N)}}{\partial z_i^{(N)}}$

- For layer $k = N - 1 \; downto \; 0$
  - For $i = 1 \dots D_k$

    - $\dfrac{\partial Div}{\partial y_i^{(k)}} = \sum_j w_{ij}^{(k+1)} \dfrac{\partial Div}{\partial z_j^{(k+1)}}$

    - $\dfrac{\partial Div}{\partial z_i^{(k)}} = \dfrac{\partial Div}{\partial y_i^{(k)}} f_k' \left( z_i^{(k)} \right)$

    - $\dfrac{\partial Div}{\partial w_{ji}^{(k+1)}} = y_j^{(k)} \dfrac{\partial Div}{\partial z_i^{(k+1)}}$    for $j = 1 \dots D_{k+1}$

# Backward Pass

- Output layer (N) :
  - For $i = 1 \ldots D_N$

    - $\dfrac{\partial Div}{\partial y_i} = \dfrac{\partial Div(Y,d)}{\partial y_i^{(N)}}$

    - $\dfrac{\partial Div}{\partial z_i^{(N)}} = \dfrac{\partial Div}{\partial y_i^{(N)}} \dfrac{\partial y_i^{(N)}}{\partial z_i^{(N)}}$

- For layer $k = N - 1 \; downto \; 0$
  - For $i = 1 \ldots D_k$

    - $\dfrac{\partial Div}{\partial y_i^{(k)}} = \sum_j w_{ij}^{(k+1)} \dfrac{\partial Div}{\partial z_j^{(k+1)}}$

    - $\dfrac{\partial Div}{\partial z_i^{(k)}} = \dfrac{\partial Div}{\partial y_i^{(k)}} f_k'\left(z_i^{(k)}\right)$

    - $\dfrac{\partial Div}{\partial w_{ji}^{(k+1)}} = y_j^{(k)} \dfrac{\partial Div}{\partial z_i^{(k+1)}}$   for $j = 1 \ldots D_{k+1}$

Called "Backpropagation" because the derivative of the loss is propagated "backwards" through the network

Very analogous to the forward pass:

Backward weighted combination of next layer

Backward equivalent of activation

128

# For comparison: the forward pass again

- Input: $D$ dimensional vector $\mathbf{x} = [x_j, \ j = 1 \dots D]$

- Set:
  - $D_0 = D,$ is the width of the $0^{\text{th}}$ (input) layer
  - $y_j^{(0)} = x_j, \ j = 1 \dots D; \qquad y_0^{(k=1\dots N)} = x_0 = 1$

- For layer $k = 1 \dots N$
  - For $j = 1 \dots D_k$
    - $z_j^{(k)} = \sum_{i=0}^{N_k} w_{i,j}^{(k)} y_i^{(k-1)}$
    - $y_j^{(k)} = f_k\left(z_j^{(k)}\right)$

- Output:
  - $Y = y_j^{(N)}, j = 1 .. D_N$

# Special cases



- Have assumed so far that
    1. The computation of the output of one neuron does not directly affect computation of other neurons in the same (or previous) layers
    2. Outputs of neurons only combine through weighted addition
    3. Activations are actually differentiable
    – All of these conditions are frequently not applicable
- Will not dwell on the topic in class, but explained in slides
    – Will appear in quiz.  Please read the slides

# Special Case 1. Vector activations



- Vector activations: all outputs are functions of all inputs

# Special Case 1. Vector activations



Scalar activation: Modifying a $z_i$ only changes corresponding $y_i$

$$y_i^{(k)} = f\left(z_i^{(k)}\right)$$

Vector activation: Modifying a $z_i$ potentially changes all, $y_1 \dots y_M$

$$\begin{bmatrix} y_1^{(k)} \\ y_2^{(k)} \\ \vdots \\ y_M^{(k)} \end{bmatrix} = f\left(\begin{bmatrix} z_1^{(k)} \\ z_2^{(k)} \\ \vdots \\ z_D^{(k)} \end{bmatrix}\right)$$

132

# "Influence" diagram



Scalar activation: Each $z_i$ influences *one* $y_i$

Vector activation: Each $z_i$ influences all, $y_1 \dots y_M$

# The number of outputs



- Note: The number of outputs ($y^{(k)}$) need not be the same as the number of inputs ($z^{(k)}$)
  - May be more or fewer

134

# Scalar Activation: Derivative rule



$$\frac{\partial Div}{\partial z_i^{(k)}} = \frac{\partial Div}{\partial y_i^{(k)}} \frac{dy_i^{(k)}}{dz_i^{(k)}}$$

- In the case of *scalar* activation functions, the derivative of the error w.r.t to the input to the unit is a simple product of derivatives

# Derivatives of vector activation



$$\frac{\partial Div}{\partial z_i^{(k)}} = \sum_j \frac{\partial Div}{\partial y_j^{(k)}} \frac{\partial y_j^{(k)}}{\partial z_i^{(k)}}$$

Note: derivatives of scalar activations are just a special case of vector activations:

$$\frac{\partial y_j^{(k)}}{\partial z_i^{(k)}} = 0 \ \ for \ \ i \ \neq j$$

- For *vector* activations the derivative of the error w.r.t. to any input is a sum of partial derivatives
  - Regardless of the number of outputs $y_j^{(k)}$

# Special cases

- Examples of vector activations and other special cases on slides
  - Please look up
  - Will appear in quiz!

# Overall Approach

- For each data instance
  - **Forward pass**: Pass instance forward through the net. Store all intermediate outputs of all computation
  - **Backward pass**: Sweep backward through the net, iteratively compute all derivatives w.r.t weights
- Actual loss is the sum of the divergence over all training instances

$$\mathbf{Loss} = \frac{1}{|\{X\}|} \sum_X Div(Y(X), d(X))$$

- Actual gradient is the sum or average of the derivatives computed for each training instance

$$\nabla_W \mathbf{Loss} = \frac{1}{|\{X\}|} \sum_X \nabla_W Div(Y(X), d(X)) \qquad W \leftarrow W - \eta \nabla_W \mathbf{Loss}^{\mathrm{T}}$$

# Training by BackProp

- Initialize all weights $(\boldsymbol{W}^{(1)}, \boldsymbol{W}^{(2)}, \ldots, \boldsymbol{W}^{(K)})$

- Do:

    - **Initialize** $Err = 0$; For all $i, j, k$, initialize $\dfrac{dErr}{dw_{i,j}^{(k)}} = 0$

    - For all $t = 1:T$ (Loop over training instances)
        - **Forward pass:** Compute
            - Output $\boldsymbol{Y_t}$
            - $Err += \boldsymbol{Div(Y_t, d_t)}$
        - **Backward pass:** For all $i, j, k$:
            - Compute $\dfrac{d\boldsymbol{Div(Y_t, d_t)}}{dw_{i,j}^{(k)}}$
            - Compute $\dfrac{dErr}{dw_{i,j}^{(k)}} += \dfrac{d\boldsymbol{Div(Y_t, d_t)}}{dw_{i,j}^{(k)}}$

    - For all $i, j, k$, update:

    $$w_{i,j}^{(k)} = w_{i,j}^{(k)} - \frac{\eta}{T} \frac{dErr}{dw_{i,j}^{(k)}}$$

- Until $Err$ has converged

# Vector formulation

- For layered networks it is generally simpler to think of the process in terms of vector operations
  - Simpler arithmetic
  - Fast matrix libraries make operations *much* faster

- We can restate the entire process in vector terms
  - On slides, please read
  - This is what is *actually* used in any real system
  - Will appear in quiz

# Vector formulation



$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

$$\mathbf{z}_k = \begin{bmatrix} z_1^{(k)} \\ z_2^{(k)} \\ \vdots \\ z_{D_k}^{(k)} \end{bmatrix}$$

$$\mathbf{y}_k = \begin{bmatrix} y_1^{(k)} \\ y_2^{(k)} \\ \vdots \\ y_{D_k}^{(k)} \end{bmatrix}$$

$$\mathbf{W}_k = \begin{bmatrix} w_{11}^{(k)} & w_{21}^{(k)} & \vdots & w_{D_{k-1}1}^{(k)} \\ w_{12}^{(k)} & w_{22}^{(k)} & \vdots & w_{D_{k-1}2}^{(k)} \\ \cdots & \cdots & \ddots & \vdots \\ w_{1D_k}^{(k)} & w_{2D_k}^{(k)} & \cdots & w_{D_{k-1}D_k}^{(k)} \end{bmatrix}$$

$$\mathbf{b}_k = \begin{bmatrix} b_1^{(k)} \\ b_2^{(k)} \\ \vdots \\ b_{D_{k+1}}^{(k)} \end{bmatrix}$$

- Arrange all inputs to the network in a vector $\mathbf{x}$
- Arrange the *inputs* to neurons of the kth layer as a vector $\mathbf{z}_k$
- Arrange the outputs of neurons in the kth layer as a vector $\mathbf{y}_k$
- Arrange the weights to any layer as a matrix $\mathbf{W}_k$
  - Similarly with biases

# Vector formulation



$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

$$\mathbf{z}_k = \begin{bmatrix} z_1^{(k)} \\ z_2^{(k)} \\ \vdots \\ z_{D_k}^{(k)} \end{bmatrix}$$
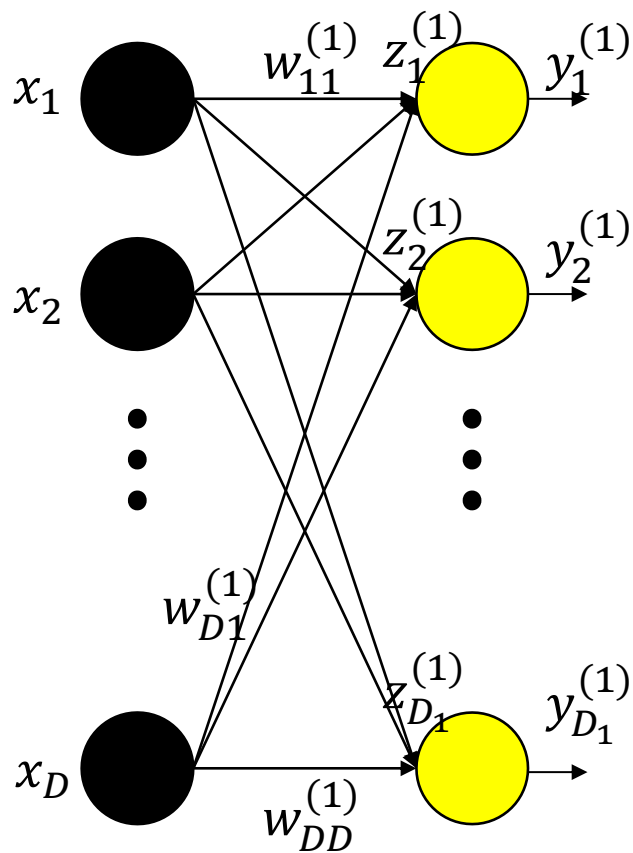
$$\mathbf{y}_k = \begin{bmatrix} y_1^{(k)} \\ y_2^{(k)} \\ \vdots \\ y_{D_k}^{(k)} \end{bmatrix}$$

$$\mathbf{W}_k = \begin{bmatrix} w_{11}^{(k)} & w_{21}^{(k)} & \vdots & w_{D_{k-1}1}^{(k)} \\ w_{12}^{(k)} & w_{22}^{(k)} & \vdots & w_{D_{k-1}2}^{(k)} \\ \ldots & \ldots & \ddots & \vdots \\ w_{1D_k}^{(k)} & w_{2D_k}^{(k)} & \ldots & w_{D_{k-1}D_k}^{(k)} \end{bmatrix}$$

$$\mathbf{b}_k = \begin{bmatrix} b_1^{(k)} \\ b_2^{(k)} \\ \vdots \\ b_{D_{k+1}}^{(k)} \end{bmatrix}$$

- The computation of a single layer is easily expressed in matrix notation as (setting $\mathbf{y_0} = \mathbf{x}$):

$$\mathbf{z}_k = \mathbf{W}_k \mathbf{y}_{k-1} + \mathbf{b}_k$$

$$\mathbf{y}_k = f_k(\mathbf{z}_k)$$

# The forward pass: Evaluating the network

x

$y_0 = x$

# The forward pass

**x**

**$W_1, b_1$**    **$z_1$**

$$z_1 = W_1 x + b_1$$

158

# The forward pass



$$\mathbf{y}_1 = \boldsymbol{f}_1(\mathbf{z}_1)$$

The Complete computation

$$\mathbf{y}_1 = f_1(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$

# The forward pass



$$\mathbf{z}_2 = \mathbf{W}_2\mathbf{y}_1 + \mathbf{b}_2$$

The Complete computation

$$\mathbf{y}_1 = f_1(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$

# The forward pass



$$\mathbf{y}_2 = \boldsymbol{f}_2(\mathbf{z}_2)$$

The Complete computation

$$\mathbf{y}_2 = f_2(\mathbf{W}_2 f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2)$$

# The forward pass



The Complete computation

$$\mathbf{y}_2 = f_2(\mathbf{W}_2 f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2)$$

$$\mathbf{z}_N = \mathbf{W}_N \mathbf{y}_{N-1} + \mathbf{b}_N$$

162

# The forward pass



The Complete computation

$$Y = f_N(\mathbf{W}_N f_{N-1}(\dots f_2(\mathbf{W}_2 f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2)\dots) + \mathbf{b}_N)$$

# Forward pass



## Forward pass:

Initialize $\qquad \boxed{\mathbf{y}_0 = \mathbf{x}}$

For k = 1 to N: $\boxed{\mathbf{z}_k = \mathbf{W}_k \mathbf{y}_{k-1} + \mathbf{b}_k} \quad \boxed{\mathbf{y}_k = \boldsymbol{f}_k(\mathbf{z}_k)}$

Output $\qquad \boxed{\text{Y} = \mathbf{y}_N}$

# The Forward Pass

- Set $\mathbf{y}_0 = \mathbf{x}$

- For layer k = 1 to N:
  - Recursion:
$$\mathbf{z}_k = \mathbf{W}_k \mathbf{y}_{k-1} + \mathbf{b}_k$$
$$\mathbf{y}_k = \boldsymbol{f}_k(\mathbf{z}_k)$$
- Output:
$$\mathbf{Y} = \mathbf{y}_N$$

# The backward pass



- The network is a nested function

$$Y = f_N(\mathbf{W}_N f_{N-1}(... f_2(\mathbf{W}_2 f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) ...) + \mathbf{b}_N)$$

- The error for any $\mathbf{x}$ is also a nested function

$$Div(Y, d) = Div(f_N(\mathbf{W}_N f_{N-1}(... f_2(\mathbf{W}_2 f_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) ...) + \mathbf{b}_N), d)$$

# Calculus recap 2: The Jacobian

- The derivative of a vector function w.r.t. vector input is called a *Jacobian*

- It is the matrix of partial derivatives given below

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix} = f\left( \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_D \end{bmatrix} \right)$$

Using vector notation

$$\mathbf{y} = f(\mathbf{z})$$

$$J_{\mathbf{y}}(\mathbf{z}) = \begin{bmatrix} \dfrac{\partial y_1}{\partial z_1} & \dfrac{\partial y_1}{\partial z_2} & \cdots & \dfrac{\partial y_1}{\partial z_D} \\[2mm] \dfrac{\partial y_2}{\partial z_1} & \dfrac{\partial y_2}{\partial z_2} & \cdots & \dfrac{\partial y_2}{\partial z_D} \\[2mm] \cdots & \cdots & \ddots & \cdots \\[2mm] \dfrac{\partial y_M}{\partial z_1} & \dfrac{\partial y_M}{\partial z_2} & \cdots & \dfrac{\partial y_M}{\partial z_D} \end{bmatrix}$$

Check: $\quad \Delta \mathbf{y} = J_y(\mathbf{z}) \Delta \mathbf{z}$

# Jacobians can describe the derivatives of neural activations w.r.t their input



$$J_{\boldsymbol{y}}(\mathbf{z}) = \begin{bmatrix} \dfrac{dy_1}{dz_1} & 0 & \cdots & 0 \\ 0 & \dfrac{dy_2}{dz_2} & \cdots & 0 \\ \cdots & \cdots & \ddots & \cdots \\ 0 & 0 & \cdots & \dfrac{dy_D}{dz_D} \end{bmatrix}$$

- **For Scalar activations**
  - Number of outputs is identical to the number of inputs
- Jacobian is a diagonal matrix
  - Diagonal entries are individual derivatives of outputs w.r.t inputs
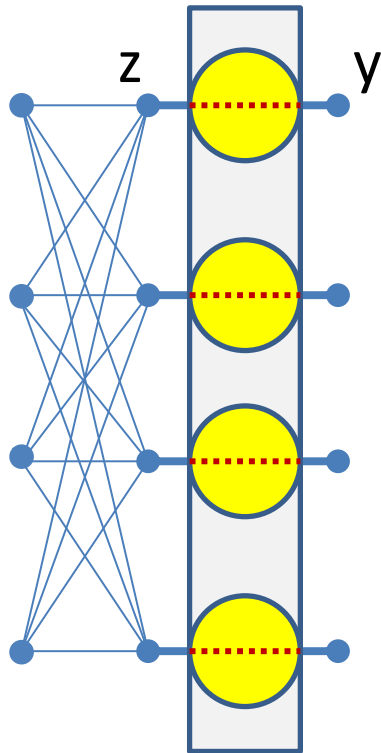  - Not showing the superscript "(k)" in equations for brevity

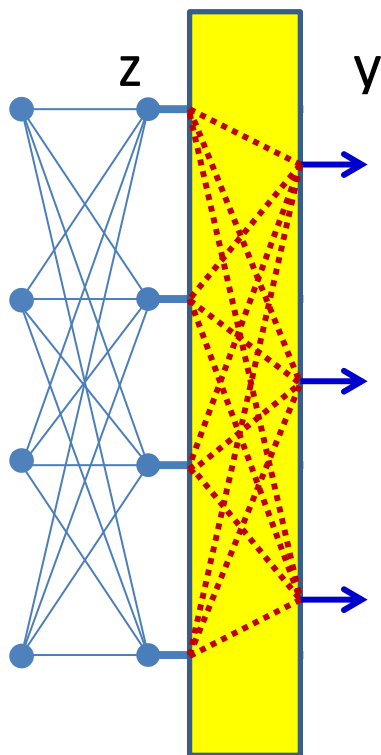# Jacobians can describe the derivatives of neural activations w.r.t their input



$$y_i = f(z_i)$$

$$J_y(\mathbf{z}) = \begin{bmatrix} f'(y_1) & 0 & \cdots & 0 \\ 0 & f'(y_2) & \cdots & 0 \\ \cdots & \cdots & \ddots & \cdots \\ 0 & 0 & \cdots & f'(y_M) \end{bmatrix}$$

- **For scalar activations (shorthand notation):**
  - Jacobian is a diagonal matrix
  - Diagonal entries are individual derivatives of outputs w.r.t inputs

# For *Vector* activations



$$J_y(\mathbf{z}) = \begin{bmatrix} \dfrac{\partial y_1}{\partial z_1} & \dfrac{\partial y_1}{\partial z_2} & \cdots & \dfrac{\partial y_1}{\partial z_D} \\[2ex] \dfrac{\partial y_2}{\partial z_1} & \dfrac{\partial y_2}{\partial z_2} & \cdots & \dfrac{\partial y_2}{\partial z_D} \\[2ex] \cdots & \cdots & \ddots & \cdots \\[2ex] \dfrac{\partial y_M}{\partial z_1} & \dfrac{\partial y_M}{\partial z_2} & \cdots & \dfrac{\partial y_M}{\partial z_D} \end{bmatrix}$$

- Jacobian is a full matrix
  - Entries are partial derivatives of individual outputs w.r.t individual inputs

# Special case: Affine functions
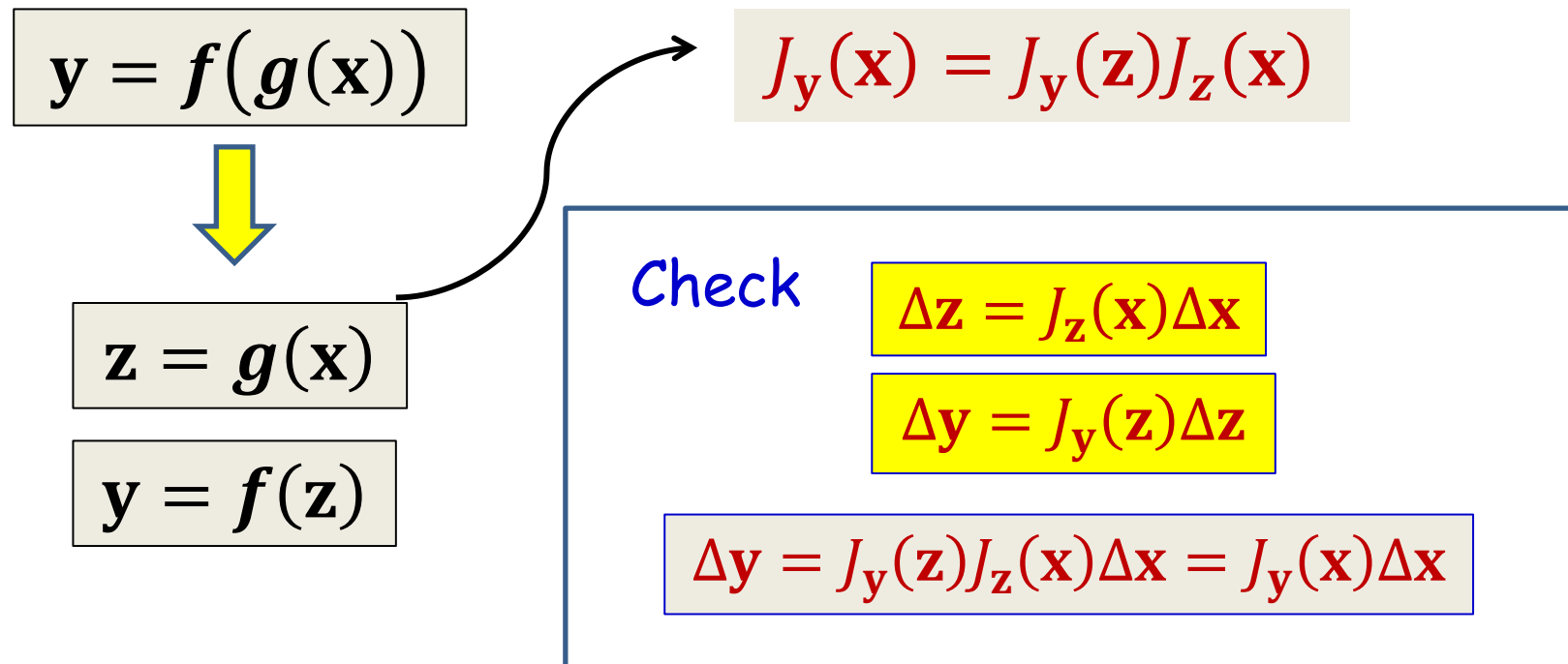
$$\mathbf{z} = \mathbf{W}\mathbf{y} + \mathbf{b}$$



$$J_{\mathbf{z}}(\mathbf{y}) = \mathbf{W}$$

- Matrix **W** and bias **b** operating on vector **y** to produce vector **z**
- The Jacobian of **z** w.r.t **y** is simply the matrix **W**

# Vector derivatives: Chain rule

- We can define a chain rule for Jacobians
- **For vector functions of vector inputs:**

$$y = f(g(\mathbf{x}))$$

$$J_{\mathbf{y}}(\mathbf{x}) = J_{\mathbf{y}}(\mathbf{z})J_{\mathbf{z}}(\mathbf{x})$$

$$\mathbf{z} = g(\mathbf{x})$$

$$y = f(\mathbf{z})$$

Check

$$\Delta\mathbf{z} = J_{\mathbf{z}}(\mathbf{x})\Delta\mathbf{x}$$

$$\Delta\mathbf{y} = J_{\mathbf{y}}(\mathbf{z})\Delta\mathbf{z}$$

$$\Delta\mathbf{y} = J_{\mathbf{y}}(\mathbf{z})J_{\mathbf{z}}(\mathbf{x})\Delta\mathbf{x} = J_{\mathbf{y}}(\mathbf{x})\Delta\mathbf{x}$$

Note the order: The derivative of the outer function comes first

172

# Vector derivatives: Chain rule

- *The chain rule can combine Jacobians and Gradients*
- **For *scalar* functions of vector inputs ($g()$ is vector):**

$$D = f\big(g(\mathbf{x})\big)$$

$$\nabla_{\mathbf{x}} D = \nabla_{\mathbf{z}}(D) J_{\mathbf{z}}(\mathbf{x})$$
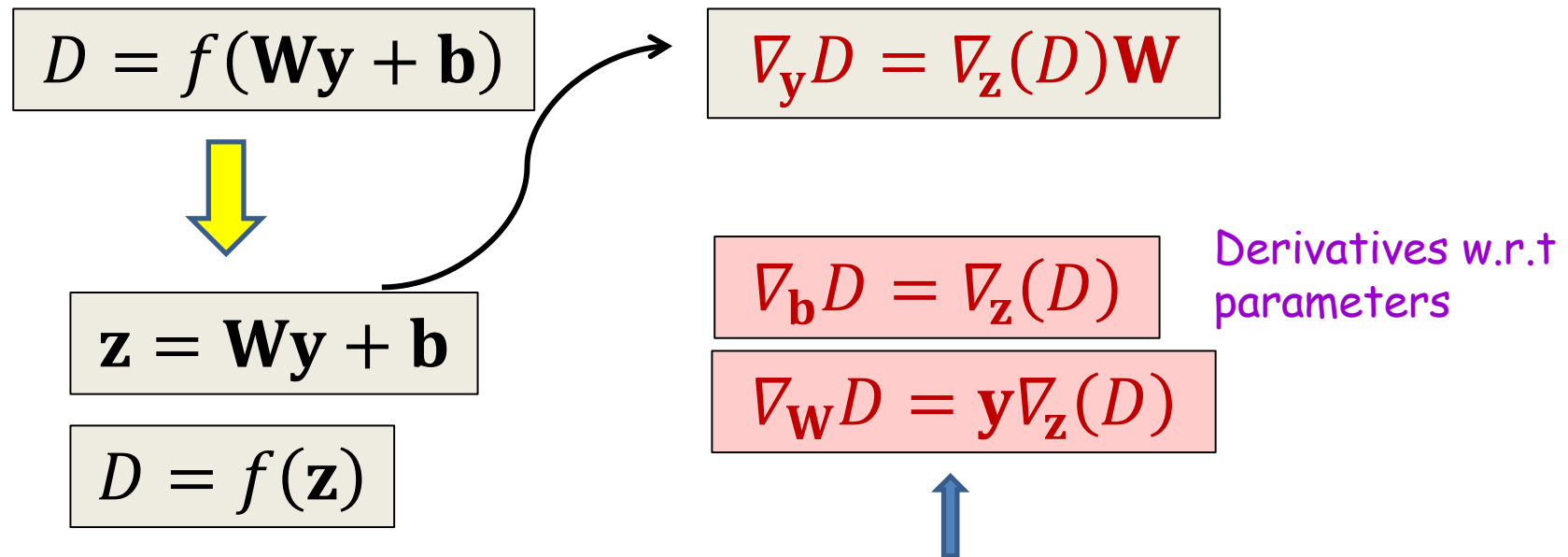
$$\mathbf{z} = g(\mathbf{x})$$

$$D = f(\mathbf{z})$$

Check

$$\Delta \mathbf{z} = J_{\mathbf{z}}(\mathbf{x}) \Delta \mathbf{x}$$

$$\Delta D = \nabla_{\mathbf{z}}(D) \Delta \mathbf{z}$$

$$\Delta D = \nabla_{\mathbf{z}}(D) J_{\mathbf{z}}(\mathbf{x}) \Delta \mathbf{x} = \nabla_{\mathbf{x}} D \Delta \mathbf{x}$$

Note the order: The derivative of the outer function comes first

# Special Case

- Scalar functions of Affine functions

$$D = f(\mathbf{W}\mathbf{y} + \mathbf{b})$$

$$\nabla_{\mathbf{y}} D = \nabla_{\mathbf{z}}(D)\mathbf{W}$$

$$\mathbf{z} = \mathbf{W}\mathbf{y} + \mathbf{b}$$

$$D = f(\mathbf{z})$$

$$\nabla_{\mathbf{b}} D = \nabla_{\mathbf{z}}(D)$$

Derivatives w.r.t parameters

$$\nabla_{\mathbf{W}} D = \mathbf{y}\nabla_{\mathbf{z}}(D)$$

Note reversal of order. This is in fact a simplification of a product of tensor terms that occur in the *right* order
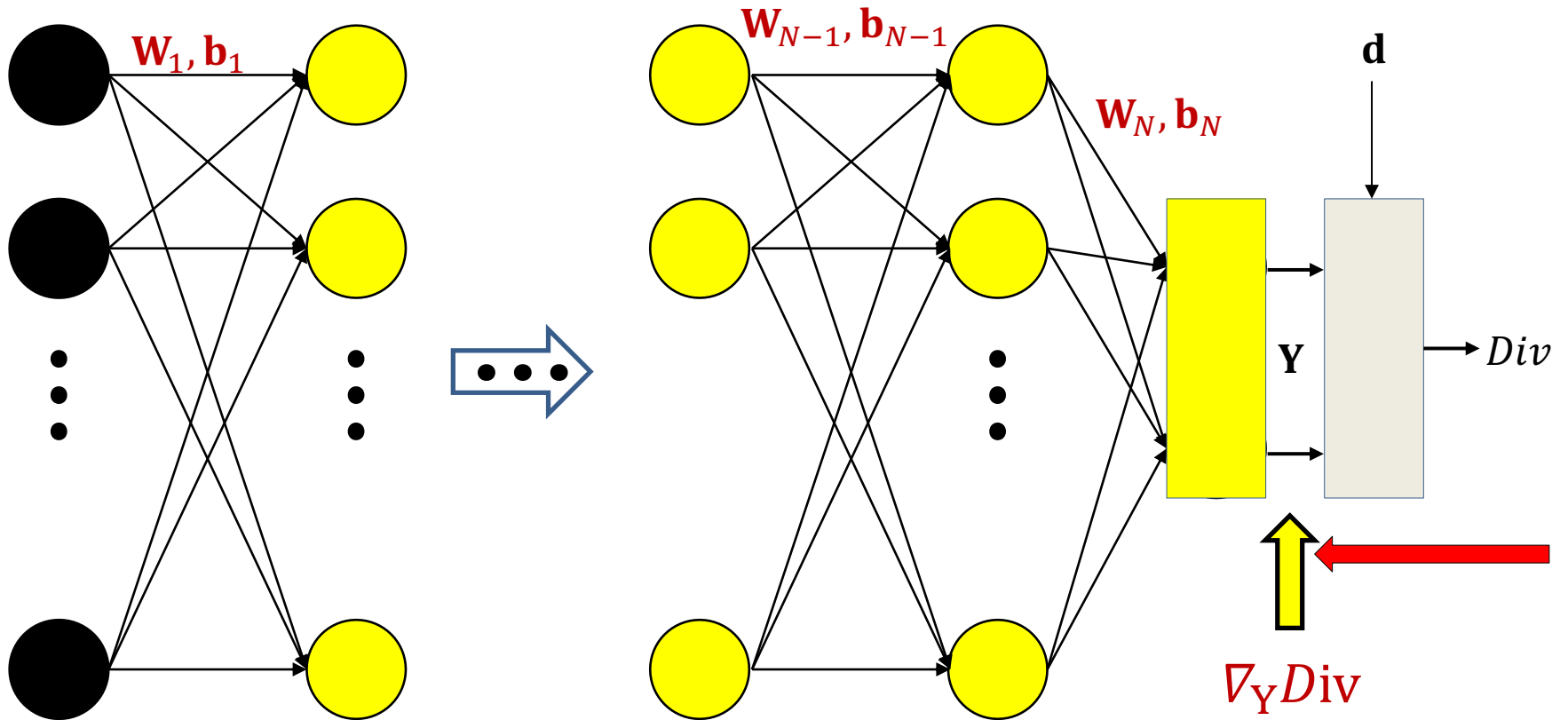
174

# The backward pass



In the following slides we will also be using the notation $\nabla_{\mathbf{z}}\mathbf{Y}$ to represent the Jacobian $J_{\mathbf{Y}}(\mathbf{z})$ to explicitly illustrate the chain rule
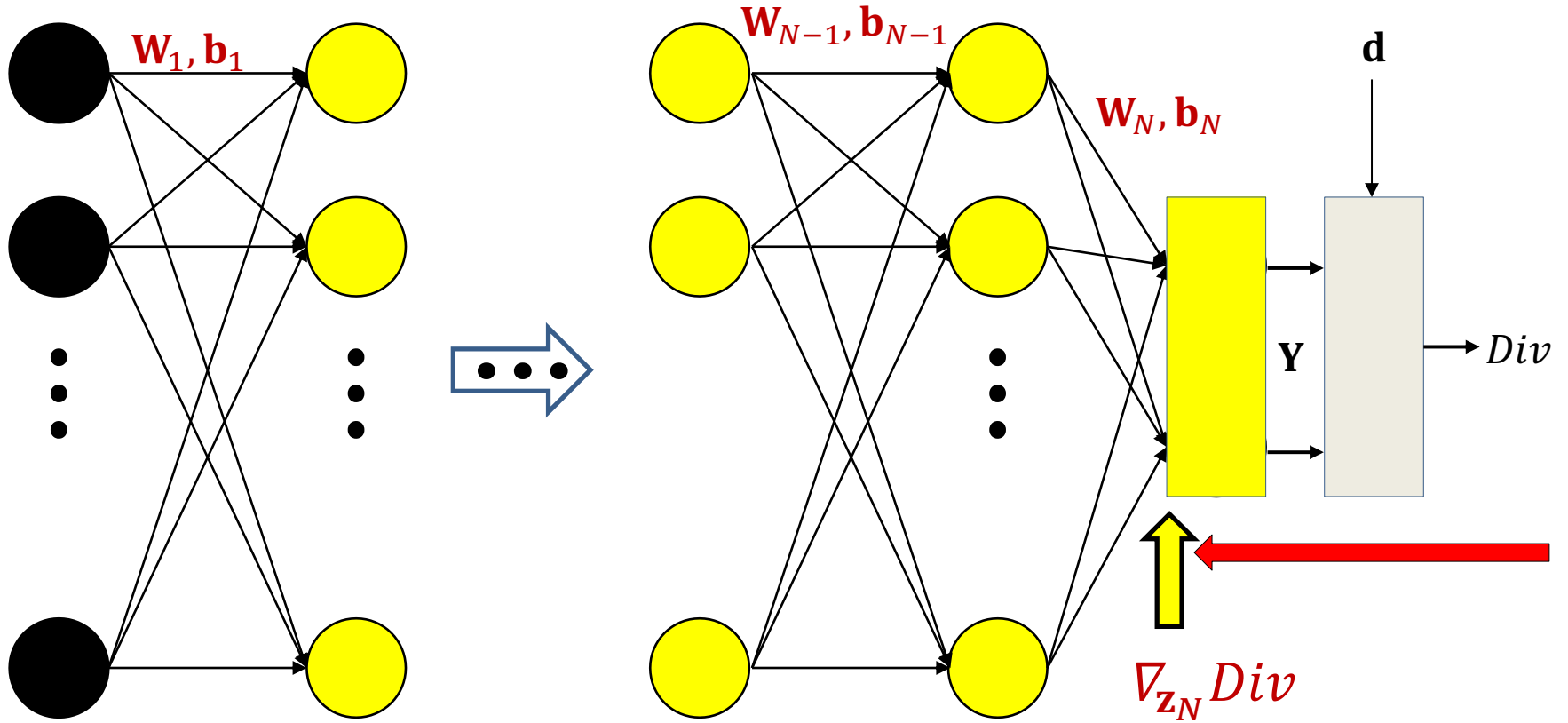
In general $\nabla_{\mathbf{a}}\mathbf{b}$ represents a derivative of $\mathbf{b}$ w.r.t. $\mathbf{a}$ and could be a gradient (for scalar $\mathbf{b}$) Or a Jacobian (for vector $\mathbf{b}$)

# The backward pass



First compute the gradient of the divergence w.r.t. Y.
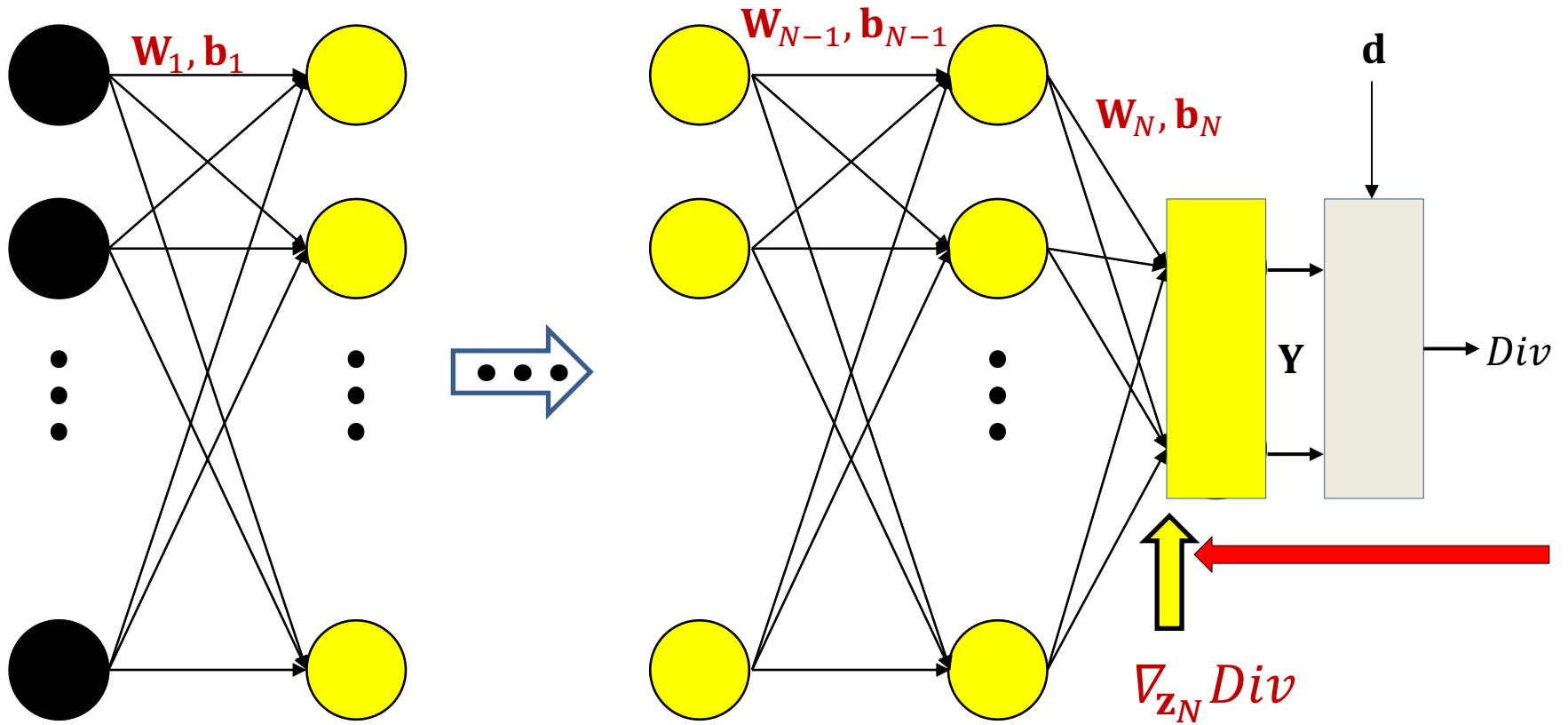The actual gradient depends on the divergence function.

# The backward pass



$$\nabla_{\mathbf{z}_N} Div = \nabla_{Y} Div . \nabla_{\mathbf{z}_N} \mathbf{Y}$$

Already computed      New term

# The backward pass



$$\nabla_{\mathbf{z}_N} Div = \nabla_{\mathbf{Y}} Div \, J_{\mathbf{Y}}(\mathbf{z}_N)$$

Already computed          New term

# The backward pass



$$\nabla_{\mathbf{y}_{N-1}} Div = \nabla_{\mathbf{z}_N} Div . \nabla_{\mathbf{y}_{N-1}} \mathbf{z}_N$$

Already computed    New term

# The backward pass



$$\nabla_{\mathbf{y}_{N-1}} Div = \nabla_{\mathbf{z}_N} Div \; \mathbf{W}_N$$

Already computed    New term

# The backward pass



$$\nabla_{\mathbf{y}_{N-1}} Div = \nabla_{\mathbf{z}_N} Div \ \mathbf{W}_N$$

$$\nabla_{\mathbf{y}_{N-1}} Div$$

$$\nabla_{\mathbf{W}_N} Div = \mathbf{y}_{N-1} \nabla_{\mathbf{z}_N} Div$$

$$\nabla_{\mathbf{b}_N} Div = \nabla_{\mathbf{z}_N} Div$$

# The backward pass



$$\nabla_{\mathbf{z}_{N-1}} Div = \underline{\nabla_{\mathbf{y}_{N-1}} Div}.\ \underline{\nabla_{\mathbf{z}_{N-1}} \mathbf{y}_{N-1}}$$

Already computed          New term

# The backward pass



$$\nabla_{\mathbf{z}_{N-1}} Div = \nabla_{\mathbf{y}_{N-1}} Div \, J_{\mathbf{y}_{N-1}}(\mathbf{z}_{N-1})$$

The Jacobian will be a diagonal matrix for scalar activations

$\nabla_{\mathbf{z}_{N-1}} Div$

# The backward pass



$$\nabla_{\mathbf{y}_{N-2}} Div = \nabla_{\mathbf{z}_{N-1}} Div . \nabla_{\mathbf{y}_{N-2}} \mathbf{z}_{N-1}$$

# The backward pass



$$\nabla_{\mathbf{y}_{N-2}} Div = \nabla_{\mathbf{z}_{N-1}} Div \ \mathbf{W}_{N-1}$$

# The backward pass



$$\nabla_{\mathbf{y}_{N-2}} Div$$

$$\nabla_{\mathbf{y}_{N-2}} Div = \nabla_{\mathbf{z}_{N-1}} Div \ \mathbf{W}_{N-1}$$

$$\nabla_{\mathbf{W}_{N-1}} Div = \mathbf{y}_{N-2} \nabla_{\mathbf{z}_{N-1}} Div$$

$$\nabla_{\mathbf{b}_{N-1}} Div = \nabla_{\mathbf{z}_{N-1}} Div$$

# The backward pass



$$\nabla_{\mathbf{z}_1} Div = \nabla_{\mathbf{y}_1} Div \, J_{\mathbf{y}_1}(\mathbf{z}_1)$$

# The backward pass



$$\nabla_{\mathbf{W}_1} Div = \mathbf{x}\nabla_{\mathbf{z}_1} Div$$

$$\nabla_{\mathbf{b}_1} Div = \nabla_{\mathbf{z}_1} Div$$

In some problems we will also want to compute the derivative w.r.t. the input

# The Backward Pass

- Set $\mathbf{y}_N = Y$, $\mathbf{y}_0 = \mathbf{x}$
- Initialize: Compute $\nabla_{\mathbf{y}_N} Div = \nabla_Y Div$

- For layer k = N downto 1:
  - Compute $J_{\mathbf{y}_k}(\mathbf{z}_k)$
    - Will require intermediate values computed in the forward pass
  - Recursion:
  $$\nabla_{\mathbf{z}_k} Div = \nabla_{\mathbf{y}_k} Div \, J_{\mathbf{y}_k}(\mathbf{z}_k)$$
  $$\nabla_{\mathbf{y}_{k-1}} Div = \nabla_{\mathbf{z}_k} Div \, \mathbf{W}_k$$
  - Gradient computation:
  $$\nabla_{\mathbf{W}_k} Div = \mathbf{y}_{k-1} \nabla_{\mathbf{z}_k} Div$$
  $$\nabla_{\mathbf{b}_k} Div = \nabla_{\mathbf{z}_k} Div$$

# The Backward Pass

- Set $\mathbf{y}_N = Y$, $\mathbf{y}_0 = \mathbf{x}$
- Initialize: Compute $\nabla_{\mathbf{y}_N} Div = \nabla_Y Div$

- For layer k = N downto 1:
  - Compute $J_{\mathbf{y}_k}(\mathbf{z}_k)$
    - Will require intermediate values computed in the forward pass
  - Recursion:    <mark>Note analogy to forward pass</mark>

$$\nabla_{\mathbf{z}_k} Div = \nabla_{\mathbf{y}_k} Div \, J_{\mathbf{y}_k}(\mathbf{z}_k)$$
$$\nabla_{\mathbf{y}_{k-1}} Div = \nabla_{\mathbf{z}_k} Div \, \mathbf{W}_k$$

  - Gradient computation:

$$\nabla_{\mathbf{W}_k} Div = \mathbf{y}_{k-1} \nabla_{\mathbf{z}_k} Div$$
$$\nabla_{\mathbf{b}_k} Div = \nabla_{\mathbf{z}_k} Div$$

# For comparison: The Forward Pass

- Set $\mathbf{y}_0 = \mathbf{x}$

- For layer k = 1 to N:
  - Recursion:
  $$\mathbf{z}_k = \mathbf{W}_k \mathbf{y}_{k-1} + \mathbf{b}_k$$
  $$\mathbf{y}_k = \boldsymbol{f}_k(\mathbf{z}_k)$$
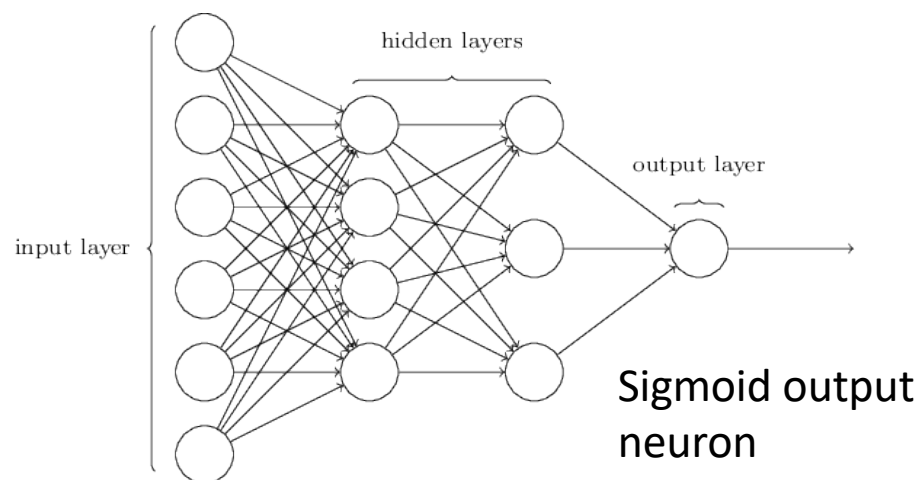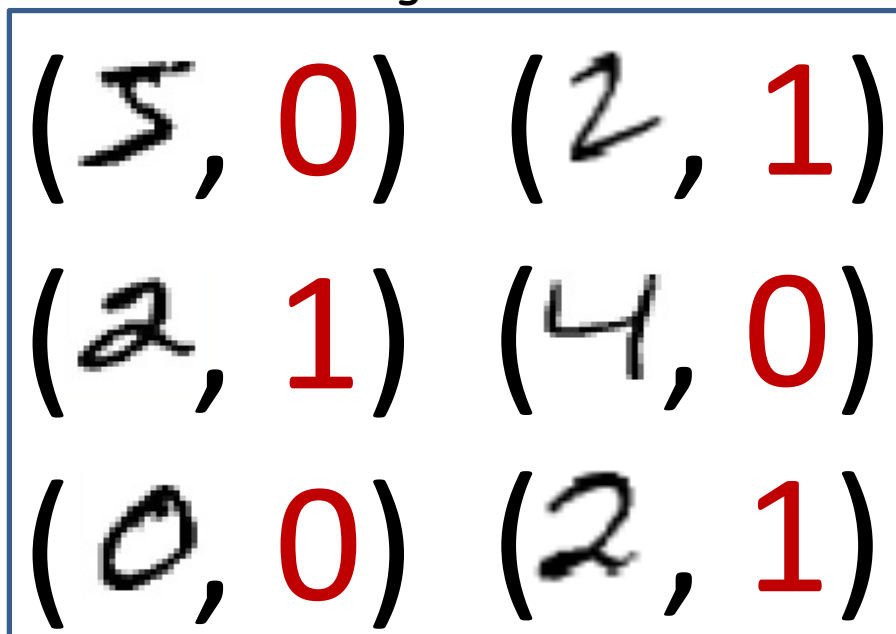- Output:
  $$\mathbf{Y} = \mathbf{y}_N$$

# Neural network training algorithm

- Initialize all weights and biases $(\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots, \mathbf{W}_N, \mathbf{b}_N)$

- Do:

  - $Err = 0$

  - For all $k$, initialize $\nabla_{\mathbf{W}_k} Err = 0$, $\nabla_{\mathbf{b}_k} Err = 0$

  - For all $t = 1:T$

    - Forward pass : Compute
      - Output $\boldsymbol{Y(X_t)}$
      - Divergence $\boldsymbol{Div(Y_t, d_t)}$
      - $Err \mathrel{+}= \boldsymbol{Div(Y_t, d_t)}$

    - Backward pass: For all $k$ compute:
      - $\nabla_{\mathbf{y}_k} Div = \nabla_{\mathbf{z}_{k+1}} Div \, \mathbf{W}_k$
      - $\nabla_{\mathbf{z}_k} Div = \nabla_{\mathbf{y}_k} Div \, J_{\mathbf{y}_k}(\mathbf{z}_k)$
      - $\nabla_{\mathbf{W}_k} \boldsymbol{Div(Y_t, d_t)}$; $\nabla_{\mathbf{b}_k} \boldsymbol{Div(Y_t, d_t)}$
      - $\nabla_{\mathbf{W}_k} Err \mathrel{+}= \nabla_{\mathbf{W}_k} \boldsymbol{Div(Y_t, d_t)}$; $\nabla_{\mathbf{b}_k} Err \mathrel{+}= \nabla_{\mathbf{b}_k} \boldsymbol{Div(Y_t, d_t)}$

  - For all $k$, update:

    $$\mathbf{W}_k = \mathbf{W}_k - \frac{\eta}{T}\left(\nabla_{\mathbf{W}_k} Err\right)^T; \qquad \mathbf{b}_k = \mathbf{b}_k - \frac{\eta}{T}\left(\nabla_{\mathbf{W}_k} Err\right)^T$$

- Until $Err$ has converged

# Setting up for digit recognition

### Training data

$$(5, 0)\ (2, 1)$$
$$(2, 1)\ (4, 0)$$
$$(0, 0)\ (2, 1)$$



input layer

hidden layers
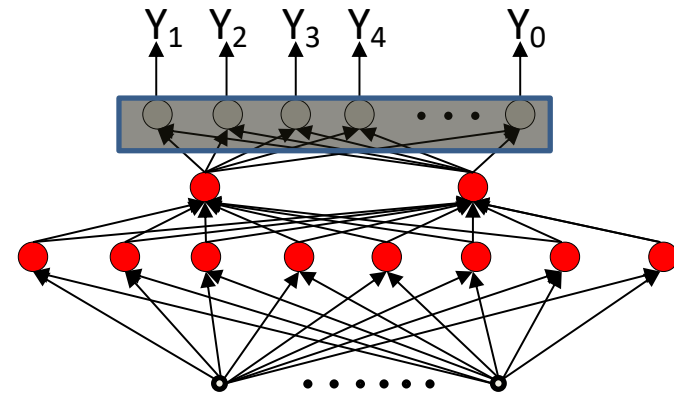
output layer

Sigmoid output neuron

- Simple Problem: Recognizing "2" or "not 2"
- Single output with sigmoid activation
  - $Y \in (0,1)$
  - $d$ is either 0 or 1
- Use KL divergence
- Backpropagation to learn network parameters

# Recognizing the digit

Training data



(5, 0)   (2, 1)
(2, 1)   (4, 0)
(0, 0)   (2, 1)

- More complex problem: Recognizing digit
- Network with 10 (or 11) outputs
  - First ten outputs correspond to the ten digits
    - Optional 11th is for none of the above
- Softmax output layer:
  - Ideal output: One of the outputs goes to 1, the others go to 0
- Backpropagation with KL divergence to learn network

194

# Issues

- Convergence: How well does it learn
  - And how can we improve it
- How well will it generalize (outside training data)
- What does the output really mean?
- *Etc..*

# Next up

- Convergence and generalization