

Backpropagation algorithm:

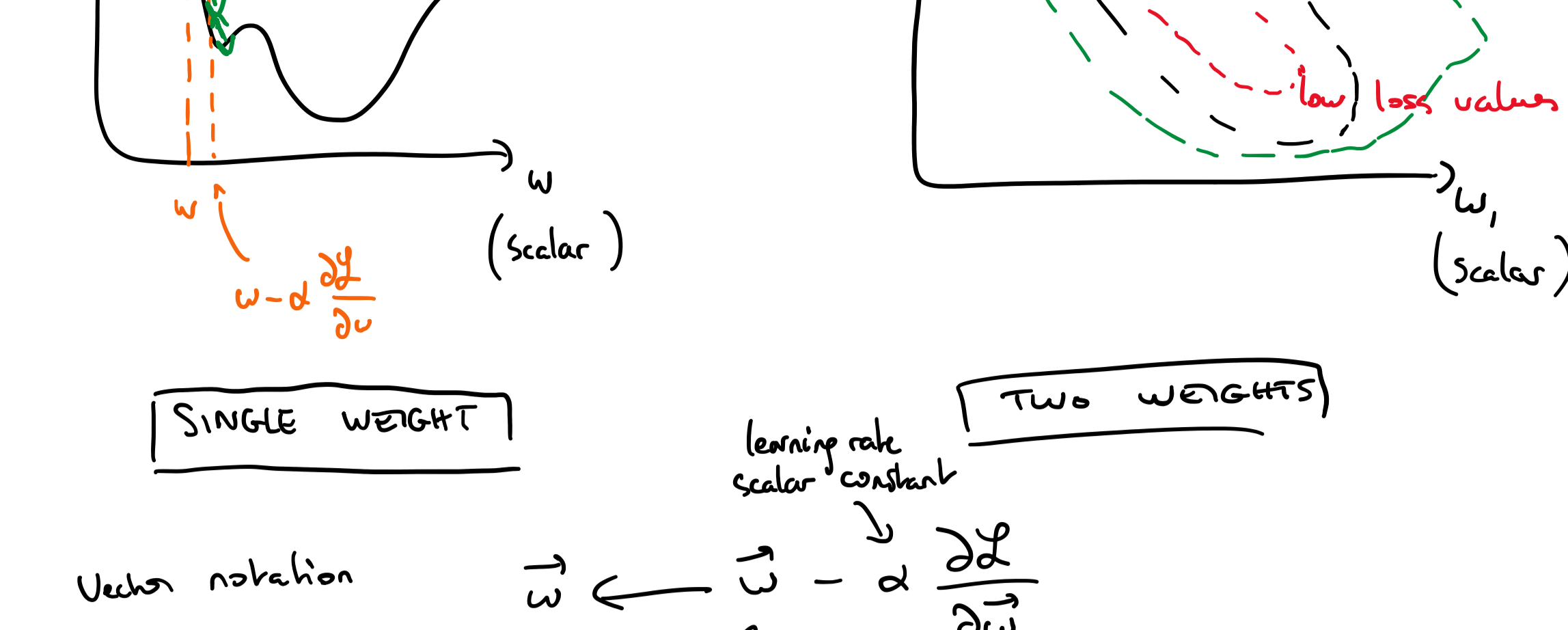
Find the values of a deep neural network's (DNN) weights given: - a training set
- an objective function / loss / cost.

Gradient descent updates the parameters in the direction of steepest descent (opposite to the gradient)

$$w \leftarrow w - \alpha \frac{\partial \mathcal{L}}{\partial w}$$

weight learning rate scalar constant partial derivative of the loss w.r.t to weight w

opposite direction to the gradient



SINGLE WEIGHT **TWO WEIGHTS**

learning rate scalar constant

Vector notation: $\vec{w} \leftarrow \vec{w} - \alpha \frac{\partial \mathcal{L}}{\partial \vec{w}}$

d-dimensional vector d-dimensional vector

assuming the model has J weights.

Chain rule

If we have $f(x)$ and $x(t)$ which are both univariate functions,

$$\frac{d}{dt}(f(x(t))) = \frac{df}{dx} \frac{dx}{dt}$$

Noticing example:

- for simplifying (of exposition):
- single training example
 - single input feature (input is a scalar)
 - activation function: logistic function
 - loss: least squares

Model: $z = wx + b$

scalar scalar scalar

input is scalar

$$y = \sigma(z) = \frac{1}{1+e^{-z}}$$

$$\mathcal{L} = \frac{1}{2} (y-t)^2$$

Add a term to our loss called a regularizer (a way to get a "simpler" explanation for the data)

Regularizer: $\frac{\lambda}{2} w^2$ (encourage weight to be close to 0)

hyperparameter (the larger λ is, the more we give importance to the regularizer during optimization)

$$\mathcal{L} = \frac{1}{2} (y-t)^2 + \frac{\lambda}{2} w^2$$

With $R = \frac{1}{2} w^2$, we write:

$$\mathcal{L}_{reg} = \frac{1}{2} (y-t)^2 + \lambda R$$

Goal: to train this model, we need to compute:

$$\frac{\partial \mathcal{L}_{reg}}{\partial w} \quad \frac{\partial \mathcal{L}_{reg}}{\partial b}$$

Naive approach based on calculus:

$$\frac{\partial \mathcal{L}_{reg}}{\partial w} = \frac{\partial}{\partial w} \left(\frac{1}{2} (\sigma(wx+b) - t)^2 + \frac{\lambda}{2} w^2 \right)$$

by definition

$$= \frac{1}{2} \frac{\partial}{\partial w} (\sigma(wx+b) - t)^2 + \lambda w$$

apply the chain rule

$$= \frac{1}{2} \cancel{\sigma'} (\sigma(wx+b) - t) \frac{\partial}{\partial w} (\sigma(wx+b) - t) + \lambda w$$

constant does not depend on w

apply the chain rule

$$= (\sigma(wx+b) - t) \sigma'(wx+b) \frac{\partial}{\partial w} (wx+b) + \lambda w$$

constant does not depend on w

$$\frac{\partial \mathcal{L}_{reg}}{\partial w} = (\sigma(wx+b) - t) \sigma'(wx+b) x + \lambda w$$

$$\frac{\partial \mathcal{L}_{reg}}{\partial b} = \frac{\partial}{\partial b} \left(\frac{1}{2} (\sigma(wx+b) - t)^2 + \frac{\lambda}{2} w^2 \right)$$

independent of b \Rightarrow derivative $\frac{\partial}{\partial b} = 0$

$$= \frac{1}{2} \cancel{\sigma'} (\sigma(wx+b) - t) \frac{\partial}{\partial b} (\sigma(wx+b) - t) + 0$$

independent of b

$$\frac{\partial \mathcal{L}_{reg}}{\partial b} = (\sigma(wx+b) - t) \sigma'(wx+b)$$

$\frac{\partial wx+b}{\partial b} = 1$ & applied the chain rule.

Several drawbacks

1. These calculations are cumbersome.
2. Calculations involved a lot of redundant work.
3. Final expressions of $\frac{\partial \mathcal{L}_{reg}}{\partial w}$ and $\frac{\partial \mathcal{L}_{reg}}{\partial b}$ involved repeated terms.

Multivariable chain rule

$$\frac{df(x(t), y(t))}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

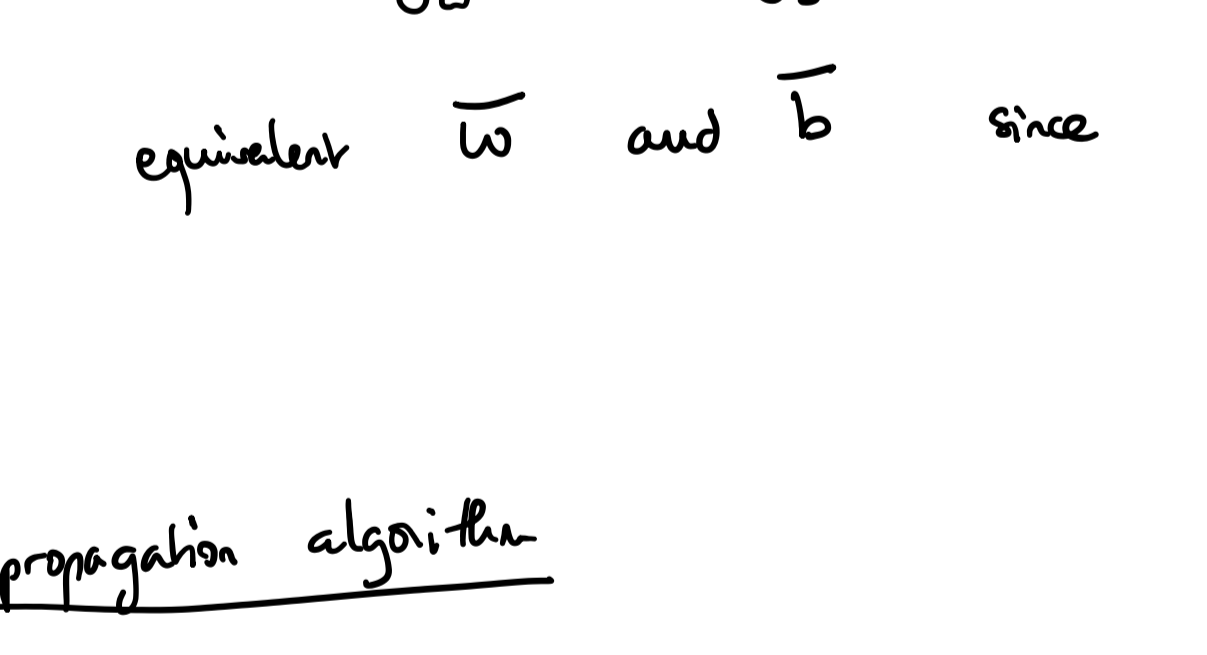
f depends on x, which depends on t f depends on y, which depends on t

Notation: $\bar{r} = \frac{\partial \mathcal{L}}{\partial r}$ any loss we would like to minimize (can be a regularized loss)

We can rewrite the chain rule: $\bar{t} = \bar{x} \frac{dx}{dt} + \bar{y} \frac{dy}{dt}$ (if $f = \mathcal{L}$)

Backpropagation:

Computation graph



$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2} (y-t)^2$$

$$R = \frac{1}{2} w^2$$

$$\mathcal{L}_{reg} = \mathcal{L} + \lambda R$$

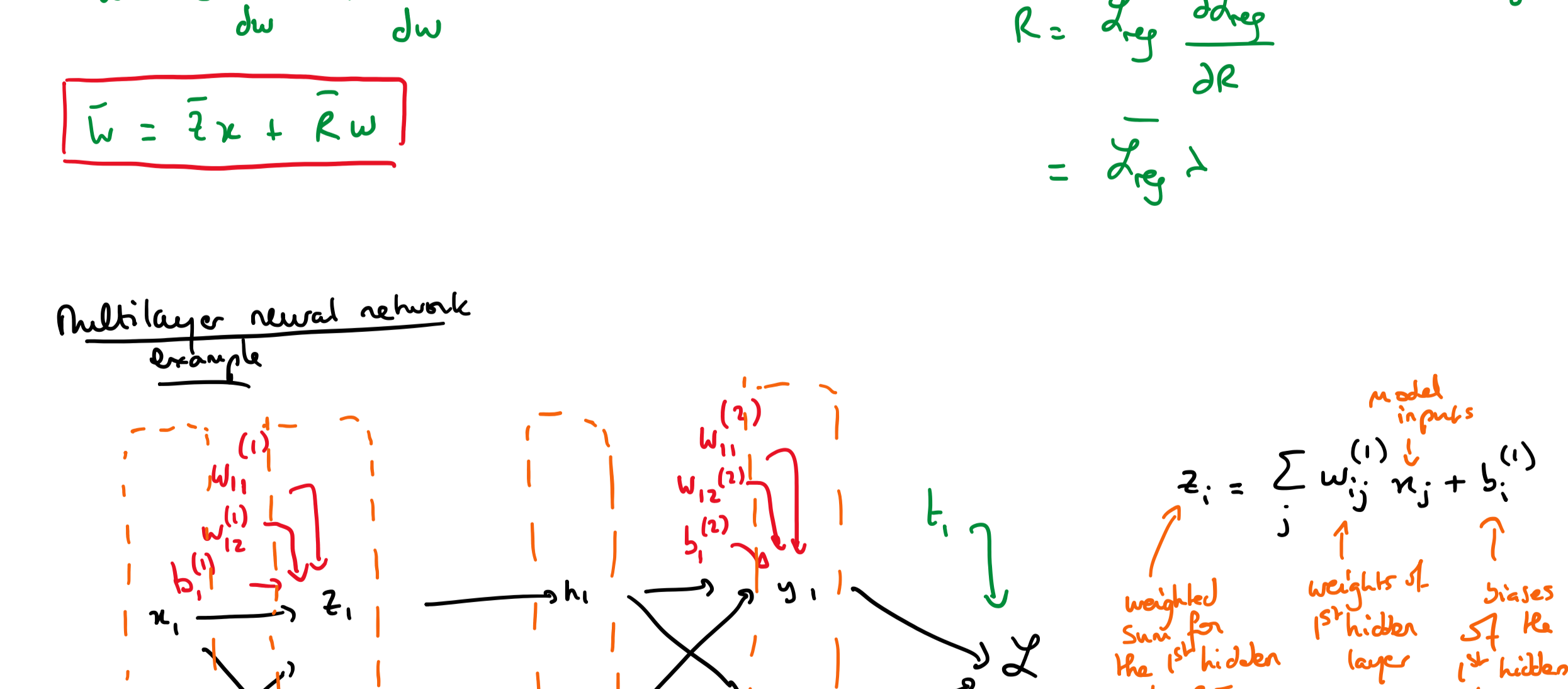
Goal: Compute $\frac{\partial \mathcal{L}_{reg}}{\partial w}$ and $\frac{\partial \mathcal{L}_{reg}}{\partial b}$

equivalent \bar{w} and \bar{b} since $\bar{w} = \frac{\partial \mathcal{L}_{reg}}{\partial w}$ and $\bar{b} = \frac{\partial \mathcal{L}_{reg}}{\partial b}$

Backpropagation algorithm

- Forward pass:** For $i \in 1..N$ Compute v_i as a function of $P_0(v_i)$ (parents of v_i)
- Backward pass:** For i in $N-1, \dots, 1$ $\bar{v}_i = \sum_{j \in Ch(v_i)} \bar{v}_j \frac{\partial v_j}{\partial v_i}$ (application of the chain rule)
- $\bar{v}_N = 1$ $\bar{v}_w = \frac{\partial \mathcal{L}_{reg}}{\partial w} = 1$
- v_i are the nodes in the computation graph. v_i is the quantity we'd like to get derivatives of, let's our loss (here, regularized loss)

Here v_1, v_2, \dots, v_N are the nodes of the computational graph in topological order (parents come before children).



Multi-layer neural network example

