

Lecture 2: Backpropagation

Peter Bloem
Deep Learning 2020

dlvu.github.io



Today's lecture will be entirely devoted to the backpropagation algorithm. The heart of all deep learning.

THE PLAN

- part 1:** review
- part 2:** scalar backpropagation
- part 3:** tensor backpropagation
- part 4:** automatic differentiation

2



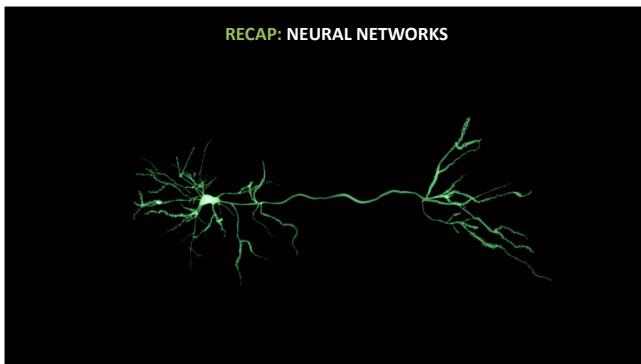
In **the first part**, we will review the basics of neural networks, and describe backpropagation in a scalar setting. That is, we will treat each individual element as a single number, and simply loop over all these numbers to do backpropagation over the whole network. This simplifies the derivation, but it is ultimately a slow algorithm with a complicated notation.

In **the second part**, we translate neural networks to operations on vectors, matrices and tensors (the higher-dimensional analogue of a matrix). This allows us to simplify our notation, and more importantly, massively speed up the computation of neural networks. Backpropagation on tensors is a little more difficult to do than backpropagation on scalars.

In **the third part**, we will make the final leap from manually worked out and implemented backpropagation system to full-fledged automatic differentiation: we will show you how to build a system that take care of the gradient computation entirely by itself. This is the technology behind software like pytorch and tensorflow.

PART ONE: REVIEW

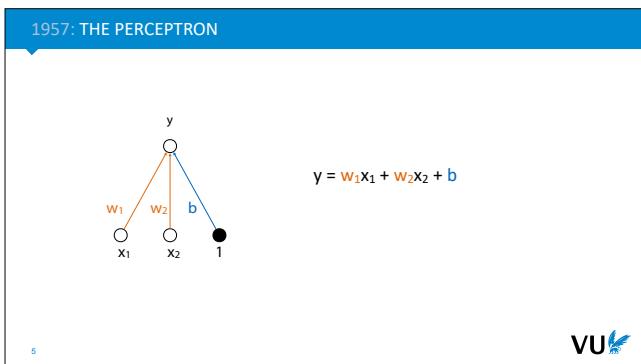




We'll start with a quick recap of the basic principles behind neural networks. The name neural network is a bit of a historical artifact.

In the very early days of AI (the late 1950s), researchers decided to take a simple approach to AI. They started with a single brain cell: a neuron. A neuron receives multiple different signals from other cells through connections called **dendrites**. It processes these in a relatively simple way, deriving a single new signal, which it sends out through its single **axon**. The axon branches out so that the single signal can reach other cells.

image source: <http://www.sciencealert.com/scientists-build-an-artificial-neuron-that-fully-mimics-a-human-brain-cell>

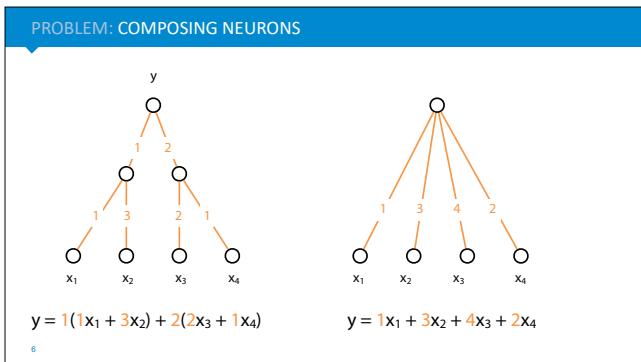


This principle needed to be radically simplified to work with computers of that age, but doing so yielded one of the first successful machine learning systems: the perceptron (also seen in the [video](#) in the first lecture).

The perceptron had a number of inputs (the features in modern parlance), each of which was multiplied by a **weight**. These result was summed, together with a **bias** parameter, and the sign of this result was used for classification or regression.

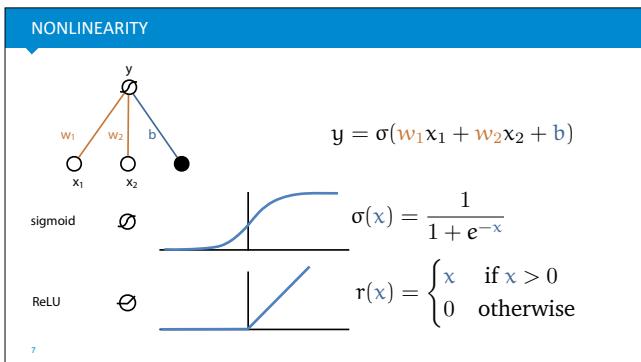
Note that the **intercept** can be represented as just another input that we just fix to always be 1. This is called a **bias node**.

Of course, there is nothing new here: this is just a basic linear regression or linear classification model. The real power of the brain is in chaining multiple neurons together in a network.



This is where the perceptron turns out to be too simple an abstraction. Because composing perceptrons (making the output of one perceptron the input of another) doesn't make it more powerful. All you end out with is something that is equivalent to another linear model. We're not creating models that can learning non-linear functions.

We've removed the bias node here for clarity, but that doesn't affect our conclusions: any composition of affine functions is itself an affine function.



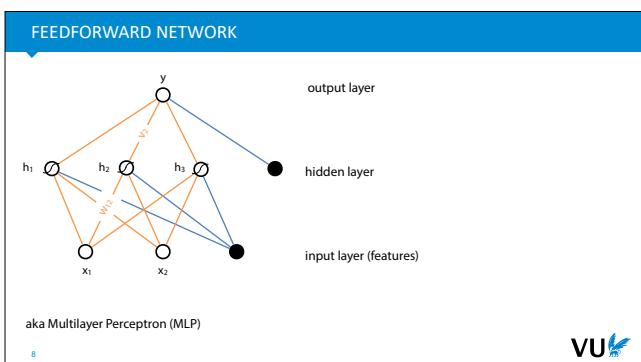
The simplest solution is to apply a nonlinear function to each neuron, called the **activation function**. This is a scalar function we apply to the output of a perceptron after all the weighted inputs have been combined.

One popular option (especially in the early days) is the **logistic sigmoid**, which we've seen already. Applying a sigmoid means that the sum of the inputs can range from negative infinity to positive infinity, but the output is always in the interval $[0, 1]$.

Another, more recent nonlinearity is the linear rectifier, or ReLU nonlinearity. This function just sets every negative input to zero, and keeps everything else the same.

Not using an activation function is also called using a linear activation.

We've seen one interpretation of the sigmoid function already: it turns every perceptron into a logistic classification unit. For now, we won't look much more into the meaning of the nonlinearities: we know we need them to make a stack of neuron more expressive than a single neuron, and by trial and error, we've found a few that work well for us.



Using these nonlinearities, we can arrange single neurons into **neural networks**. Any arrangement makes a neural network, but for ease of training, this arrangement was the most popular for a long time. It's called the **feedforward network** or **multilayer perceptron**. We arrange a layer of hidden units in the middle, each of which acts as a perceptron with a nonlinearity, connecting to all input nodes. Then we have one or more output nodes, connecting to all hidden layers. Crucially:

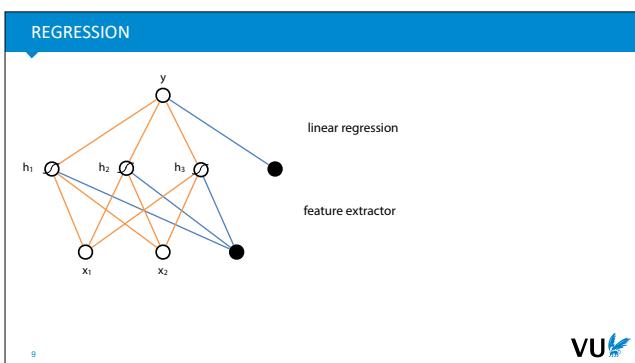
- There are no cycles, the network feeds forward from input to output.
- Nodes in the same layer are not connected to each other, or to any other layer that the

previous one.

- Each layer is fully connected to the previous layer, every node in one layer connects to every node in the layer before it.

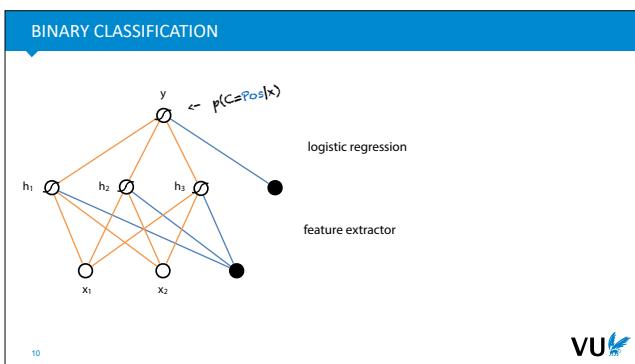
In the 80s and 90s they usually had just one hidden layer, because we hadn't figured out how to train deeper networks.

Every orange and blue line in this picture represents one parameter of the model.

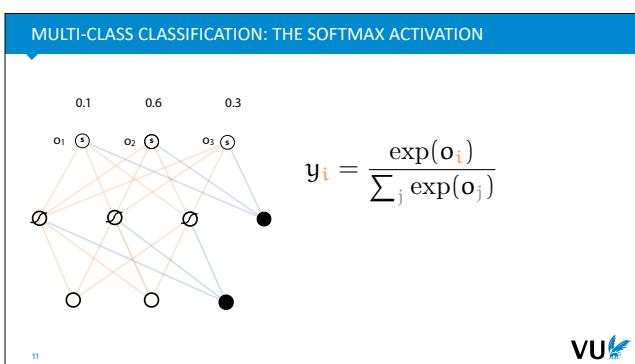


If we want to train a regression model, we put non-linearities on the hidden nodes, and no activation on the output node. That way, the output can range from negative to positive infinity.

We can think of the first layer as learning some nonlinear transformation of the features, and the second layer as performing linear regression on the features.



If we have a classification problem with two classes (a binary classification problem), we can place a sigmoid activation on the output layer, so that the output is between 0 and 1. We can then interpret this as the probability that the input has the **positive** class.

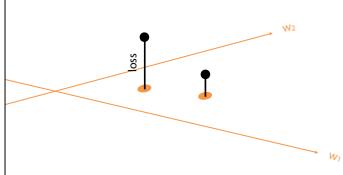


For multi-class classification, we can use the **softmax activation**. We create a single output node per class, and ensure that they sum to one. We can then interpret the output of the network as **class probabilities**.

This activation is a little unusual in that it's not strictly element-wise: to compute the value of one output node, it looks at the inputs of all the other output nodes. To compute it we simply take the exponent of each output node (to ensure that they are all positive) and then divides each by the total (to ensure that they sum to one).

After the softmax we can interpret the output of node y_3 as the probability that x has class 3, and train with cross entropy loss.

HOW DO WE FIND GOOD WEIGHTS?



$$\text{model}_{\theta}(x) = y$$
$$\text{loss}_{x,t}(\theta) = \|\text{model}_{\theta}(x) - t\|$$



12

To find good weights we first define a loss function. This is a function of a particular model (represented by its **the weights**) to a scalar value. **The better our model, the lower the loss.** If we imagine a model with just two weights, the **model space** forms a plane. For every point in this plane, our loss function defines a scalar loss, which we can draw above the plane as a surface: the **loss surface** (sometimes also called, more poetically, the **loss landscape**).

Make sure you understand the difference between the model, a function from the inputs x to the outputs y in which the parameters θ act as constants, and the loss, a function from the parameters to a loss value, in which the data acts as constants.

The symbol θ is a common notation referring to the set of all weights (sometimes combined into a vector, sometimes just a set).

$$\arg \min_{\theta} \text{loss}_{\text{data}}(\theta)$$



SOME COMMON LOSS FUNCTIONS

classification	squared errors	$\ y - t\ $ with $y = \text{model}_0(x)$
	absolute errors	$\ y - t\ _1 = \sum_i \text{abs}(y_i - t_i)$
	binary cross-entropy	$-\log p_0(t)$ with $t \in \{0, 1\}$
	cross-entropy	$-\log p_0(t)$ with $t \in \{0, \dots, K\}$
	hinge loss	$\max(0, 1 - ty)$ with $t \in \{-1, 1\}$



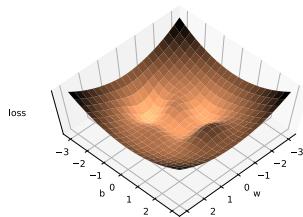
14

Here are some common loss functions for situations where we have examples (t) of what the model output (y) should be for a given input (x).

The squared error losses are derived from basic regression. The (binary) cross entropy comes from logistic regression (as shown last lecture) and the hinge loss comes from support vector machine classification. You can find their derivations in most machine learning books/courses. We won't elaborate on them here, except to note that in all cases the loss is lower if the model output (y) is closer to the example output (t).

The loss can be computed for a single example or for multiple examples. **In almost all cases, the loss for multiple examples is just the sum over all their individual losses.**

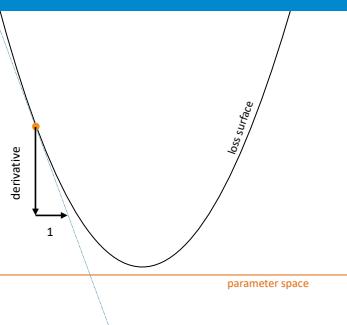
LOSS SURFACE



15

We want to follow the loss surface down to the lowest point.

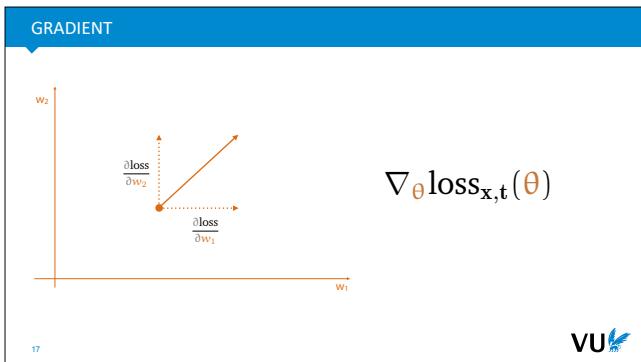
DERIVATIVE



16

In one dimension, we know that the derivative of a function (like the loss) tells us how much a function increases or decreases in we take a step of size 1 to the right.

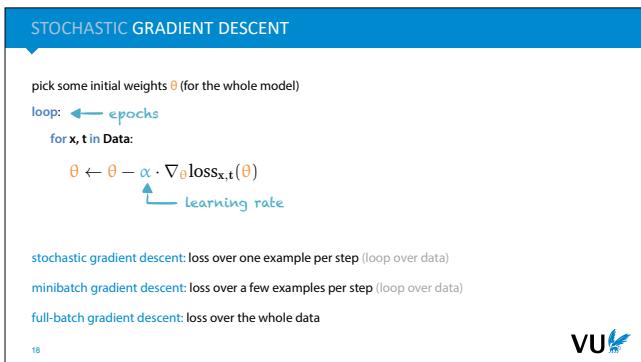
To be more precise, it tells us how much the best linear approximation (the [tangent line](#)) increases or decreases.



If our input space has multiple dimensions, like our loss surface, we can simply take the derivative with respect to each input, separately, treating the others as constants. This is called a **partial derivative**. The collection of all possible partial derivatives is called the **gradient**.

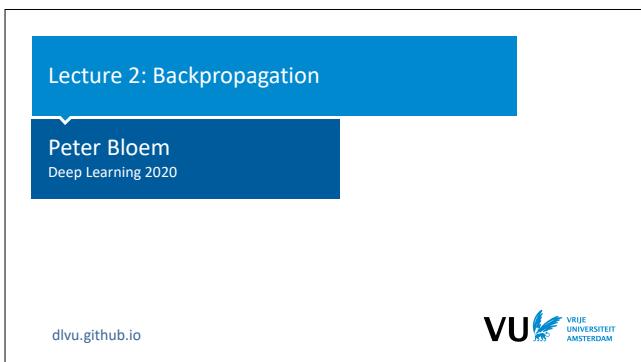
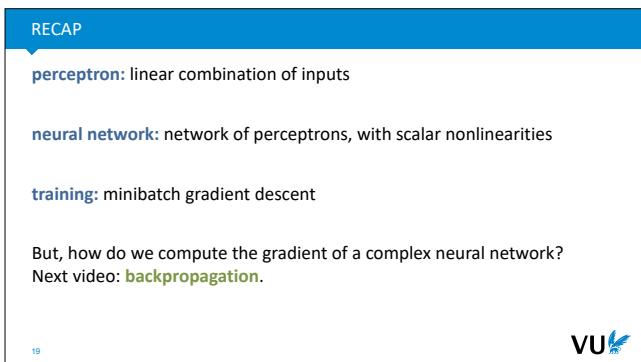
If we interpret the gradient as a vector, it points in the direction in which the function grows the fastest. Taking a step in the opposite direction means we are walking down the loss surface.

The symbol for the gradient is a downward pointing triangle called a nabla. The subscript indicates the variable over which we are taking the derivatives. Note that in this case we are treating



This is the idea behind the gradient descent algorithm. We compute the gradient, take a small step in the opposite direction and repeat. The reason we take small steps is that the gradient is only the direction of steepest ascent locally, the further we move from our current position the worse an approximation the tangent hyperplane will be for the function that we are actually trying to follow.

In deep learning, we almost always use **minibatch gradient descent**, but there are some examples.

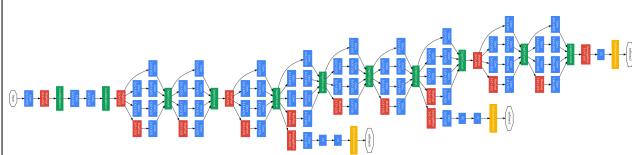


PART TWO: SCALAR BACKPROPAGATION

How do we work out the gradient for a neural network?



Working out a gradient is usually done by hand, with pen and paper. This function is then transferred to code, and used in a gradient descent loop. But the more complicated our model becomes, the more complex it becomes to work out a complete formulation of the gradient.



For very simple models like linear regression or logistic regression, we can just work out the gradient by hand, with pen and paper.

However, here is a diagram of the sort of network we'll be encountering (the GoogLeNet). We can't work out a complete gradient for this kind of architecture by hand. We need help.

GRADIENT COMPUTATION: THE SYMBOLIC APPROACH

WolframAlpha computational intelligence.

derivative of $l(w) = (d + z * 1/(1+e^a - (c + v*(b + 1/(1+e^a - (x*w))))))$

Extended Keyboard Upload

Examples

Random

Derivative: Approximate form Step-by-step solution

$$\frac{\partial}{\partial w} l(w) = d + \frac{z}{1 + e^{-(c+v(b+1/(1+e^{-xw})))}} = \frac{\partial \text{Hold}\left[\frac{z}{e^{-v(b+1/(1+e^{-xw}))}+1}+d\right]}{\partial w}$$

Alternate form:

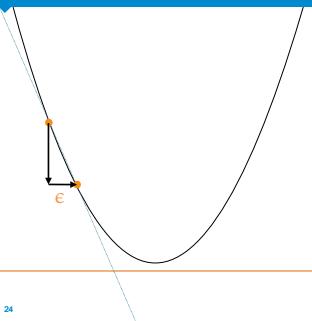
$$l'(w) = \frac{v x z e^{w x}}{(e^{w x} + 1)^2 \left(e^{v(b+1/(e^{-w x}+1))+c} + 1\right)} - \frac{v x z e^{w x}}{(e^{w x} + 1)^2 \left(e^{v(b+1/(e^{-w x}+1))+c} + 1\right)^2}$$

Of course, working out derivatives is a pretty mechanical process. We could easily take all the rules we know, and put them into some algorithm. This is called symbolic differentiation, and it's what systems like Wolfram Alpha do for us.

Unfortunately, as you can see here, the gradients it returns get pretty horrendous the deeper the neural network gets. This approach becomes impractical very quickly.

Note that in symbolic differentiation we get an answer that is **independent of the input**. We get a function that we can then feed an input to.

GRADIENT COMPUTATION: THE NUMERIC APPROACH



$$\text{deriv} \approx \frac{f(x + \epsilon) - f(x)}{\epsilon}$$

the "method of finite differences"



Another approach is to compute the gradient numerically. For instance by the method of finite differences: we take a small step ϵ and, see how much the function changes. The amount of change divided by the step size is a good estimate for the gradient if ϵ is small enough.

Numeric approaches are sometimes used, but it's very expensive to make them accurate if you have a large number of parameters.

Note that in the numeric approach, you only get an answer for a particular input. If you want to compute the gradient at some other point in space, you have to compute again.

BACKPROPAGATION: THE MIDDLE GROUND

Work out parts of the derivative **symbolically**
chain these together in a **numeric computation**.

secret ingredient: **the chain rule**.

25



Backpropagation is a kind of middle ground between symbolic and numeric gradient descent.

THE CHAIN RULE

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

shorthand:

$$x \rightarrow g \rightarrow f \quad \frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x}$$

26



The single thing that makes backpropagation possible is *the chain rule of differentiation*. If we want the derivative of a function which is the composition of two other functions, in this case **f** and **g**, we can take the derivative of **f** with respect to the output of **g** and multiply it by the derivative of **g** with respect to the input **x**.

Since we'll be using the chain rule *a lot*, we'll introduce a simple shorthand to make it a little easier to parse. We draw a little diagram of which function feeds into which. This means we know what the argument of each function is, so we can remove the arguments from our notation.

INTUITION FOR THE CHAIN RULE

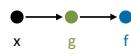
slope of **f** over **g**

$$f(g) = s_f g + b_f$$

$$g(x) = s_g x + b_g$$

with $s_f = \frac{\partial f}{\partial g}$

$$s_g = \frac{\partial g}{\partial x}$$



$$\begin{aligned} f(g(x)) &= s_f(s_g x + b_g) + b_f \\ &= s_f s_g x + s_f b_g + s_g \end{aligned}$$

slope of **f** over **x** constant

27



Since the chain rule is the heart of backpropagation, and backpropagation is the heart of deep learning, we should probably take some time to see why the chain rule is true at all.

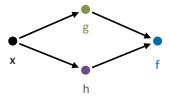
If we imagine that **f** and **g** are linear functions, it's pretty straightforward to show that this is true. They may not be, of course, but the nice thing about calculus is that locally, we can treat them as linear functions (if they are differentiable). In an infinitesimally small neighbourhood **f** and **g** are exactly linear.

In that case **f** and **g** both have a *slope*, s_f and s_g , which is simply the derivative at **x**. They also have additive constants b_f and b_g , which we're not interested in. For these linear approximations, we can easily work out what the function $f(g(x))$ looks like and what its slope is. It's the slope of **f** times the slope of **g** (plus some constant value that doesn't depend on **x**). This is exactly what the chain rule tells us.

If you want an example: imagine you have one investment scheme that doubles your money, and one that triples it: then, if you take the outcome of the first and put it into the second, you are multiplying your money by 6.

Note that this doesn't quite count as a proof, but it's hopefully enough to give you some intuition for why the chain rule holds.

MULTIVARIATE CHAIN RULE



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} + \frac{\partial f}{\partial h} \frac{\partial h}{\partial x}$$

28



Since we'll be looking at some pretty elaborate computation graphs, we'll need to be able to deal with this situation as well. We have a computation graph, as before, but f depends on x through **two different** operations. How do we take the derivative of f over x ?

The multivariate chain rule tells us that we can simply apply the chain rule along g , taking h as a constant, and sum it with the chain rule along h taking g as a constant.

INTUITION FOR THE MULTIVARIATE CHAIN RULE

$$\begin{aligned} f(g, h) &= s_1 g + s_2 h + b_f \\ g(x) &= s_g x + b_g \quad \text{with } s_1 = \frac{\partial f}{\partial g} \quad s_2 = \frac{\partial f}{\partial h} \quad s_g = \frac{\partial g}{\partial x} \quad s_h = \frac{\partial h}{\partial x} \\ h(x) &= s_h x + b_h \end{aligned}$$

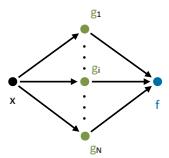
$$\begin{aligned} f(g(x), h(x)) &= s_1(s_g x + b_g) + s_2(s_h x + b_h) + b_f \\ &= s_1 s_g x + s_1 b_g + s_2 s_h x + s_2 b_h + b_f \\ &= (\underbrace{s_1 s_g + s_1 s_h}_{\text{slope of } f \text{ over } x} x + \underbrace{s_2 b_g + s_2 b_h + b_f}_{\text{constant}}) \end{aligned}$$

29



We can see why this holds in the same way as before. The short story, since all functions can be taken to be linear, their slopes distribute out into a sum

MULTIVARIATE CHAIN RULE



$$\frac{\partial f}{\partial x} = \sum_i \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial x}$$

30



If we have more than two paths from the input to the output, we simply sum over all of them.

EXAMPLE

$$f(x) = \frac{2}{\sin(e^{-x})}$$

31



ITERATING THE CHAIN RULE

$$f(x) = \frac{2}{\sin(e^{-x})}$$

operations:

$$d(c) = \frac{2}{c}$$

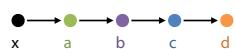
$$c(b) = \sin b$$

$$b(a) = e^a$$

$$a(x) = -x$$

$$f(x) = d(c(b(a(x))))$$

computation graph:



32

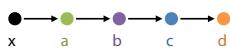
Here's an example of how the chain rule can help us to automate complicated derivatives. We start with the function on the top left. We break up its form into a series of smaller operations. The entire function f is then just a chain of these small operations chained together. We can draw this in a diagram as we did before.

We will call this kind of diagram a **computation graph**.

Normally, we wouldn't break a function up in such small operations. This is just a simple example to illustrate the principle.

ITERATING THE CHAIN RULE

$$\begin{aligned} \frac{\partial f}{\partial x} &= \frac{\partial d}{\partial x} \\ &= \frac{\partial d}{\partial c} \frac{\partial c}{\partial x} \\ &= \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial x} \\ &= \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x} \end{aligned}$$



33

Now, to work out the derivative of f , we can *iterate the chain rule*. We apply it again and again, until the derivative of f over x is expressed as a long product of derivatives of operation outputs over their inputs.

GLOBAL AND LOCAL DERIVATIVES

$$\frac{\partial f}{\partial x} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

global derivative local derivatives

We call the larger derivative of f over x the **global derivative**. And we call the individual factors, the derivatives of the operation output wrt to their inputs, the local derivatives.



34

BACKPROPAGATION

The **BACKPROPAGATION** algorithm:

- break your computation up into a sequence of operations what counts as an operation is up to you
- work out the **local derivatives** symbolically.
- compute the **global derivative** numerically by computing the local derivatives and multiplying them



35

WORK OUT THE LOCAL DERIVATIVES SYMBOLICALLY

$$f(x) = \frac{2}{\sin(e^{-x})}$$

operations:

$$d(c) = \frac{2}{c}$$

$$c(b) = \sin b$$

$$b(a) = e^a$$

$$a(x) = -x$$

$$\frac{\partial f}{\partial x} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

$$= -\frac{2}{c^2} \cdot \cos b \cdot e^a \cdot -1$$

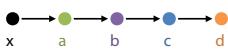


36

For each local derivative, we work out, on pen and paper the symbolic derivative. Note that we could fill in the a , b and c in the result, but we don't. We simply leave them as is, and continue with the numerical part of the algorithm.

COMPUTE A FORWARD PASS ($X = -4.499$)

$$f(-4.499) = 2$$



$$d = \frac{2}{c} = 2$$

$$c = \sin b = 1$$

$$b = e^a = 90$$

$$a = -x = 4.499$$



37

We start by computing what is known as a forward pass. We pick some input, in this case $x = -4.499$, and we feed it through the computation graph. Crucially, we save all the intermediate values, as well as the output.

Note that at this point, we are no longer computing solutions in general. We are computing our function for a specific input. We'll be computing the gradient for this specific input as well.

COMPUTE A FORWARD PASS ($X = -4.499$)

$$f(-4.499) = 2$$

$$\frac{\partial f}{\partial x} = -\frac{2}{c^2} \cdot \cos b \cdot e^a \cdot -1$$

$$d = \frac{2}{c} = 2$$

$$= -\frac{2}{1^2} \cdot \cos 90 \cdot e^{4.499} \cdot -1$$

$$c = \sin b = 1$$

$$= -2 \cdot 0 \cdot 90 \cdot -1 = 0$$

$$b = e^a = 90$$

$$a = -x = 4.499$$



38

Keeping all intermediate values from the forward pass in memory, we go back to our symbolic expression of the derivative. Here, we fill in the intermediate values a , b and c . After we do this, we can finish the multiplication numerically, giving us a numeric value of the gradient of f at $x = -4.499$.

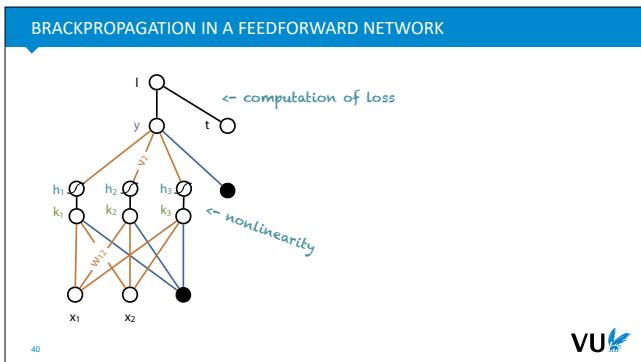
BACKPROPAGATION

The BACKPROPAGATION algorithm:

- break your computation up into a sequence of operations
what counts as an operation is up to you
- work out the **local derivatives** symbolically.
- compute the **global derivative** numerically
by computing the local derivatives and multiplying them
- Much more accurate than finite differences
only source of inaccuracy is the numeric computation of the operations.
- Much faster than symbolic differentiation
The backward pass has (broadly) the same complexity as the forward.



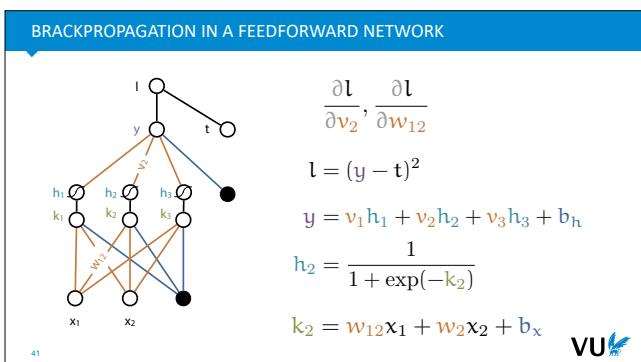
39



To explain how backpropagation works in a neural network, we extend our diagram a little bit, to make it closer to the actual computation graph we'll be using.

First, we separate the hidden node into the result of the linear operation k_i and the application of the nonlinearity h_i . Second, since we're interested in the derivative of the loss rather than the output of the network, we extend the network with one more step: the computation of the loss (over one example to keep things simple). In this final step, the output y of the network is compared to the target value t from the data, producing a loss value.

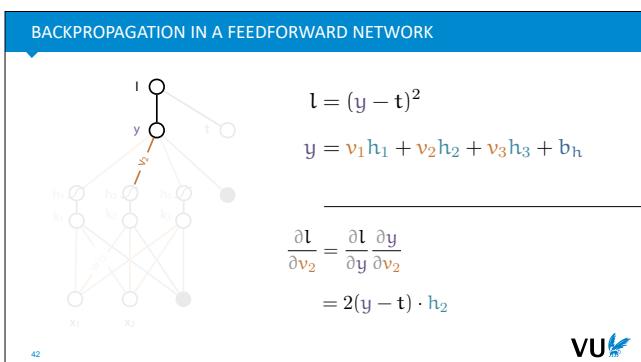
You can think of t as another input node, like x_1 and x_2 , but one to which the model doesn't have access.



We want to work out the gradient of the loss over the parameters. We'll isolate two parameters, v_2 in the second layer, and w_{12} in the first, and see how backpropagation operates.

First, we have to break the computation of the loss into operations. If we take the graph on the left to be our computation graph, then we end up with the operations of the right.

To simplify things, we'll compute the loss over only one example. We'll also use a simple squared error loss.



For the derivative with respect to v_2 , we'll only need these two operations. Anything below doesn't affect the result.

To work out the derivative we apply the chain rule, and work out the local derivatives symbolically.

BACKPROPAGATION IN A FEEDFORWARD NETWORK

$$\begin{aligned}\frac{\partial l}{\partial v_2} &= \frac{\partial l}{\partial y} \frac{\partial y}{\partial v_2} \\ &= 2(y - t) \cdot h_2 \\ &= 2(10.1 - 12.1) \cdot .99 = -3.96 \\ v_2 &\leftarrow v_2 - \alpha \cdot -3.96\end{aligned}$$

VU

We then do a forward pass with some values. We get an output of 10.1, which should have been 12.1, so our loss is 4. We keep all intermediate values in memory.

We then take our product of local derivatives, fill in the numeric values from the forward pass, and compute the derivative over v_2 .

When we apply this derivative in a gradient descent update, v_2 changes as shown below.

BACKPROPAGATION IN A FEEDFORWARD NETWORK

$$\begin{aligned}l &= (y - t)^2 \\ y &= v_1 h_1 + v_2 h_2 + v_3 h_3 + b_y \\ h_2 &= \frac{1}{1 + \exp(-k_2)} \\ k_2 &= w_{12} x_1 + w_{22} x_2 + b_k \\ \frac{\partial l}{\partial w_{12}} &= \frac{\partial l}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial k_2} \frac{\partial k_2}{\partial w_{12}} \\ &= 2(y - t) \cdot v_2 \cdot h_2(1 - h_2) \cdot x_1 \\ \text{derivative of sigmoid: } \sigma(x) &= \sigma(x)(1 - \sigma(x))\end{aligned}$$

VU

Let's try something a bit earlier in the network: the weight w_{12} . We add two operations, apply the chain rule and work out the local derivatives.

BACKPROPAGATION

$$\begin{aligned}\frac{\partial l}{\partial y} \\ \frac{\partial l}{\partial y} \frac{\partial y}{\partial h_2} &= \frac{\partial l}{\partial h_2} \\ \frac{\partial l}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial k_2} &= \frac{\partial l}{\partial k_2} \\ \frac{\partial l}{\partial y} \frac{\partial y}{\partial h_2} \frac{\partial h_2}{\partial k_2} \frac{\partial k_2}{\partial w_{12}} &= \frac{\partial l}{\partial w_{12}}\end{aligned}$$

VU

Note that when we're computing the derivative for w_{12} , we are also, along the way computing the derivatives for y , h_2 and k_2 .

This useful when it comes to implementing backpropagation. We can walk backward from the computation graph and compute the derivative of the loss for every node. For the nodes below, we just multiply the local gradient. This means we can very efficiently compute any derivatives we need.

We will show this more precisely in the last part of this lecture.

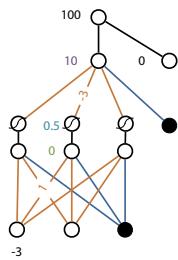
BACKPROPAGATION

The **BACKPROPAGATION** algorithm:

- break your computation up into a sequence of operations what counts as an operation is up to you
- work out the **local derivatives** symbolically.
- compute the **global derivative** numerically by computing the local derivatives and multiplying them
- Walk backward from the loss, *accumulating* the derivatives.

VU

BUILDING SOME INTUITION



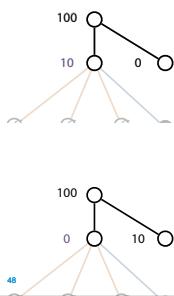
VU

47

To finish up, let's see if we can build a little intuition for what all these accumulated derivatives mean.

Here is a forward pass for some weights and some inputs. Backpropagation starts with the loss, and walks down the network, figuring out at each step how every value contributed to the result of the forward pass. Every value that contributed positively to a positive loss should be lowered, every value that contributed positively to a negative loss should be increased, and so on.

IMAGINE THAT y IS A PARAMETER



$$y \leftarrow y - \alpha \cdot 2(y - t) \\ = y - \alpha \cdot 20$$

$$y \leftarrow y - \alpha \cdot 2(y - t) \\ = y + \alpha \cdot 20$$

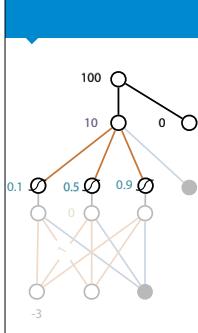
VU

48

We'll start with the first value below the loss: y . Of course, this isn't *parameter* of the network, but let's imagine for a moment that it is. What would the gradient descent update rule look like if we try to update y ?

Even though we can't change y directly, this is the effect we want to achieve: we want to change the parameters we *can* change so that we achieve this change in y .

Note that the error (the derivative of the loss) takes **the sign** into account. If our model output is 0, and our target is 10, the loss is the same, but the error is *minus* 20, and the update rule for y



$$v_2 \leftarrow v_2 - \alpha \cdot 2(y - t)h_2 \\ v_1 \leftarrow v_1 - \alpha \cdot 20 \cdot 0.1 \\ v_2 \leftarrow v_2 - \alpha \cdot 20 \cdot 0.5 \\ v_3 \leftarrow v_3 - \alpha \cdot 20 \cdot 0.9$$

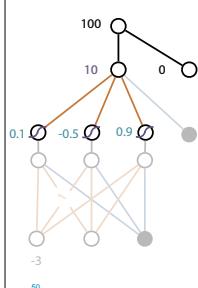
VU

49

Instead of changing y , we have to change the values that influenced y . Here we see what that looks like for the **weights** of the second layer.

Note that the current value of the weight doesn't factor into the update. Only how much influence the weight had on the value of y in the forward pass. The higher the activation of the **source node**, the more the weight gets adjusted.

NEGATIVE ACTIVATION



$$\tanh \theta$$

$$v_1 \leftarrow v_1 - \alpha \cdot 20 \cdot 0.1 \\ v_2 \leftarrow v_2 + \alpha \cdot 20 \cdot 0.5 \\ v_3 \leftarrow v_3 - \alpha \cdot 20 \cdot 0.9$$

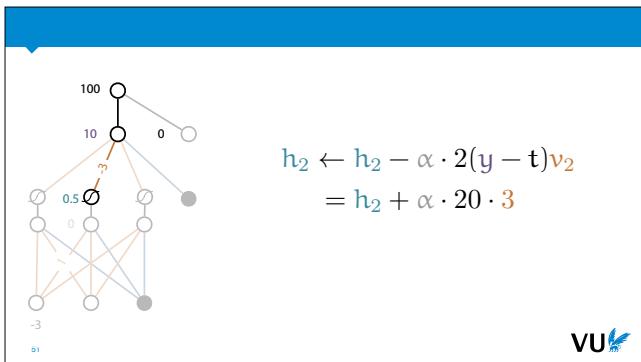
VU

50

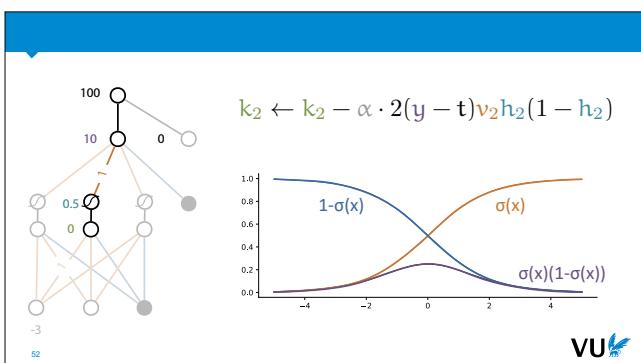
The sigmoid activation we've used so far allows only positive values to emerge from the hidden layer. If we switch to an activation that also allows negative activations (like a linear activation or a **tanh** activation), we see that backpropagation very naturally takes the sign into account.

In this case, we want to update in such a way that y decreases, but we note that the weight v_2 is multiplied by a *negative* value. This means that (for this instance) v_2 contributes negatively to the loss, and its value should be increased.

Note that the sign of v_2 itself doesn't matter. Whether it's positive or negative, its value should increase.



If we could change h_2 directly, this is how we'd do it. We take the value of v_2 to be a constant. We want to decrease the output of the network. Since v_2 makes a *negative* contribution to the loss, we can achieve this by *increasing* the activation of the source node of v_2 .

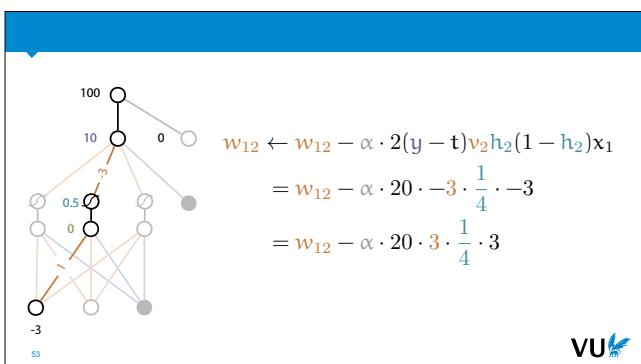


Moving down to k_2 , remember that the derivative of the sigmoid is the output of the sigmoid times 1 minus that output.

We see here, that in the extreme regimes, the sigmoid is *resistant to change*. The closer to 1 or 0 we get the smaller the weight update becomes.

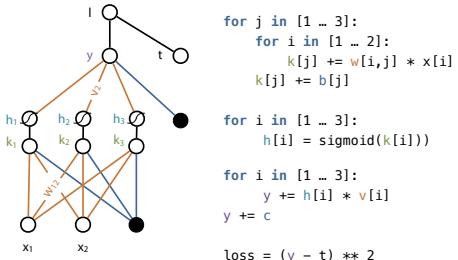
This is actually a great downside of the sigmoid activation, and one of the big reasons it was eventually replaced by the ReLU as the default choice for hidden units. We'll come back to this in later lectures.

Nevertheless, this update rule tells us what the change is to k_2 that we *want to achieve* by changing the gradients we can actually change (in this case *the weights* of layer 1).



Finally, we come to the weights of the first layer. As before, we want the output of the network to *decrease*. To achieve this, we want h_2 to *increase* (because v_2 is negative). However, the input x_1 is negative, so we should decrease w_{12} to increase h_2 . This is all beautifully captured by the chain rule: the two negatives of x_1 and v_2 cancel out and we get a positive value which we subtract from w_{12} .

FORWARD PASS IN PSEUDOCODE

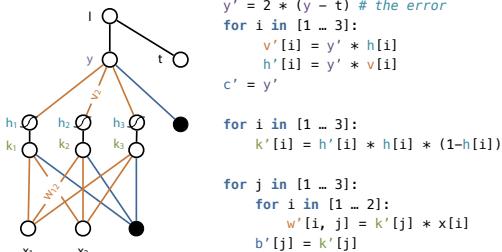


54



Assume that **k** and **y** are initialized with 0s.

BACKWARD PASS IN PSEUDOCODE



55



Note that we don't implement the derivations from slide 39 directly. Instead, we work backwards down the neural network: computing the derivative of each node as we go by taking the derivative of the loss over the outputs and multiplying it by the local derivative.

RECAP

Backpropagation: method for computing derivatives.

Combined with **gradient descent** to train NNs.

Middle ground between symbolic and numeric computation.

- Break computation up into operations.
- Work out **local derivatives** symbolically.
- Work out **global derivative** numerically.

56



Lecture 2: Backpropagation

Peter Bloem
Deep Learning 2020

PART THREE: TENSOR BACKPROPAGATION

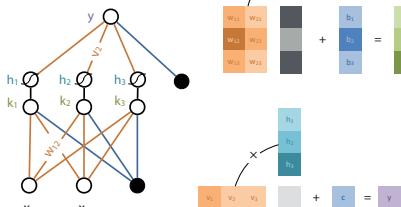
Expressing backpropagation in vector, matrix and tensor operations.



We've seen what neural networks are, how to train them by gradient descent and how to compute that gradient by backpropagation.

In order to scale up to larger and more complex structures, we need two make our computations as efficient as possible, and we'll also need to simplify our notation. There's one insight that we are going to get a lot of mileage out of.

IT'S ALL JUST LINEAR ALGEBRA



When we look at the computation of a neural network, we can see that most operations can be expressed very naturally in those of linear algebra.

The multiplication of weights by their inputs is a multiplication of **a matrix of weights** by a vector of inputs.

The addition of a bias is the addition of **a vector of bias parameters**.

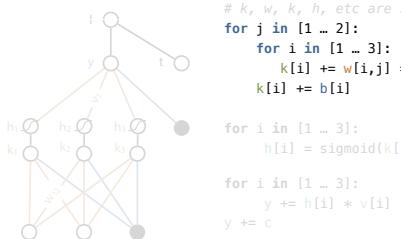
The nonlinearities can be implemented as simple element-wise operations.

This perspective buys us two things. First...

$$f(x) = \mathbf{V} \sigma(\mathbf{W}x + \mathbf{b}) + \mathbf{c}$$

Our notation becomes extremely simple. We can describe the whole operation of a neural network with one easily parseable equation. This expressiveness will be sorely needed when we move to more complicated networks.

MATRIX MULTIPLICATION



```
# k, w, k, h, etc are zero'd arrays
for j in [1 .. 2]:
    for i in [1 .. 3]:
        k[i] += w[i,j] * x[i]
        k[i] += b[i]

for i in [1 .. 3]:
    h[i] = sigmoid(k[i])

for i in [1 .. 3]:
    y += h[i] * v[i]
y += c

loss = (y - c) ** 2
```

The second reason is that the biggest computational bottleneck in a neural network is the matrix multiplication. This operation is more than quadratic while everything else is linear.

Matrix multiplication (and other tensor operations like it) can be parallelized and implemented efficiently but it's a lot of work. Ideally, we'd like to let somebody else do all that work (like the implementers of numpy) and then just call their code to do the matrix multiplications.

This is especially important if you have access to a GPU. A matrix multiplication can easily be offloaded to the GPU for much faster processing, but for a loop over an array, there's no benefit.

This is called **vectorizing**: taking a piece of code written in for loops, and getting rid of the loops by expressing the function as a sequence of linear algebra operations.

VECTORIZING

Express everything as operations on vectors, matrices and tensors.

Get rid of all the loops

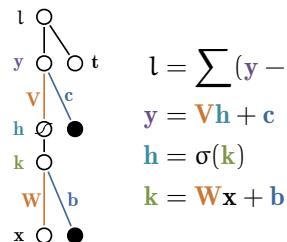
Makes the notation *simpler*.

Makes the execution *faster*.

62



FORWARD



$$\begin{aligned} l &= \sum (y - t)^2 \\ y &= Vh + c \\ h &= \sigma(k) \\ k &= Wx + b \end{aligned}$$

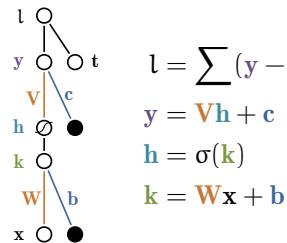
$$\begin{aligned} k &= w \cdot \text{dot}(x) + b \\ h &= \text{sigmoid}(k) \\ y &= v \cdot \text{dot}(h) + c \\ l &= (y - t) ** 2 \end{aligned}$$

63



Here's what the forward pass looks like in pseudocode. When you implement this in numpy it'll look almost the same.

BUT WHAT ABOUT THE BACKWARD PASS?



$$\begin{aligned} l &= \sum (y - t)^2 \\ y &= Vh + c \\ h &= \sigma(k) \\ k &= Wx + b \end{aligned}$$

$$\frac{\partial l}{\partial W} = ? \quad \frac{\partial l}{\partial y} \frac{\partial y}{\partial h} \frac{\partial h}{\partial k} \frac{\partial k}{\partial W}$$

64



Of course, we lose a lot of the benefit of vectorizing if the backward pass isn't also expressed in terms of matrix multiplications. How do we vectorize the backward pass?

On the left, we see the forward pass of our loss computation as a set of operations on vectors and matrices. (We've generalized things by giving the network a vector output y and vector target t).

To generalize backpropagation to this view, we might ask if something similar to the chain rule exists for vectors and matrices. Firstly, can we define something analogous to the derivative of one matrix over another, and secondly, can we break this apart into a product of local derivatives,

possibly giving us a sequence of matrix multiplications?

The answer is yes, there are many ways of applying calculus to vectors and matrices, and there are many chain rules available in these domains. However, things can quickly get a little hairy, so we need to tread carefully.

GRADIENTS, JACOBIANS, ETC

$f(\mathbf{A}) = \mathbf{B}$, $\frac{\partial \mathbf{B}}{\partial \mathbf{A}}$: derivatives of every element of \mathbf{A} over every element of \mathbf{B}

for instance:

$$f \left(\begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \right) = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

$$J_f = \begin{pmatrix} \frac{\partial b_1}{\partial a_1} & \frac{\partial b_1}{\partial a_2} \\ \frac{\partial b_2}{\partial a_1} & \frac{\partial b_2}{\partial a_2} \\ \frac{\partial b_1}{\partial a_3} & \frac{\partial b_2}{\partial a_3} \end{pmatrix}$$

function returns a
scalar vector matrix
arrange derivatives in a:
scalar scalar vector matrix
inputs a
vector vector matrix ?
matrix matrix ? ?

65

The derivatives of high-dimensional objects are easily defined. We simply take the derivative of every input over every output, and we arrange all possibilities into some logical shape. For instance, if we have a vector-to-vector function, the natural analog of the derivative is a matrix with all the partial derivatives in it.

However, once we get to matrix/vector operations or matrix/matrix operations, the only way to logically arrange every input with every output is a tensor of higher dimension than 2.

NB: We don't normally apply the differential symbol to non-scalars like this. We'll introduce better notation later.

$\frac{\partial l}{\partial W} = \frac{\partial l}{\partial y} \frac{\partial y}{\partial h} \frac{\partial h}{\partial k} \frac{\partial k}{\partial W}$
uh oh

VU

66

As we see, even if we could come up with this kind of chain rule, one of the local derivatives is now a vector over a matrix. The result could only be represented in a 3-tensor. There are two problems with this:

- If the layer has n inputs and n outputs, we are computing n^3 derivatives, even though we were only ultimately interested in n^2 of them (the derivatives of W). In the scalar algorithm we only ever had two nested loops (an n^2 complexity), and we only ever stored one gradient for one node in the computation graph. Now we suddenly have n^3 complexity and n^3 memory use. We were supposed to make things faster.
- We can easily represent a 3-tensor, but there's no obvious, default way to multiply a 3-tensor with a matrix, or with a vector (in fact there are many different ways). The multiplication of the chain rule becomes very complicated this way.

In short, we need a more careful approach.

SIMPLIFYING ASSUMPTIONS

1) The computation graph *always* has a single, scalar output: l

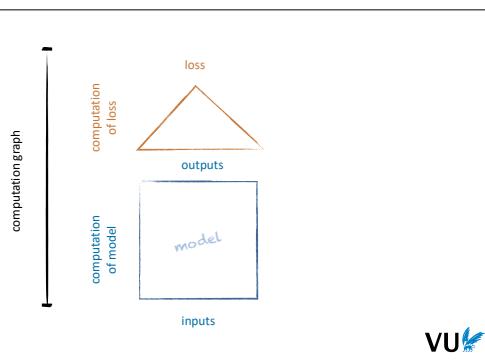
2) We are only ever interested in the derivative of l .

$$\frac{\partial l}{\partial \mathbf{W}} \quad \begin{matrix} \leftarrow \text{scalar} \\ \leftarrow \text{tensor of any dimension} \end{matrix}$$

	scalar	vector	matrix
scalar	scalar	vector	matrix
vector	vector	matrix	?
matrix	matrix	?	?
3-tensor	3-tensor	?	?

67

To work out how to do this we make these following simplifying assumptions.



Note that this doesn't mean we can only ever train neural networks with a single scalar output. Our neural networks can have any number of outputs of any shape and size. However, the *loss* we define over those outputs needs to map them all to a single scalar value. The computation graph is always the model, plus the computation of the loss.

THE GRADIENT

We will call $\frac{\partial l}{\partial \mathbf{W}}$ the gradient of l (with respect to \mathbf{W})

Commonly written as $\nabla_{\mathbf{W}} l$

Nonstandard assumption: $\nabla_{\mathbf{W}} l$ has the same *shape* as \mathbf{W}

$$[\nabla_{\mathbf{W}} l]_{ijk} = \frac{\partial l}{\partial \mathbf{W}_{ijk}}$$



We call this derivative the gradient. This is a common term, but we will deviate from the standard approach in one respect.

Normally, the gradient defined as a row vector, so that it can operate on a space of column vectors by matrix multiplication.

In our case, we are never interested in matrix multiplying the gradient by anything. We only ever want to sum the gradient of l wrt \mathbf{W} with the original matrix \mathbf{W} in a gradient update step. For this reason we define the gradient as having the same shape as the tensor with respect to which we are taking the derivative.

In the example shown, \mathbf{W} is a 3-tensor. The gradient of l wrt \mathbf{W} has the same shape as \mathbf{W} , and at element (i, j, k) it holds the scalar derivative of l wrt \mathbf{W}_{ijk} .

With these rules, we can use tensors of any shape and dimension and always have a well defined gradient.

NEW NOTATION: THE GRADIENT FOR \mathbf{W}

$$\nabla_{\mathbf{W}} \mathbf{l} \leftarrow \text{always the same} = \mathbf{W}^{\nabla}$$

$\nabla_{\mathbf{W}} \mathbf{l}$ ← actual object of interest

$$\mathbf{b}^{\nabla} = \nabla_{\mathbf{b}} \mathbf{l} = \left[\frac{\partial \mathbf{l}}{\partial \mathbf{b}_1} \dots \frac{\partial \mathbf{l}}{\partial \mathbf{b}_n} \right]^T$$

$$\mathbf{y}^{\nabla} = \nabla_{\mathbf{y}} \mathbf{l} = \frac{\partial \mathbf{l}}{\partial \mathbf{y}}$$



The standard gradient notation isn't very suitable for our purposes. It puts the loss front and center, but that will be the same for all our gradients. The object that we're actually interested in is relegated to a subscript. Also, it isn't very clear from the notation what the shape is of the tensor that we're looking at.

We can think of the nabla as an operator like a transposition taking an inverse. It turns a matrix into another matrix, a vector into another vector and so on.

For this reason, we'll introduce a new notation. This isn't standard, so don't expect to see it anywhere else, but it will help to clarify our mathematics a lot as we go forward.

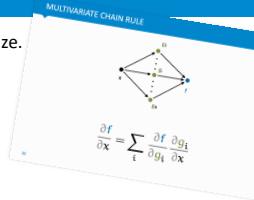
We've put the \mathbf{W} front and center, so it's clear that the result of taking the gradient is also a matrix, and we've removed the loss, since we've assumed that we are always taking the gradient of the loss.

The notation works the same for vectors and even for scalars. When we refer to a single element of the gradient, we will follow our convention for matrices, and use the non-bold version of its letter.

This is the gradient of \mathbf{l} . Since \mathbf{l} never changes, we'll refer to this as "the gradient for \mathbf{W} ".

TENSOR BACKPROPAGATION

Work out *scalar* derivatives first, then vectorize.
Use the multivariate chain rule to derive the scalar derivative.



Apply the chain rule step by step.

Start at the loss and work backward, *accumulating* the gradients.

71



This gives us a good way of thinking about the gradients, but we still don't have a chain rule to base out back propagation on.

The main trick we will use is to stick to scalar derivatives as much as possible. Once we have worked out the derivative in purely scalar terms, we will find a way to vectorize their computation.

GRADIENT FOR \mathbf{y}



$$\begin{aligned} \mathbf{y}_i^{\nabla} &= \frac{\partial \mathbf{l}}{\partial \mathbf{y}_i} \\ &= \frac{\partial \sum_k (y_k - t_k)^2}{\partial y_i} = \sum_k \frac{\partial (y_k - t_k)^2}{\partial y_i} \\ &= \sum_k 2(y_k - t_k) \frac{\partial y_k - t_k}{\partial y_i} = 2(y_i - t_i) \end{aligned}$$

$$\mathbf{y}^{\nabla} = 2 \begin{pmatrix} y_1 - t_1 \\ \vdots \\ y_n - t_n \end{pmatrix} = 2 \cdot (\mathbf{y} - \mathbf{t})$$



We start simple: what is the gradient for \mathbf{y} ? This is a vector, because \mathbf{y} is a vector. Let's first work out the derivative of the i -th element of \mathbf{y} . This is purely a scalar derivative so we can simply use the rules we already know. We get $2(y_i - t_i)$ for that particular derivative.

Then, we just re-arrange all the derivatives for y_i into a vector, which gives us the complete gradient for \mathbf{y} .

The final step requires a little creativity: we need to figure out how to compute this vector using only basic linear algebra operations on the given vectors and matrices. In this case it's not so complicated: we get the gradient for \mathbf{y} by element-

wise subtracting \mathbf{t} from \mathbf{y} and multiplying by 2.

We haven't needed any chain rule yet, because our computation graph has only one edge.

GRADIENT FOR \mathbf{V}

$$\begin{aligned} \mathbf{V}_{ij} &= \frac{\partial \mathbf{L}}{\partial \mathbf{V}_{ij}} \\ &= \sum_k \frac{\partial \mathbf{L}}{\partial y_k} \frac{\partial y_k}{\partial \mathbf{V}_{ij}} = \sum_k y_k^v \frac{\partial y_k}{\partial \mathbf{V}_{ij}} \\ &= \sum_k y_k^v \frac{\partial [\mathbf{V}\mathbf{h} + \mathbf{c}_k]}{\partial \mathbf{V}_{ij}} = \sum_k y_k^v \frac{\partial [\mathbf{V}_{ki} \mathbf{h}_i]_k + c_k}{\partial \mathbf{V}_{ij}} \\ &= \sum_k y_k^v \frac{\partial \sum_l \mathbf{V}_{kl} \mathbf{h}_i + c_k}{\partial \mathbf{V}_{ij}} \\ &= \sum_{kl} y_k^v \frac{\partial \mathbf{V}_{kl} \mathbf{h}_i}{\partial \mathbf{V}_{ij}} = y_i^v \frac{\partial \mathbf{V}_{ij} \mathbf{h}_i}{\partial \mathbf{V}_{ij}} \\ &= y_i^v \mathbf{h}_i \end{aligned}$$

$$\mathbf{V}^\nabla = \begin{pmatrix} \dots & y_i^v \mathbf{h}_i & \dots \end{pmatrix} = \mathbf{y}^\nabla \mathbf{h}^T$$

73

Let's move one step down and work out the gradients for \mathbf{V} .

RECIPE

To work out a gradient \mathbf{X}^∇ for some \mathbf{X} :

- Write down a **scalar derivative** of the loss over *one element* of \mathbf{X} .
- Use the **multivariate chain rule** to sum over all outputs.
- Work out the scalar derivative.
- **Vectorize** the computation of \mathbf{X}^∇ in terms of the original inputs.

74



GRADIENT FOR \mathbf{h}

$$\begin{aligned} \mathbf{h}_i^\nabla &= \frac{\partial \mathbf{L}}{\partial \mathbf{h}_i} \\ &= \sum_k \frac{\partial \mathbf{L}}{\partial y_k} \frac{\partial y_k}{\partial \mathbf{h}_i} = \sum_k y_k^v \frac{\partial y_k}{\partial \mathbf{h}_i} \\ &= \sum_k y_k^v \frac{\partial [\mathbf{V}\mathbf{h} + \mathbf{c}_k]}{\partial \mathbf{h}_i} = \sum_k y_k^v \frac{\partial \sum_l \mathbf{V}_{kl} \mathbf{h}_l + c_k}{\partial \mathbf{h}_i} \\ &= \sum_k y_k^v \frac{\partial \mathbf{V}_{ki} \mathbf{h}_i + c_k}{\partial \mathbf{h}_i} = \sum_k y_k^v \frac{\partial \mathbf{V}_{ki} \mathbf{h}_i}{\partial \mathbf{h}_i} \\ &= \sum_k y_k^v \mathbf{V}_{ki} \end{aligned}$$

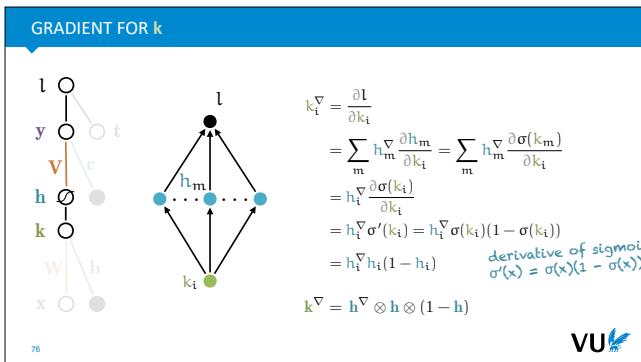
$$\mathbf{h}^\nabla = \begin{pmatrix} \sum_k y_k^v \mathbf{V}_{ki} \end{pmatrix} = (\mathbf{y}^\nabla \mathbf{I}^T \otimes \mathbf{V}) \mathbf{I}$$

or in numpy:
`(gy[:, None] * v).sum(axis=1)`

75

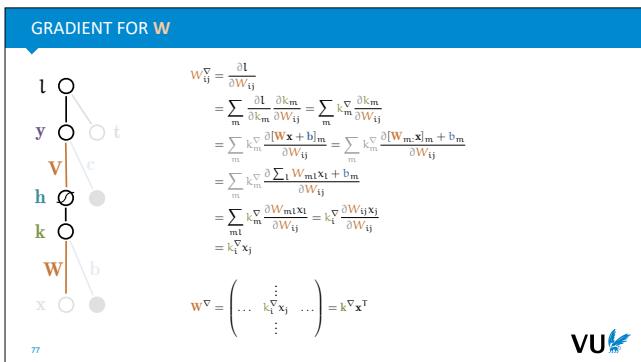
Since this is an important principle to grasp, let's keep going until we get to the other set of parameters, \mathbf{W} . We'll leave the biases as an exercise.

Most of the derivation is the same as before. However, sin

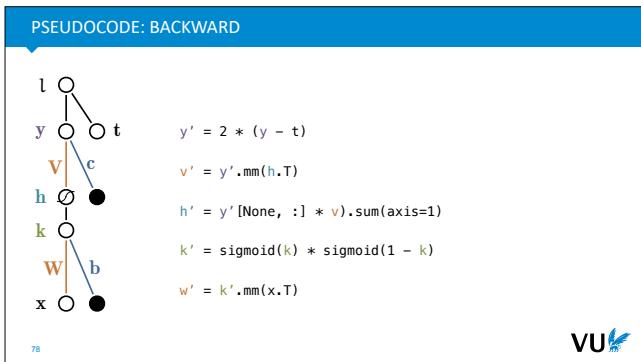


Now, at this point, when we analyze \mathbf{k} , note that we already have the gradient over \mathbf{h} . This means that we no longer have to apply the chain rule to anything above \mathbf{h} . We can draw this scalar computation graph, and work out the local gradient for \mathbf{k} in terms for the gradient for \mathbf{h} .

Given that, working out the gradient for \mathbf{k} is relatively easy, since the operation from \mathbf{k} to \mathbf{h} is an element-wise one.



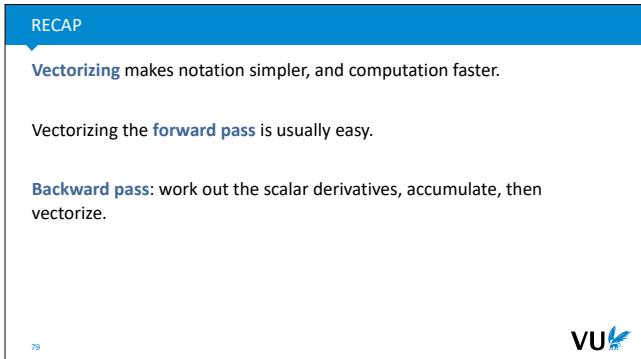
Finally, the gradient for \mathbf{W} . The situation here is exactly the same as we saw earlier for \mathbf{V} (matrix in, vector out, matrix multiplication), so we should expect the derivation to have the same form (and indeed it does).



And here's the backward pass that we just derived.

We've left the derivatives of the bias parameters out. You'll have to work these out to implement the first homework exercise.

(note that ' cannot be part of a variable name in actual python. you'll have to come up with some other name)



Lecture 2: Backpropagation

Peter Bloem
Deep Learning 2020

dlvu.github.io

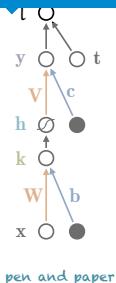


PART FOUR: AUTOMATIC DIFFERENTIATION

Letting the computer do all the work.

We've simplified the computation of derivatives a lot. All we've had to do is work out local derivatives and chain them together. We can go one step further: if we let the computer keep track of our computation, and provide some backwards functions for basic operations, the computer can work out the whole backpropagation algorithm for us. This is called **automatic differentiation** (or sometimes **autograd**).

THE STORY SO FAR



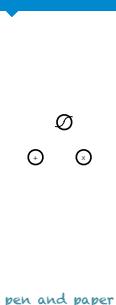
pen and paper
82

```
v' = y'.mm(h.T)  
h' = y'[None, :] * v.sum(axis=1)  
k' = sigmoid(k) * sigmoid(1 - k)  
w' = k'.mm(x.T)
```

in the computer



THE IDEAL

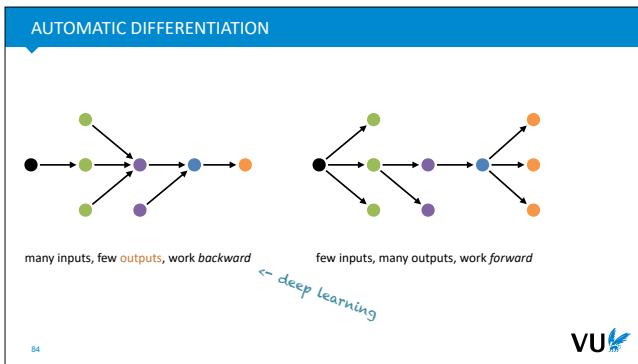


pen and paper

```
k = w.dot(x) + b  
h = sigmoid(k)  
y = v.dot(h) + c  
l = (y - t) ** 2  
l.backward() # start backprop
```

in the computer

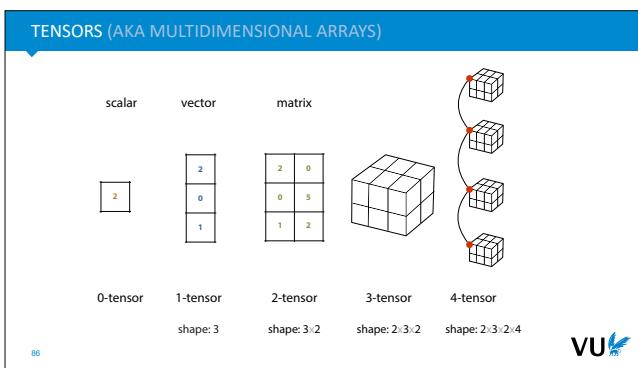
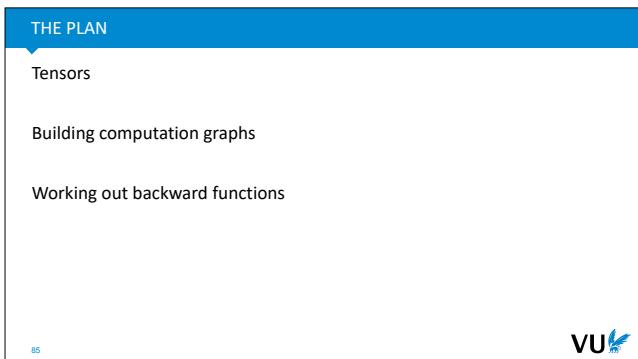
This is what we want to achieve. We work out on pen and paper the local derivatives of various modules, and then we chain these modules together in a computation graph *in our code*. The computer keeps the graph in memory, and can automatically work out the backpropagation.



This kind of algorithm is called automatic differentiation. What we've been doing so far is called **backward, or reverse mode automatic differentiation**. This is efficient if you have few output nodes. Note that if we had two output nodes, we'd need to do a separate backward pass for each.

If you have few inputs, it's more efficient to start with the inputs and apply the chain rule working forward. This is called forward mode automatic differentiation.

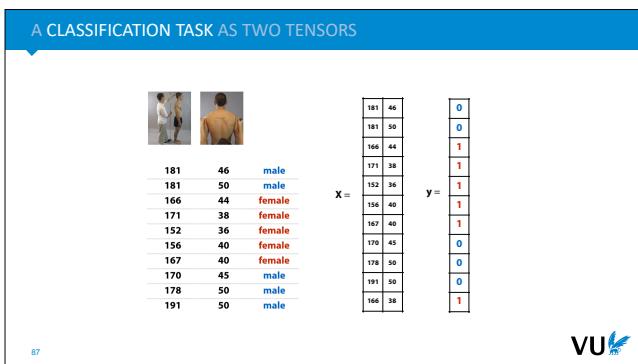
Since we assume we only have one output node (where the gradient computation is concerned), we will always use reverse mode.



The basic datastructure of our system will be the tensor. A tensor is a generic name for family of datastructures that includes a scalar, a vector, a matrix and so on.

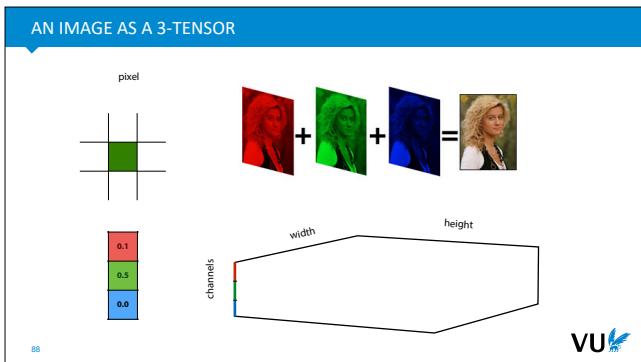
There is no good way to visualize a 4-dimensional structure. We will occasionally use this form to indicate that there is a fourth dimension along which we can index the tensor.

We will assume that whatever data we work with (images, text, sounds), will in some way be encode into one or more tensors, before it is fed into our system. Let's look at some examples.



A simple dataset, with numeric features can simply be represented as a matrix. For the labels, we usually create a separate corresponding vector for the labels.

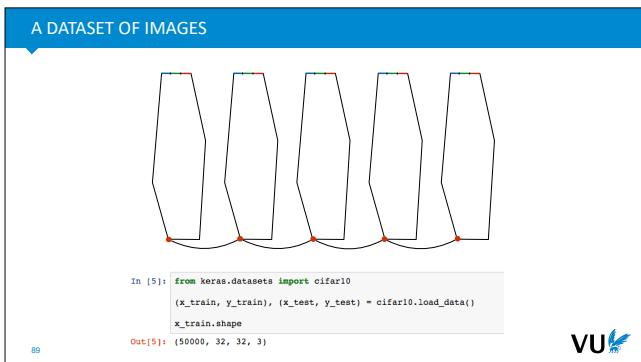
Any categoric features or labels should be converted to numeric features (normally by one-hot coding).



Images can be represented as 3-tensors. In an RGB image, the color of a single pixel is represented using three values between 0 and 1 (how red it is, how green it is and how blue it is). This means that an RGB image can be thought of as a stack of three color channels, represented by matrices.

This stack forms a 3-tensor.

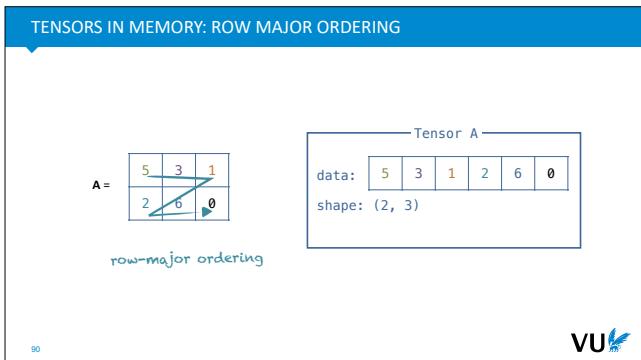
source: [http://www.adsell.com/
scanning101.html](http://www.adsell.com/scanning101.html)



If we have a dataset of images, we can represent this as a 4-tensor, with dimensions indexing the instances, their width, their height and their color channels respectively. Below is a snippet of code showing that when you load the CIFAR10 image data in Keras, you do indeed get a 4-tensor like this.

There is no agreed standard ordering for the dimensions in a batch of images. Tensorflow and Keras use (batch, height, width, channels), whereas Pytorch uses (batch, channels, height, width).

(You can remember the latter with the mnemonic “**bachelor chow**”.)

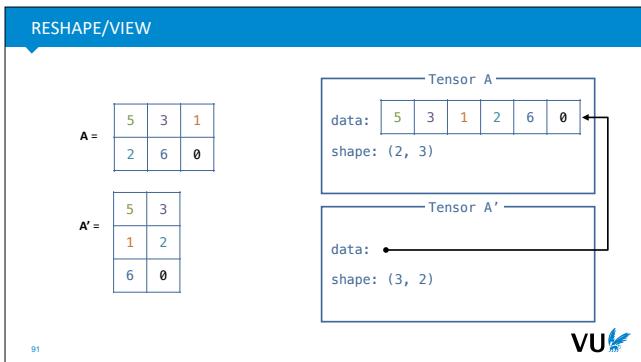


It's important to realize that even though we think of tensors as multidimensional arrays, in memory, they are necessarily laid out as a single line of numbers. The Tensor object knows the shape that these numbers should take, and uses that to compute whatever we need it to compute.

We can do this in two ways: scanning along the rows first and then the columns is called **row-major ordering**. This is what numpy and pytorch both do. The other option is (unsurprisingly) called column major ordering. In higher dimensional tensors, the dimensions further to the right are always scanned before the dimensions to their left.

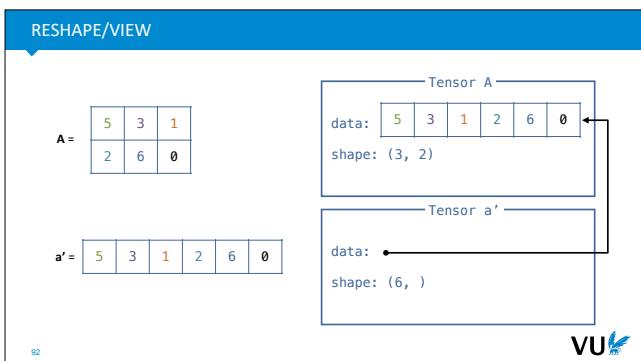
Imagine looping over all elements in a tensor with shape (a, b, c, d) in four nested loops. If you want to loop over the elements in the order they are in memory, then the loop over d would be the innermost loop, then c, the b and then a.

This allows us to perform some operations very cheaply, **but we have to be careful**.



Here is one such cheap operation: a reshape. This changes the shape of the tensor, but not the data.

To do this cheaply, we can create a new tensor object with the same data as the old one, but with a different shape.



EXAMPLE: NORMALIZE COLORS

```
images = load_images(...)
n, h, w, c = images.size()

images = images.reshape(n*h*w, c)
images = images / images.max(dim=0)
images = images.reshape(n, h, w, c)

# this makes no sense, but it is allowed.
images = images.reshape(n*h, w*c)
images = images.reshape(h, n, c, w)
```

93



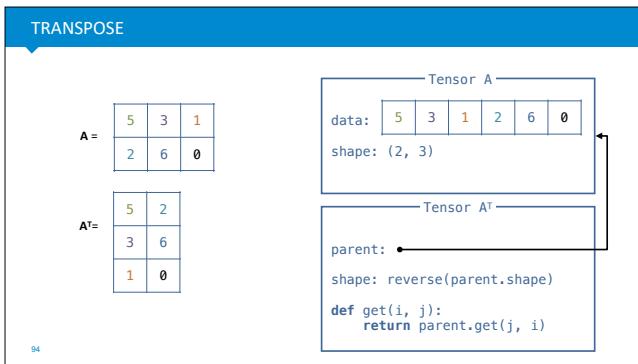
Imagine that we wanted to scale a dataset of images in such way that over the whole dataset, each color channel has a maximal value of 1 (independent of the other channels).

To do this, we need a 3-vector of the maxima per color channel. Pytorch only allows us to take the maximum in one direction, but we can reshape the array so that all the directions we're not interested in are collapsed into one dimension.

Afterwards, we can reverse the reshape to get our image dataset back.

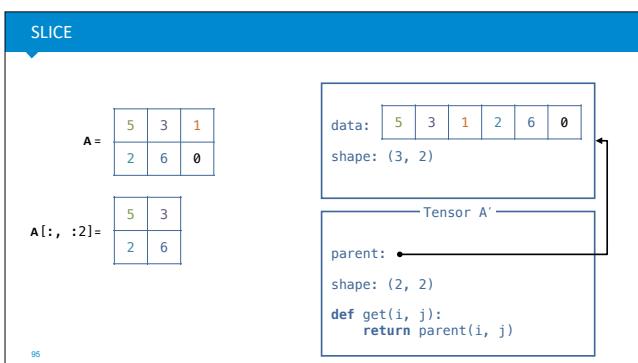
We have to be careful: the last three lines in this slide form a perfectly valid reshape, but the `c` dimension in the result does *not* contain our color values.

In general, you're fine if you **collapse dimensions that are next to each other**, and uncollapse them *in the same order*.

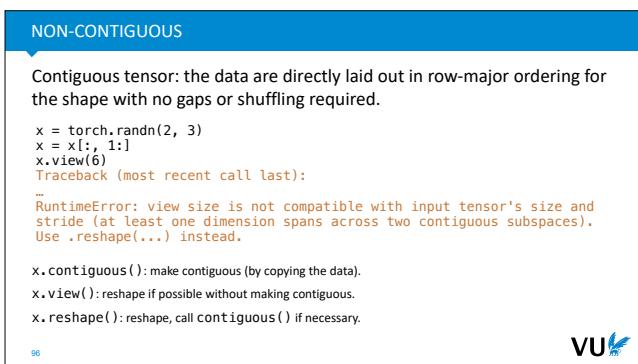


A transpose operation can also be achieved cheaply. In this case, we just keep a reference to the old tensor, and whenever a user requests the element (i, j) we return (j, i) from the original tensor.

NB: This isn't precisely how pytorch does this, but the effect is the same: we get a transposed tensor in constant time, by viewing the same data in memory in a different way.

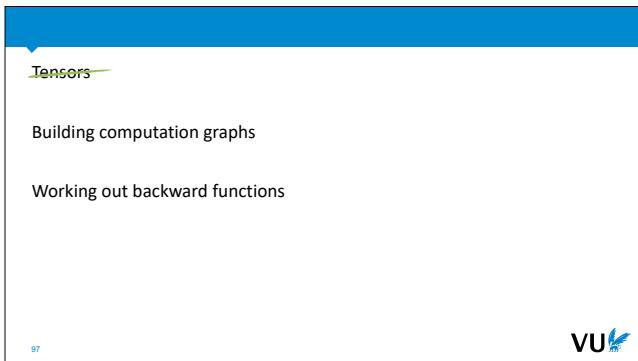


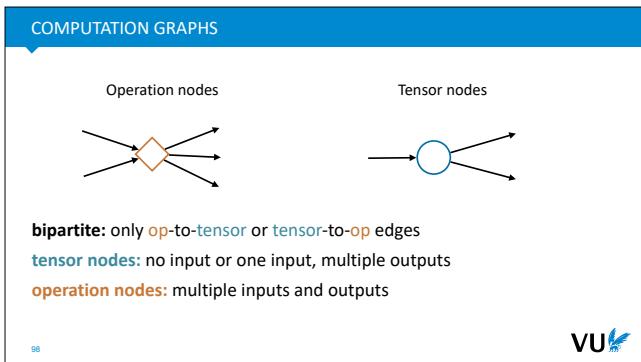
Even slicing can be accomplished by referring to the original data.



For some operations, however, the data needs to be contiguous. That is, the tensor data in memory needs to be one uninterrupted string of data in row major ordering with no gaps. If this isn't the case, pytorch will throw an exception like this.

You can fix this by calling `.contiguous()` on the tensor. The price you pay is the linear time and space complexity of copying the data.

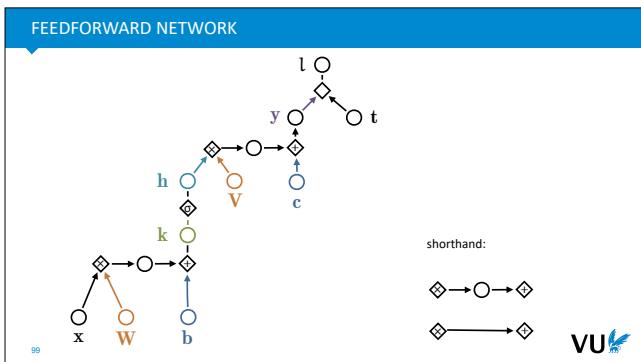




We'll need to be a little more precise in our notations. From now on we'll draw computation graphs like this, making the operation explicit.

(There doesn't seem to be a standard notation. This will work for our purposes).

In tensorflow operations are called *ops*, and in pytorch they're called *functions*.



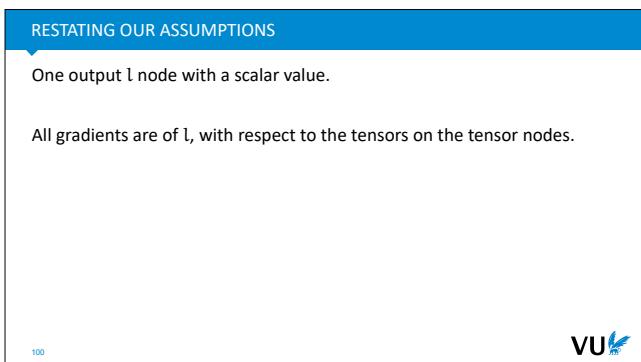
As an example, here is our MLP expressed in our new notation.

Just as before, it's up to us how *granular* we make the computation. We can wrap the whole computation of the loss in a single operation, but we can also separate the matrix multiplication by the **weights** and the addition of the **bias vector** into separate operations.

Note that the *granularity* with which we represent a computation in a computation graph is up to us. We could, for instance also choose to combine the multiplication and the addition into a single operation.

If there's a uniform direction to the computation, we'll leave out the arrowheads for clarity, but you should think of a computation graph as a directional one. All connections have a definite direction.

When we draw intricate computation graphs, we may connect an op node to another opnode, without explicitly drawing the tensor node in between. This should be understood as a shorthand for a proper bipartite



We hold on to the same assumptions we had before. Without these two assumptions things would become a lot more complicated.

IMPLEMENTATIONS

```
TensorNode: —> O
  value: <Tensor>
  gradient: <Tensor>
  source: <OpNode>
```

```
OpNode: > O
  inputs: List<TensorNode>
  outputs: List<TensorNode>
  op: <Op>
```

101



To store a computation graph in memory, we need three objects.

A TensorNode object holds the tensor value at that node. It holds the gradient of the tensor at that node (to be filled by the backpropagation algorithm) and it holds a pointer to the Operation Node that produced it (which can be null for leaf nodes).

An OpNode object represents an instance of a particular

DEFINING AN OPERATION

```
Op: stores anything we need for the backward pass
  f(A) = B
  forward(context, inputs): forward_f : A → B
    # given the inputs, compute the outputs
    ...
  backward(context, outputs_gradient): backward_f : B^∇ → A^∇
    # given the gradient of the loss wrt to the outputs
    # compute the gradient of the loss wrt to the inputs
    ...
```

102



We also need to implement the computation of the operation itself. This is done in two functions (these are usually class functions or static functions).

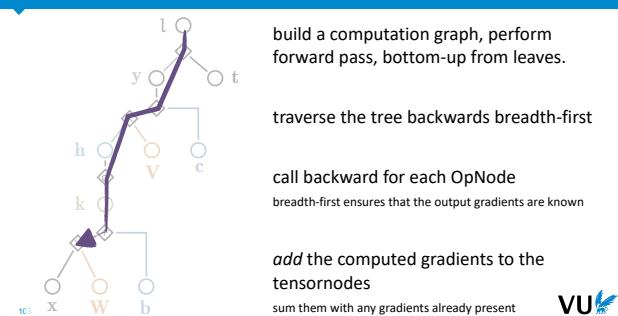
The function forward computes the outputs given the inputs (just like any function in any programming language).

The function backward takes the gradient for the outputs (the gradient of the loss wrt to the outputs) and produces the gradient for the inputs.

Both functions are also given a context object. This is a datastructure (a dictionary or a list) to which the forward can add any value which it needs to save for the backward.

Note that the backward function does not compute the local derivative: it computes the accumulated gradient of the loss over the inputs (given the accumulated gradient of the loss over the outputs).

BACKPROPAGATION



This is the algorithm for backpropagation in this setting. We chain together these modules into a computation graph, and perform a forward pass to compute a loss. We then walk backward from the loss node breadth-first to compute the gradients for each tensor node. Because we walk breadth first, we ensure that we've always computed the gradients for the outputs of each Operation node we encounter already, and we can simply call its backward to compute the gradients for its inputs.

Note the last point: some tensor nodes will be inputs to multiple operation nodes. By the multivariate chain rule, we should sum the gradients they get from all the OpNodes they feed into. We can achieve this easily by just initializing

the gradients to zero and adding the gradients we compute to any that are already there.

BUILDING COMPUTATION GRAPHS IN CODE

Lazy execution: build your graph, compile it, feed data through.

Eager execution: perform forward pass, keep track of computation graph.

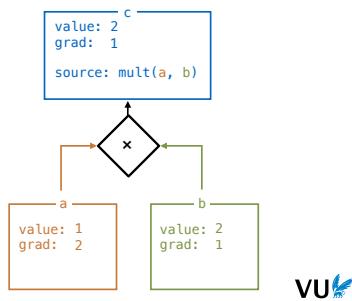
104



There are two common strategies for constructing the computation graph: lazy and eager execution.

EXAMPLE: LAZY EXECUTION

```
a = TensorNode()  
b = TensorNode()  
  
c = a * b  
  
m = Model(  
    in=(a, b),  
    loss=c)  
  
m.train(1, 2))
```



105



Here's one example of how a lazy execution API might look.

Note that when we're building the graph, we're not telling it which values to put at each node. We're just defining the *shape* of the computation, but not performing the computation itself.

When we create the model, we define which nodes are the input nodes, and which node is the loss node. We then provide the input values and perform the forward and backward passes.

BUILDING THE COMPUTATION GRAPH: LAZY EXECUTION

Tensorflow 1.x default, Keras default

Define the computation graph.

- Compile it.
- Iterate backward/forward over the data

Fast. Many possibilities for optimization. Easy to serialise models. Easy to make training parallel.

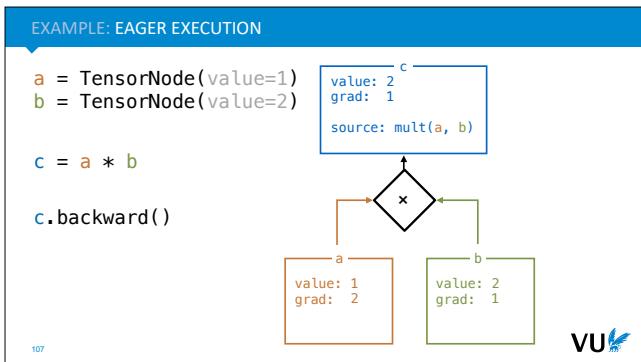
Difficult to debug. Model must remain static during training.

106



In lazy execution, which was the norm during the early years of deep learning, we build the computation graph, but we don't yet specify the values on the tensor nodes. When the computation graph is finished, we define the data, and we feed it through.

This is fast since the deep learning system can optimize the graph structure during compilation, but it makes models hard to debug: if something goes wrong during training, you probably made a mistake while defining your graph, but you will only get an error while passing data through it. The resulting stack trace will never point to the part of your code where you made the mistake.

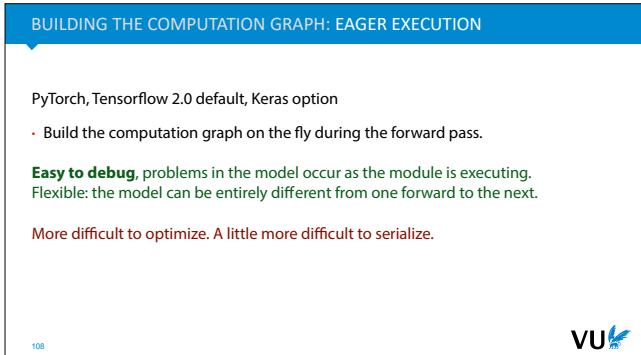


In eager mode deep learning systems, we create a node in our computation graph (a `TensorNode`) by specifying what data it should contain. The result is a tensor object that stores both the data, and the gradient over that data (which will be filled later).

Here we create the variables `a` and `b`. If we now apply an operation to these, for instance to multiply their values, the result is another variable `c`. Languages like python allow us to *overload* the `*` operator it looks like we're just computing multiplication, but behind the scenes, we are creating a computation graph that records all the computations we've done.

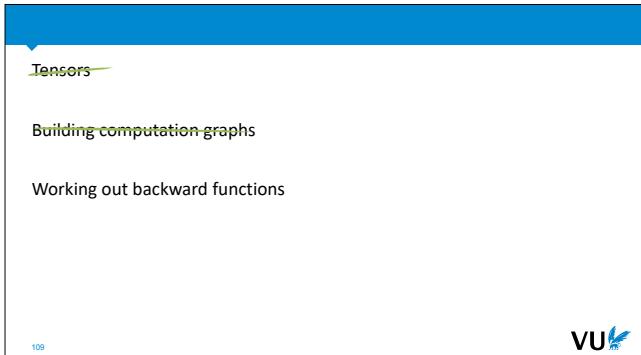
We compute the data stored in `c` by running the computation, but we also store references to the variables that were used to create `c`, and the operation that created it.

Using this graph, we can perform the back propagation from a given node that we designate as *the loss node*. We work our way down the graph computing the derivative of each variable with respect to `c`. At the start the `TensorNodes` do not have their grad's filled in, but at the end of the backward, all gradients have been computed.



In eager execution, we simply execute all operations immediately on the data, and collect the computation graph on the fly. We then execute the backward and ditch the graph we've collected.

This makes debugging much easier, and allows for more flexibility in what we can do in the forward pass. It can, however be a little difficult to wrap your head around. Since we'll be using pytorch later in the course, we'll show you how eager execution works step by step.



WORKING OUT THE BACKWARD FUNCTION

```

class Plus(Op):
    def forward(context, a, b):
        # a, b are matrices of the same size
        return a + b

    def backward(context, goutput):
        return goutput, goutput

```

$$\begin{aligned}
A_{ij}^{\nabla} &= \frac{\partial l}{\partial A_{ij}} \\
&= \sum_{kl} \frac{\partial l}{\partial S_{kl}} \frac{\partial S_{kl}}{\partial A_{ij}} = \sum_{kl} S_{kl}^{\nabla} \frac{\partial S_{kl}}{\partial A_{ij}} \\
&= \sum_{kl} S_{kl}^{\nabla} \frac{\partial [A + B]_{kl}}{\partial A_{ij}} = \sum_{kl} S_{kl}^{\nabla} \frac{\partial A_{kl} + B_{kl}}{\partial A_{ij}} \\
&= S_{ij}^{\nabla} \frac{\partial A_{ij}}{\partial A_{ij}} = S_{ij}^{\nabla}
\end{aligned}$$

$A^{\nabla} = S^{\nabla}$ $B^{\nabla} = S^{\nabla}$

VU

The final ingredient we need is a large collection of operations with backward functions worked out. We'll show how to do this for a few examples.

First, an operation that sums two matrices element-wise.

The recipe is the same as we saw in the last part:

- 1) Write out the scalar derivative for a single element.
- 2) Use the multivariate chain rule to sum over all outputs.
- 3) Vectorize the result.

Note that when we draw the computation graph, we can think of everything that happens between S and l as a single module: we are already given the gradient of l over S , so it doesn't matter if it's one operation or a million.

Again, we can draw the computation graph at any granularity we like: very fine individual operations like summing and multiplying or very coarse-grained operations like entire NNs.

For this operation, the context object is not needed, we can perform the backward pass without remembering anything about the forward pass.

BACKWARD: A FEW MORE EXAMPLES

- Sigmoid
- Row-wise sum
- Expand

VU

To finish up, we'll show you the implementation of some more operations. You'll be asked to do a few of these in the second homework exercise.

SIGMOID

```

class Sigmoid(Op):
    def forward(context, x):
        # x is a tensor of any shape
        sigx = 1 / (1 + (- x).exp())
        context['sigx'] = sigx
        return sigx

    def backward(context, goutput):
        sigx = context['sigx']
        return goutput * sigx * (1 - sigx)

```

$$\begin{aligned}
X_{ijk}^{\nabla} &= \sum_{abc} Y_{abc}^{\nabla} \frac{\partial Y_{abc}}{\partial X_{ijk}} \\
&= \sum_{abc} Y_{abc}^{\nabla} \frac{\partial \sigma(X_{abc})}{\partial X_{ijk}} \\
&= Y_{ijk}^{\nabla} \frac{\partial \sigma(X_{ijk})}{\partial X_{ijk}} \\
&= Y_{ijk}^{\nabla} \sigma(X_{ijk}) (1 - \sigma(X_{ijk}))
\end{aligned}$$

$X^{\nabla} = Y^{\nabla} \otimes \sigma(X) \otimes (1 - \sigma(X))$

VU

This is a pretty simple derivation, but it shows two things:

- 1) We can easily do a backward over functions that output high dimensional tensors, but we should sum over all dimensions of the output when we apply the multivariate chain rule.
- 2) The backward function illustrates the utility of the **context object**. To work out the backward, we need to know the original input. We could just have stored that for every operation, but as we see here, we've done part of the computation already (an in other backwards the inputs aren't necessary). If we give the operation control over what information it

wants to store for the backward pass, we are maximally flexible.

ROW-WISE SUM

```
class RowSum(Op):
    def forward(context, x):
        # x is a matrix
        sumd = x.sum(axis=1)
        context['m'] = x.shape[1]
        return sumd

    def backward(context, gy):
        n, m = gy.shape[0], context['m']
        return gy[None, :].expand(n, m)
```

113

$$\begin{array}{c} \text{Diagram: } \text{x} \text{ (blue circle)} \xrightarrow{\diamond} \text{y} \text{ (orange circle)} \xrightarrow{\diamond} \text{l} \text{ (white circle)} \\ \text{Equation: } \mathbf{x}_{ij}^{\nabla} = \sum_k \mathbf{y}_k^{\nabla} \frac{\partial \mathbf{y}_k}{\partial \mathbf{x}_{ij}} = \sum_k \mathbf{y}_k^{\nabla} \frac{\partial \sum_l \mathbf{x}_{kl}}{\partial \mathbf{x}_{ij}} \\ = \sum_{kl} \mathbf{y}_k^{\nabla} \frac{\partial \mathbf{x}_{kl}}{\partial \mathbf{x}_{ij}} = \mathbf{y}_i^{\nabla} \frac{\partial \mathbf{x}_{ij}}{\partial \mathbf{x}_{ij}} \\ = \mathbf{y}_i^{\nabla} \\ \mathbf{x}^{\nabla} = \mathbf{y}^{\nabla} \mathbf{1}^T \end{array}$$



EXPAND

```
class Expand(Op):
    def forward(context, x, size):
        # x is a scalar
        return np.full(x, size=size)

    def backward(context, gy):
        return gy.sum(), None
```

114

$$\begin{array}{c} \text{Diagram: } \text{x} \text{ (blue circle)} \xrightarrow{\diamond} \text{Y} \text{ (orange circle)} \xrightarrow{\diamond} \text{l} \text{ (white circle)} \\ \text{size } \text{ (blue circle)} \\ \text{Equation: } \mathbf{x}^{\nabla} = \sum_{ab} \mathbf{Y}_{ab}^{\nabla} \frac{\partial \mathbf{Y}_{ab}}{\partial \mathbf{x}} \\ = \sum_{ab} \mathbf{Y}_{ab}^{\nabla} \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \\ = \sum_{ab} \mathbf{Y}_{ab}^{\nabla} \end{array}$$



Note that the output is a vector, because we've summed out one dimension. Therefore, when we apply the multivariate chain rule, we sum over only one dimension.

This is a sum like the earlier example, so we don't need to save any tensor values from the forward pass. However, we do need to remember the size of the dimension we summed out. Therefore, we use the context to store just one object.

The `expand` function we use here is not available in numpy, but it does exist in pytorch. In numpy we can use `repeat`.

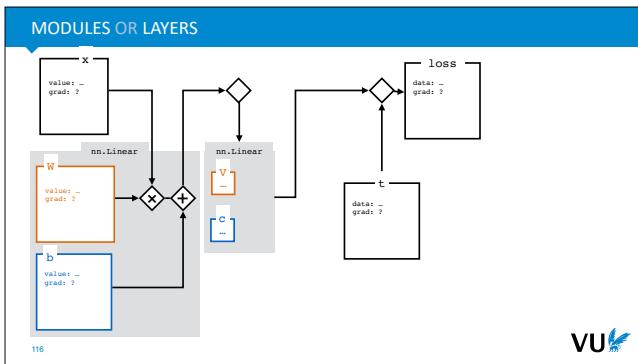
Tensors

[Building computation graphs](#)

[Working out backward functions](#)

115

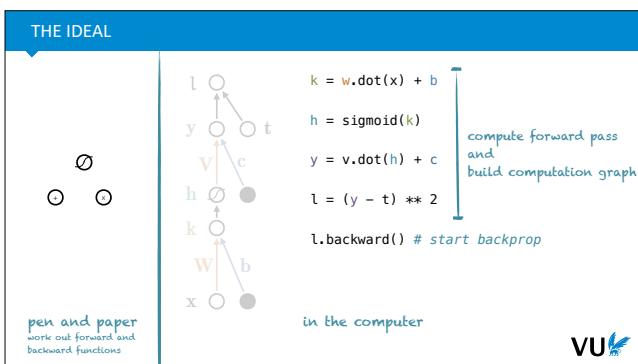




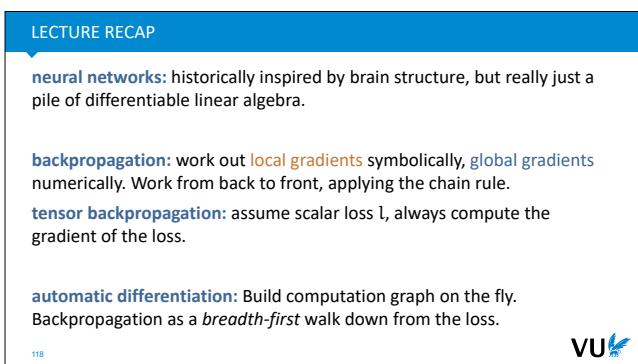
Most deep learning frameworks also have a way of combining model parameters and computation into a single unit, often called a **module** or a **layer**.

In this case a Linear module (as it is called in Pytorch) takes care of implementing the computation of a single layer of a neural network (sans activation) and of remembering the weights and the bias.

Modules have a **forward** function which performs the computation of the layer, but they *don't* have a **backward** function. They just chain together operations, and the backpropagation calls the backwards on the operations. In other words, it doesn't matter to the backpropagation whether we use a module or apply all the operations by hand.



Which completes the picture we wanted to create: a system where all we have to do is perform some computations on some tensors, just like we would do in numpy, and all the building of computation graphs and all the backpropagating is handled for us automatically.



So

THANK YOU FOR YOUR ATTENTION

divu@peterbloem.nl

119