

CS11-747 Neural Networks for NLP

Building a Neural Network Toolkit for NLP

minnn

Graham Neubig



Carnegie Mellon University

Language Technologies Institute

Site

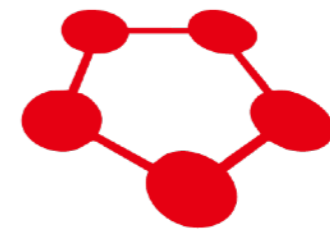
<https://phontron.com/class/nn4nlp2021/>

Neural Network Frameworks

theano

dy/net

Caffe



Chainer

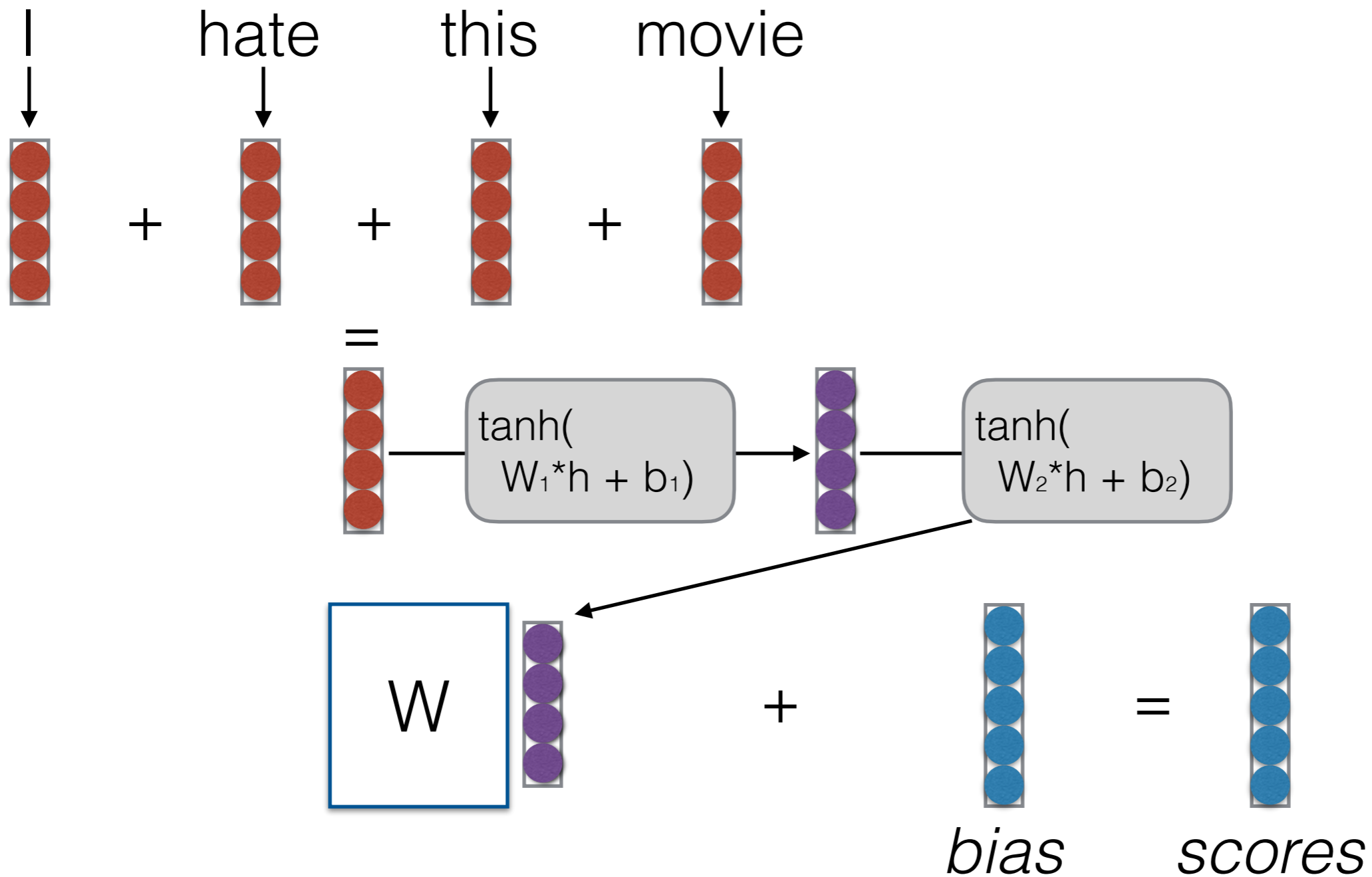


PYTORCH



minnn

Example App: Deep CBOW Model



Algorithm Sketch for NN App Code

- Create a model
- For each example
 - **create a graph** that represents the computation you want
 - **calculate the result** of that computation
 - if training
 - perform **back propagation**
 - **update** parameters

Numerical Computation Backend

- Most neural network libraries use a backend for numerical computation
- **PyTorch/Tensorflow:** MKL, CUDNN, custom-written kernels
- **minnn:** numpy/CuPy

```
import numpy as np
```

```
a = [[1, 0], [0, 1]]
```

```
b = [[4, 1], [2, 2]]
```

```
np.dot(a, b)
```

```
array([[4, 1],  
       [2, 2]])
```

- Many many different operations
- CuPy is a clone of NumPy that works on GPU

Tensors

- An n-dimensional array

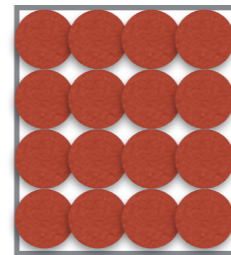
Scalar



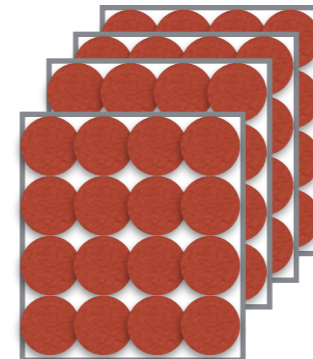
Vector



Matrix



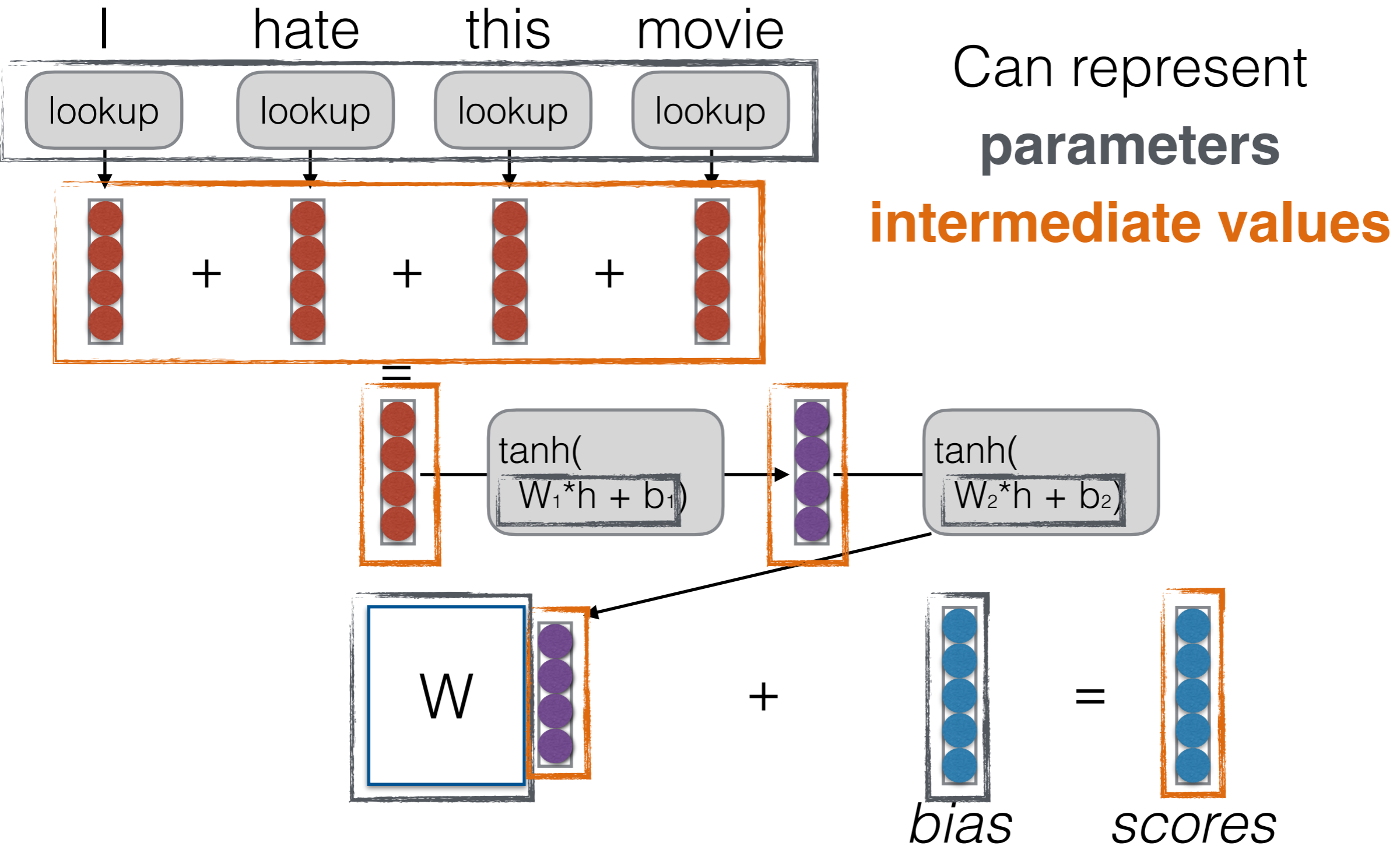
3-dim Tensor



...

- Widely used in neural networks
- Implementation in minnn saves both values and gradients

Tensors in Neural Networks



Tensor Data Structure Definition

Tensor

class Tensor:

def `__init__`(self, data: xp.ndarray):

self.data: xp.ndarray = data

gradient, should be the same size as data

self.grad: Union[Dict[int, xp.ndarray], xp.ndarray] = **None**

generated from which operation?

self.op: Op = **None**

Algorithm Sketch

- Create a model
- For each example
 - **create a graph** that represents the computation you want
 - **calculate the result** of that computation
 - if training
 - perform **back propagation**
 - **update** parameters

Example Model Creation (in App Code)

```
# Define the model
EMB_SIZE = args.emb_size
HID_SIZE = args.hid_size
HID_LAY = args.hid_layer

W_emb = model.add_parameters((nwords, EMB_SIZE))

W_h = [model.add_parameters(
    (HID_SIZE, EMB_SIZE if lay == 0 else HID_SIZE),
    initializer='xavier_uniform')
    for lay in range(HID_LAY)]

W_sm = model.add_parameters((ntags, HID_SIZE),
    initializer='xavier_uniform')
```

Model Class, Adding Parameters

```
# Model: collection of parameters
```

```
class Model:
```

```
    def __init__(self):
```

```
        self.params: List[Parameter] = []
```

```
    def add_parameters(self, shape,
```

```
                        initializer='normal',
```

```
                        **initializer_kwargs):
```

```
        init_f = getattr(Initializer, initializer)
```

```
        data = init_f(shape, **initializer_kwargs)
```

```
        param = Parameter(data)
```

```
        self.params.append(param)
```

```
    return param
```

Parameter Initialization

- Neural nets must have weights that are not identical to learn non-identical features
- **Uniform Initialization:** Initialize weights in some range, such as $[-0.1, 0.1]$ for example
 - *Problem!* Depending on the size of the net, inputs to downstream nodes may be very large
- **Glorot (Xavier) Initialization, He Initialization:** Initialize based on the size of the matrix

Glorot Init: $\sqrt{\frac{6}{d_{in} + d_{out}}}$

NN App Algorithm Sketch

- Create a model
 - For each example
- Greedy Computation
(cf Lazy Computation)

- **create a graph** that represents the computation you want
- **calculate the result** of that computation
- if training
 - perform **back propagation**
 - **update** parameters

Example Graph Creation (in App Code)

```
mn.reset_computation_graph()

emb = mn.lookup(W_emb, words)
h = mn.sum(emb, axis=0)
for W_h_i, b_h_i in zip(W_h, b_h):
    h = mn.tanh(mn.dot(W_h_i, h) + b_h_i)
return mn.dot(W_sm, h) + b_sm
```

Computation Graph

```
class ComputationGraph:
    # global cg
    _cg: 'ComputationGraph' = None

    @classmethod
    def get_cg(cls, reset=False):
        if ComputationGraph._cg is None or reset:
            ComputationGraph._cg = ComputationGraph()
        return ComputationGraph._cg

    def __init__(self):
        self.ops: List[Op] = []

    def reg_op(self, op: Op):
        assert op.idx is None
        op.idx = len(self.ops)
        self.ops.append(op)
```

Operations

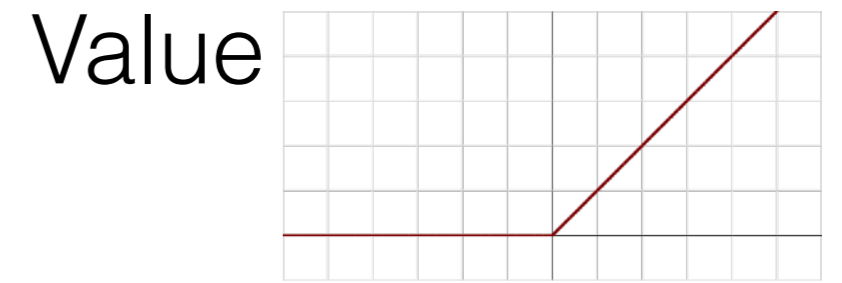
- Operations must know:
- **Forward:** how to calculate their value given input

$$f(\mathbf{u})$$

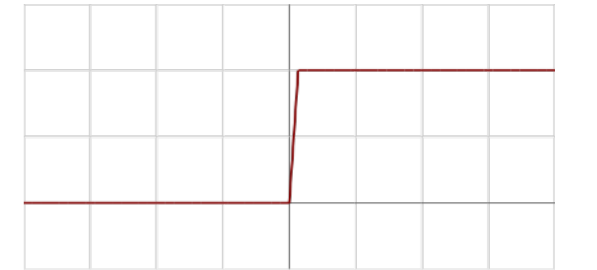
- **Backward:** how to calculate their derivative given following derivative

$$\frac{\partial f(\mathbf{u})}{\partial \mathbf{u}} \quad \frac{\partial \mathcal{F}}{\partial f(\mathbf{u})}$$

Example Op: Relu



Gradient



```
class OpRelu(Op):
    def forward(self, t: Tensor):
        arr_relu = t.data
        arr_relu[arr_relu < 0.0] = 0.0
        t_relu = Tensor(arr_relu)
        self.store_ctx(t=t, t_relu=t_relu, arr_relu=arr_relu)
        return t_relu

    def backward(self):
        t, t_relu, arr_relu = self.get_ctx('t', 't_relu', 'arr_relu')
        if t_relu.grad is not None:
            grad_t = xp.where(arr_relu > 0.0, 1.0, 0.0) * t_relu.grad
            t.accumulate_grad(grad_t)

def relu(param): return OpRelu().full_forward(param)
```

NN App Algorithm Sketch

- Create a model
- For each example
 - **create a graph** that represents the computation you want
 - **calculate the result** of that computation
 - if training
 - perform **back propagation**
 - **update** parameters

Backward Code

```
def backward(t: Tensor, alpha=1.):  
    # first put grad to the start one  
    t.accumulate_grad(alpha)  
    # locate the op  
    op = t.op  
    # backward the whole graph!!  
    cg = ComputationGraph.get_cg()  
    for idx in reversed(range(op.idx+1)):  
        cg.ops[idx].backward()
```

NN App Algorithm Sketch

- Create a model
- For each example
 - **create a graph** that represents the computation you want
 - **calculate the result** of that computation
 - if training
 - perform **back propagation**
 - **update** parameters

Remember: Many Different Update Rules

- **Simple SGD:** update with only gradients
- **Momentum:** update w/ running average of gradient
- **Adagrad:** update downweighting high-variance values
- **Adam:** update w/ running average of gradient, downweighting by running average of variance

SGD Update Rule

```
class SGDTrainer(Trainer):
    def __init__(self, model: Model, lr=0.1):
        super().__init__(model)
        self.lr = lr

    def update(self):
        lr = self.lr
        for p in self.model.params:
            if p.grad is not None:
                if isinstance(p.grad, dict): # sparsely update to save time!
                    self.update_sparse(p, p.grad, lr)
                else:
                    self.update_dense(p, p.grad, lr)
            # clean grad
            p.grad = None

    def update_dense(self, p: Parameter, g: xp.ndarray, lr: float):
        p.data -= lr * g

    def update_sparse(self, p: Parameter,
                      gs: Dict[int, xp.ndarray], lr: float):
        for widx, arr in gs.items():
            p.data[widx] -= lr * arr
```

Still Some Things Left!

- We've left off the details of some underlying parts.
- What about more operations?
- What about more optimizers?
- **Challenge:** can you make a more sophisticated model?

<https://github.com/neubig/minnn-assignment/>

Questions?