# Machine Learning - Exercise 4

## Companion Slides

Ali Athar   Sabarinath Mahadevan
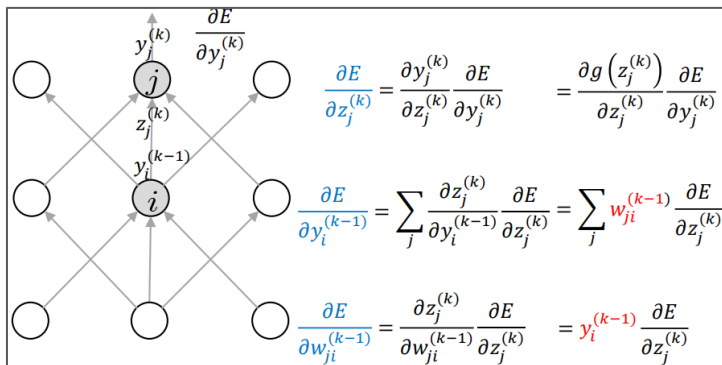
December 6, 2018
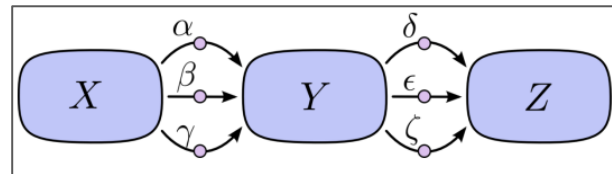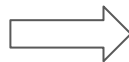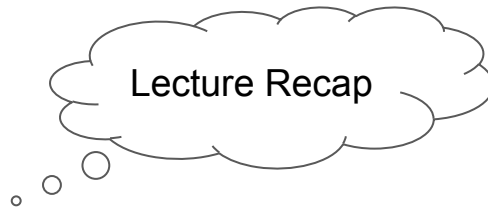
# Exercise Goal



Backpropagation for fixed network



Lecture Recap

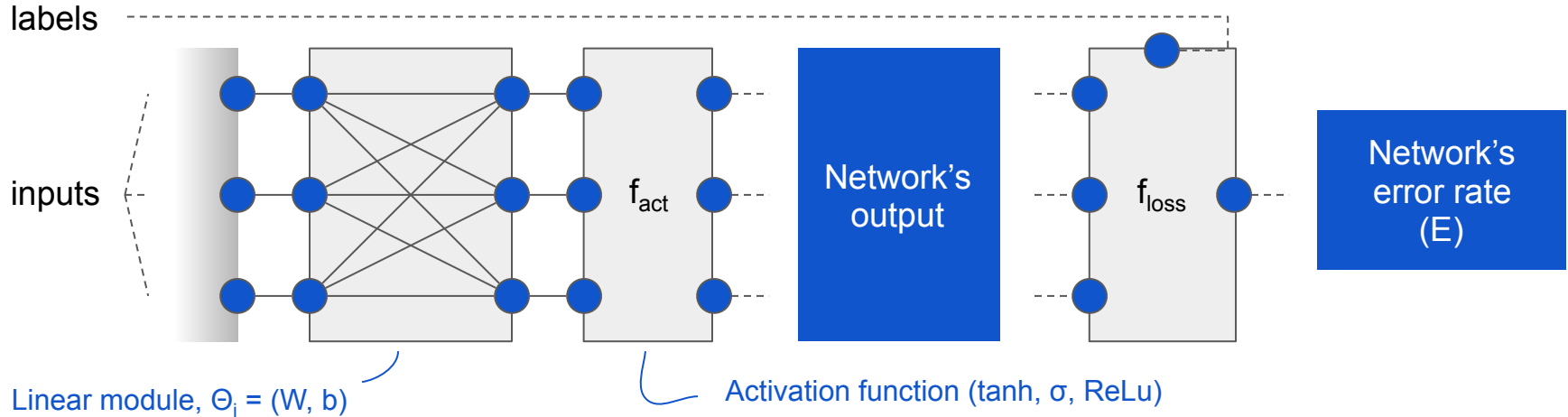General backpropagation with computational graphs

This exercise is about

- Understanding backpropagation, deriving formulas, optimizing them
- Implement simple neural network framework yourself
- Digit recognition

# Recap: Neural Networks



labels

inputs

$f_{act}$

Network's output

$f_{loss}$

Network's error rate (E)

Linear module, $\Theta_i = (W, b)$

Activation function (tanh, σ, ReLu)

Parameters

- Training data (inputs) $X = \{x_i\}_{i=1..N}$ with $x_i \in \mathbb{I}$, N the batch size

- Training labels $T = \{t_i\}_{i=1..N}$ with $x_i \in \mathbb{O}$

- Network is a parametrized, (sub-)differentiable function $F(X,\Theta) : \mathbb{I} \times \mathbb{P} \rightarrow \mathbb{O}$

    - e.g., $\mathbb{O} = \mathbb{R}^{Dim}$ (regression), $\mathbb{O} = [0,1]^{Dim}$ (prob. classification)

- Loss (criterion) $\mathbf{L}(T, F(X,\Theta)) : \mathbb{O} \times \mathbb{O} \rightarrow \mathbb{R}$, put on top of output to measure performance

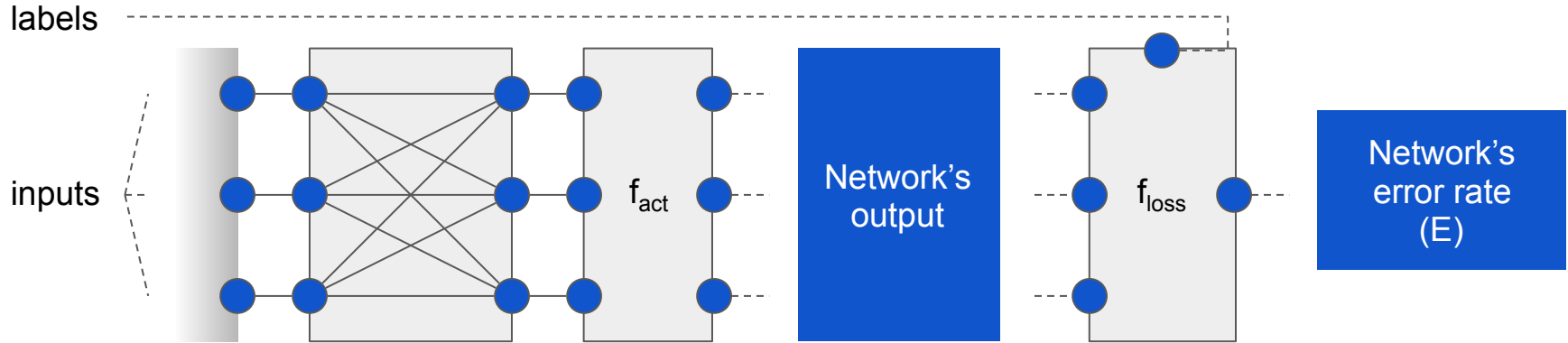    - find optimal parameters: $\Theta^* = \text{argmin}_\Theta \; \mathbf{L}(T, F(X,\Theta))$

# Recap: Backpropagation



▸ Optimize towards lower error rate, i.e., lower E
  ▸ Take derivative of E with respect to each modules parameters, follow gradient
  ▸ Example: Gradient Descent: $\Theta = \Theta - \lambda * D_{\Theta}(E(x))$
  ▸ $D_{\Theta}(E(x)) = D_{\Theta}(E)$ for brevity
▸ How to calculate $D_{\Theta}(E)$
  ▸ Reverse order of modules
  ▸ Module gets $D_{out}(E)$, calculates $D_{\Theta}(E)$, passes $D_{in}(E)$ to next module

Derivative w.r.t. modules parameters $\Theta$ at point x

Learning rate

# Example: Linear/Fully Connected Module



x

$y = W^T \cdot x + b$

$\Theta = (W, b)$

**Element-wise:**

$$y = \left[ \sum_{i=1}^{N_{in}} (w_{ij} \cdot x_i) + b_j \right]_j$$

⚠️ **Without batching**

Given: Derivative with respect to output $\frac{dE(x)}{dy}$

Calculate:

▶ Derivatives with respect to parameters $\Theta$

$$\frac{dE(x)}{dw_{ij}} = \frac{dE(x)}{dy} \cdot \frac{dy}{dw_{ij}} = \frac{dE(x)}{dy_j} \cdot x_i = \left[ x \cdot \frac{dE}{dy} \right]_{ij}$$

$$\frac{dE(x)}{db} = \frac{dE(x)}{dy} \cdot \frac{dy}{db} = \frac{dE(x)}{dy} \cdot I = \frac{dE(x)}{dy}$$

# Example: Linear/Fully Connected Module



x - - - - y = W$^T$·x + b

Θ = (W, b)

Element-wise:

$$y = \left[ \sum_{i=1}^{N_{in}} (w_{ij} \cdot x_i) + b_j \right]_j$$

⚠ ! Without batching

Given: Derivative with respect to output $\frac{dE(x)}{dy}$

Calculate:

▸ Derivative with respect to input

$$\frac{dE(x)}{dx} = \frac{dE(x)}{dy} \cdot \frac{dy}{dx} = \frac{dE(x)}{dy} \cdot W^T$$

# Example: Linear/Fully Connected Module



$x$ ----- $y = W^T \cdot x + b$

$\Theta = (W, b)$

Without batching

Putting it together:

```
fprop(x):
    cache.x = x
    return WT*x + b
```
run training data through
(forwards)

```
bprop(dE):
    dW = cache.x * dE
    db = dE
    return dE * W
```
run gradients through
(backwards)

```
update(rate):
    W = W – rate*dW
    b = b – rate*dbT
```
update the parameters
(grad. descent)

# Mini-Batching

▶ Batch learning
  ▸ All training samples processed at once, parameters updated once at the end
  ▸ Stable, well understood, many acceleration techniques, but slow

▶ Stochastic learning
  ▸ Each training sample separately, parameters updated at each step
  ▸ Noisy (though may lead to better results), fast

▶ Mini-batching
  ▸ Middle ground, batches of data processed, bundled updates
  ▸ Combine advantages, reduce drawbacks

▶ Example
  ▸ Linear Module $f$ with input dimension $N_{in}$ and output dimension $N_{out}$, batch size $n$

$$f(x) = W^T \cdot \begin{bmatrix} | & | & & | \\ x^1 & x^2 & \dots & x^n \\ | & | & & | \end{bmatrix} + \begin{bmatrix} | & | & & | \\ b & b & \dots & b \\ | & | & & | \end{bmatrix}$$

broadcast (i.e., repeat) $b$

mini-batch matrix

# Batching Update Rule

- (Mini-)Batch learning
    - Multiple samples processed at once
    - Calculate gradient for each sample, but don't update the parameters
    - After processing the batch, update using a sum of all gradients
    - Learning rate has to be adapted, e.g., divide E by batch size

    - Example: Gradient Descent
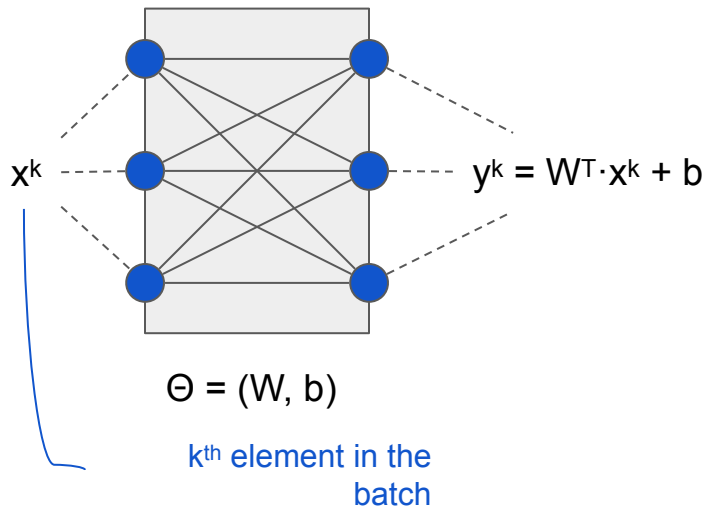    $$\Theta = \Theta - \lambda \cdot \sum_{k=1}^{N} D_\Theta(E(x^k))$$

    Derivative of E w.r.t parameters $\Theta$ at point $x^k$

    - To make things easier, we write
    $$\sum_{k=1}^{N} D_\Theta(E(x^k)) = D_\Theta(E(x)) = D_\Theta(E)$$

# Example: Linear/Fully Connected Module - Batching



$x^k$     $y^k = W^T \cdot x^k + b$

$\Theta = (W, b)$

$k^{th}$ element in the batch

Given: Derivatives with respect to outputs $\frac{dE(x^k)}{dy^k}$

Plural!

Calculate:

▸ Derivatives with respect to parameters $\Theta$

$$\frac{dE(x)}{dw_{ij}} = \sum_{k=1}^{N} \left( \frac{dE(x^k)}{dy^k} \cdot \frac{dy^k}{dw_{ij}} \right) = \left[ \sum_{k=1}^{N} \left( x^k \cdot \frac{dE(x^k)}{dy^k} \right) \right]_{ij}$$

$$= \left[ x \cdot \frac{dE(x)}{dy} \right]_{ij}$$

$$\frac{dE(x)}{db} = \sum_{k=1}^{N} \left( \frac{dE(x^k)}{dy^k} \cdot \frac{dy^k}{db} \right) = \sum_{k=1}^{N} \left( \frac{dE(x^k)}{dy^k} \right)$$

$$= sum_{col-wise} \left( \frac{dE(x)}{dy} \right)$$

Deriv. w.r.t outputs assumed to be given row-wise:

$$\frac{dE(x)}{dy} = \begin{bmatrix} - & \frac{dE(x^1)}{dy^1} & - \\ & \ldots & \\ - & \frac{dE(x^n)}{dy^n} & - \end{bmatrix}$$

Visual Computing Institute

RWTH AACHEN UNIVERSITY

# Example: Linear/Fully Connected Module - Batching



$x^k$

$y^k = W^T \cdot x^k + b$

$\Theta = (W, b)$

$k^{th}$ element in the batch

Given: Derivatives with respect to outputs $\dfrac{dE(x^k)}{dy^k}$

Plural!

Calculate:

▸ Derivatives with respect to inputs

$$\frac{dE(x)}{dx} = \left[ \frac{dE(x^k)}{dx^k} \right]_k = \left[ \frac{dE(x^k)}{dy^k} \cdot \frac{dy^k}{dx^k} \right]_k$$

$$= \left[ \frac{dE(x^k)}{dy^k} \cdot W^T \right]_k = \frac{dE(x)}{dy} \cdot W^T$$

Deriv. w.r.t outputs assumed to be given row-wise:

$$\frac{dE(x)}{dy} = \begin{bmatrix} - & \frac{dE(x^1)}{dy^1} & - \\ & \ldots & \\ - & \frac{dE(x^n)}{dy^n} & - \end{bmatrix}$$

# Example: Training a Network

```
1.   network = [module_1, module_2, …, module_n], loss = f_loss
2.
3.   for X, T in batched(inputs, labels) do
4.           z = X
5.           for module in net do
6.                   z = module.fprop(z)
7.           end for
8.           E = loss.fprop(z,T)
9.           dz = loss.bprop(1/batchSize)        // Normalization for batch size
10.          for module in reversed(net) do
11.                  dz = module.brop(dz)
12.          end for
13.          for module in net do
14.                  module.update(rate)
15.          end for
16.  end for
```

# Debugging Tip: Gradient Checking

Check the Jacobian *J* from *bprop* with numerical differentiation

▸ Numerical approach: Column-wise (here for the first column)

$$x_+ = (x_1 + h_1, x_2, ..., x_n)$$
$$x_- = (x_1 - h_1, x_2, ..., x_n)$$

$$J_{-,1} = \frac{fprop(x_+) - fprop(x_-)}{2h_1}$$

▸ Backprop: Row-wise (here for the first row)

$$fprop(x)$$
$$J_{1,-} = bprop(1, 0, ..., 0)$$

▸ Advice

  ▸ Use (small) random *x*
  $$h_i = \sqrt{\epsilon} * max(x_i, 1)$$

# Expected Results/Tips for MNIST

- [Linear(28x28, 10), Softmax]
  - should give ± 750 errors
- [Linear(28x28, 200), tanh, Linear(200,10), Softmax]
  - should give ± 250 errors
- Typical learning rates
  - $\lambda \in [0.1, 0.01]$
- Typical batch sizes
  - $N_B \in [100, 1000]$
- Weight initialization
  - $W \in \mathbb{R}^{M \times N}$
  - $W \sim N \ (0, \sqrt{\frac{2}{M+N}})$, i.e., sampled from normal distribution around 0 with deviation $\sqrt{\frac{2}{M+N}}$
  - $b = 0$
- Pre-process the data
  - Dividide values by 255 (= max pixel value)

Visual Computing Institute

RWTH AACHEN UNIVERSITY