# Midterm and automatic differentiation (aka autograd)

## Lecture 06 — CS 577 Deep Learning

Instructor: Yutong Wang

Computer Science
Illinois Institute of Technology

September 25, 2024

# Topics

- The midterm
- Attention models
- Tricks for doing homework 3.
  - `np.matmul` with multidim arrays
  - Batched matrix multiplication
  - Batched diagonalization
  - Batched outer product
  - `np.moveaxis`
- Automatic differentiation

Midterm ↑

exam ↓

# More on homework 3: Next-word-prediction

$\varphi$ word embedding   $\varphi(over)$   $\varphi(the) \in \mathbb{R}^d$

The quick brown fox jumps over the ____.

$$\mathbf{X}^{(i)} = \begin{bmatrix} \mathbf{x}^{(i,1)} & \mathbf{x}^{(i,2)} & \cdots & \mathbf{x}^{(i,C-1)} & \mathbf{x}^{(i,C)} \end{bmatrix} \in \mathbb{R}^{d \times C} \quad \text{is a } d \times C \text{ matrix}$$

a single data point   (also called a prompt)

the

$$y^{(i)} \in \{1, \ldots, K\} \quad \text{is set of all candidate words}$$

Vocabulary

# More on homework 3: Next-word-prediction

The quick brown fox jumps over the ____.

$$\mathbf{X}^{(i)} = \begin{bmatrix} \mathbf{x}^{(i,1)} & \mathbf{x}^{(i,2)} & \cdots & \mathbf{x}^{(i,C-1)} & \mathbf{x}^{(i,C)} \end{bmatrix} \in \mathbb{R}^{d \times C} \quad \text{is a } d \times C \text{ matrix}$$

the

should have the most important

empty

$$y^{(i)} \in \{1, \ldots, K\} \quad \text{is set of all candidate words}$$

# More on homework 3: Next-word-prediction

The quick brown fox jumps over the ____.

$$\mathbf{X}^{(i)} = \begin{bmatrix} \mathbf{x}^{(i,1)} & \mathbf{x}^{(i,2)} & \cdots & \mathbf{x}^{(i,C-1)} & \mathbf{x}^{(i,C)} \end{bmatrix} \in \mathbb{R}^{d \times C} \quad \text{is a } d \times C \text{ matrix}$$

- unrealistic
- good for learning + writing paper

$y^{(i)} \in \{\pm 1\}$     we are keeping it simple

# Notations

- $C$ as `n_context` — dim of word embed 700-ish
- $d$ as `n_features` — training
- $n$ as `n_samples`
- $q$ as `n_reduced`, where $q < d$ — projection onto small space of dim $q$
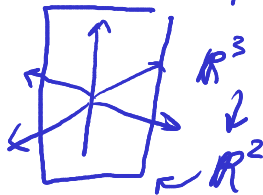
Let

$$\theta = [W^{(1)}, W^{(2)}, w^{(3)}]$$

where

$$W^{(1)} \text{ and } W^{(2)} \in \mathbb{R}^{q \times d}$$

and

$$w^{(3)} \in \mathbb{R}^{d}$$

and

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \,\text{softmax}\left( X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)} \right).$$
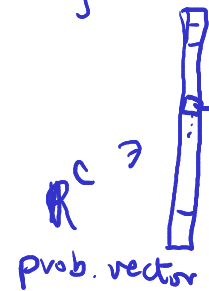
attention

$\mathbb{R}^3$

$\mathbb{R}^2$

subspace of some concept.

# ("Single-layer" and "single-headed") Attention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \, \text{softmax}\left( X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)} \right).$$
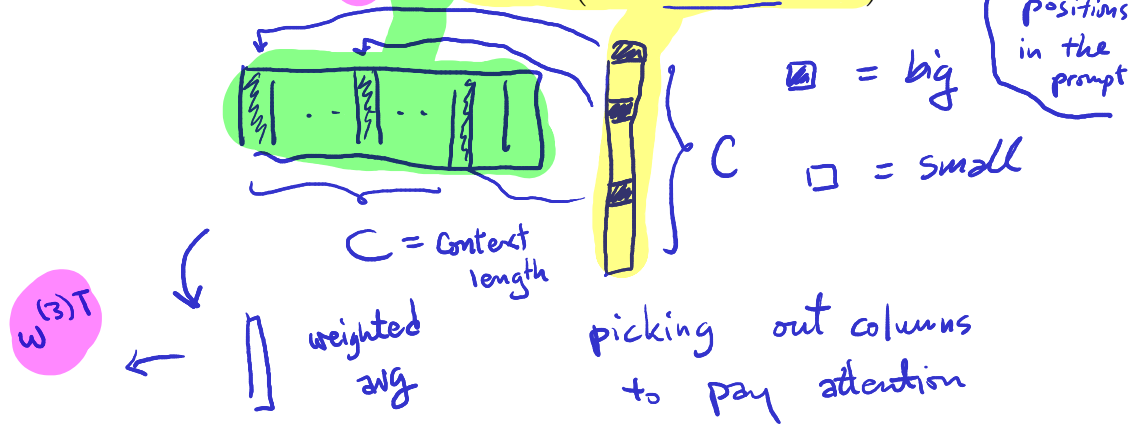
$$[x^{(i,1)} \quad \cdots\cdots \quad x^{(i,C)}]$$

$\propto \exp\left( x^{(i,j)\top} W^{(2)\top} W^{(1)} x^{(i,C)} \right)$

prop to

$j$

$\mathbb{R}^C \ni$

prob. vector

$\sum$ to 1 and positive

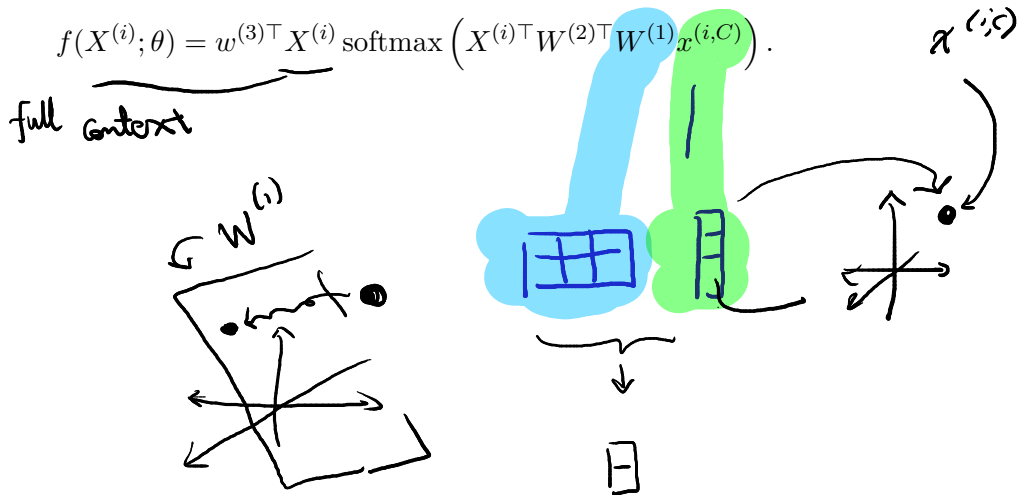Similarity between $j$-th word emb & $C$-th word emb

◦ Note softmax playing totally different role compared to CE.
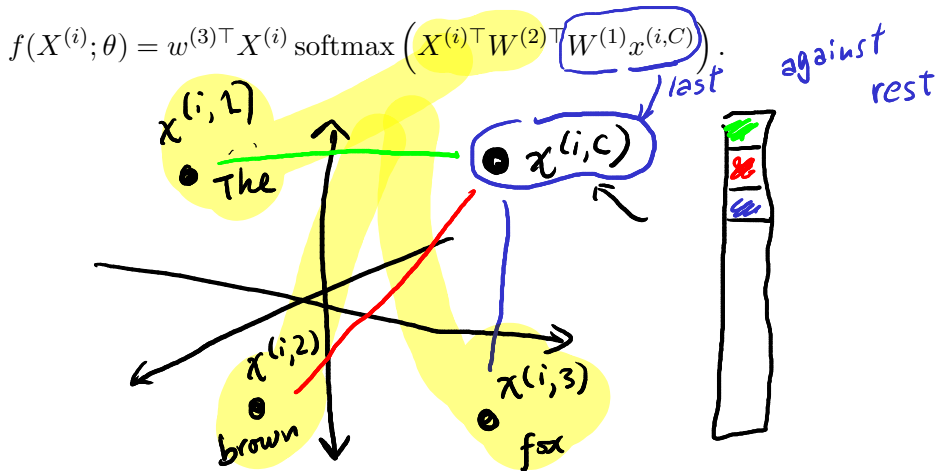
# ("Single-layer" and "single-headed") Attention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left(X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)}\right).$$

Learning to pick positions in the prompt

▨ = big

□ = small

$C$

$C$ = context length

$w^{(3)T}$

weighted avg

picking out columns to pay attention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left( X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)} \right).$$

full context

$x^{(i,c)}$

$W^{(1)}$

# ("Single-layer" and "single-headed") Attention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left(X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)}\right).$$
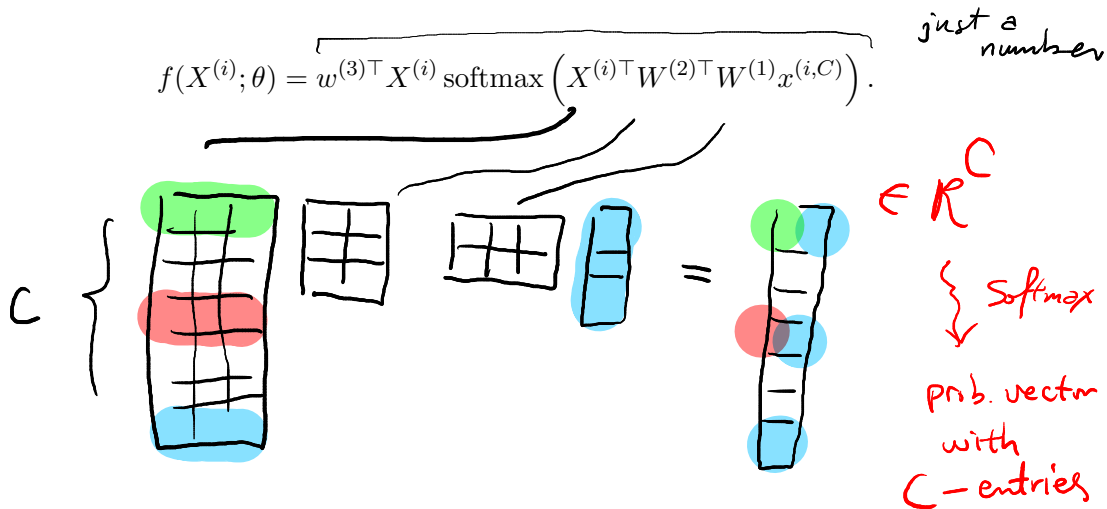
# ("Single-layer" and "single-headed") Attention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left( X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)} \right).$$

just a number



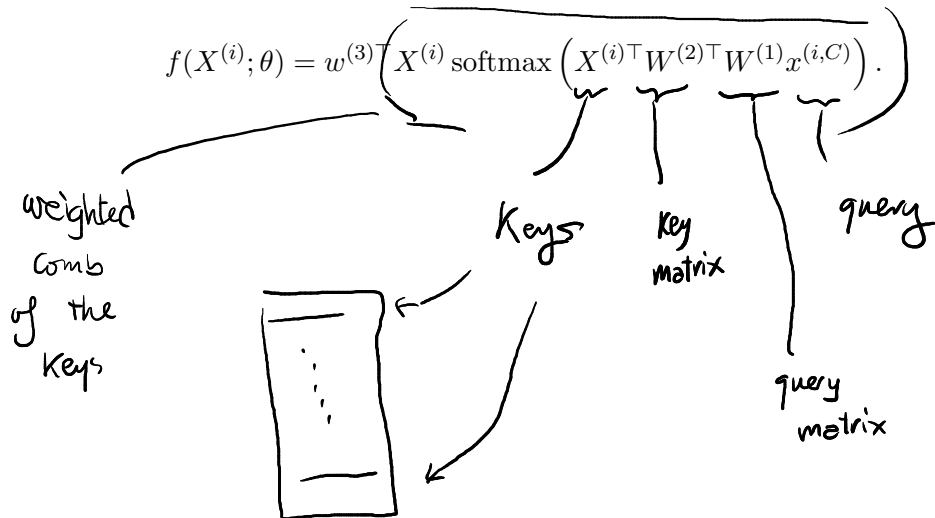$\in \mathbb{R}^C$

{ softmax

prob. vector with C-entries

# ("Single-layer" and "single-headed") Attention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left( X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)} \right).$$

# ("Single-layer" and "single-headed") Attention

$$f(X^{(i)}; \theta) = w^{(3)\top} \left( X^{(i)} \operatorname{softmax} \left( X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)} \right) \right).$$

weighted
comb
of the
keys

Keys

key
matrix

query

query
matrix

# ("Single-layer" and "single-headed") Attention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left(X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)}\right).$$
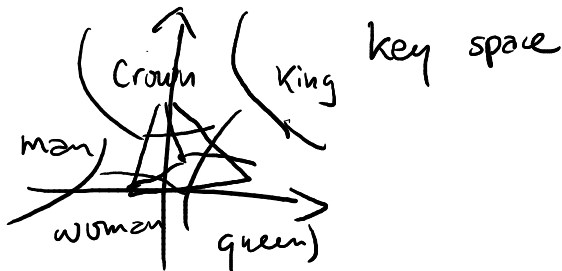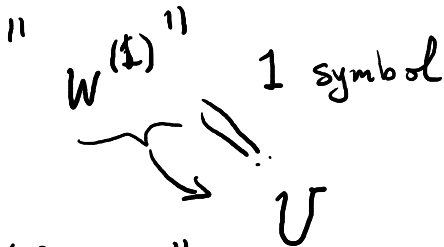
$$f(X^{(i)}; \theta) = \underline{w^{(3)\top}} X^{(i)} \, \mathrm{softmax}\left( X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)} \right).$$

"

$W^{(1)}$ " 1 symbol

$U$

Reason for "(1), (2)" is to reuse W for weights

# ("Single-layer" and "single-headed") Attention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left(X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)}\right).$$ ← typo

bold **W** capital stands multi-dim arr
including matrix & tensor

bold **w** lower case = vector | In this class

non-bold $w$ l.c. = scalar | Goodfellow convention

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left(X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)}\right).$$
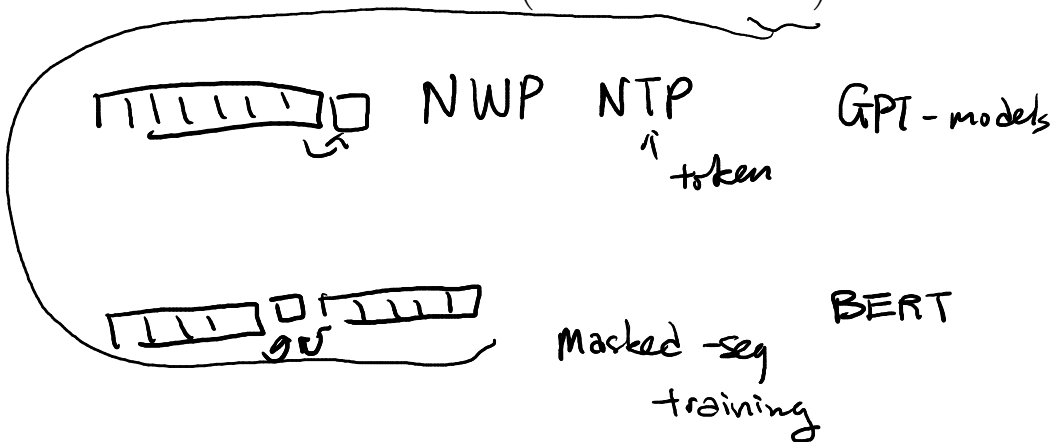
$\searrow$ sign

$C$ — dot products

$\{\pm 1\}$

$+1 = \text{foo}$

$-1 = \text{bar}$

$\{1, \cdots\cdots, K\}$ $\searrow$ argmax

$1 = \text{the}$

$2 = \text{an}$

$3 = \text{and}$

$$f(X^{(i)}; \theta) = w^{(3)\top} X^{(i)} \operatorname{softmax}\left(X^{(i)\top} W^{(2)\top} W^{(1)} x^{(i,C)}\right).$$



NWP NTP

token

GPT - models

Masked - seq
training

BERT

# Tricks for doing homework 3.

- `np.matmul` with multidim arrays
- Batched matrix multiplication
- Batched diagonalization
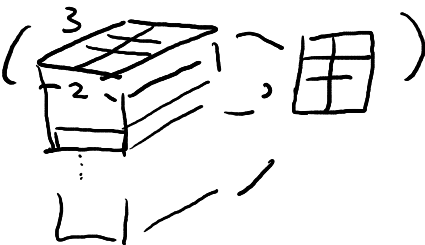- Batched outer product
- `np.moveaxis`

# `np.matmul` with multidim arrays

X

rank 3 tensor

$W^{(2)}$  rank 2 tensor  aka matrix

```
1  # IN (n_samples, n_context, n_features) @ (n_features, n_reduced)
2  Keys = np.matmul(X, W2)
3  # OUT  (n_samples, n_context, n_reduced)
```

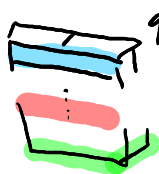np. matmul ( [cube diagram 3, 2, 1] , [grid] )

a stack of matmul

# Batched ~~matrix~~ multiplication

*matrix -vector*

Keys

queries.shape

```
1  queries.shape == (n_samples, n_reduced)
2
3  # (n_samples, n_context, n_reduced)
4  np.sum(
5      # (n_samples, n_context, n_reduced) * (n_samples, 1, n_reduced)
6      Keys * np.expand_dims(queries, axis=1),
7      # --> (n_samples, n_context, n_reduced)
8      axis=2)
9  # --> (n_samples, n_context)
```

queries

sum(axis =2)

match the
0th dim
and do
stackwise matmul

keys ²    c

key    queries

# Batched diagonalization



```
1  softmaxKq.shape == (n_samples, n_context)
2
3  # (n_samples, 1, n_context) * (n_context, n_context)
4  softmaxKq_diag = # YOUR CODE HERE
5  # --> (n_context, n_context, n_context)
```
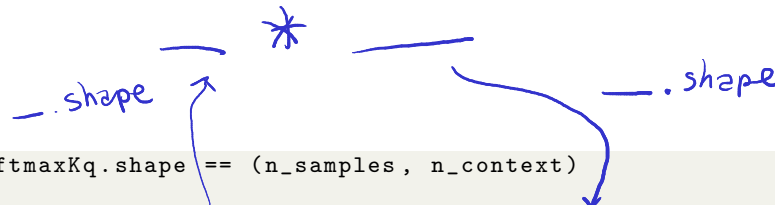
# Batched diagonalization

_.shape

*

_.shape

```
1  softmaxKq.shape == (n_samples, n_context)
2
3  # (n_samples, 1, n_context) * (n_context, n_context)
4  softmaxKq_diag = # YOUR CODE HERE
5  # --> (n_context, n_context, n_context)
```

# Batched outer product

softmaxKq

(n samples, ncontext)

$j$-th entry

$\Pi_i P_j$

$P_c$

$P_i$

$P_1 \ldots P_c$

$i$-th entry

```python
# (n_samples, 1, n_context) * (n_samples, n_context, 1)
softmaxKq_outer = # YOUR CODE HERE
# --> (n_samples, n_context, n_context)
```

# Batched outer product



expand_dims( axis = 1 )

flip

*

```
1 # (n_samples, 1, n_context) * (n_samples, n_context, 1)
2 softmaxKq_outer = # YOUR CODE HERE
3 # --> (n_samples, n_context, n_context)
```

n-context

n-context

*

# `np.moveaxis`

*X^(i)* stacked     not touching D

```
1  X.shape == (n_samples, n_context, n_features)
2
3  # (n_samples, n_features, n_context) @ (n_samples, n_context, n_context)
4  A = np.moveaxis(X,1,2) @ D
5  # --> (n_samples, n_features, n_context)
```

Want batchwise matmul

↳ transpose each stack

# Welcome to the second half the course

- Understand how automatic differentiation works under the hood
- What is overfitting, when it happens, and what can you do about it
- Transfer learning
- Convolution layers, attention
- How to compute things in deep learning really really fast.

# Automatic differentiation

Let's compute

$$f(x, y) = \log(x^2 + 1) \exp(y) + y^2 x$$

$$x, y \in \mathbb{R}$$

using automatic differentiation.

**Eval** $f(x, y) = \log(x^2 + 1)\exp(y) + y^2 x$ **at** $(x, y) = (2, 3)$

$x$

$1$

$y$

$y^2$

$x^2$

$x^2 + 1$

$\exp(y)$

$y^2 \cdot x$

$\log(x^2 + 1)$

$\log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$\log(x^2 + 1) \cdot \exp(y)$

**Eval** $f(x, y) = \log(x^2 + 1) \exp(y) + y^2 x$ **at** $(x, y) = (2, 3)$

**Eval** $f(x, y) = \log(x^2 + 1)\exp(y) + y^2 x$ **at** $(x, y) = (2, 3)$



$x$

$1$

$y$

$y^2$

$x^2$

$x^2 + 1$   5

$\log(5)$

$\exp(y)$

$y^2 \cdot x$

$\log(x^2 + 1)$

$\log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$\log(x^2 + 1) \cdot \exp(y)$

**Eval** $f(x, y) = \log(x^2 + 1)\exp(y) + y^2 x$ **at** $(x, y) = (2, 3)$

**Eval** $f(x, y) = \log(x^2 + 1)\exp(y) + y^2 x$ **at** $(x, y) = (2, 3)$

**Eval** $f(x, y) = \log(x^2 + 1)\exp(y) + y^2 x$ **at** $(x, y) = (2, 3)$



$x$

$1$

$y$

$y^2$

$x^2$

$\exp(y)$

$y^2 \cdot x$

DFS search

$x^2 + 1$

$\log(x^2 + 1)$

$\log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$\log(x^2 + 1) \cdot \exp(y)$

**Eval** $\nabla f(x, y) = \log(x^2 + 1) \exp(y) + y^2 x$ **at** $(2, 3)$

# Eval $\nabla f(x,y) = \log(x^2 + 1)\exp(y) + y^2 x$ at $(2,3)$



$\mathcal{I}(c) = c$

$\mathcal{I}(a+b)$

$\dfrac{\partial \mathcal{I}}{\partial c}(a+b)\nabla a$

$x$

$y$

$1$

pow

$y^2$

pow

$x^2$

exp

mul

add

$\exp(y)$

$y^2 \cdot x$

$x^2 + 1$

log

add

$\log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$\log(x^2 + 1)$

mul

$a$

$b$

$\log(x^2 + 1) \cdot \exp(y)$

# Eval $\nabla f(x, y) = \log(x^2 + 1)\exp(y) + y^2 x$ at $(2, 3)$



$x$

$1$

$y$ $\xrightarrow{\texttt{pow}}$ $y^2$

$x \xrightarrow{\texttt{pow}} x^2$

$y \xrightarrow{\texttt{exp}} \exp(y)$

$y^2 \xrightarrow{\texttt{mul}} y^2 \cdot x$

$x^2 \xrightarrow{\texttt{add}} x^2 + 1$

$x^2 + 1 \xrightarrow{\texttt{log}} \log(x^2 + 1)$

$\log(x^2 + 1) \xrightarrow{\texttt{mul}} \log(x^2 + 1) \cdot \exp(y)$

$\log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$\nabla f$

$f(a, b) = I(a, b)$

$\nabla_a f = \frac{\partial I}{\partial c}(a + b)\nabla_a$

$\underbrace{\quad}_{1}$

$\nabla_b$

$\nabla_a$

**Eval $\nabla f(x, y) = \log(x^2 + 1)\exp(y) + y^2 x$ at $(2, 3)$**



$x$

$1$

$y$ — pow → $y^2$

$x$ — pow → $x^2$

$x^2$ — add → $x^2 + 1$

$1$ — add → $x^2 + 1$

$y$ — exp → $\exp(y)$

$y^2$ — mul → $y^2 \cdot x$

$x^2 + 1$ — log → $\log(x^2 + 1)$

$\log(x^2 + 1)$ — mul → $\log(x^2 + 1) \cdot \exp(y)$

$\exp(y)$ — mul → $\log(x^2 + 1) \cdot \exp(y)$

$\log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$ — add

$y^2 \cdot x$

$\frac{\partial I}{\partial a}(a, b)$

$\frac{\partial I}{\partial c}\frac{\partial}{\partial a}(a, b)$

$I(a, b)$

$\underbrace{\log(x^2 + 1)}_{a} \cdot \underbrace{\exp(y)}_{b}$

$\frac{\partial I}{\partial c} \cdot b$

**Eval** $\nabla f(x, y) = \log(x^2 + 1) \exp(y) + y^2 x$ **at** $(2, 3)$



$x$

$1$

$y$

pow

$y^2$

pow

$x^2$

add

$x^2 + 1$

exp

$\exp(y)$

log

$\log(x^2 + 1)$

mul

$y^2 \cdot x$

mul

$\dfrac{\partial I}{\partial b}(a \cdot b)$

$= \dfrac{\partial I}{\partial c}(a \cdot b)$

$\dfrac{\partial}{\partial b}(a \cdot b)$

add

$\log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$\dfrac{\partial I}{\partial b}$

$I(a, b)$

$\log(x^2 + 1) \cdot \exp(y)$

$a \quad b$

$\dfrac{\partial I}{\partial a}$

# An autograd "module"

```
1  class ag: # AutoGrad
2      class Scalar: # Scalars with grads
3          def __init__(self, value, op="", _backward= lambda : None,
       inputs=[], label=""):
4
5              self.value = float(value)
6              self.grad = 0.0
7
8              # ... lines skipped
```

# An autograd "module"

```python
class ag: # AutoGrad
    class Scalar: # Scalars with grads
        def __init__(self,  value, op="", _backward= lambda : None,
    inputs=[], label=""):

            self.value = float(value)
            self.grad = 0.0

            # ... lines skipped

        def __add__(self, other): # ...
        def __mul__(self, other): # ...
        def __pow__(self, exponent): # ...

    def exp(input): # ...
    def log(input): # ...
```
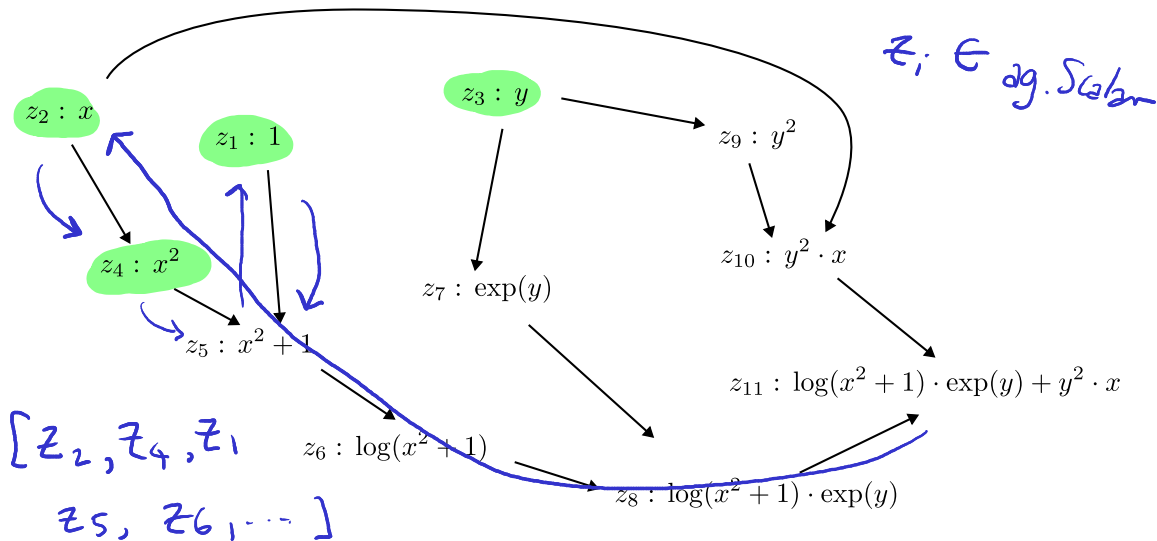
# Implement $f$ in our framework

```
1 x = ag.Scalar(2, label="z2:x")    # wrap x
2 y = ag.Scalar(3, label="z3:y")    # wrap y
3
4 # implement log(x^2+1)*exp(y) + y^2*x
5 def f(x,y):
6     z1 = ag.Scalar(1,label= "z1:1")    # constant 1
7     z2 = x
8     z3 = y
9
10    z4 = z2**2
11    z4.label = "z4:x^2"
12
13    # ... lines skipped
```

along left first

$z_i \in$ ag. Scalar

$z_2 : x$

$z_1 : 1$

$z_3 : y$

$z_9 : y^2$

$z_4 : x^2$

$z_7 : \exp(y)$

$z_{10} : y^2 \cdot x$

$z_5 : x^2 + 1$

$z_6 : \log(x^2 + 1)$

$z_{11} : \log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$z_8 : \log(x^2 + 1) \cdot \exp(y)$

$[z_2, z_4, z_1$

$z_5, z_6, \cdots]$

# Forward pass — DFS — computation graph



$z_2 : x$

$z_1 : 1$

$z_3 : y$

$z_9 : y^2$

$z_4 : x^2$

$z_5 : x^2 + 1$

$z_7 : \exp(y)$

$z_{10} : y^2 \cdot x$

$z_{11}$ inputs
$= [z_8, z_{10}]$

$z_6 : \log(x^2 + 1)$

$z_{11} : \log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$z_8 : \log(x^2 + 1) \cdot \exp(y)$

$y^2$ 1

$z_2 : x$ 2

exp(3) $\frac{1}{x^2+1}$ · 2x

$z_1 : 1$ 1

1 exp(3) $\frac{1}{x^2+1}$

$z_3 : y$ 3

x 2y + 1·log(5) exp(3)

$z_9 : y^2$ 9

x

call backward here

1

4

$z_4 : x^2$

1 exp(3) $\frac{1}{x^2+1}$

5

exp(3)

$z_7 : \exp(y)$

$z_{10} : y^2 \cdot x$ 18

grad = 1

$z_5 : x^2 + 1$

1 exp(3) · $\frac{1}{x^2+1}$

log(5)

1·log(5)

$z_{11} : \log(x^2 + 1) \cdot \exp(y) + y^2 \cdot x$

$z_6 : \log(x^2 + 1)$

1 · exp(3)

1

$z_8 : \log(x^2 + 1) \cdot \exp(y)$

log(5) exp(3)

## Exercise 2.A

Implement `ag.Scalar.topological_sort`
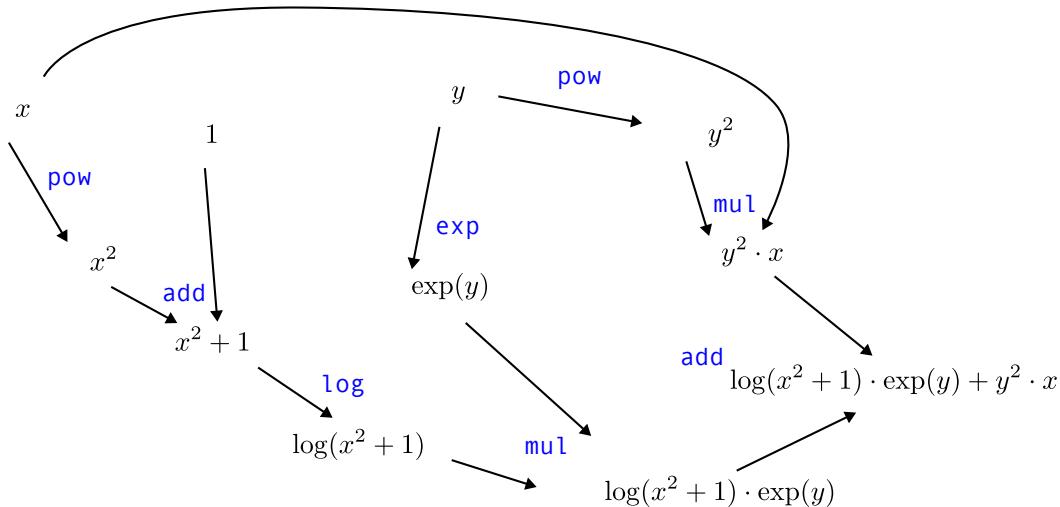
```
1        def topological_sort(self):
2            ## EXERCISE 2.A
3            ## YOUR CODE HERE
```

```
1 Check your answer
2 [z2:x,
3  z4:x^2     pow(2),
4  z1:1,
5  z5:x^2 + 1     add,
6  z6:log(x^2 + 1)     log,
7  z3:y,
8  z7:exp(y)     exp,
9  z8:log(x^2 + 1) * exp(y)     mul,
10 z9:y^2     pow(2),
11 z10:y^2 * x     mul,
12 z11:log(x^2 + 1) * exp(y) + y^2 * x     add]
```

*← leaves first* .

*← root last*

# Eval $\nabla f(x, y) = \log(x^2 + 1)\exp(y) + y^2 x$ at $(2, 3)$

# Backward and add

```python
def backward(self):
    self.grad = 1.0
    topo_order = self.topological_sort()
    for node in reversed(topo_order):
        node._backward()

def __add__(self, other):
    assert isinstance(other, ag.Scalar)
    output = ag.Scalar(self.value + other.value,
                       inputs=[self, other], op="add")
    def _backward():
        #
        self.grad += output.grad
        other.grad += output.grad

    output._backward = _backward
    return output
```

# Exercise 2.B

Implement `ag.Scalar.topological_sort`

```
1         def __mul__(self, other):
2             # ...
3             def _backward():
4                 ## EXERCISE 2.B
5                 ## YOUR CODE HERE
6                 return None
7 # ... and also the backward function in
8         def __pow__(self, exponent): # exponent is just a python float
9 # ... and in
10    def exp(input):
```

self * other

x ** exponent

```
1 Check your answer
2 f(x,y) 50.3264246157732
3 x.grad 25.068429538550134
4 y.grad 44.3264246157732
```

# Exercise 2.C

Implement `ag.Scalar.relu`

```
def relu(input):
    output = ag.Scalar(max(0, input.value), inputs=[input], op="relu")
    def _backward():
        ## EXERCISE 2.C
        ## YOUR CODE HERE
# ...
```

```
# Check your answer with the plot
xs_raw = [i - 5 for i in range(10)]
xs = [ag.Scalar(i) for i in xs_raw]
ys = [ag.relu(x)**3 for x in xs]
[y.backward() for y in ys]
grads = [x.grad for x in xs]
plt.plot(xs_raw, grads, label='autograd')
plt.plot(xs_raw, [3*max(0,x_raw)**2 for x_raw in xs_raw], label='manual grad')
```

# What can we do with `ag.Scalar`?

For the remaining of this class, we will build a "knock-off" PyTorch consisting of:

- A `Model` class that encapsulates the parameters and the forward pass
- A `Loss` class for calculating Mean Squared Error (MSE)
- An `Optimizer` class for performing gradient descent
- A training loop

We will fit a 1-hidden layer neural network

$$f(x; \mathbf{w}_1, \mathbf{b}_1, \mathbf{w}_2, b_2) = \mathbf{w}_2^\top \mathrm{relu}(\mathbf{w}_1 x + \mathbf{b}_1) + b_2$$

where $\mathbf{w}_1$ and $\mathbf{w}_2$ are both vectors of `n_hidden` dimensions.

# The `Model` class

```
1  class Model:
2      def __init__(self, n_hidden, rng_seed=42):
3          np.random.seed(rng_seed)
4
5          w1np = np.random.randn(n_hidden)
6          # ...
7          b2np = np.random.randn(1)
8
9          self.w1 = [ag.Scalar(val) for val in w1np]
10         # ...
11         self.b2 = [ag.Scalar(val) for val in b2np]
12
13         self.parameters = self.w1 + self.b1 + self.w2 + self.b2
14
15     def forward(self, x): # ...
```

# The `forward` function

```
1    def forward(self, x):
2        # "upgrade" x into ag.Scalars
3        x_scalar = [ag.Scalar(val) for val in x]
4        n_samples = len(x_scalar)
5
6        # calculate the forward
7
8        ## YOUR CODE HERE
9        return [ag.Scalar(0.0) for i in range(n_samples)]
```

# The Loss class

```
1 class Loss :
2     def mse ( self , predictions , targets ) :
3         # mean squared error
4         assert len ( predictions ) == len ( targets )
5         n_samples = len ( predictions )
6         loss = ag . Scalar (0.0)
7
8         # YOUR CODE HERE
9
10        return loss
```
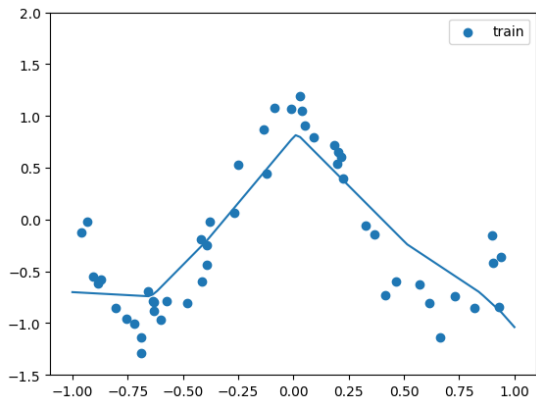
# The Optimizer class

```python
class Optimizer:
    def __init__(self, parameters, lr=0.01):
        self.parameters = parameters
        self.lr = lr

    def zero_grad(self):
        # YOUR CODE HERE
        pass

    def step(self):
        # YOUR CODE HERE
        pass
```

# Training Loop

```
1  model = Model(n_hidden=20)
2  loss_fn = Loss()
3  optimizer = Optimizer(model.parameters, lr=0.1)
4
5  for epoch in range(100):
6      optimizer.zero_grad()
7      output = model.forward(xnp)
8      loss = loss_fn.mse(output, ynp)
9      loss.backward()
10     optimizer.step()
11     if epoch % 10 == 0:
12         print(f"Iteration {epoch}, Loss: {loss.value}")
```

# Exercise 3

# References I