

# 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

## More on homework 3: Next-word-prediction

---

The quick brown fox jumps over the -----.

$\mathbf{X}^{(i)} = [\mathbf{x}^{(i,1)} \quad \mathbf{x}^{(i,2)} \quad \dots \quad \mathbf{x}^{(i,C-1)} \quad \mathbf{x}^{(i,C)}] \in \mathbb{R}^{d \times C}$  is a  $d \times C$  matrix

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

## More on homework 3: Next-word-prediction

---

The quick brown fox jumps over the -----.

$\mathbf{X}^{(i)} = [\mathbf{x}^{(i,1)} \quad \mathbf{x}^{(i,2)} \quad \dots \quad \mathbf{x}^{(i,C-1)} \quad \mathbf{x}^{(i,C)}] \in \mathbb{R}^{d \times C}$  is a  $d \times C$  matrix

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

# Notations

---

- $C$  as `n_context`
- $d$  as `n_features`
- $n$  as `n_samples`
- $q$  as `n_reduced`, where  $q < d$

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

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

# (“Single-layer” and “single-headed”) Attention

---

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

# 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

---

```
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)
```



# Batched matrix multiplication

---

```
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)
```

# 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 outer product

---

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

## np.moveaxis

---

```
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)
```

# 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

---

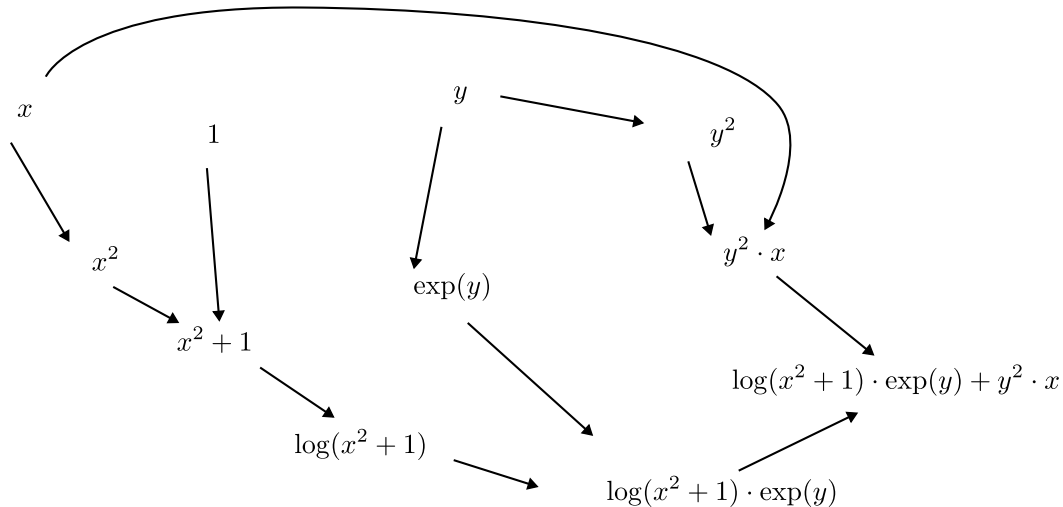
Let's compute

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

using automatic differentiation.

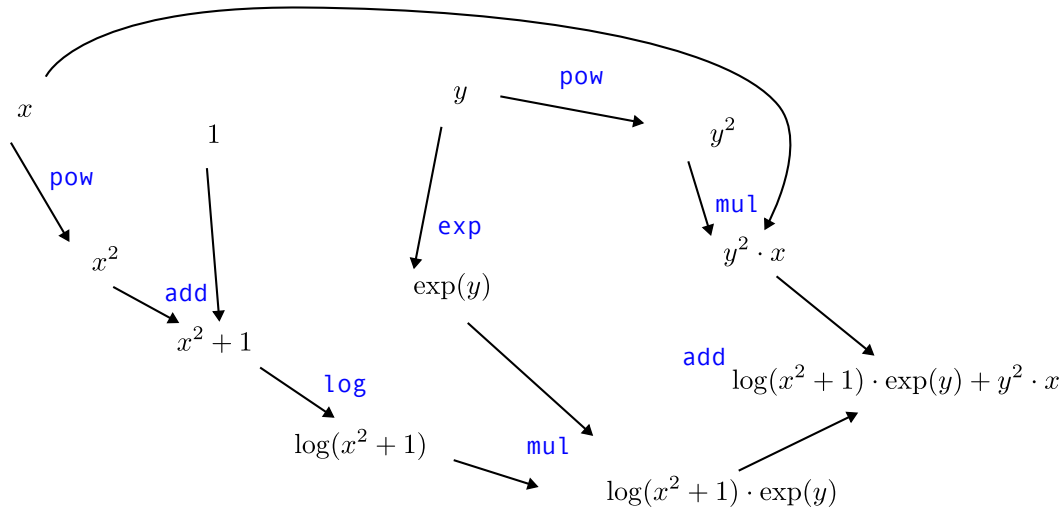
**Eval**  $f(x, y) = \log(x^2 + 1) \exp(y) + y^2 x$  **at**  $(x, y) = (2, 3)$

---



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

---





# An autograd “module”

---

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

# An autograd “module”

---

```
1 class ag: # AutoGrad
2     class Scalar: # Scalars with grads
3         def __init__(self, value, op="", _backward= lambda : None,
4             inputs=[], label=""):
5             self.value = float(value)
6             self.grad = 0.0
7
8             # ... lines skipped
9
10        def __add__(self, other): # ...
11        def __mul__(self, other): # ...
12        def __pow__(self, exponent): # ...
13
14    def exp(input): # ...
15    def log(input): # ...
```

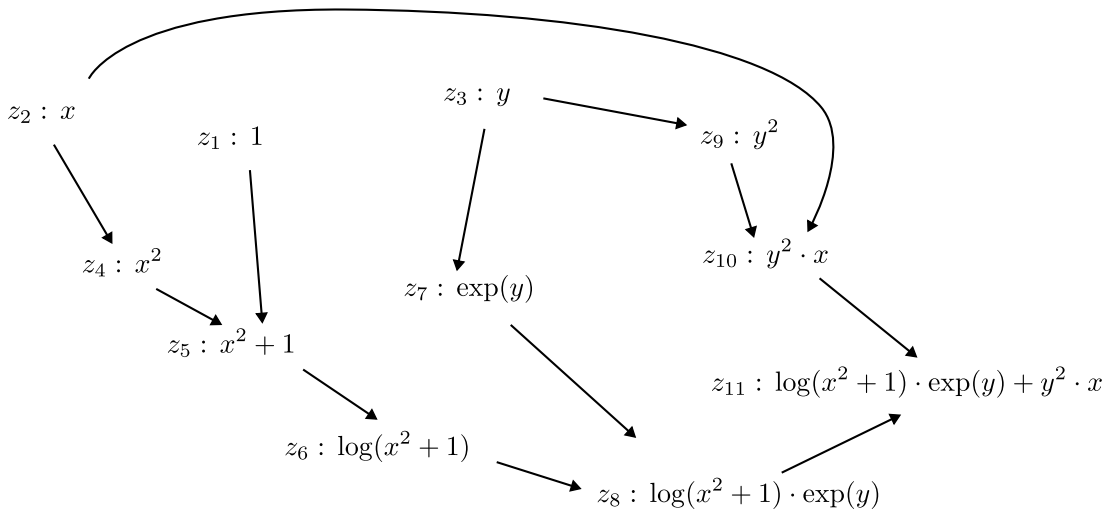
# Implement $f$ in our framework

---

```
1 x = ag.Scalar(2, label="z2:x")
2 y = ag.Scalar(3, label="z3: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")
7     z2 = x
8     z3 = y
9
10    z4 = z2**2
11    z4.label = "z4:x^2"
12
13    # ... lines skipped
```

# Forward pass — DFS — computation graph

---



## Exercise 2.A

---

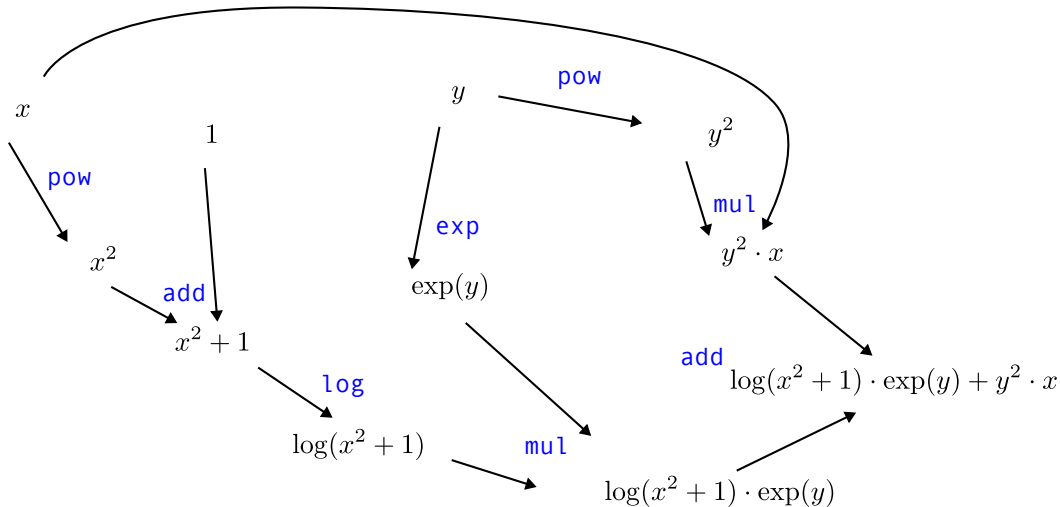
Implement `ag.Scalar.topological_sort`

```
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]
```

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

---



# Backward and add

---

```
1  def backward(self):
2      self.grad = 1.0
3      topo_order = self.topological_sort()
4      for node in reversed(topo_order):
5          node._backward()
6
7  def __add__(self, other):
8      assert isinstance(other, ag.Scalar)
9      output = ag.Scalar(self.value + other.value,
10                          inputs=[self, other], op="add")
11
12     def _backward():
13         # pass
14         self.grad += output.grad
15         other.grad += output.grad
16
17     output._backward = _backward
18     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):
```

```
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`

```
1  def relu(input):
2      output = ag.Scalar(max(0, input.value), inputs=[input], op="relu")
3
4      def _backward():
5          ## EXERCISE 2.C
6          ## YOUR CODE HERE
7
8  # ...
9
10 # Check your answer with the plot
11 xs_raw = [i - 5 for i in range(10)]
12 xs = [ag.Scalar(i) for i in xs_raw]
13 ys = [ag.relu(x)**3 for x in xs]
14 [y.backward() for y in ys]
15 grads = [x.grad for x in xs]
16 plt.plot(xs_raw, grads, label='autograd')
17 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

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

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

---

```
1 class Optimizer:
2     def __init__(self, parameters, lr=0.01):
3         self.parameters = parameters
4         self.lr = lr
5
6     def zero_grad(self):
7         # YOUR CODE HERE
8         pass
9
10    def step(self):
11        # YOUR CODE HERE
12        pass
```

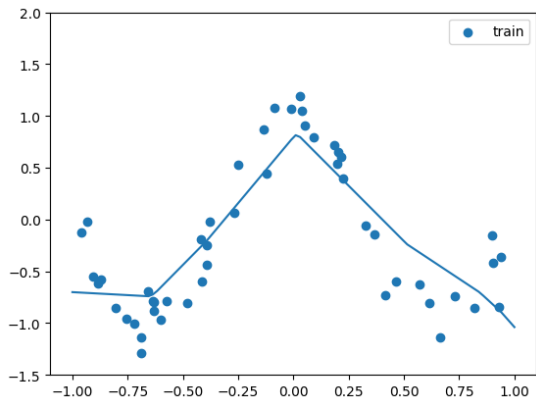
# Training Loop

---

```
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

---