Topics:

- Backpropagation
- Matrix/Linear Algebra view

CS 4644-DL / 7643-A ZSOLT KIRA

- Assignment 1 out!
 - Due June 5th (with grace period June 7th)
 - Start now, start now, start now!
 - Start now, start now, start now!
 - Start now, start now, start now!
- Resources:
 - These lectures
 - Matrix calculus for deep learning
 - <u>Gradients notes</u> and <u>MLP/ReLU Jacobian notes</u>.
 - **Topic OH:** Assignment 1 and matrix calculus (TBA)
- Piazza: Project teaming thread
 - Project Proposal: June 15th, Project Check-in: July 1st.
 - Project proposal overview during my OH (Thursday 12:30pm ET, recorded)

The main goal of the project proposal is to:

- 1. Clearly state the problem you aim to investigate.
- 2. Explain why it is interesting or important (the motivation).
- 3. Describe how you plan to address it (methods, experiments, datasets).
- 4. Provide a short overview of related work and background reading.

Your proposal is a chance for us to give you feedback on your project direction, so don't worry if all your ideas are not fully concrete yet. Think of it as a condensed version of the "Introduction" and "Related Work" sections of a typical Computer Vision, NLP, or Robotics conference paper.

Some important dates: the project proposal is due on Feb 14th (11:59pm EST)

Your proposal should address the following questions and points:

- 1. What: What problem are you investigating?
- 2. So What: Why is it interesting or significant?
- 3. Now What: What course of action or methodology do you propose?
- Related Work & Background: Summarize the readings or papers you have examined that provide context.

Additionally, please include initial thoughts on the following, even if they are not fully concrete:

- Uniqueness: How does your approach differ from (or build upon) prior works?
- Benchmark/Dataset/Setup: What benchmark, dataset, or experimental setup will you use?
- Algorithms: What class of algorithm(s) will you propose or implement?
- Evaluation: How will you measure performance (quantitative) and analyze results (qualitative/visualization)? What kinds of analysis do you expect to run?
- Risks/Roadblocks: What potential challenges do you foresee, and how might you address them?

You are free to choose any topic you would like. Below are some suggestions if you are having trouble:

FAQs:

I don't have a group. Can I do my project solo? No, you will be penalized for not gaining approval for a solo group.

It's the week/day/hour before the proposal deadline, and I couldn't find a group. Can I do my project solo?

This is insufficient justification. You will be penalized for not gaining approval for a solo group.

I don't have a group. When should I start panicking?

We are releasing this info now, but note the proposal is not due until Feb 14. Expect churn in groups until early Feb (students will drop, find other topics of interest, etc.). Ideally, you want to start locking down groups after assignment 2 is due and begin experimenting with ideas before finalizing a project topic.

Can we use existing implementations?

Yes, you can use existing implementations (do not forget to cite). However, remember that your project must be comprehensive per the project guidelines; don't provide just the results and shallow analysis. If you use existing models/code, this is even more important, as we can't judge how much you have learned from that. The key is to iterate between running, making hypotheses/claims about what's happening (with evidence, which you should show), and documenting the decisions you made in response to improving things. Also, be cautious about unofficial implementations. Often, they have bugs; we typically try to avoid using implementations that don't reproduce the results from the paper. You can often tell as the good ones have tables comparing their performance to the paper in the README.

Changing topic after the proposal deadline?

It depends. It is a higher risk further in the term. Feel free to attend office hours to ask or make a private post if you encounter this. Ultimately, it's up to you and your group. The proposal doesn't lock you into anything; however, the expectations for the final report do not change, even if you changed topics the week before.

Can we use pre-trained models?

You are allowed to use pre-trained models. However, just running open-source code is not sufficient. What we're looking for in the project is making novel contributions in the field (can be novel in both algorithms and applications). In that case, using a pre-trained model and fine-tuning on a dataset is insufficient by itself. However, using a pre-trained model and doing some non-trivial DL as part of the pipeline is acceptable. Do not forget to cite.

Can we use libraries other than Pytorch? Yes.

What makes a good project?

Usually, projects where the dataset is readily available and existing benchmarks are solid options. However, you are free to experiment.

Projects where you collect your own data take longer, require more infrastructure, and generally are not good candidates. We don't discourage passion projects, so if you have something in mind, go for it! Just be forewarned that you should likely start looking into the infrastructure now to ensure you have both the right data and enough data amenable to DL. An additional challenge here is that you will have to create baseline models (naive or otherwise) to show DL shows improvement over something.

Deep reinforcement learning (RL) projects are tough for the same reason. In deep RL, you'd be using a simulator that generates the data you train on. Often, these simulators require significant debugging to work, so be forewarned. Again, if you're passionate, don't let this stop you from pursuing projects in this space.

What type of resources are required to train deep neural networks (DNNs) for the project? Depending on the dataset, it can take hours or days for models to train on non-trivial datasets. If you haven't used GPUs or cloud computing before and plan to use them (often necessary for the Facebook projects or non-trivial datasets), plan on accounting for time to learn how to use these technologies.

An Example Project Proposal:

A sample will be posted on Canvas. We will make a post with a sample final project report after the proposal deadline.

CS 7643 - Deep Learning

Project Proposal

Learning from Demonstrations using Parameter-efficient Finetuning of Language Models

Team Name: Precision Pilots

Team Members: Neelabh Sinha, Sniodha Verma, Ananva Sharma

January 26, 2025

1 Project Summarv

The knowledge of Language Models (LMs) as next-word predictors have been proven to be very powerful recently. However, there is still a challenge in aligning their responses to user objectives. The ideal way is to finetune them for a downstream task, but that is inefficient and resource consuming due to the scale of these models. Recently, incontext learning, which is prompting LMs with in-context examples to get results for a task under zero-shot settings, has also produced great results. It stands out for its efficiency; however, for long sequence tasks, it encounters limitations due to restrictive context size of LMs, i.e., a maximum cap of sequence length under which they produce good results. When providing many in-context examples, it is natural that the prompt size for long-sequence tasks can exceed this limit. We aim to bridge this gap between finetuning and in-context learning by building a technique that is efficient and improves outputs without expensive finetuning/training these large LMs on large datasets.

2 Approach

To address the above problem, we propose Parameter-Efficient Fine-Tuning (PEFT) of LMs on these downstream tasks with prompts containing in-context examples. This involves exploring efficient methodologies to enhance their performance in downstream tasks. These adapters are methods/architectures to use with Transformer-based models that can improve performance on downstream tasks by training none or significantly less number of parameters. This direction is promising, primarily because of two reasons:

- 1. It fine-tunes a much smaller number of parameters, solving the limitation of inefficiency of end-to-end supervised finetuning.
- 2. It can compensate for context size limitation with in-context learning as the finetuning will be done using the specific downstream tasks

As in the paper [16], we will focus on two tasks to start: Natural Language Inference and Paraphrase Identification. The hypothesis is that the finetuning with prompts containing in-context examples will improve the context-size limitation of these models and the examples themselves will provide the LM with the information it needs to perform a new task.

We will start from pre-trained OPT [24] models from smallest to largest scale and experiment with different kinds of PEFT techniques. We have categorized those techniques into multiple types as listed in Table 1 and will try to implement at least one from each category to do a comparative study of how they generalize. Some of them are also illustrated in Figure 1.

The pre-trained models are available on HuggingFace1

Туре	PEFT Techniques	
Prompt-based Techniques	Prompt-tuning [12], Prefix-tuning [13], SPOT [20], IPT [17]	
Representation-based Techniques	LoRA [11], KronA [4]	
Series-based Techniques	Series Adapters [9], Adamix [21], Sparse Adapters [8], LeTS [5]	
Parallel Techniques	Parallel Adapters [6]	
Hybrid Techniques	MAM Adapters [7], Unipelt [15], Compacter [14], S4 [1]	

Table 1: Existing Techniques for PEFT





We will also study how these adapters behave with increasing scale of LMs. Since OPT provide range of scaled LMs from 125 million to 30 billion, we can do this easily. At last, we will analyze the in-domain and out-of-domain accuracy along with memory and execution time for all the models across the range of experiments designed above. We will compare our methods against end-to-end supervised finetuning and in-context learning as baselines. Tentative steps are described as follows -

- 1. Devise methods to get in-context examples: We will explore using cosine similarity or SBERT similarity, as well as look into more sophisticated methods.
- 2. Pick a Language Model and Adapter; We will use one of the OPT variants with an adapter from Table 1. The choice will be based on types, usability, etc.
- 3. Finetune: We can either do vanilla finetuning [10] of pattern-based finetuning [16] based on experimentation results that we get.
- 4. Evaluate: We will be evaluating our models through the following metrics:
 - In-Distribution and Out-of-Distribution Accuracies
 - Execution Time
- Memory

CS 7643 - Deep Learning

5. Experiment and Ablation Study: Observe the change in performance with:

Scale of LM

- · across different adapters and their variants
- · Baseline comparisons with SFT and ICL

3 Related Work

Supervised finetuning of LMs fine-tune a pretrained language model on a specific task using a randomized classification head on top. A variation of it is pattern-based finetuning [16] where there is a verbalizer [18] that maps LM output logits to desired answer, and rest of the finetuning is similar. In-context learning [2], on the other hand, is a method that doesn't update model's weights but provides examples similar to the desired task to instruct LMs how to respond to a given task

We are exploring PEFT, which is a technique to improve LM performance by not modifying its original weights but adding additional elements and finetune those to allow the LM to adapt to these new tasks. There are multiple such methods that are presented and cited in Table 1.

4 Datasets

As defined before, we will focus on two tasks; Natural Language Inference (NLI) and Paraphrase Identification, For NLI, we can use the MNLI [23] and RTE [3] dataset. And for paraphrase identification, we can use the QQP [19]

CS 7643 - Deep Learning

Project Proposal

dataset. NLIV2 [22] is also a great dataset which contains many of such tasks with given in-context examples. So, this will remove our need for pre-processing and finding adequate examples and prompts.

Project Proposal

We can also look at some generation tasks for our proposed methodology later on, based on timeline. Text summarization is one good example. We can use the NLIV2 [22] for that also, as it contains summarization examples from multiple other datasets.

All these datasets are publicly available through the website of the authors of the paper, or HuggingFace2,

5 Group Members

Neelabh Sinha, Snigdha Verma, Ananva Sharma

6 Ethical Implications

Ethical considerations accompany the initiative to fine-tune language models through parameter-efficient fine-tuning (PEFT). Efficient fine-tuning could lessen obstacles to unethical use, such as customizing language models to produce misleading, harmful, or biased content. Additionally, biases inherent in pre-trained language models could be exacerbated by fine-tuning, potentially resulting in discriminatory or unfair results in later applications. Privacy issues might emerge if sensitive data unintentionally becomes part of the fine-tuning datasets, risking exposure through model outputs. Although PEFT seeks to reduce resource consumption compared to traditional methods. even training smaller adapters demands considerable computational power, posing environmental and economic challenges. Furthermore, tasks like Natural Language Inference and Paraphrase Identification could spread misinformation if outputs aren't carefully scrutinized or if biases in datasets are perpetuated. Advances in PEFT may disproportionately benefit organizations with robust computational access, potentially widening gaps between wellresourced and less-resourced communities. To counteract these problems, we commit to conducting comprehensive audits of training datasets, incorporating fairness and robustness into our evaluations, documenting our models' limitations and intended purposes, and adhering to ethical standards to promote responsible development and use of our approaches

References

- [1] Jiaao Chen, Aston Zhang, Xingjian Shi, Mu Li, Alex Smola, and Divi Yang. Parameter-efficient fine-tuning design spaces. In The Eleventh International Conference on Learning Representations, 2023.
- [2] Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment challenge. In Joaquin Quiñonero-Candela, Ido Dagan, Bernardo Magnini, and Florence d'Alché Buc, editors, Machine Learning Challenges. Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Tectual Entailment, pages 177-190, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- [3] Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment challenge. In Joaquin Quiñonero-Candela, Ido Dagan, Bernardo Magnini, and Florence d'Alché Buc, editors, Machine Learning Challenges, Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Tectual Entailment, pages 177-190, Berlin, Heidelberg, 2006, Springer Berlin Heidelberg,
- [4] Ali Edalati, Marzieh Tahaei, Ivan Kobyzev, Vahid Partovi Nia, James J. Clark, and Mehdi Rezagholizadeh. Krona: Parameter efficient tuning with kronecker adapter, 2022,
- 15] Cheng Fu, Hanxian Huang, Xinvun Chen, Yuandong Tian, and Jishen Zhao, Learn-to-share; A hardwarefriendly transfer learning framework exploiting computation and parameter sharing. In Marina Meila and Tong Zhang, editors, Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pages 3469-3479. PMLR, 18-24 Jul 2021.
- [6] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. In International Conference on Learning Representations, 2022.
- [7] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning, 2022.

²https://huggingface.co/docs/datasets/en/index

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Stretch pixels into column

Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n





We can find the steepest descent direction by computing the derivative (gradient):

$$f'(a) = \lim_{h \to 0} \frac{f(a+h) - f(a)}{h}$$

- Steepest descent direction is the negative gradient
- Intuitively: Measures how the function changes as the argument a changes by a small step size
 - As step size goes to zero
- In Machine Learning: Want to know how the loss function changes as weights are varied

Derivatives

 Can consider each parameter separately by taking partial derivative of loss function with respect to that parameter





The same two-layered neural network corresponds to adding another weight matrix

 We will prefer the linear algebra view, but use some terminology from neural networks (& biology)



Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n





Large (deep) networks can be built by adding more and more layers

Three-layered neural networks can represent **any function**

 The number of nodes could grow unreasonably (exponential or worse) with respect to the complexity of the function

We will show them **without edges**:





 $f(x, W_1, W_2, W_3) = \sigma(W_2 \sigma(W_1 x))$

Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n



Adding More Layers!

- We are learning complex models with significant amount of parameters (millions or billions)
- How do we compute the gradients of the loss (at the end) with respect to internal parameters?
- Intuitively, want to understand how small changes in weight deep inside are propagated to affect the loss function at the end







To develop a general algorithm for this, we will view the function as a **computation graph**

Graph can be any **directed acyclic** graph (DAG)

 Modules must be differentiable to support gradient computations for gradient descent

A **training algorithm** will then process this graph, **one module at a time**





















Note that we must store the **intermediate outputs of all layers**!

This is because we will need them to compute the gradients (the gradient equations will have terms with the output values in them)





Step 2: Compute Gradients wrt parameters: Backward Pass







Step 2: Compute Gradients wrt parameters: Backward Pass







Step 2: Compute Gradients wrt parameters: Backward Pass









Computing the Gradients of Loss



Step 1: Compute Loss on Mini-Batch: Forward Pass
Step 2: Compute Gradients wrt parameters: Backward Pass
Step 3: Use gradient to update all parameters at the end



$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$$

Backpropagation is the application of gradient descent to a computation graph via the chain rule!





Backpropagation: a simple example



Figure adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n



Deep Learning = Differentiable Programming

- Computation = Graph
 - Input = Data + Parameters
 - Output = Loss
 - Scheduling = Topological ordering
- What do we need to do?
 - Generic code for representing the graph of modules
 - Specify modules (both forward and backward function)
 - Backpropagation implementation on the graph



Modularized implementation: forward / backward API

cl



Graph (or Net) object (rough psuedo code)

ass Co	<pre>mputationalGraph(object):</pre>
#	
def	<pre>forward(inputs):</pre>
ŝ	# 1. [pass inputs to input gates]
	# 2. forward the computational graph:
	<pre>for gate in self.graph.nodes_topologically_sorted():</pre>
	gate.forward()
	return loss # the final gate in the graph outputs the loss
def	backward():
	<pre>for gate in reversed(self.graph.nodes_topologically_sorted()):</pre>
	<pre>gate.backward() # little piece of backprop (chain rule applied)</pre>
	return inputs_gradients





Modularized implementation: forward / backward API





Modularized implementation: forward / backward API



<pre>.ass MultiplyGate(object):</pre>			
<pre>def forward(x,y):</pre>			
$z = x^*y$			
<pre>self.x = x # must keep these around!</pre>			
self.y = y			
return z			
<pre>def backward(dz):</pre>			
dx = self.y * dz # [dz/dx * dL/dz]			
dy = self.x * dz # [dz/dy * dL/dz]			
return [dx, dy]			





Example: Caffe layers

Example.			cudnn_lcn_layer.cpp	dismantle layer headers	a year ago
ranch: master - caffe / src / caff	e / layers /	Create new file Upload files Find file History	Cudnn_lcn_layer.cu	dismantle layer headers	a year ago
			Cudnn_Irn_layer.cpp	dismantle layer headers	a year ago
shelhamer committed on GitHub Merge pull request #4630 from BIGene/load_hdf5_fix Latest commit e687a71 21 days ago			cudnn_Irn_layer.cu	dismantle layer headers	a year ago
			cudnn_pooling_layer.cpp	dismantle layer headers	a year ago
absval_layer.cpp	dismantle layer headers	a year ago	cudnn_pooling_layer.cu	dismantle layer headers	a year ago
absval_layer.cu	dismantle layer headers	a year ago	Cudnn_relu_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
accuracy_layer.cpp	dismantle layer headers	a year ago	Cudnn_relu_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
argmax_layer.cpp	dismantle layer headers	a year ago	cudnn_sigmoid_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
base_conv_layer.cpp	enable dilated deconvolution	a year ago	cudnn_sigmoid_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
base_data_layer.cpp	Using default from proto for prefetch	3 months ago	cudnn_softmax_layer.cpp	dismantle layer headers	a year ago
base_data_layer.cu	Switched multi-GPU to NCCL	3 months ago	cudnn_softmax_layer.cu	dismantle layer headers	a year ago
batch_norm_layer.cpp	Add missing spaces besides equal signs in batch_norm	_layer.cpp 4 months ago	cudnn_tanh_layer.cpp	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
batch_norm_layer.cu	dismantle layer headers	a year ago	cudnn_tanh_layer.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago
batch_reindex_layer.cpp	dismantle layer headers	a year ago	data_layer.cpp	Switched multi-GPU to NCCL	3 months ago
batch_reindex_layer.cu	dismantle layer headers	a year ago	deconv_layer.cpp	enable dilated deconvolution	a year ago
bias_layer.cpp	Remove incorrect cast of gemm int arg to Dtype in Bias	Layer a year ago	deconv_layer.cu	dismantle layer headers	a year ago
bias_layer.cu	Separation and generalization of ChannelwiseAffineLay	er into BiasLayer a year ago	dropout_layer.cpp	supporting N-D Blobs in Dropout layer Reshape	a year ago
bnll_layer.cpp	dismantle layer headers	a year ago	dropout_layer.cu	dismantle layer headers	a year ago
bnll_layer.cu	dismantle layer headers	a year ago	dummy_data_layer.cpp	dismantle layer headers	a year ago
concat_layer.cpp	dismantle layer headers	a year ago	eltwise_layer.cpp	dismantle layer headers	a year ago
concat_layer.cu	dismantle layer headers	a year ago	eltwise_layer.cu	dismantle layer headers	a year ago
contrastive_loss_layer.cpp	dismantle layer headers	a year ago	elu_layer.cpp	ELU layer with basic tests	a year ago
contrastive_loss_layer.cu	dismantle layer headers	a year ago	elu_layer.cu	ELU layer with basic tests	a year ago
conv_layer.cpp	add support for 2D dilated convolution	a year ago	embed_layer.cpp	dismantle layer headers	a year ago
conv_layer.cu	dismantle-layer headers under RSD 2 Clause	a year ago	embed_layer.cu	dismantle layer headers	a year ago
crop laver.cpp	remove redundant operations in Crop laver (#5138)	2 months ago	euclidean_loss_layer.cpp	dismantle layer headers	a year ago
crop_laver.cu	remove redundant operations in Crop layer (#5138)	2 months ago	euclidean_loss_layer.cu	dismantle layer headers	a year ago
cudnn conv laver.cop	dismantle laver headers	a vear ago	exp_layer.cpp	Solving issue with exp layer with base e	a year ago
cudnn conv laver.cu	Add cuDNN v5 support, drop cuDNN v3 support	11 months ago	exp_layer.cu	dismantle layer headers	a year ago
sonsone					









- Neural networks involves composing simple functions into a computation graph
- Optimization (updating weights) of this graph is through backpropagation
 - Recursive algorithm: Gradient descent (partial derivatives) plus chain rule
- Remaining questions:
 - How does this work with vectors, matrices, tensors?
 - Across a composed function?
 - How can we implement this algorithmically to make these calculations automatic? **Automatic Differentiation**





Linear **Algebra View: Vector and Matrix Sizes**







Sizes: $[c \times (m + 1)]$ $[(m + 1) \times 1]$

Where *c* is number of classes

m is dimensionality of input



Closer Look at a Linear Classifier



Conventions:

Size of derivatives for scalars, vectors, and matrices: Assume we have scalar $s \in \mathbb{R}^1$, vector $v \in \mathbb{R}^m$, i.e. $v = [v_1, v_2, ..., v_m]^T$ and matrix $M \in \mathbb{R}^{m_1 \times m_2}$



Dimensionality of Derivatives



Conventions:

- Size of derivatives for scalars, vectors, and matrices: Assume we have scalar $s \in \mathbb{R}^1$, vector $v \in \mathbb{R}^m$, i.e. $v = [v_1, v_2, ..., v_m]^T$ and matrix $M \in \mathbb{R}^{m_1 \times m_2}$
- What is the size of $\frac{\partial v}{\partial s}$? $\mathbb{R}^{m \times 1}$ (column vector of size m)
- What is the size of $\frac{\partial s}{\partial v}$? $\mathbb{R}^{1 \times m}$ (row vector of size m)

$$\left[\frac{\partial s}{\partial v_1} \frac{\partial s}{\partial v_1} \cdots \frac{\partial s}{\partial v_m}\right]$$

$$\begin{bmatrix} \frac{\partial v_1}{\partial s} \\ \frac{\partial v_2}{\partial s} \\ \vdots \\ \frac{\partial v_m}{\partial s} \end{bmatrix}$$







This matrix of partial derivatives is called a Jacobian

(Note this is slightly different convention than on Wikipedia). Also, computationally other conventions are used.

Dimensionality of Derivatives



Conventions:

• What is the size of $\frac{\partial s}{\partial M}$? A matrix:



(Note this is slightly different convention than on Wikipedia). Also, computationally other conventions are used.





Example 1: $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x \\ x^2 \end{bmatrix} \qquad \frac{\partial}{\partial}$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} 1\\2x \end{bmatrix}$$

 $\frac{\partial(\sum_k w_k x_k)}{\partial x_i} = w_i$

Example 2:

$$y = w^{T}x = \sum_{k} w_{k}x_{k}$$
$$\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x_{1}}, \dots, \frac{\partial y}{\partial x_{m}}\right]$$
$$= [w_{1}, \dots, w_{m}] \quad \text{because}$$
$$= w^{T}$$







Example 4:

$$\frac{\partial (wAw)}{\partial w} = 2w^T A \text{ (assuming A is symmetric)}$$




- What is the size of $\frac{\partial L}{\partial W}$?
 - Remember that loss is a scalar and W is a matrix:

 $\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b3 \end{bmatrix}$ Jacobian is also a matrix: W $\begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} & \cdots & \frac{\partial L}{\partial w_{1m}} & \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial w_{21}} & \cdots & \cdots & \frac{\partial L}{\partial w_{2m}} & \frac{\partial L}{\partial b_2} \\ \cdots & \cdots & \cdots & \frac{\partial L}{\partial w_{3m}} & \frac{\partial L}{\partial b_3} \end{bmatrix}$

Dimensionality of Derivatives in ML



Batches of data are **matrices** or **tensors** (multidimensional matrices)

Examples:

- Each instance is a vector of size m, our batch is of size $[B \times m]$
- Each instance is a matrix (e.g. grayscale image) of size $W \times H$, our batch is $[B \times W \times H]$
- Each instance is a multi-channel matrix (e.g. color image with R,B,G channels) of size $C \times W \times H$, our batch is $[B \times C \times W \times H]$

Jacobians become tensors which is complicated

- Instead, flatten input to a vector and get a vector of derivatives!
- This can also be done for partial derivatives between two vectors, two matrices, or two tensors



Jacobians of Batches





Fully Connected (FC) Layer: Forward Function







Fully Connected (FC) Layer









We can employ **any differentiable** (or piecewise differentiable) function

A common choice is the **Rectified** Linear Unit

- Provides non-linearity but better gradient flow than sigmoid
- Performed element-wise

How many parameters for this layer?







Full Jacobian of ReLU layer is **large** (output dim x input dim)

- But again it is sparse
- Only diagonal values non-zero because it is element-wise
- An output value affected only by corresponding input value

Max function **funnels gradients through selected max**

Gradient will be zero if input
 <= 0













For element-wise ops, jacobian is **sparse**: off-diagonal entries always zero! Never **explicitly** form Jacobian -- instead use elementwise multiplication



- Neural networks involves composing simple functions into a computation graph
- Optimization (updating weights) of this graph is through backpropagation
 - Recursive algorithm: Gradient descent (partial derivatives) plus chain rule
- Remaining questions:
 - How does this work with vectors, matrices, tensors?
 - Across a composed function? Next!
 - How can we implement this algorithmically to make these calculations automatic? **Automatic Differentiation**





Composition of Functions: $f(g(x)) = (f \circ g)(x)$

A complex function (e.g. defined by a neural network):

$$f(x) = g_{\ell} (g_{\ell-1}(\dots g_1(x)))$$
$$f(x) = g_{\ell} \circ g_{\ell-1} \dots \circ g_1(x)$$

(Many of these will be parameterized)

(Note you might find the opposite notation as well!)

Composition of Functions & Chain Rule



$$\begin{array}{ccc} \mathbf{x} \in \mathbb{R}^1 & \longrightarrow & \mathbf{z} \in \mathbb{R}^1 & \longrightarrow & \mathbf{y} \in \mathbb{R}^1 \\ & & g_1 & & g_2 \\ & & & y = g_2 \big(g_1(x) \big) \end{array}$$

$$=\frac{\partial \partial y}{\partial \partial x} \frac{\partial z}{\partial x}$$

Scalar Multiplication









Matrix Multiplication





$$\left[\begin{array}{c} \frac{\partial y_i}{\partial x_j} \\ \end{array}\right] = \left[\begin{array}{c} \frac{\partial y_i}{\partial z_k} \\ \frac{\partial z_k}{\partial x_j} \\ \end{array}\right]$$

$$\frac{\partial y_i}{\partial x_j} = \sum_k \frac{\partial y_i}{\partial z_k} * \frac{\partial z_k}{\partial x_j}$$











$$h^0 \in \mathbb{R}^d \longrightarrow h^1 \in \mathbb{R}^d \longrightarrow \dots \longrightarrow h^l \in \mathbb{R}^d$$







 $h^0 \in \mathbb{R}^d \longrightarrow h^1 \in \mathbb{R}^d \longrightarrow \dots \longrightarrow h^l \in \mathbb{R}^d \longrightarrow L \in \mathbb{R}^1$

Which directions is more efficient to multiply?







$$\bar{L} = 1$$

$$\bar{p} = \frac{\partial L}{\partial p} = -\frac{1}{p}$$

where $p = \sigma(w^T x)$ and $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\bar{u} = \frac{\partial L}{\partial u} = \frac{\partial L}{\partial p} \quad \frac{\partial p}{\partial u} = \bar{p} \ \sigma(1 - \sigma)$$

$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u} \quad \frac{\partial u}{\partial w} = \bar{u}x^T$$

We can do this in a combined way to see all terms together:

$$\overline{w} = \frac{\partial L}{\partial p} \frac{\partial p}{\partial u} \frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T$$
$$= -\left(1 - \sigma(w^T x)\right) x^T$$

This effectively shows gradient flow along path from L to w





The chain rule can be computed as a **series of scalar, vector, and matrix linear algebra operations**



Extremely efficient in graphics processing units (GPUs)

$\overline{w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T$				
		[]	[] [1v1]
	1x1	1X1	171	1xd









We can do this in a combined way to see all terms together:

$$\begin{split} \overline{w} &= \frac{\partial L}{\partial p} \, \frac{\partial p}{\partial u} \, \frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T \\ &= -\left(1 - \sigma(w^T x)\right) x^T \end{split}$$

This effectively shows gradient flow along path from L to w

Computation Graph / Global View of Chain Rule





Computational / Tensor View



Graph View



Backpropagation View (Recursive Algorithm)

Different Views of Equivalent Ideas



- **Backpropagation:** Recursive, modular algorithm for chain rule + gradient descent
- When we move to vectors and matrices:
 - Composition of functions (scalar)
 - Composition of functions (vectors/matrices)
 - Jacobian view of chain rule
 - Can view entire set of calculations as linear algebra operations (matrix-vector or matrix-matrix multiplication)
- Automatic differentiation:
 - Reduction of modules to simple operations we know (simple multiplication, etc.)
 - Automatically build computation graph in background as write code
 - Automatically compute gradients via backward pass





Automatic Differentiation



Deep Learning = Differentiable Programming

- Computation = Graph
 - Input = Data + Parameters
 - Output = Loss
 - Scheduling = Topological ordering
- What do we need to do?
 - Generic code for representing the graph of modules
 - Specify modules (both forward and backward function)



Modularized implementation: forward / backward API





Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Backpropagation does not really spell out how to **efficiently** carry out the necessary computations

But the idea can be applied to **any directed acyclic graph** (DAG)

 Graph represents an ordering constraining which paths must be calculated first

Given an ordering, we can then iterate from the last module backwards, **applying the chain rule**

- We will store, for each node, its local gradient function/computation for efficiency
- We will do this automatically by computing backwards function for primitives and as you write code, express the function with them

This is called reverse-mode **automatic differentiation**







Computation = Graph

- Input = Data + Parameters
- Output = Loss
- Scheduling = Topological ordering

Auto-Diff

 A family of algorithms for implementing chain-rule on computation graphs







Automatic differentiation:

- Carries out this procedure for us on arbitrary graphs
- Knows derivatives of primitive functions
- As a result, we just define these (forward) functions and don't even need to specify the gradient (backward) functions!

$$\bar{L} = 1$$

$$\bar{p} = \frac{\partial L}{\partial p} = -\frac{1}{p}$$

where $p = \sigma(w^T x)$ and $\sigma(x) = \frac{1}{1+e^{-x}}$

$$\bar{u} = \frac{\partial L}{\partial u} = \frac{\partial L}{\partial p} \quad \frac{\partial p}{\partial u} = \bar{p} \ \sigma(1 - \sigma)$$

$$\bar{w} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial u} \quad \frac{\partial u}{\partial w} = \bar{u}x^T$$

We can do this in a combined way to see all terms together:

$$\overline{w} = \frac{\partial L}{\partial p} \frac{\partial p}{\partial u} \frac{\partial u}{\partial w} = -\frac{1}{\sigma(w^T x)} \sigma(w^T x) (1 - \sigma(w^T x)) x^T$$
$$= -\left(1 - \sigma(w^T x)\right) x^T$$

This effectively shows gradient flow along path from L to w





- Key idea is to explicitly store computation graph in memory and corresponding gradient functions
- Nodes broken down to basic primitive computations (addition, multiplication, log, etc.) for which corresponding derivative is known









We want to find the **partial derivative of output f** (output) with respect to **all intermediate variables**

Assign intermediate variables

Simplify notation: Denote bar as: $\overline{a_3} = \frac{\partial f}{\partial a_3}$

Start at end and move backward







Example





$$\overline{a_1} = \frac{\partial f}{\partial a_1} = \frac{\partial f}{\partial a_3} \quad \frac{\partial a_3}{\partial a_1} = \frac{\partial f}{\partial a_3} \quad \frac{\partial (a_1 + a_2)}{\partial a_1} = \frac{\partial f}{\partial a_3} \quad \mathbf{1} = \overline{a_3}$$
$$\overline{a_2} = \frac{\partial f}{\partial a_2} = \frac{\partial f}{\partial a_3} \quad \frac{\partial a_3}{\partial a_2} = \overline{a_3}$$

Addition operation distributes gradients along all paths!







Multiplication operation is a gradient switcher (multiplies it by the values of the other term)

$$\overline{x_2} = \frac{\partial f}{\partial a_2} \quad \frac{\partial a_2}{\partial x_2} = \frac{\partial f}{\partial a_2} \quad \frac{\partial (x_1 x_2)}{\partial x_2} = \overline{a_2} x_1$$

$$\overline{x_1} = \frac{\partial f}{\partial a_2} \quad \frac{\partial a_2}{\partial x_1} = \overline{a_2} x_2$$

Patterns of Gradient Flow: Multiplication



Several other patterns as well, e.g.:

Max operation **selects** which path to push the gradients through

- Gradient flows along the path that was "selected" to be max
- This information must be recorded in the forward pass



The flow of gradients is one of the most important aspects in deep neural networks

If gradients do not flow backwards properly, learning slows or stops!

Patterns of Gradient Flow: Other



A graph is created on the fly

from torch.autograd import Variable

```
x = Variable(torch.randn(1, 20))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 20))
```

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
```



(Note above)





Back-propagation uses the dynamically built graph

from torch.autograd import Variable

x = Variable(torch.randn(1, 20))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 20))

```
i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()
```

next_h.backward(torch.ones(1, 20))



From pytorch.org

Computation Graphs in PyTorch






Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Neural Turing Machine



Figure reproduced with permission from a Twitter post by Andrej Karpathy.



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

- Computation graphs are not limited to mathematical functions!
- Can have control flows (if statements, loops) and backpropagate through algorithms!
- Can be done dynamically so that gradients are computed, then nodes are added, repeat
- Differentiable programming



Adapted from figure by Andrej Karpathy

Power of Automatic Differentiation



Backpropagation, and automatic differentiation, allows us to optimize **any** function composed of differentiable blocks

- No need to modify the learning algorithm!
- The complexity of the function is only limited by computation and memory







A network with two or more hidden layers is often considered a **deep** model

Depth is important:

- Structure the model to represent an inherently compositional world
- Theoretical evidence that it leads to parameter efficiency
- Gentle dimensionality reduction (if done right)







There are still many design decisions that must be made:

- Architecture
- Data Considerations
- Training and Optimization
- Machine Learning Considerations







