Topics:

- Autodiff
- Optimization

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 (linked today)
- Piazza: Project teaming thread
 - Project Proposal: June 15th, Project Check-in: July 1st.
 - Project proposal reviewed last time

https://faculty.cc.gate	E					
CIOS ᅇ InfoReady	💞 GradeChange 🛛 😭 Wo	orkday 🥌 IC Docs	🕒 Deliverables 🗴	Kira FY2024.xlsx 🛛 👯 N	MediaSpace 💁 Slate	Reader 🛛 🗮 VLM-25
CS4644: https: CS7643: https:	://www.gradescope. ://www.gradescope.	com/courses/103 com/courses/103	7086 7087			
Today) < > Ma	ay 2025 👻		Ę		Month -
SUN 27	MON 28	TUE 29	WED 30	THU May 1	FRI 2 • 11am Zach's OH	SAT 3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20 • 11:30am Kausar	21	22 • 12:30pm Instruc	23 • 11:30am Kausar	24
25	26 • 2:30pm Mili's O	27 • 11:30am Kausar • 3:30pm Yipu's C	28 • 3:30pm Mili's O	29 • 12:30pm Instruc • 3:30pm Mili's O	30 • 11:30am Kausar • 3:30pm Yipu's C	31
CS7643/46	44 Sum25	0) Eastern Time	aw York		Ge	oole Calendar

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





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









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!



Adapted from figure by Marc'Aurelio Ranzato, Yann LeCun







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**





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	IXI	IVI	Ixa			









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







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





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!





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







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








Architectural Considerations



Determining what modules to use, and how to connect them is part of the **architectural design**

- Guided by the type of data used and its characteristics
 - Understanding your data is always the first step!
- Lots of data types (modalities) already have good architectures
 - Start with what others have discovered!
- The flow of gradients is one of the key principles to use when analyzing layers







- Combination of linear and non-linear layers
- Combination of only linear layers has same representational power as one linear layer
- Non-linear layers are crucial
 - Composition of non-linear layers enables complex transformations of the data

 $w_1^T(w_2^T(w_3^Tx)) = w_4^Tx$







Linear and Non-Linear Modules

Several aspects that we can **analyze**:

- Min/Max
- Correspondence between input & output statistics
- Gradients
 - At initialization (e.g. small values)
 - At extremes
- Computational complexity







- Min: 0, Max: 1
- Output always positive
- Saturates at both ends
- Gradients
 - Vanishes at both end
 - Always positive
- Computation: Exponential term



$$h^{\ell} = \sigma (h^{\ell-1})$$
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

 $1 + e^{-x}$







Centered

- Saturates at both ends
- Gradients
 - Vanishes at both end
 - Always positive
- Still somewhat computationally heavy



$$h^{\ell} = tanh(h^{\ell-1})$$





- Min: 0, Max: Infinity
- Output always positive
- **No saturation** on positive end!
- Gradients
 - 0 if $x \le 0$ (dead ReLU)
 - Constant otherwise (does not vanish)
- Cheap to compute (max)



$$h^{\ell} = max(0, h^{\ell-1})$$



Rectified Linear Unit

- Min: -Infinity, Max: Infinity
- Learnable parameter!
- No saturation
- Gradients
 - No dead neuron
- Still cheap to compute

Leaky ReLU



$$h^{\ell} = max(\alpha h^{\ell-1}, h^{\ell-1})$$



- Activation functions is still area of research!
 - Though many don't catch on

 In Transformer architectures, other activations such as GeLU is common



From "Gaussian Error Linear Units (GELUs)", Hendrycks & Gimpel

Variations: ELU, GeLU, etc.



Selecting a Non-Linearity

Which **non-linearity** should you select?

- Unfortunately, no one activation function is best for all applications
- ReLU is most common starting point
 - Sometimes leaky ReLU can make a big difference
- Sigmoid is typically avoided unless clamping to values from [0,1] is needed



Demo

<u>http://playground.tensorflow.org</u>





Initialization



Initializing the Parameters

The parameters of our model must be initialized to something

- Initialization is extremely important!
 - Determines how statistics of outputs (given inputs) behave
 - Determines how well gradients flow in the beginning of training (important)
 - Could limit use of full capacity of the model if done improperly
- Initialization that is close to a good (local) minima will converge faster and to a better solution





Initializing values to a constant value leads to a degenerate solution!

- What happens to the weight updates?
- Each node has the same input from previous layers so gradients will be the same
- As a results, all weights will be updated to the same exact values

input layer 1 hidden layer 2

 $w_i = c \quad \forall i$





Common approach is small normally distributed random numbers

- E.g. $N(\mu, \sigma)$ where $\mu = 0, \sigma = 0.01$
- Small weights are preferred since no feature/input has prior importance
- Keeps the model within the linear region of most activation functions









Deeper networks (with many layers) are more sensitive to initialization

- With a deep network, activations (outputs of nodes) get smaller
 - Standard deviation reduces significantly
- Leads to small updates smaller values multiplied by upstream gradients



Distribution of activation values of a network with tanh nonlinearities, for increasingly deep layers

From "Understanding the difficulty of training deep feedforward neural networks." AISTATS, 2010.



Limitation of Small Weights



Ideally, we'd like to maintain the variance at the output to be similar to that of input!

This condition leads to a simple initialization rule, sampling from uniform distribution:

Uniform
$$\left(-\frac{\sqrt{6}}{n_j+n_{j+1}},+\frac{\sqrt{6}}{n_j+n_{j+1}}\right)$$

Where n_j is fan-in

 (number of input nodes)
 and n_{j+1} is fan-out
 (number of output nodes)



Distribution of activation values of a network with tanh nonlinearities, for increasingly deep layers

From "Understanding the difficulty of training deep feedforward neural networks." AISTATS, **2010.**



Xavier Initialization

In practice, simpler versions perform empirically well:

$$N(0,1) * \sqrt{\frac{1}{n_j}}$$

- This analysis holds for tanh or similar activations.
- Similar analysis for ReLU activations leads to:

$$N(0,1) * \sqrt{\frac{1}{n_j/2}}$$

"Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification", ICCV, 2015.





Summary

Key takeaway: Initialization matters!

- Determines the activation (output) statistics, and therefore gradient statistics
- If gradients are small, no learning will occur and no improvement is possible!
- Important to reason about output/gradient statistics and analyze them for new layers and architectures





Regularization



Many standard regularization methods still apply!

$$L = |y - Wx_i|^2 + \lambda |W|$$

where |W| is element-wise

Example regularizations:

- L1/L2 on weights (encourage small values)
- L2: $L = |y Wx_i|^2 + \lambda |W|^2$ (weight decay)
- Elastic L1/L2: $|y Wx_i|^2 + \alpha |W|^2 + \beta |W|$

Regularization





Problem: Network can learn to rely strong on a few features that work really well

• May cause **overfitting** if not representative of test data

From: Dropout: A Simple Way to Prevent Neural Networks from Overfitting, Srivastava et al.

Preventing Co-Adapted Features





An idea: For each node, keep its output with probability *p*

Activations of deactivated nodes are essentially zero

Choose whether to mask out a particular node each iteration

From: Dropout: A Simple Way to Prevent Neural Networks from Overfitting, Srivastava et al.

Dropout Regularization



In practice, implement with a mask calculated each iteration

 During testing, no nodes are dropped



From: Dropout: A Simple Way to Prevent Neural Networks from Overfitting, Srivastava et al.





- During training, each node has an expected *p* * *fan_in* nodes
- During test all nodes are activated
- Principle: Always try to have similar train and test-time input/output distributions!

Solution: During test time, scale outputs (or equivalently weights) by p

• i.e. $W_{test} = pW$



From: Dropout: A Simple Way to Prevent Neural Networks from Overfitting, Srivastava et al.

Inference with Dropout



Interpretation 1: The model should not rely too heavily on particular features

 If it does, it has probability 1 – p of losing that feature in an iteration



From: Dropout: A Simple Way to Prevent Neural Networks from Overfitting, Srivastava et al.





Interpretation 1: The model should not rely too heavily on particular features

 If it does, it has probability 1 – p of losing that feature in an iteration

Interpretation 2: Training 2^n networks:

- Each configuration is a network
- Most are trained with 1 or 2 minibatches of data



From: Dropout: A Simple Way to Prevent Neural Networks from Overfitting, Srivastava et al.





Optimizers



Deep learning involves complex, compositional, non-linear functions

The **loss landscape** is extremely **nonconvex** as a result

There is **little direct theory** and a **lot of intuition/rules of thumbs** instead

 Some insight can be gained via theory for simpler cases (e.g. convex settings)







It used to be thought that existence of local minima is the main issue in optimization

There are other **more impactful issues**:

- Noisy gradient estimates
- Saddle points
- Ill-conditioned loss surface



From: Identifying and attacking the saddle point problem in highdimensional non-convex optimization, Dauphi et al., 2014.





We use a subset of the data at each iteration to calculate the loss (& gradients)

- This is an unbiased estimator but can have high variance
- This results in noisy steps in gradient descent

 $L = \frac{1}{M} \sum_{i} L(f(x_i, W), y_i)$





Several **loss surface geometries** are difficult for optimization

Several types of minima: Local minima, plateaus, saddle points

Saddle points are those where the gradient of orthogonal directions are zero

 But they disagree (it's min for one, max for another)









- Gradient descent takes a step in the steepest direction (negative gradient)
- Intuitive idea: Imagine a ball rolling down loss surface, and use momentum to pass flat surfaces

 $v_i = \beta v_{i-1} + \frac{\partial L}{\partial w_{i-1}}$ Update Velocity (starts as 0, $\beta = 0.99$)

 $w_i = w_{i-1} - \alpha v_i$ Update Weights

• Generalizes SGD ($\beta = 0$)

$$w_i = w_{i-1} - \alpha \frac{\partial L}{\partial w_i}$$



Adding Momentum



Velocity term is an exponential moving average of the gradient

$$v_i = \beta v_{i-1} + \frac{\partial L}{\partial w_{i-1}}$$

$$v_{i} = \beta(\beta v_{i-2} + \frac{\partial L}{\partial w_{i-2}}) + \frac{\partial L}{\partial w_{i-1}}$$
$$= \beta^{2} v_{i-2} + \beta \frac{\partial L}{\partial w_{i-2}} + \frac{\partial L}{\partial w_{i-1}}$$

There is a general class of accelerated gradient methods, with some theoretical analysis (under assumptions)

Accelerated Descent Methods



Equivalent formulation:

$$egin{aligned} &v_i = eta v_{i-1} - lpha rac{\partial L}{\partial w_{i-1}} & ext{Update Velocity} \ & ext{(starts as 0)} \end{aligned}$$







Key idea: Rather than combining velocity with current gradient, go along velocity **first** and then calculate gradient at new point

 We know velocity is probably a reasonable direction

$$\widehat{w}_{i-1} = w_{i-1} + \beta v_{i-1}$$

$$v_i = \beta v_{i-1} + \frac{\partial L}{\partial \widehat{w}_{i-1}}$$

$$w_i = w_{i-1} - \alpha v_i$$



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



Georgia ∤ Tech ∦
Momentum

Note there are **several equivalent formulations** across deep learning frameworks!

Resource:

https://medium.com/the-artificialimpostor/sgd-implementation-inpytorch-4115bcb9f02c





• Activation Functions: Use ReLU, GeLU, etc.

• Initialization: Important for initial activation and gradient statistics

• **Normalization:** Use dynamic normalization with learnable parts

- **Optimization:** Use momentum (helps w/ local minima, etc.)
 - Next: More sophisticated gradient history/statistics in update rule





- Various mathematical ways to characterize the loss landscape
- If you liked Jacobians... meet:



 Gives us information about the curvature of the loss surface

Hessian and Loss Curvature



Condition number is the ratio of the largest and smallest eigenvalue

 Tells us how different the curvature is along different dimensions

If this is high, SGD will make **big** steps in some dimensions and **small** steps in other dimension

Second-order optimization methods divide steps by curvature, but expensive to compute







Per-Parameter Learning Rate

Idea: Have a dynamic learning rate for each weight

Several flavors of **optimization algorithms**:

- RMSProp
- Adagrad
- Adam

SGD can achieve similar results in many cases but with much more tuning



Idea: Use gradient statistics to reduce learning rate across iterations

Denominator: Sum up gradients over iterations

Directions with **high curvature will have higher gradients**, and learning rate will reduce $G_{i} = G_{i-1} + \left(\frac{\partial L}{\partial w_{i-1}}\right)^{2}$ $w_{i} = w_{i-1} - \frac{\alpha}{\sqrt{G_{i} + \epsilon}} \frac{\partial L}{\partial w_{i-1}}$

As gradients are accumulated learning rate will go to zero

Duchi, et al., "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization"





Solution: Keep a moving average of squared gradients!

Does not saturate the learning rate

$$G_{i} = \beta G_{i-1} + (1 - \beta) \left(\frac{\partial L}{\partial w_{i-1}}\right)^{2}$$

$$w_i = w_{i-1} - \frac{\alpha}{\sqrt{G_i + \epsilon}} \frac{\partial L}{\partial w_{i-1}}$$





Combines ideas from above algorithms

Maintains both first and second moment statistics for gradients

$$v_i = \beta_1 v_{i-1} + (1 - \beta_1) \left(\frac{\partial L}{\partial w_{i-1}} \right)$$

$$G_{i} = \beta_{2} G_{i-1} + (1 - \beta_{2}) \left(\frac{\partial L}{\partial w_{i-1}}\right)^{2}$$

$$w_i = w_{i-1} - \frac{\alpha v_i}{\sqrt{G_i + \epsilon}}$$

But unstable in the beginning (one or both of moments will be tiny values)

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015





Solution: Time-varying bias correction

Typically $\beta_1 = 0.9$, $\beta_2 = 0.999$

So \hat{v}_i will be small number divided by (1-0.9=0.1) resulting in more reasonable values (and \hat{G}_i larger)

$$v_{i} = \beta_{1} v_{i-1} + (1 - \beta_{1}) \left(\frac{\partial L}{\partial w_{i-1}}\right)$$
$$G_{i} = \beta_{2} G_{i-1} + (1 - \beta_{2}) \left(\frac{\partial L}{\partial w_{i-1}}\right)^{2}$$

$$\widehat{v}_{i} = \frac{v_{i}}{1 - \beta_{1}^{t}} \quad \widehat{G}_{i} = \frac{G_{i}}{1 - \beta_{2}^{t}}$$
$$w_{i} = w_{i-1} - \frac{\alpha \,\widehat{v}_{i}}{\sqrt{\widehat{G}_{i} + \epsilon}}$$





Optimizers behave differently depending on landscape

Different behaviors such as **overshooting**, **stagnating**, **etc**.

Plain SGD+Momentum can generalize better than adaptive methods, but requires more tuning

See: Luo et al., Adaptive Gradient Methods with Dynamic Bound of Learning Rate, ICLR 2019



From: https://mlfromscratch.com/optimizers-explained/#/



Georgia Tech First order optimization methods have learning rates

Theoretical results rely on **annealed** learning rate

Several schedules that are typical:

- Graduate student!
- Step scheduler
- Exponential scheduler
- Cosine scheduler



From: Leslie Smith, "Cyclical Learning Rates for Training Neural Networks"

Learning Rate Schedules

