

Revisiting Line Searches

Machine Learning Reading Group Fall 2025

Betty Shea

2025-11-18

University of British Columbia

```
for a in range(1, num_epoch+1):  
    for (inputs, targets) in train_dataloader:  
        optimizer.zero_grad()  
        outputs = model(inputs)  
        loss = criterion(outputs, targets)  
        loss.backward()  
        optimizer.step()
```

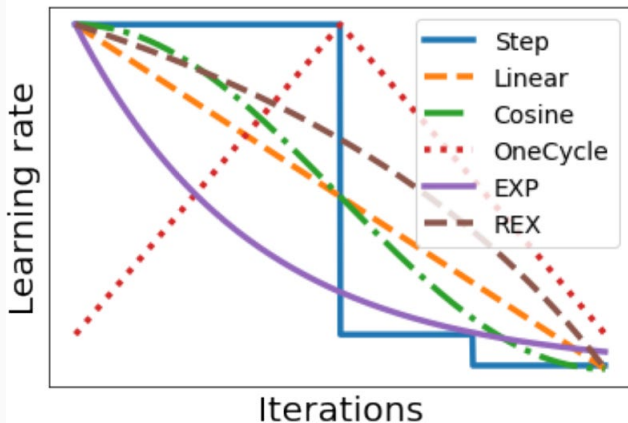
```
for a in range(1, num_epoch+1):  
    for (inputs, targets) in train_dataloader:  
        optimizer.zero_grad()  
        outputs = model(inputs)  
        loss = criterion(outputs, targets)  
        loss.backward()  
        optimizer.step()
```

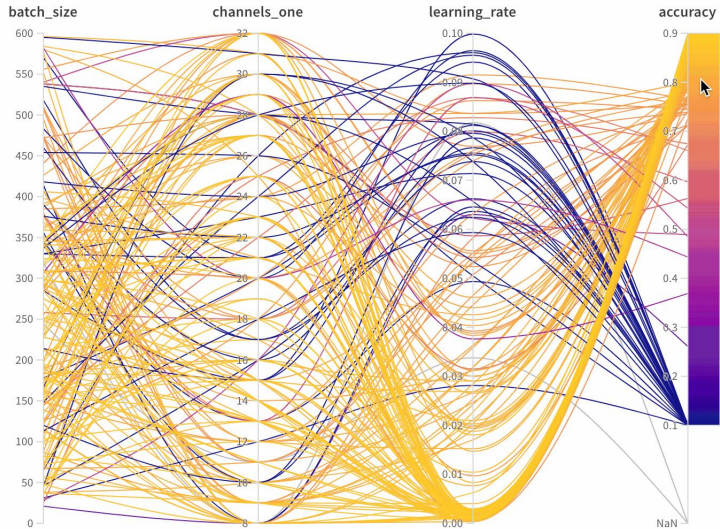
$$w_{t+1} = w_t + \alpha_t p_t$$

```
for a in range(1, num_epoch+1):  
    for (inputs, targets) in train_dataloader:  
        optimizer.zero_grad()  
        outputs = model(inputs)  
        loss = criterion(outputs, targets)  
        loss.backward()  
        optimizer.step()
```

$$w_{t+1} = w_t + \alpha_t p_t$$

$$\alpha_t = ?? \qquad \alpha_t = \alpha$$



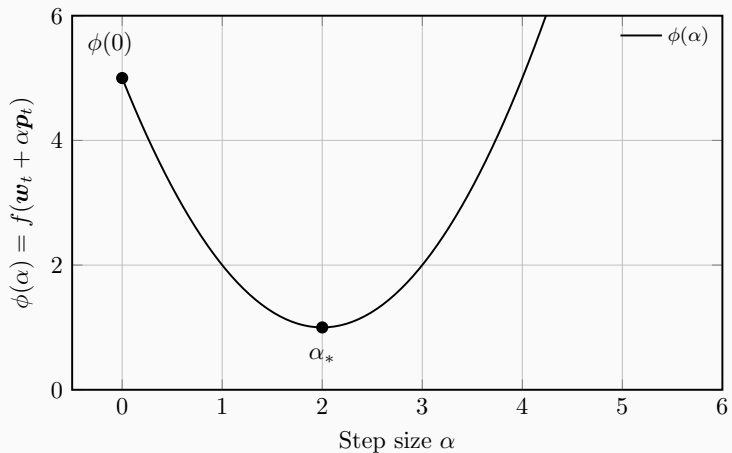


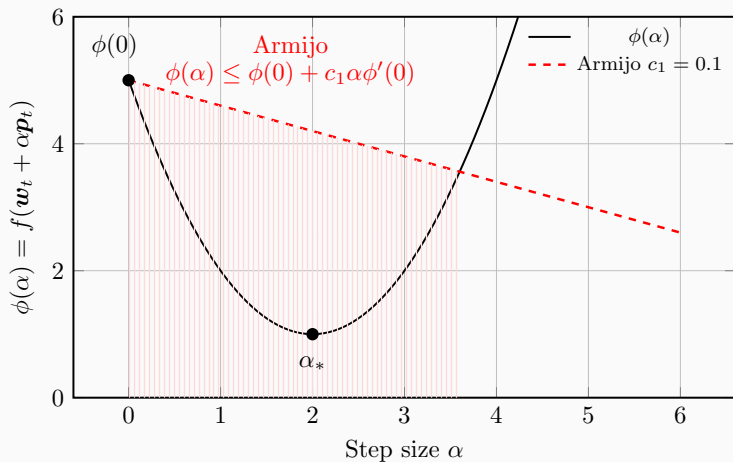
Can we set α_t on the fly? Some ideas outside of the scope of this talk:

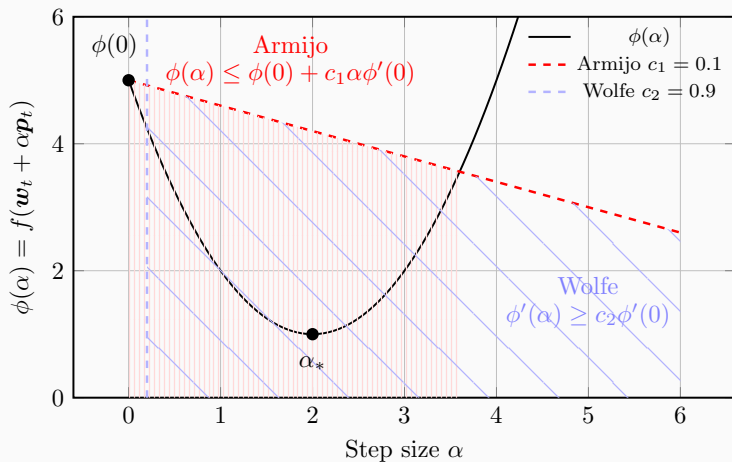
- Barzilai-Borwein
- Polyak step sizes (ask Curtis)
- Preconditioning
- Hyper-gradient descent

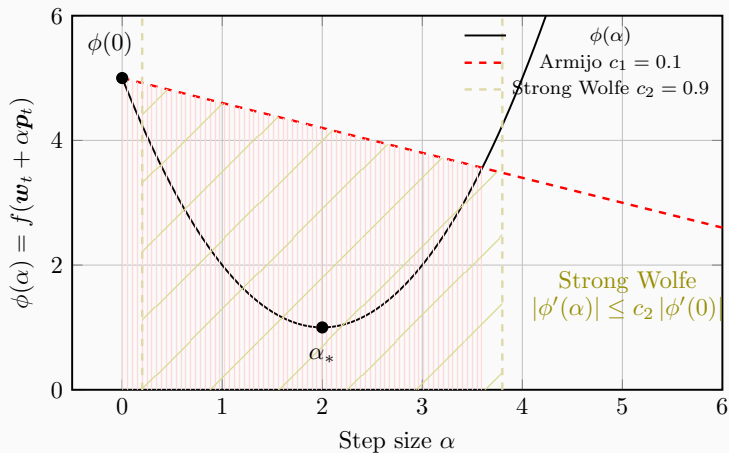
Sub-problem within every iteration of the main loop.

- Outer loop: training loop.
- Inner loop: find “good” learning rate α given fixed
 - weights w_t
 - search direction p_t
- $\phi(\alpha) = f(w_t + \alpha p_t)$
- $\phi'(0) = \nabla f(w_t)^\top p_t$, often assumed to be negative









Why aren't line searches everywhere?



1. Stochasticity
2. Ineffectiveness
3. Cost per iteration

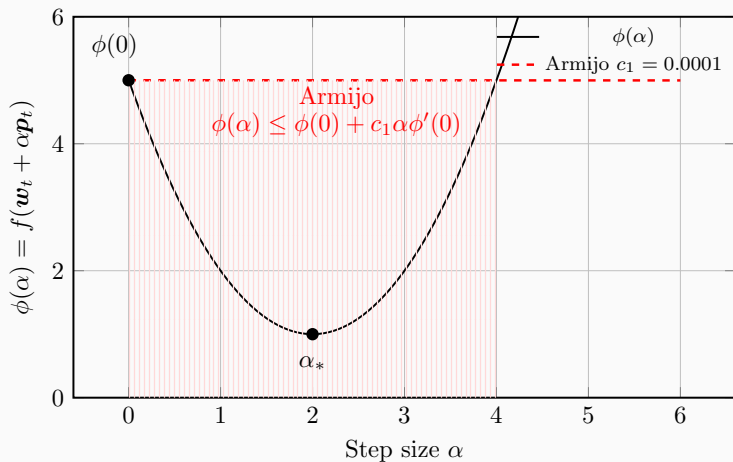
- Line searches were originally designed for the deterministic setting.

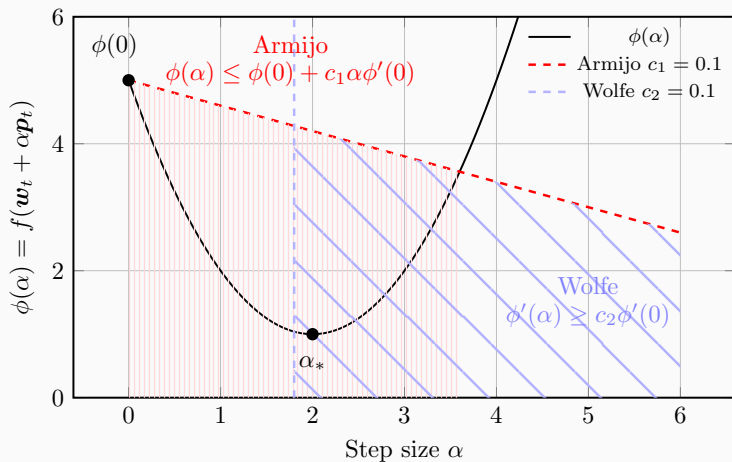
- Line searches were originally designed for the deterministic setting.
- There are theoretical results for stochastic versions.
 - Stochastic backtracking Armijo line searches [Vaswani et al. 2019, 2020, 2025]
 - Polyak Non-Monotone Stochastic backtracking Armijo line search [Galli et al. 2023] (ask Curtis)

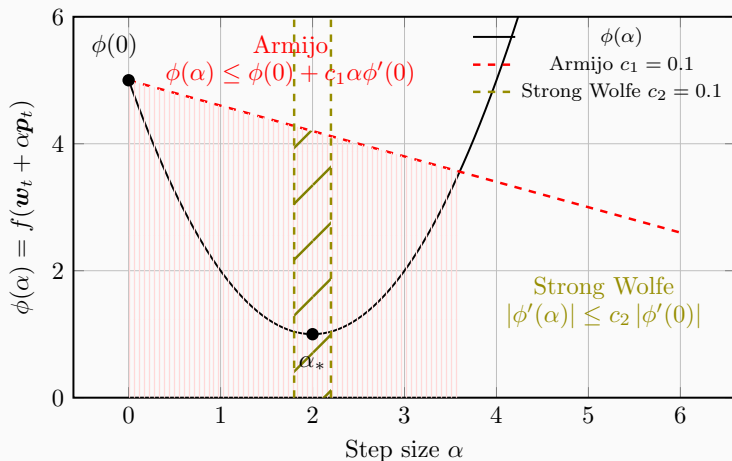
- Line searches were originally designed for the deterministic setting.
- There are theoretical results for stochastic versions.
 - Stochastic backtracking Armijo line searches [Vaswani et al. 2019, 2020, 2025]
 - Polyak Non-Monotone Stochastic backtracking Armijo line search [Galli et al. 2023] (ask Curtis)
- Some indication that noise from using mini-batches is beneficial for a line search. [Roulet et al. 2024]

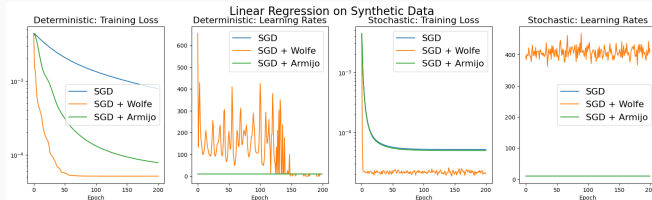
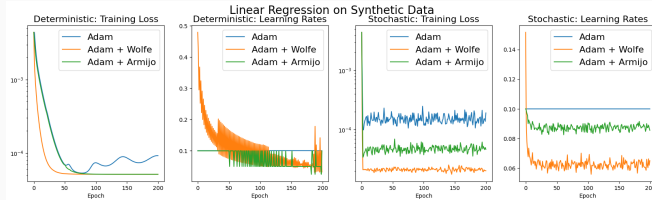
- Are we really doing a line search?

- Are we really doing a line search?
- Usually backtracking Armijo with very loose conditions
- More of a sanity check than a search

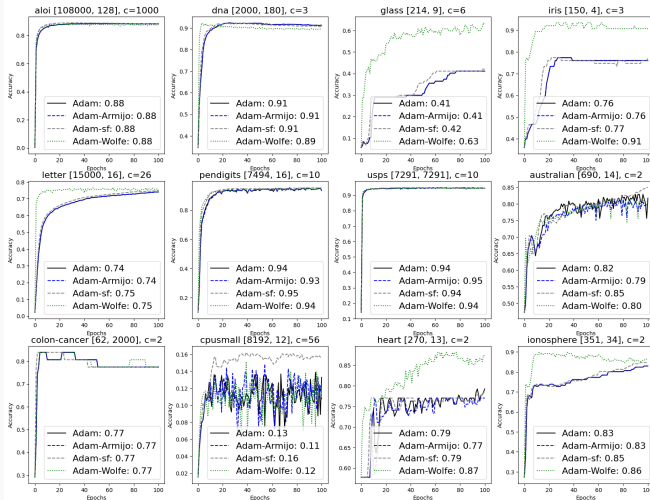




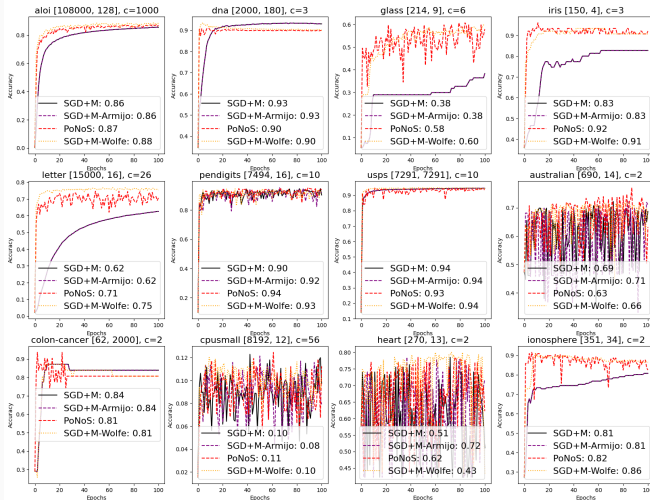




Logistic Regression Test Accuracy, Batch Size = 32



Logistic Regression Test Accuracy, Batch Size = 32



Searching is costly because of additional forward passes

- Generally, number of backtracks is closely monitored
- Wolfe is worse - requires additional backward passes too

Searching is costly because of additional forward passes

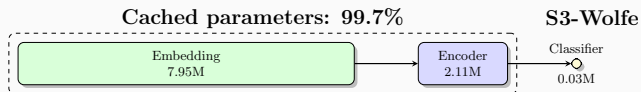
- Generally, number of backtracks is closely monitored
- Wolfe is worse - requires additional backward passes too

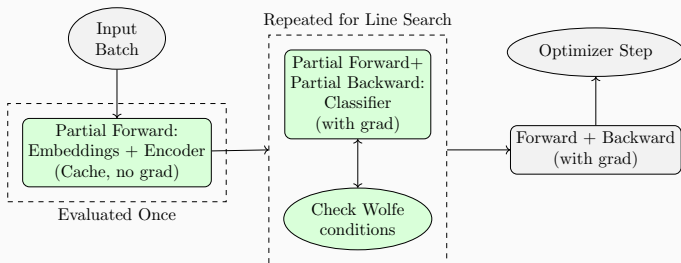
Is a line search needed on all model parameters?

Maybe not in all cases....

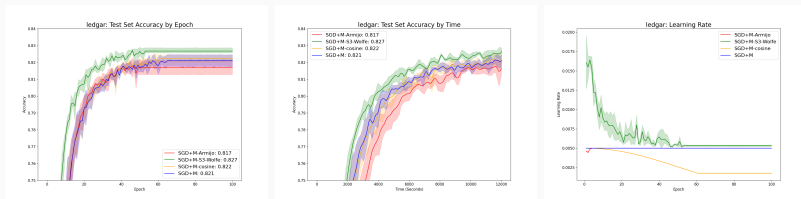
- Almost all parameters in BERT model in embedding and encoder layers
- Classifier parameters change most during fine-tuning
- Classifier parameters are a small fraction of all weights

- Line search accurately on small subset of parameters
- “Premium optimization” on selective parameters
- Similar to using parameter groups





- Two forward passes and one backward pass regardless of accuracy
- With some work, can get rid of the final forward pass



Full fine-tuning BertForSequenceClassification model initialized with pretrained legal-bert-base-uncased checkpoint on the ledger dataset.

Using a line search requires considering tradeoffs.

1. Accuracy vs cost per iteration

Using a line search requires considering tradeoffs.

1. Accuracy vs cost per iteration, or
2. Selectivity vs cost per iteration

Using a line search requires considering tradeoffs.

1. Accuracy vs cost per iteration, or
2. Selectivity vs cost per iteration

With large models, tradeoff 1 will always sacrifice accuracy.

Tradeoff 2 asks which parameters are worth spending optimization effort.

- Randomly initialized weights
 - Line search can give learning rates that are too far from the fixed learning rates used in the rest of the model
- Other criteria for search other than decreasing loss
 - Existing work on curvature aware learning rates suggest that being too greedy can hurt overall performance