**Trinity College Dublin**
Coláiste na Tríonóide, Baile Átha Cliath
The University of Dublin

# Final Assignment
# CS7DS2 Optimisation Algorithms for Data Analysis

Ujjayant Kadian (22330954)

April 29, 2025

In this report, we explore the use of Polyak's adaptive step size in mini-batch stochastic gradient descent (SGD), an approach aimed at improving convergence behavior in optimization problems. The objective is to implement and analyze this method through a series of experiments involving synthetic and text datasets using various models, including linear regression and transformers. By comparing Polyak's step size with traditional constant step size and other optimizers like Adam, we aim to understand its effectiveness in terms of convergence speed, robustness to noise, and generalization performance.

# 1 (a) Implementation of Polyak Adaptive Step Size Optimizer

This section describes the implementation of a custom PyTorch optimizer that incorporates the Polyak adaptive step size for mini-batch Stochastic Gradient Descent (SGD), following the formulation proposed by Loizou et al. (2021).

The primary objective was to modify the conventional SGD update rule to employ an adaptive step size $\alpha$ at each iteration, computed based on the current mini-batch loss and gradient. Specifically, for a mini-batch $\mathcal{N}$, the step size is calculated as:

$$\alpha = \frac{f_\mathcal{N}(\theta) - f_\mathcal{N}^*}{\nabla f_\mathcal{N}(\theta)^T \nabla f_\mathcal{N}(\theta) + \epsilon} = \frac{f_\mathcal{N}(\theta) - f_\mathcal{N}^*}{\|\nabla f_\mathcal{N}(\theta)\|^2 + \epsilon}$$

where $f_\mathcal{N}(\theta)$ denotes the mini-batch loss, $\nabla f_\mathcal{N}(\theta)$ its gradient, $f_\mathcal{N}^*$ the minimum loss (in practical scenarios, the minimum loss of a function $f$ is ideally 0, thus it is generally set to 0), and $\epsilon$ a small positive constant to ensure numerical stability (avoiding divison by 0).

A custom optimizer class, `PolyakSGD`, was implemented by subclassing `torch.optim.Optimizer`. The key components are as follows:

**Initialization (`__init__`)**: Accepts model parameters along with optional arguments for $\epsilon$ (default $10^{-8}$) and $f^*$ (default $0.0$). Basic validation ensures $\epsilon$ is strictly positive.

**Optimization Step (`step`)**: Expects a closure function that clears existing gradients, performs a forward pass to compute $f_\mathcal{N}(\theta)$, and backpropagates to compute $\nabla f_\mathcal{N}(\theta)$. The squared $L_2$ norm of the gradient is computed by summing the squares of all gradient components across parameter groups. The step size $\alpha$ is then calculated, ensuring non-negativity with a safeguard $\max(0.0, \alpha)$. Parameters are updated via:

$$\theta \leftarrow \theta - \alpha \nabla f_\mathcal{N}(\theta)$$

This update is performed within a `torch.no_grad()` context to avoid tracking in the computational graph.

**Utility Features**: A property `last_step_size` stores the most recent step size for debugging and analysis.

**Testing Methodology and Results**

To ensure the correctness of the implementation, a comprehensive set of tests was developed using `pytest`, targeting critical components individually.

**Gradient Norm Calculation (`test_gradient_norm_calculation`)**: Verifies accurate computation of the squared $L_2$ norm using manually specified gradients. For example, with gradients $[1, 2, 3]$, the expected result is $1^2 + 2^2 + 3^2 = 14$.

**Step Size Formula (`test_step_size_formula`)**: Validates correct computation of $\alpha$ based on known loss and gradient values. For instance, with loss 7.0, target 2.0, and gradient norm 25, the expected step size is approximately:

$$\frac{7.0 - 2.0}{25 + \epsilon} \approx 0.2$$

**Quadratic Optimization (`test_quadratic_optimization`)**: Assesses the optimizer's ability to minimize a simple convex function $f(\theta) = \theta^2$, ensuring convergence towards the global minimum at $\theta = 0$.

**Synthetic Linear Regression (`test_synthetic_regression`)**: Evaluates the optimizer in a realistic machine learning scenario using synthetic data and mean squared error loss. Convergence towards true model parameters is checked across epochs (similar to the next part 1 (b)).

**Input Validation (`test_invalid_eps`)**: Ensures that improper initialization with non-positive $\epsilon$ triggers appropriate error handling.

**Usage Pattern (`test_usage_example`)**: Confirms that the optimizer integrates correctly within a standard mini-batch training loop using closures.

All tests passed successfully. The results demonstrate that the `PolyakSGD` optimizer accurately computes the gradient norms and adaptive step sizes, correctly updates parameters, and operates reliably both in controlled mathematical settings and standard machine learning tasks. These outcomes provide confidence in the robustness and correctness of the implementation, justifying its use in the subsequent experimental investigations in parts 1(b) through 1(e).

# 1(b) Synthetic Data Comparison: Constant vs. Polyak Step Size

This section investigates the performance of the implemented Polyak step size for mini-batch SGD (PolyakSGD) compared to standard SGD with a constant step size (SGD). The comparison is performed on a synthetic linear regression task, analyzing the impact of varying mini-batch sizes and levels of noise in the training data.

**Experimental Setup**

- Data Generation: Synthetic data was generated using the `make_synthetic_linreg` function. This function creates $N$ data points $(X, y)$ according to the model $y = Xw_{\text{true}} + b_{\text{true}} + \epsilon$, where $X \in \mathbb{R}^{N \times d}$, $w_{\text{true}} \in \mathbb{R}^d$, $b_{\text{true}} \in \mathbb{R}$, and $\epsilon$ is Gaussian noise with a specified standard deviation (`noise_std`). For these experiments, $N = 1000$ samples and $d = 20$ dimensions were used, as defined in the configuration file (`synthetic_regression_config.yaml`). Noise levels tested were $\sigma \in \{0.1, 1.0, 5.0\}$.

- Model and Loss: A simple linear regression model (LinReg) consisting of a single `torch.nn.Linear` layer was used, trained with the Mean Squared Error loss (`torch.nn.MSELoss`)

- Optimizers:

    - **SGD:** Standard SGD with a constant learning rate. Crucially, the learning rate for SGD was optimized for each combination of noise level, batch size, and random seed using a grid search (`find_best_lr` function over `learning_rates: [0.001, 0.01, 0.1, 1.0]`) to ensure a fair comparison against the adaptive method. The best learning rate was chosen based on achieving the lowest final loss without significant oscillations.

    - **PolyakSGD:** The custom optimizer implemented in section 1(a), using $f^* = 0$ and $\epsilon = 10^{-8}$.

- Training: Models were trained for 20 epochs (`config['training']['epochs']`). Experiments were run across different mini-batch sizes: 1, 16, 64, 256, "full" (where "full" means $N = 1000$) and multiple random seeds: 0, 42, 1337 (`config['data']['seeds']`). Summary plots showing final loss vs. noise/batch size (Figures 1 and 2) present results across these seeds, while plots illustrating specific dynamics like step-size variation or convergence speed (Figures 3–5) use specific seeds (e.g., seed 42) for clarity. The core training loop is managed by the `train_model` function within the `synthetic_regression.py` script.

**Analysis of Results**

The experiments compared the final loss achieved, the convergence speed, and the behavior of the Polyak step size itself under different conditions.

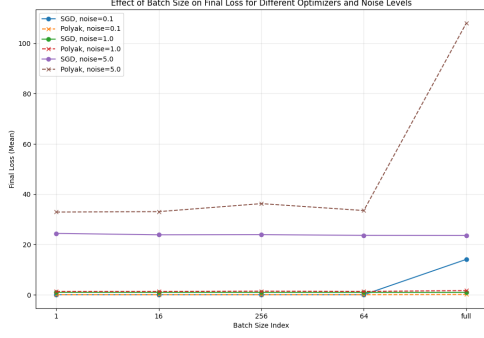**Impact of Noise and Batch Size on Final Loss:**

Figure 1: Effect of Batch Size on Final Loss for Different Optimizers and Noise Levels
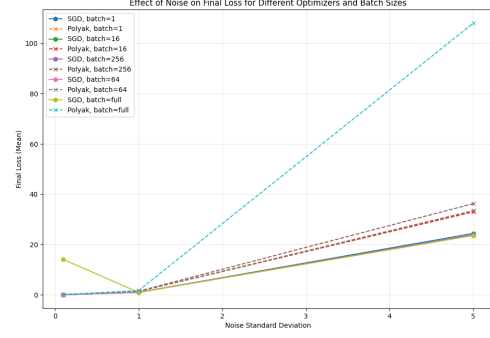
Figure 2: Effect of Noise on Final Loss for Different Optimizers and Batch Sizes

Figures 1 and 2 illustrate the interplay between noise, batch size, and the choice of optimizer on the final training loss (averaged over seeds).

1. Low Noise ($\sigma = 0.1$): Polyak SGD demonstrates remarkable stability across all batch sizes, achieving a low final loss consistently. Constant SGD (with its tuned LR) also performs well but shows a significant degradation in performance when using the full batch size. This suggests that even with a tuned learning rate, full-batch gradient descent can struggle in low-noise scenarios, potentially overshooting the minimum, whereas Polyak's adaptive nature handles this regime better.

2. Moderate Noise ($\sigma = 1.0$): Both optimizers perform well and achieve similar low final losses across all batch sizes. The adaptivity of Polyak and the careful tuning of the constant LR lead to comparable stable performance in this regime.

3. High Noise ($\sigma = 5.0$): A clear difference emerges. Constant SGD, benefiting from the optimized learning rate for each setting, maintains a relatively stable (though higher) final loss across all batch sizes. In contrast, Polyak SGD's performance degrades significantly as the batch size increases, performing particularly poorly with the full batch. This indicates that Polyak's step size calculation $\alpha = \frac{f_{\mathcal{N}} - f^*}{\|\nabla f_{\mathcal{N}}\|^2 + \epsilon}$ becomes highly sensitive to noise when averaged over larger batches. High noise can inflate the loss estimate $f_{\mathcal{N}}$ while potentially reducing the relative magnitude of the gradient norm $\|\nabla f_{\mathcal{N}}\|^2$, leading to excessively large and detrimental steps. Tuned constant SGD appears more robust in this high-noise setting.
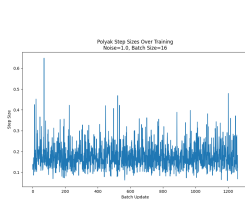
**Polyak Step Size Behavior:**



Figure 3.1: Polyak Step Size Variation (Batch=16, Noise=1.0, Seed=42)
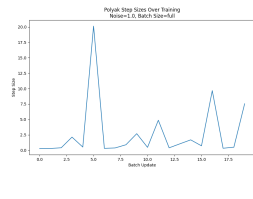


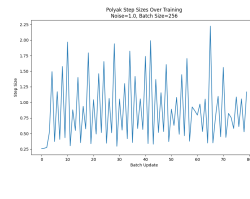Figure 3.2: Polyak Step Size Variation (Batch=Full, Noise=1.0, Seed=42)



Figure 4.1: Polyak Step Size Variation (Batch=256, Noise=1.0, Seed=42)



Figure 4.2: Polyak Step Size Variation (Batch=256, Noise=1.0, Seed=1337)

Figures 3 and 4 provide insight into how the Polyak step size $\alpha$ dynamically changes during training.

- Effect of Batch Size: As seen in Figure 3 (comparing batch=16 vs. full batch for noise=1.0, seed=42), smaller batch sizes (like 16) lead to more frequent updates but generally smaller, though still variable, step sizes. The step size reacts to the noisy gradient and loss estimates from each small batch. With the full batch, updates are infrequent (one per epoch), but the step sizes can become extremely large and erratic (spiking over 20 in the example). This suggests that even with noise averaged over the full dataset, the resulting gradient norm can sometimes become very small relative to the current loss, causing the step size formula to yield huge values.

- Interaction with SGD Noise: The step size calculation directly uses the mini-batch loss $f_{\mathcal{N}}$ and gradient norm $\|\nabla f_{\mathcal{N}}\|^2$, both of which are affected by the randomness (noise $\epsilon$ in the data and sampling randomness for batch sizes < full). High noise can lead to poor estimates, causing volatility in $\alpha$. If noise happens to create a situation where the gradient norm is small while the loss estimate remains relatively high, the step size can explode. This explains the poor performance observed in Figure 2 for Polyak with large batches and high noise.

- Behavior Near Minimum: The assignment asks why the step size might increase near the minimum. Theoretically, as $\theta$ approaches the minimum $\theta^*$, both the numerator ($f_{\mathcal{N}}(\theta) - f^*$) and the denominator ($\|\nabla f_{\mathcal{N}}(\theta)\|^2 + \epsilon$) should approach zero. If the gradient norm shrinks faster than the loss difference, the step size could indeed increase. However, the plots (e.g., Figure 3.1 for seed 42) don't show a clear increasing trend towards the end; instead, they show continued variability driven by noise. The high volatility, especially with larger batches (Figure 3.2, Figure 4), suggests that noise often dominates the step size calculation, preventing a smooth theoretical behavior near the minimum.

- Consistency Across Seeds: Figure 4 shows step size behavior for the same batch size (256) and noise (1.0) but different random seeds (42 vs 1337). While both runs exhibit large, volatile step sizes, the exact patterns differ, highlighting the sensitivity to the specific data sampling and noise realization (with Figure 4.2 being comparatively smoother around epoch 30).
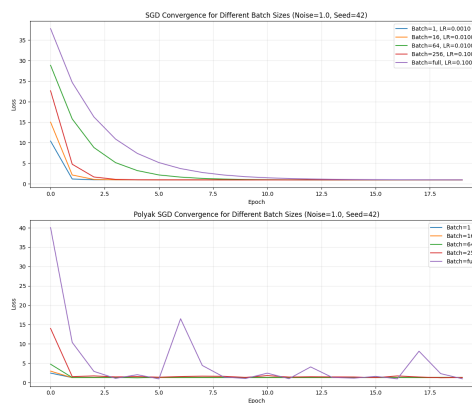
**Convergence Rate Comparison:**



Figure 5: Convergence Speed Comparison (SGD vs. Polyak, Noise=1.0, Seed=42)

Small/Medium Batches (1, 16, 64, 256): In the moderate noise case shown (seed=42), Polyak SGD generally converges faster than constant SGD (with its tuned LR), reaching the final loss plateau in fewer epochs. Its adaptive nature allows it to take more aggressive steps initially when the loss is high and gradients are informative.

Full Batch: Constant SGD shows slow but steady convergence. Polyak SGD initially decreases the loss rapidly but then exhibits extreme instability, with loss spiking erratically, consistent with the observed large step sizes (Figure 3.2). This highlights the risk associated with Polyak in the full-batch setting, even with moderate noise.

The experiments on synthetic linear regression data reveal a trade-off. Polyak SGD can offer faster convergence than even a well-tuned constant SGD, particularly with small-to-medium batch sizes and low-to-moderate noise, due to its adaptive step size. However, its performance is sensitive to noise, especially when combined with large batch sizes (including full-batch). In high-noise scenarios or full-batch settings, the step size calculation can become unstable, leading to poor convergence or divergence. Tuned constant-step SGD, while potentially converging slower, demonstrates greater robustness across different noise levels and batch sizes, provided the learning rate is appropriately selected for the specific regime. The choice between them depends on the noise characteristics of the problem and the computational budget for tuning a constant learning rate versus accepting the potential volatility of the adaptive method.

# 1(c) Comparison on Week 6 Assignment Data

This section evaluates the performance of constant step size SGD (SGD) and Polyak SGD (PolyakSGD) on the optimization problem defined in the Week 6 assignment. This problem features a non-convex loss landscape with multiple minima, providing a different challenge compared to the convex linear regression task in section 1(b).

**Experimental Setup**

- Model and Loss Function: The Week6Model directly optimizes a 2-dimensional parameter vector $x$. The loss function for a single data point $w$ (from the training data) is defined as:

$$l(x, w) = \min\left(20\|x - w - 1\|^2, \|(x - w - 1) + (9, 10)\|^2\right)$$

where $x, w \in \mathbb{R}^2$. The total loss $f(\theta)$ is the average of $l(x, w)$ over all $w$ in the training dataset (or mini-batch). This function is designed to have multiple local minima. Since the training data $w$ is generated from a Gaussian cluster centered near $(0, 0)$ (make_gaussian_cluster with std=0.25), the term $x - w - 1$ is approximately $x - (1, 1)$. The loss function thus has minima near $x \approx (1, 1)$ (where the first term is small) and near $x \approx (1, 1) - (9, 10) = (-8, -9)$ (where the second term is small). The minimum near $(-8, -9)$ is expected to be the global minimum due to the scaling factor of 20 applied to the first term.

- Data: The training data consists of $m = 25$ points generated using make_gaussian_cluster (seed=42, std=0.25), forming a tight cluster around the origin.

- Optimizers:

    - **SGD:** Constant step size SGD with a fixed learning rate lr = 0.01 (as specified in week6_experiment_config.yaml). Unlike section 1(b), the LR was not tuned per batch size for this experiment.

    - **PolyakSGD:** The adaptive step size optimizer with $f^* = 0$ and $\epsilon = 10^{-8}$.

- Training: Models were trained for 100 epochs (config['training']['epochs']). Experiments were run with mini-batch sizes: 1, 5, 10, 25 (where 25 corresponds to the full batch size for this dataset). The week6_experiment.py script orchestrates the training and logging, including tracking the parameter history (param_history) for plotting trajectories.

**Analysis of Results**
**Convergence Speed and Stability:**



Figure 6: Convergence Comparison (SGD vs. Polyak) for Week 6 Data

Figure 6 shows the loss curves over epochs for both optimizers across different batch sizes.

- **SGD:** Constant SGD exhibits smooth and stable convergence for all batch sizes. The final loss values are slightly lower for larger batch sizes (around 2.5 for batches 5, 10, 25 vs. 3.5 for batch 1), likely due to reduced gradient noise. Convergence speed is broadly similar across batch sizes.

- **PolyakSGD:** Polyak shows rapid initial convergence, often faster than SGD. However, its stability is highly dependent on the batch size. For batch size 1, it converges quickly to a low loss value (around 2.0, lower than any SGD variant) and remains stable. For larger batch sizes (5, 10, and especially 25), the loss becomes highly unstable after the initial descent, exhibiting large spikes. This instability, similar to the findings in section 1(b) but more pronounced here, suggests that the Polyak step size calculation struggles when averaging over larger batches on this non-convex landscape, potentially leading to large, detrimental steps.

5

**Optimization Trajectories and Minima:**



Figure 7a: SGD Trajectories on Loss Contour (Week 6)



Figure 7b: PolyakSGD Trajectories on Loss Contour (Week 6)

Figure 7 overlays the optimization paths on the loss function's contour plot.

- **SGD:** The trajectories for constant SGD are smooth and direct for all batch sizes. They consistently converge to the minimum near $x \approx (1, 1)$. This appears to be the local minimum closest to t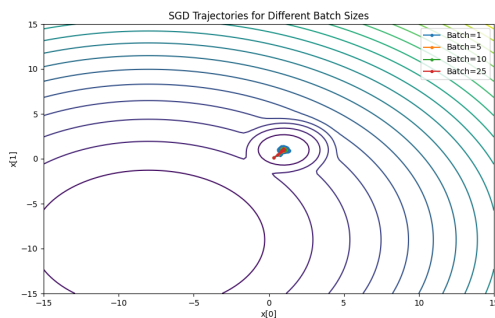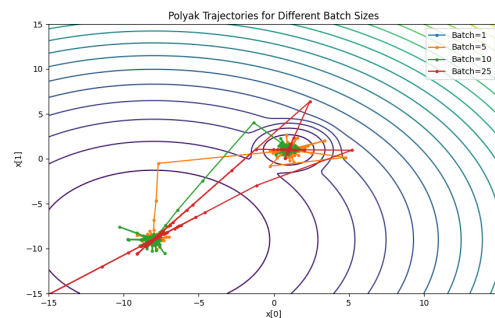he typical initialization point (random normal). The noise inherent in smaller batches (like batch=1) introduces minor oscillations but doesn't fundamentally alter the path towards this nearest minimum.

- **PolyakSGD:** The trajectories for Polyak SGD differ significantly. For batch size 1, the path is relatively stable and converges quickly and accurately to the minimum near $x \approx (-8, -9)$, which is believed to be the global minimum. For larger batch sizes (5, 10, 25), the trajectories become increasingly erratic. While they also tend towards the global minimum near $(-8, -9)$, they take large, unstable steps, often overshooting and oscillating wildly, especially with the full batch size (25). This corresponds to the loss spikes seen in Figure 6.

**Discussion:**

- **Escaping Local Minima:** Constant SGD, with the chosen fixed learning rate, consistently gets trapped in the local minimum near $(1, 1)$. It fails to escape and find the deeper minimum. Polyak SGD, particularly with small batch sizes (batch=1), successfully avoids the local minimum and converges to the global minimum near $(-8, -9)$. The adaptive, potentially larger steps taken by Polyak seem crucial in navigating the landscape and moving away from the basin of attraction of the local minimum. However, for larger batches, while Polyak still moves towards the global minimum, the steps become too large and unstable, hindering efficient convergence.

- **Convergence Points:** As noted, SGD converges to the local minimum, while Polyak (when stable) converges to the global minimum. They do not converge to the same point in this non-convex scenario.

- **Optimal Batch Sizes:** For constant SGD on this problem, larger batch sizes (5, 10, 25) yielded slightly lower final loss values compared to batch size 1, suggesting that reducing gradient noise was beneficial for converging accurately within the local minimum's basin. For Polyak SGD, the opposite was true: batch size 1 provided the best performance, achieving stable convergence to the global minimum with the lowest final loss. Larger batch sizes led to instability and degraded performance. Therefore, the optimal batch size differs significantly between the two methods for this problem. A small batch size is preferable for Polyak SGD to maintain stability while leveraging its ability to escape local minima, whereas a larger batch size might be slightly better for constant SGD if converging to the nearest minimum is acceptable.

In summary, for the non-convex Week 6 problem, Polyak SGD (specifically with batch size 1) demonstrated a clear advantage over constant SGD by finding the global minimum. However, its sensitivity to batch size was again apparent, with larger batches causing significant instability. Constant SGD was stable but failed to escape the local minimum.

# 1(d) Transformer Model Comparison: SGD vs. Polyak vs. Adam

This section extends the investigation to a more complex, real-world task: training a character-level GPT-style Transformer model, similar to the one used in the Machine Learning module's Week 9 assignment. The goal is

to compare the performance of constant step size SGD (SGD), Polyak SGD (PolyakSGD), and the widely used Adam optimizer (Adam). Performance is evaluated based on convergence speed, final loss, and overfitting characteristics.

**Experimental Setup**

- Model: A standard decoder-only Transformer model (`GPTLanguageModel` from `transformer_week9.py`) was used. The architecture includes multi-head self-attention, feed-forward layers, layer normalization, and dropout. The implementation was slightly restructured from the original Week 9 script for easier integration into the experimental framework. Key hyperparameters (embedding dimension $n\_embd = 192$, number of heads $n\_head = 4$, number of layers $n\_layer = 2$, dropout $dropout = 0.2$, block size $block\_size = 256$) were managed via the `TransformerDataManager` class, using values consistent with the Week 9 setup.

- Data: The model was trained on the `input_childSpeech_trainingSet.txt` dataset, managed by `TransformerDataManager`. This class handles tokenization (character-level), train/validation splitting (90/10 split), and batch generation (`batch_size=64`).

- Loss: Standard Cross-Entropy loss was used, suitable for language modeling.

- Optimizers:

    - **SGD:** Constant step size SGD.
    - **PolyakSGD:** The adaptive step size optimizer implemented in section 1(a) ($f^* = 0$, $\epsilon = 10^{-8}$).
    - **Adam:** The standard Adam optimizer from `torch.optim`.

- Tuning: As required by the assignment, SGD and Adam parameters were tuned before the final comparison runs to ensure they performed near their best.

    - The `tune_sgd_learning_rate` function performed a grid search over learning rates $[0.0001, 0.001, 0.01, 0.1]$ for SGD, selecting the one yielding the lowest validation loss after `tuning_iters=500`.
    - The `tune_adam_hyperparams` function performed a grid search over learning rates $[0.0001, 0.001, 0.01]$, $\beta_1$ values $[0.9, 0.95]$, and $\beta_2$ values $[0.999, 0.99]$ for Adam, again selecting the best combination based on validation loss after 500 iterations.
    - The final comparison runs used these tuned parameters (or the defaults specified in `week9_transformer_config.yaml` if tuning was disabled). Note: The specific tuned values used are logged in the experiment output files.

- Training: Models were trained for a maximum of `max_iters=2000` iterations, with evaluation performed every `eval_interval=200` iterations. Early stopping was enabled (patience=3, min_delta=0.001) to terminate training if validation loss did not improve for 3 consecutive evaluations, preventing excessive training time, especially for potentially unstable optimizers. The `week9_transformer.py` script managed the training, tuning, and result logging.

**Analysis of Results**



Figure 8: Training and Validation Loss



Figure 9: Generalization Gap



Figure 10: Polyak Step Size Evolution

**Convergence Speed and Stability:** Figure 8 shows the training and validation loss curves for the three optimizers using their (potentially tuned) parameters.

- **Adam:** Demonstrates the fastest convergence, rapidly decreasing both training and validation loss within the first 200–400 iterations to reach the lowest loss values among the three optimizers (around 0.5). It remains very stable thereafter.

- **SGD:** Converges much more slowly and steadily than Adam, taking the full 2000 iterations to reach its minimum loss (around 0.6 for training, 0.7 for validation). It exhibits stable behavior throughout training.

- **PolyakSGD:** Shows highly unstable behavior. After an initial decrease, the loss spikes dramatically around iteration 200 and continues to oscillate wildly between high ($\sim$5.5–7.0) and moderate ($\sim$2.8) loss values for the remainder of the training. It fails to converge effectively.

**Overfitting Analysis:**

Figure 9 plots the generalization gap (Validation Loss − Training Loss) over iterations. A smaller, stable gap indicates better generalization and less overfitting.

- **Adam:** Maintains a very small and stable generalization gap throughout training, increasing only slightly towards the end. This suggests Adam effectively controls overfitting for this model and dataset.

- **SGD:** Shows a larger generalization gap than Adam, which peaks around iteration 1000 before decreasing slightly. While stable, it indicates slightly more overfitting than Adam.

- **PolyakSGD:** The generalization gap mirrors its loss instability. It spikes significantly early on but then fluctuates around zero, even becoming negative at times (implying validation loss was momentarily lower than training loss, likely due to the noisy nature of the updates). When not spiking, its gap is competitive, suggesting the underlying model isn't inherently overfitting drastically, but the optimizer's instability prevents consistent performance.

**Polyak Step Size Behavior:**

Figure 10 shows the adaptive step size chosen by PolyakSGD during training.

The step sizes are extremely volatile, exhibiting sharp spikes (reaching values $> 30$) that coincide with the large loss spikes observed in Figure 8. Between spikes, the step size can be relatively small. The Exponential Moving Average (EMA) trend line shows a general increase in the smoothed step size during the latter half of training, suggesting growing instability.

This confirms that the instability observed in Polyak's loss curves is directly linked to its tendency to compute excessively large step sizes on this complex, high-dimensional optimization landscape, even with $f^* = 0$. The interaction between the mini-batch loss estimate and the gradient norm likely leads to near-zero denominators in the step size formula intermittently.

**Discussion Summary:**

On the Transformer language modeling task, Adam (after tuning) clearly provided the best performance, achieving rapid convergence to the lowest validation loss with minimal overfitting. Constant SGD (after tuning) was stable and generalized reasonably well but converged significantly slower than Adam. Polyak SGD proved highly unstable for this problem. While its adaptive nature might theoretically be beneficial, the step size calculation led to extreme volatility and prevented effective convergence.

The large spikes in step size suggest that the assumptions underlying the Polyak step rule (particularly regarding the relationship between loss value and gradient norm) may not hold well in this high-dimensional, non-convex setting, or that the mini-batch estimates are too noisy.

Evaluating overfitting was difficult for Polyak due to its instability, but when it wasn't spiking, its generalization gap was comparable to the others. Overall, for this task, the tuned Adam optimizer is the superior choice, while Polyak SGD is unsuitable due to instability.

# 1(e) Brief Investigation of Modified Adam with Polyak Step Size

Given the instability of Polyak SGD, especially under noise, this section briefly explores modifying Adam to incorporate the Polyak step size, aiming to leverage Adam's momentum for stability while retaining Polyak's adaptivity. The investigation uses the synthetic linear regression task from section 1(b).

**Implementation (PolyakAdam)**

A custom `PolyakAdam` optimizer was implemented (see Appendix). It combines standard Adam components (momentum, adaptive scaling via $\beta_1$, $\beta_2$) with the Polyak step calculation: $\alpha_{\text{Polyak}} = \max\left(0.0, \frac{f_\mathcal{N} - f^*}{\|\nabla f_\mathcal{N}\|^2 + \epsilon}\right)$.

Instead of using $\alpha_{\text{Polyak}}$ directly, it employs a hybrid update: aggressive boosted steps early on, followed by standard Adam updates where the step size is scaled by a factor influenced by $\alpha_{\text{Polyak}}$ (controlled by `polyak_factor`, capped to prevent excessive steps). An AMSGrad option was included for potential robustness (theoretically prevents overly large steps late in training) but examined only briefly.

**Experimental Setup**

PolyakAdam (with/without AMSGrad, using default parameters: `alpha=0.003`, `betas=(0.9, 0.999)`, `polyak_factor=0.3`) was tested on the synthetic linear regression task ($N = 1000$, $d = 20$, MSE loss), focusing on the high noise case ($\sigma = 5.0$). It was compared against SGD (tuned learning rate) and PolyakSGD.

**Results and Discussion**



Figure 11: Step Sizes Over Training (Polyak Variants, Noise=0.1, Batch=256)



Figure 12: Convergence Comparison (All Optimizers, Noise=5.0, Seed=42)

**Convergence:** As shown in Figure 12, `PolyakAdam` significantly improves stability compared to `PolyakSGD` under high noise. It effectively reduces loss, often outperforming tuned SGD, especially with smaller batches. The extreme instability of `PolyakSGD` is largely suppressed. The AMSGrad variant offered only marginal additional stability.

**Step Size Behavior:** Figure 11 (using noise=0.1 for clarity) reveals that while `PolyakSGD`'s step size oscillates, the effective step size in `PolyakAdam` (and `PolyakAdam+AMSGrad`) tends to escalate and become volatile later in training. This suggests the interaction between Adam's state and the Polyak boost factor creates a different instability mechanism, potentially amplifying steps as Adam's internal denominator shrinks.

Incorporating the Polyak step calculation into Adam (`PolyakAdam`) successfully mitigates the instability of pure `PolyakSGD` in high-noise scenarios, leading to better convergence. However, this hybrid approach introduces its own dynamic where step sizes can escalate later in training. While promising, `PolyakAdam` would require careful tuning of its specific hyperparameters (`alpha`, `polyak_factor`) and further study to be considered a robust alternative.

# 2. Optimization Concepts

This section addresses fundamental concepts in optimization related to step size selection and constraint handling.

## (a) Line Search in Gradient Descent

Line search is a strategy used in iterative optimization methods like gradient descent to determine an appropriate step size $\alpha_k$ at each iteration $k$. Instead of using a fixed step size, line search aims to find a step size that provides sufficient decrease in the objective function $f(\theta)$ along the chosen search direction $p_k$ (which is typically the negative gradient, $p_k = -\nabla f(\theta_k)$, for gradient descent). The update rule is $\theta_{k+1} = \theta_k + \alpha_k p_k$.

The core idea is to approximately solve a one-dimensional optimization problem at each iteration:

$$\min_{\alpha > 0} f(\theta_k + \alpha p_k)$$

In practice, finding the exact minimum is often too expensive. Instead, inexact line search methods are commonly used. These methods seek an $\alpha_k$ that satisfies certain conditions, such as the Wolfe conditions

or Armijo condition, which guarantee sufficient decrease in the function value and prevent excessively small steps. For example, the Armijo condition requires:

$$f(\theta_k + \alpha_k p_k) \leq f(\theta_k) + c_1 \alpha_k \nabla f(\theta_k)^T p_k$$

for some constant $c_1 \in (0, 1)$. A common approach is *backtracking line search*, where one starts with an initial guess for $\alpha$ (e.g., $\alpha = 1$) and iteratively reduces it (e.g., $\alpha \leftarrow \beta\alpha$ for $\beta \in (0, 1)$) until the condition is met.

**Advantage:**

*Improved Convergence/Stability*: Line search can lead to more robust and sometimes faster convergence compared to a poorly chosen fixed step size, as it adapts the step length to the local geometry of the loss landscape. It helps avoid overshooting the minimum or taking excessively small steps.

**Disadvantage:**

*Computational Cost*: Performing the line search procedure at each iteration requires additional function (and sometimes gradient) evaluations, increasing the computational cost per iteration compared to using a fixed step size.

**Use with Stochastic Gradient Descent (SGD):**

Line search is generally not used with standard mini-batch SGD. The primary reasons are:

- **Noisy Gradients:** The gradient $\nabla f_{\mathcal{N}}(\theta)$ computed on a mini-batch $\mathcal{N}$ is only a noisy estimate of the true gradient $\nabla f(\theta)$. Performing an accurate line search along this noisy direction is computationally expensive and may not be beneficial, as the "optimal" step size for the mini-batch loss might not be good for the true loss.

- **Cost per Iteration:** SGD relies on very fast iterations. The overhead of function evaluations required for line search would significantly slow down the training process, negating the main advantage of SGD's rapid updates. The inherent noise in SGD often provides enough exploration, and simple step size decay schedules are typically preferred.

## (b) Change of Variables for Constraints

A change of variables (or reparameterization) is a technique to handle constraints on decision variables $\theta$ by transforming them into unconstrained variables $\phi$. We define a mapping $\theta = g(\phi)$ such that for any value of the unconstrained variable $\phi$, the resulting $\theta$ automatically satisfies the desired constraints. The optimization is then performed with respect to $\phi$ using an unconstrained optimization algorithm on the transformed objective function $\hat{f}(\phi) = f(g(\phi))$. After finding the optimal $\phi^*$, the optimal constrained variables are recovered as $\theta^* = g(\phi^*)$.

**Example:**

Suppose we want to minimize a function $f(\theta_1, \theta_2)$ subject to the constraint that $\theta_1 > 0$. We can introduce an unconstrained variable $\phi_1$ and define the transformation:

$$\theta_1 = e^{\phi_1}, \quad \theta_2 = \phi_2$$

Here, $\theta_1 = e^{\phi_1}$ ensures that $\theta_1$ is always positive, regardless of the value of the unconstrained variable $\phi_1 \in \mathbb{R}$. The variable $\theta_2$ is already unconstrained, so we can set $\theta_2 = \phi_2$.

We then perform unconstrained optimization on the function $\hat{f}(\phi_1, \phi_2) = f(e^{\phi_1}, \phi_2)$ with respect to $\phi_1$ and $\phi_2$. If the unconstrained optimization yields $\phi_1^*$ and $\phi_2^*$, the solution for the original constrained problem is $\theta_1^* = e^{\phi_1^*}$ and $\theta_2^* = \phi_2^*$. This approach effectively incorporates the constraint into the definition of the variables.

## (c) Penalty Method for Constraints

The penalty method transforms a constrained optimization problem into an unconstrained one by adding a penalty term to the objective function. This penalty term is designed to be large when the constraints are violated and small (or zero) when they are satisfied. The optimizer is thus discouraged from exploring regions that violate the constraints.

For an equality constraint $h(\theta) = 0$, a common penalty is $\mu h(\theta)^2$. For an inequality constraint $g(\theta) \leq 0$, a common penalty is $\mu(\max(0, g(\theta)))^2$. Here, $\mu > 0$ is a penalty parameter. The modified objective function becomes:

$$\hat{f}(\theta) = f(\theta) + \text{PenaltyTerm}(\theta)$$

Typically, the optimization is performed iteratively, gradually increasing the penalty parameter $\mu$ to enforce the constraints more strictly as the optimum is approached.

**Example:**

Suppose we want to minimize $f(\theta)$ subject to the constraint $\theta = 5$. We can define a penalty function and create an unconstrained objective:

$$\hat{f}(\theta) = f(\theta) + \mu(\theta - 5)^2$$

where $\mu > 0$ is the penalty parameter. An unconstrained optimization algorithm applied to $\hat{f}(\theta)$ will try to minimize both $f(\theta)$ and the penalty term $\mu(\theta - 5)^2$. As $\mu$ increases, the optimizer is forced to find solutions where $(\theta - 5)^2$ is very small, thus driving $\theta$ towards 5 and satisfying the constraint.

# A   Code

## Project Structure

```
final-assignment/
        configs/                        # Configuration files for experiments
                synthetic_regression_config.yaml
                week6_experiment_config.yaml
                week9_transformer_config.yaml
        experiments/             # Experiment scripts
                synthetic_regression.py
                week6_experiment.py
                week9_transformer.py
        results/                        # Results from experiments
                logs/
                        synthetic_regression/
                        transformer/
                        week6/
        src/                             # Source code
                datasets.py               # Dataset implementations
                models/                   # Model implementations
                        linear_regression.py
                        transformer_data.py
                        transformer_week9.py
                        week6_model.py
                optim/                    # Optimizer implementations
                        polyak_adam.py
                        polyak_sgd.py
                        sgd.py
                transformer-datasets/ # Datasets for transformer models
                utils.py                  # Utility functions
        tests/                   # Unit tests
        requirements.txt         # Python dependencies
```

## Configs

**synthetic_regression_config.yaml**

```yaml
# Configuration for synthetic regression experiments

# Data generation parameters
data:
  n_samples: 1000   # Number of data points
  dimension: 20     # Feature dimension
  noise_std_values: [0.1, 1.0, 5.0]  # Noise standard deviations to test
  seeds: [0, 42, 1337]  # Random seeds for reproducibility

# Training parameters
training:
  epochs: 20   # Number of training epochs
  batch_sizes: [1, 16, 64, 256, "full"]  # Batch sizes to test (full = use all data)

# Optimizer parameters
optimizer:
  # Learning rates for grid search (constant step size)
  learning_rates: [0.001, 0.01, 0.1, 1.0]

  # Polyak step size parameters
  polyak:
    eps: 1.0e-8  # Small constant for numerical stability
```

```
    f_star: 0.0   # Known minimum value of the loss function

  # PolyakAdam parameters
  polyak_adam:
    alpha: 0.003   # Base learning rate
    polyak_factor: 0.3  # Factor to control the influence of Polyak step size
    amsgrad: false  # Whether to use the AMSGrad variant

# Logging parameters
logging:
  log_dir: "results/logs/synthetic_regression"
  save_interval: 1 # Save logs every N epochs
```

**week6_experiment_config.yaml**

```
# Week 6 Experiment Configuration

# Data configuration
data:
  m: 25                    # Number of data points
  std: 0.25                # Standard deviation of Gaussian noise
  seed: 42                 # Random seed for reproducibility

# Training configuration
training:
  batch_sizes: [1, 5, 10, 25] # Batch sizes to experiment with
  epochs: 100                   # Number of training epochs

# Optimizer configuration
optimizer:
  sgd:
    lr: 0.01               # Constant step size for SGD
  polyak:
    eps: 1.0e-8            # Small epsilon to avoid division by zero
    f_star: 0.0           # Optimal function value (usually 0)

# Logging configuration
logging:
  log_dir: "results/logs/week6" # Directory to save results
```

**week9_transformer_config.yaml**

```
# Configuration for Week 9 Transformer Experiment
# Comparing different optimizers for transformer training

data:
  data_dir: "src/transformer-datasets"
  dataset_file: "input_childSpeech_trainingSet.txt"
  test_file: null  # Set to file path if you want to evaluate on test data
  batch_size: 64
  block_size: 256

training:
  max_iters: 2000
  eval_interval: 200
  seed: 42
  tuning_iters: 500  # Fewer iterations for hyperparameter tuning
  enable_tuning: false
  early_stopping:
    enabled: true
    patience: 3  # Number of evaluations with no improvement after which training will
        stop
    min_delta: 0.001  # Minimum change in validation loss to qualify as improvement
    min_epochs: 1  # Minimum number of epochs to train regardless of early stopping
        condition

tuning:
  sgd:
    lr_values: [0.0001, 0.001, 0.01, 0.1]
  adam:
    lr_values: [0.0001, 0.001, 0.01]
    beta1_values: [0.9, 0.95]
```

```
    beta2_values: [0.999, 0.99]

optimizers:
  Adam:
    lr: 0.001
    betas: [0.9, 0.999]
  SGD:
    lr: 0.1
  Polyak:
    eps: 1.0e-8
    f_star: 0.0

logging:
  log_dir: "results/logs/transformer"
  save_models: true
```

## datasets.py

```python
import numpy as np
import torch
from typing import Tuple, Dict, List, Callable, Optional
from src.utils import set_seed, create_random_tensor, get_device
import os


def make_synthetic_linreg(N: int, d: int, noise_std: float, seed: int) ->
    Tuple[torch.Tensor, torch.Tensor, torch.Tensor, float]:
    """
    Generate synthetic data for linear regression with the model: y =  w x  + b +   ,
    where    is Gaussian noise with standard deviation noise_std.

    Args:
        N: Number of data points to generate
        d: Dimensionality of the feature space
        noise_std: Standard deviation of the Gaussian noise
        seed: Random seed for reproducibility

    Returns:
        Tuple containing:
            - X: Input features of shape (N, d)
            - y: Target values of shape (N,)
            - w: True weight vector of shape (d,)
            - b: True bias term (scalar)
    """
    # Set seed for reproducibility
    set_seed(seed)

    # Generate input features X ~ N(0, I)
    X = create_random_tensor((N, d), seed=seed)

    # Generate true parameters w ~ N(0, I) and b ~ N(0, 1)
    # We use a fixed seed + 1 to ensure w and b are different from X but consistent
        across runs
    w = create_random_tensor((d,), seed=seed+1)
    b = torch.randn(1).item()  # Scalar bias term

    # Generate target values: y = Xw + b + noise
    y_clean = torch.matmul(X, w) + b

    # Add Gaussian noise with standard deviation noise_std
    # Use seed+2 to ensure noise is different but consistent
    noise = create_random_tensor((N,), seed=seed+2) * noise_std
    y = y_clean + noise

    return X, y, w, b


def make_gaussian_cluster(m: int = 25, std: float = 0.25, seed: int = 42) ->
    torch.Tensor:
    """
    Generate data points from a 2D Gaussian distribution centered at origin. (Week 6
        Training Data)

    Original function
```

```python
    def generate_trainingdata(m=25):
        return np.array([0,0])+0.25*np.random.randn(m,2)

    Args:
        m: Number of data points to generate
        std: Standard deviation of the Gaussian noise
        seed: Random seed for reproducibility

    Returns:
        X: Data points of shape (m, 2)
    """
    # Set seed for reproducibility
    set_seed(seed)

    # Generate points from a 2D Gaussian centered at origin
    X = create_random_tensor((m, 2), seed=seed) * std

    return X


def load_text_data(filepath: str) -> str:
    """
    Load text data from a file.

    Args:
        filepath: Path to the text file

    Returns:
        The content of the text file as a string
    """
    with open(filepath, 'r', encoding='utf-8') as f:
        return f.read()


def create_char_mappings(text: str) -> Tuple[Dict[str, int], Dict[int, str]]:
    """
    Create character-to-index and index-to-character mappings.

    Args:
        text: Text data to create mappings from

    Returns:
        Tuple containing:
            - stoi: Dictionary mapping characters to indices
            - itos: Dictionary mapping indices to characters
    """
    # Get unique characters in sorted order
    chars = sorted(list(set(text)))

    # Create mappings
    stoi = {ch: i for i, ch in enumerate(chars)}
    itos = {i: ch for i, ch in enumerate(chars)}

    return stoi, itos


def prepare_transformer_data(filepath: str,
                             train_ratio: float = 0.9,
                             seed: int = 1337) -> Tuple[torch.Tensor, torch.Tensor,
                                 Dict[str, int], Dict[int, str]]:
    """
    Prepare data for transformer model training and evaluation from a single text file.

    Args:
        filepath: Path to the text file
        train_ratio: Ratio of data to use for training (rest for validation)
        seed: Random seed for reproducibility

    Returns:
        Tuple containing:
            - train_data: Training data tensor
            - val_data: Validation data tensor
            - stoi: Character to index mapping
            - itos: Index to character mapping
    """
    # Set seed for reproducibility
```

```python
    set_seed(seed)

    # Load text data
    text_data = load_text_data(filepath)

    # Create character mappings
    stoi, itos = create_char_mappings(text_data)

    # Create encoder function
    encode = lambda s: [stoi[c] for c in s]

    # Encode data
    data = torch.tensor(encode(text_data), dtype=torch.long)

    # Split into train and validation sets
    n = int(train_ratio * len(data))
    train_data = data[:n]
    val_data = data[n:]

    return train_data, val_data, stoi, itos


def encode_text(text: str, stoi: Dict[str, int], unk_token: int = 0) -> List[int]:
    """
    Encode text using given character to index mapping.

    Args:
        text: Text to encode
        stoi: Character to index mapping dictionary
        unk_token: Index to use for unknown characters

    Returns:
        List of encoded indices
    """
    return [stoi.get(c, unk_token) for c in text]


def load_and_encode_text(filepath: str, stoi: Dict[str, int], unk_token: int = 0) ->
    torch.Tensor:
    """
    Load and encode text from a file using existing character mapping.
    Can be used for loading test data with existing vocabulary.

    Args:
        filepath: Path to the text file
        stoi: Character to index mapping
        unk_token: Index to use for unknown characters

    Returns:
        Tensor of encoded indices
    """
    # Verify file exists
    if not os.path.exists(filepath):
        raise FileNotFoundError(f"Text file not found: {filepath}")

    # Load text data
    text_data = load_text_data(filepath)

    # Encode using provided character mappings
    encoded_data = encode_text(text_data, stoi, unk_token)

    # Convert to tensor
    return torch.tensor(encoded_data, dtype=torch.long)


def get_transformer_batch(data: torch.Tensor,
                          block_size: int,
                          batch_size: int,
                          device: Optional[torch.device] = None) -> Tuple[torch.Tensor,
                              torch.Tensor]:
    """
    Generate a batch of data for transformer training or evaluation.
    Works for any dataset (train, validation, or test).

    Args:
        data: Encoded text data tensor
```

```
            block_size: Maximum context length
            batch_size: Number of sequences in a batch
            device: Device to place tensors on (defaults to utils.get_device())

        Returns:
            Tuple containing:
                - x: Input tensor of shape (batch_size, block_size)
                - y: Target tensor of shape (batch_size, block_size)
        """
        if device is None:
            device = get_device()

        # Generate random starting indices
        ix = torch.randint(len(data) - block_size, (batch_size,))

        # Create input sequences
        x = torch.stack([data[i:i+block_size] for i in ix])

        # Create target sequences (shifted by one position)
        y = torch.stack([data[i+1:i+block_size+1] for i in ix])

        # Move to device
        x, y = x.to(device), y.to(device)

        return x, y
```

## utils.py

```python
import random
import numpy as np
import torch
from torch.utils.data import DataLoader, TensorDataset
from typing import Optional, Dict, List, Tuple, Callable, Any


def set_seed(seed: int) -> None:
    """
    Set random seeds for reproducibility across all libraries.

    Args:
        seed: Integer seed value to use
    """
    random.seed(seed)
    np.random.seed(seed)
    torch.manual_seed(seed)
    if torch.cuda.is_available():
        torch.cuda.manual_seed(seed)
        torch.cuda.manual_seed_all(seed)  # For multi-GPU setups

    # Additional configurations for reproducibility
    torch.backends.cudnn.deterministic = True
    torch.backends.cudnn.benchmark = False


def get_device() -> torch.device:
    """
    Get the appropriate device (CUDA or CPU) for PyTorch operations.

    Returns:
        torch.device: The device to use for computations
    """
    return torch.device("cuda" if torch.cuda.is_available() else "cpu")


def create_random_tensor(shape: tuple,
                         seed: Optional[int] = None,
                         device: Optional[torch.device] = None) -> torch.Tensor:
    """
    Create a random tensor with a specific shape and seed.

    Args:
        shape: Tuple specifying the shape of the tensor
        seed: Optional seed for reproducibility
        device: Optional device to place the tensor on
```

```python
    Returns:
        Random tensor of the specified shape
    """
    if seed is not None:
        set_seed(seed)

    if device is None:
        device = get_device()

    return torch.randn(shape, device=device)


def get_data_loader(X: torch.Tensor, y: Optional[torch.Tensor] = None, batch_size: int =
    8) -> DataLoader:
    """
    Create a DataLoader for the given data.

    Args:
        X: Input features tensor of shape (N, d)
        y: Optional target values tensor of shape (N,). If None, creates a dataset with
            only X
        batch_size: Size of mini-batches

    Returns:
        DataLoader that yields batches of (X,) or (X, y) pairs depending on if y is
            provided
    """
    # Create dataset based on whether y is provided
    if y is not None:
        # Ensure y has the right shape for the model
        if y.dim() == 1:
            y = y.view(-1, 1)
        dataset = TensorDataset(X, y)
    else:
        # Create dataset with just X
        dataset = TensorDataset(X)

    # Create data loader
    return DataLoader(
        dataset=dataset,
        batch_size=batch_size,
        shuffle=True
    )


def create_text_tensor_dataset(data: torch.Tensor,
                               block_size: int,
                               device: Optional[torch.device] = None) -> TensorDataset:
    """
    Create a TensorDataset for text data using sliding window approach.

    Args:
        data: Encoded text data tensor
        block_size: Context length for each sample
        device: Device to place tensors on

    Returns:
        TensorDataset with input-target pairs for language modeling
    """
    if device is None:
        device = get_device()

    # Create input-target pairs
    inputs = []
    targets = []

    for i in range(len(data) - block_size):
        # Input is a sequence of block_size tokens
        inputs.append(data[i:i+block_size])
        # Target is the sequence shifted by one position
        targets.append(data[i+1:i+block_size+1])

    # Convert to tensors
    x = torch.stack(inputs).to(device)
    y = torch.stack(targets).to(device)
```

```python
    return TensorDataset(x, y)


def get_text_data_loader(data: torch.Tensor,
                         block_size: int,
                         batch_size: int,
                         shuffle: bool = True,
                         device: Optional[torch.device] = None) -> DataLoader:
    """
    Create a DataLoader for text data.

    Args:
        data: Encoded text data tensor
        block_size: Context length for each sample
        batch_size: Size of mini-batches
        shuffle: Whether to shuffle the data
        device: Device to place tensors on

    Returns:
        DataLoader for language modeling
    """
    dataset = create_text_tensor_dataset(data, block_size, device)

    return DataLoader(
        dataset=dataset,
        batch_size=batch_size,
        shuffle=shuffle
    )


def decode_text(indices: List[int], itos: Dict[int, str]) -> str:
    """
    Decode a list of token indices back to a string.

    Args:
        indices: List of integer token indices
        itos: Dictionary mapping indices to characters

    Returns:
        Decoded string
    """
    return ''.join([itos[idx] for idx in indices])


def analyze_text_data(text: str) -> Dict[str, Any]:
    """
    Analyze a text dataset and return statistics.

    Args:
        text: Text data to analyze

    Returns:
        Dictionary containing various statistics about the text
    """
    # Character-level statistics
    chars = sorted(list(set(text)))
    char_vocab_size = len(chars)
    total_chars = len(text)

    # Word-level statistics
    words = text.split()
    unique_words = set(words)
    word_count = len(words)
    word_vocab_size = len(unique_words)

    # Return statistics dictionary
    return {
        "char_vocab_size": char_vocab_size,
        "total_chars": total_chars,
        "word_vocab_size": word_vocab_size,
        "word_count": word_count,
        "sample": text[:500] if len(text) > 500 else text,
        "unique_chars": chars
    }
```

## experiments

### synthetic_regression.py

```python
import os
import yaml
import argparse
import numpy as np
import torch
import json
import matplotlib.pyplot as plt
from datetime import datetime
from typing import Dict, List, Union, Tuple, Optional, Any
from pathlib import Path
import shutil

from src.utils import set_seed, get_device, get_data_loader
from src.datasets import make_synthetic_linreg
from src.models.linear_regression import LinReg, get_mse_loss
from src.optim.sgd import SGD
from src.optim.polyak_sgd import PolyakSGD
from src.optim.polyak_adam import PolyakAdam


def parse_args() -> argparse.Namespace:
    """Parse command line arguments."""
    parser = argparse.ArgumentParser(description="Run synthetic regression experiments")
    parser.add_argument(
        "--config",
        type=str,
        default="configs/synthetic_regression_config.yaml",
        help="Path to configuration YAML file"
    )
    parser.add_argument(
        "--best_lr_path",
        type=str,
        default=None,
        help="Path to a previously saved best_learning_rates.json file to skip learning
            rate search"
    )
    parser.add_argument(
        "--demo_mode",
        action="store_true",
        help="Run only for specific noise and seed across different batch sizes for
            demonstration"
    )
    parser.add_argument(
        "--demo_noise",
        type=float,
        default=1.0,
        help="Noise level to use in demo mode"
    )
    parser.add_argument(
        "--demo_seed",
        type=int,
        default=42,
        help="Seed to use in demo mode"
    )
    return parser.parse_args()


def load_config(config_path: str) -> Dict[str, Any]:
    """Load configuration from YAML file."""
    with open(config_path, "r") as f:
        config = yaml.safe_load(f)
    return config


def train_model(
    model: torch.nn.Module,
    optimizer: torch.optim.Optimizer,
    data_loader: torch.utils.data.DataLoader,
    criterion: torch.nn.Module,
    device: torch.device,
    epochs: int,
    track_step_size: bool = False
```

```python
) -> Tuple[List[float], List[float]]:
    """
    Train a model with the given optimizer and data loader.

    Args:
        model: PyTorch model to train
        optimizer: Optimizer to use for training
        data_loader: DataLoader providing training data
        criterion: Loss function
        device: Device to run training on
        epochs: Number of epochs to train for
        track_step_size: Whether to track step sizes (for Polyak-type optimizers)

    Returns:
        Tuple of (training losses, step sizes if tracking step sizes)
    """
    model.train()
    losses = []
    step_sizes = [] if track_step_size else None

    for epoch in range(epochs):
        epoch_loss = 0.0
        batches = 0

        for X_batch, y_batch in data_loader:
            X_batch, y_batch = X_batch.to(device), y_batch.to(device)

            def closure():
                optimizer.zero_grad()
                output = model(X_batch)
                loss = criterion(output, y_batch)
                loss.backward()
                return loss

            if track_step_size:
                loss = optimizer.step(closure)
                step_sizes.append(optimizer.last_step_size)
            else:
                loss = closure()
                optimizer.step()

            epoch_loss += loss.item()
            batches += 1

        avg_epoch_loss = epoch_loss / batches
        losses.append(avg_epoch_loss)

    return losses, step_sizes


def find_best_lr(
    config: Dict[str, Any],
    noise_std: float,
    batch_size: Union[int, str],
    seed: int,
    device: torch.device,
    log_dir: str
) -> Optional[float]:
    """
    Find the best learning rate for Constant Step Size SGD for the given configuration.

    Args:
        config: Configuration dictionary
        noise_std: Noise standard deviation for data generation
        batch_size: Batch size for training (int or "full")
        seed: Random seed
        device: Device to run training on
        log_dir: Directory to save logs

    Returns:
        Best learning rate or None if no suitable learning rate found
    """
    print(f"Finding best LR for noise_std={noise_std}, batch_size={batch_size}, 
        seed={seed}")

    # Generate synthetic data
```

```python
n_samples = config["data"]["n_samples"]
dimension = config["data"]["dimension"]
epochs = config["training"]["epochs"]

X, y, true_w, true_b = make_synthetic_linreg(
    N=n_samples,
    d=dimension,
    noise_std=noise_std,
    seed=seed
)

# Use full batch if specified
actual_batch_size = n_samples if batch_size == "full" else batch_size
data_loader = get_data_loader(X, y, actual_batch_size)

# Initialize results dictionary
results = {}
best_lr = None
best_final_loss = float('inf')

# Create log subdirectory
lr_log_dir = os.path.join(log_dir,
    f"lr_search_noise{noise_std}_batch{batch_size}_seed{seed}")
os.makedirs(lr_log_dir, exist_ok=True)

# Try different learning rates
for lr in config["optimizer"]["learning_rates"]:
    # Create model and optimizer
    model = LinReg(dimension).to(device)
    criterion = get_mse_loss()
    optimizer = SGD(model.parameters(), lr=lr)

    # Train model
    try:
        losses, _ = train_model(
            model=model,
            optimizer=optimizer,
            data_loader=data_loader,
            criterion=criterion,
            device=device,
            epochs=epochs
        )

        # Check if training diverged (NaN or very high loss)
        if np.isnan(losses[-1]) or losses[-1] > 1e6:
            print(f"  LR {lr} diverged, final loss: {losses[-1]}")
            continue

        # Check if loss decreased
        loss_decreased = losses[0] > losses[-1]

        # Calculate loss oscillation
        oscillation = 0
        if len(losses) > 2:
            # Calculate average oscillation in the second half of training
            half_idx = len(losses) // 2
            for i in range(half_idx, len(losses) - 1):
                oscillation += abs(losses[i+1] - losses[i])
            oscillation /= (len(losses) - half_idx - 1) if len(losses) - half_idx - \
                1 > 0 else 1

        # Save results
        results[lr] = {
            "final_loss": losses[-1],
            "loss_decreased": loss_decreased,
            "oscillation": oscillation,
            "losses": losses
        }

        # Update best learning rate if this one is better
        if loss_decreased and losses[-1] < best_final_loss and oscillation < 0.1 * \
            losses[-1]:
            best_lr = lr
            best_final_loss = losses[-1]
```

```python
                print(f"  LR {lr} - Final loss: {losses[-1]:.6f}, Oscillation: "
                    {oscillation:.6f}")

            except Exception as e:
                print(f"  Error with LR {lr}: {str(e)}")

    # Save results to file
    with open(os.path.join(lr_log_dir, "lr_search_results.json"), "w") as f:
        # Convert to serializable format
        serializable_results = {str(k): {
            "final_loss": v["final_loss"],
            "loss_decreased": v["loss_decreased"],
            "oscillation": v["oscillation"],
            "losses": v["losses"]
        } for k, v in results.items()}
        json.dump(serializable_results, f, indent=2)

    # Plot loss curves for different learning rates
    plt.figure(figsize=(10, 6))
    for lr, result in results.items():
        plt.plot(result["losses"], label=f"LR={lr}")
    plt.xlabel("Epoch")
    plt.ylabel("Loss")
    plt.title(f"Loss vs. Epoch for Different Learning Rates\nNoise={noise_std}, Batch "
        Size={batch_size}")
    plt.legend()
    plt.savefig(os.path.join(lr_log_dir, "lr_comparison.png"))
    plt.close()

    # If no good learning rate found, use the one with lowest final loss
    if best_lr is None and results:
        best_lr = min(results.keys(), key=lambda lr: results[lr]["final_loss"])
        print(f"No ideal learning rate found. Using LR={best_lr} with lowest final "
            loss.")

    return best_lr


def compare_optimizers(
    config: Dict[str, Any],
    noise_std: float,
    batch_size: Union[int, str],
    seed: int,
    best_lr: float,
    device: torch.device,
    log_dir: str
) -> Dict[str, Any]:
    """
    Compare Constant Step Size SGD with Polyak Step Size SGD and PolyakAdam.

    Args:
        config: Configuration dictionary
        noise_std: Noise standard deviation for data generation
        batch_size: Batch size for training (int or "full")
        seed: Random seed
        best_lr: Best learning rate for Constant Step Size SGD
        device: Device to run training on
        log_dir: Directory to save logs

    Returns:
        Dictionary with results
    """
    print(f"Comparing optimizers: noise_std={noise_std}, batch_size={batch_size}, "
        seed={seed}")

    # Generate synthetic data
    n_samples = config["data"]["n_samples"]
    dimension = config["data"]["dimension"]
    epochs = config["training"]["epochs"]

    X, y, true_w, true_b = make_synthetic_linreg(
        N=n_samples,
        d=dimension,
        noise_std=noise_std,
        seed=seed
    )
```

```python
    # Use full batch if specified
    actual_batch_size = n_samples if batch_size == "full" else batch_size
    data_loader = get_data_loader(X, y, actual_batch_size)

    # Directory for comparison results
    comp_log_dir = os.path.join(
        log_dir,
        f"comparison_noise{noise_std}_batch{batch_size}_seed{seed}"
    )
    os.makedirs(comp_log_dir, exist_ok=True)

    # Train with Constant Step Size SGD
    model_sgd = LinReg(dimension).to(device)
    criterion = get_mse_loss()
    optimizer_sgd = SGD(model_sgd.parameters(), lr=best_lr)

    sgd_losses, _ = train_model(
        model=model_sgd,
        optimizer=optimizer_sgd,
        data_loader=data_loader,
        criterion=criterion,
        device=device,
        epochs=epochs
    )

    # Train with Polyak Step Size SGD
    model_polyak = LinReg(dimension).to(device)
    optimizer_polyak = PolyakSGD(
        model_polyak.parameters(),
        eps=config["optimizer"]["polyak"]["eps"],
        f_star=config["optimizer"]["polyak"]["f_star"]
    )

    polyak_losses, polyak_step_sizes = train_model(
        model=model_polyak,
        optimizer=optimizer_polyak,
        data_loader=data_loader,
        criterion=criterion,
        device=device,
        epochs=epochs,
        track_step_size=True
    )

    # Train with PolyakAdam (without AMSGrad)
    model_polyak_adam = LinReg(dimension).to(device)
    optimizer_polyak_adam = PolyakAdam(
        model_polyak_adam.parameters(),
        eps=config["optimizer"]["polyak"]["eps"],
        f_star=config["optimizer"]["polyak"]["f_star"],
        alpha=config["optimizer"]["polyak_adam"]["alpha"] if "polyak_adam" in
            config["optimizer"] else 0.003,
        polyak_factor=config["optimizer"]["polyak_adam"]["polyak_factor"] if
            "polyak_adam" in config["optimizer"] else 0.3,
        amsgrad=config["optimizer"]["polyak_adam"]["amsgrad"] if "polyak_adam" in
            config["optimizer"] and "amsgrad" in config["optimizer"]["polyak_adam"] else
            False
    )

    polyak_adam_losses, polyak_adam_step_sizes = train_model(
        model=model_polyak_adam,
        optimizer=optimizer_polyak_adam,
        data_loader=data_loader,
        criterion=criterion,
        device=device,
        epochs=epochs,
        track_step_size=True
    )

    # Train with PolyakAdam with AMSGrad enabled
    model_polyak_adam_ams = LinReg(dimension).to(device)
    optimizer_polyak_adam_ams = PolyakAdam(
        model_polyak_adam_ams.parameters(),
        eps=config["optimizer"]["polyak"]["eps"],
        f_star=config["optimizer"]["polyak"]["f_star"],
```

```python
        alpha=config["optimizer"]["polyak_adam"]["alpha"] if "polyak_adam" in
            config["optimizer"] else 0.003,
        polyak_factor=config["optimizer"]["polyak_adam"]["polyak_factor"] if
            "polyak_adam" in config["optimizer"] else 0.3,
        amsgrad=True  # Always enable AMSGrad for this optimizer
    )

    polyak_adam_ams_losses, polyak_adam_ams_step_sizes = train_model(
        model=model_polyak_adam_ams,
        optimizer=optimizer_polyak_adam_ams,
        data_loader=data_loader,
        criterion=criterion,
        device=device,
        epochs=epochs,
        track_step_size=True
    )

    # Save results
    results = {
        "sgd": {
            "losses": sgd_losses,
            "lr": best_lr
        },
        "polyak": {
            "losses": polyak_losses,
            "step_sizes": polyak_step_sizes
        },
        "polyak_adam": {
            "losses": polyak_adam_losses,
            "step_sizes": polyak_adam_step_sizes
        },
        "polyak_adam_ams": {
            "losses": polyak_adam_ams_losses,
            "step_sizes": polyak_adam_ams_step_sizes
        },
        "config": {
            "noise_std": noise_std,
            "batch_size": batch_size,
            "seed": seed,
            "epochs": epochs
        }
    }

    # Save results to file
    with open(os.path.join(comp_log_dir, "comparison_results.json"), "w") as f:
        json.dump(results, f, indent=2)

    # Plot loss comparison
    plt.figure(figsize=(10, 6))
    plt.plot(sgd_losses, label=f"Constant Step Size (LR={best_lr})")
    plt.plot(polyak_losses, label="Polyak Step Size")
    plt.plot(polyak_adam_losses, label="PolyakAdam")
    plt.plot(polyak_adam_ams_losses, label="PolyakAdam+AMSGrad")
    plt.xlabel("Epoch")
    plt.ylabel("Loss")
    plt.title(f"Loss Comparison: SGD vs Polyak vs PolyakAdam\nNoise={noise_std}, Batch
        Size={batch_size}")
    plt.legend()
    plt.savefig(os.path.join(comp_log_dir, "loss_comparison.png"))
    plt.close()

    # Plot Polyak step sizes
    plt.figure(figsize=(10, 6))
    plt.plot(polyak_step_sizes, label="Polyak SGD")
    plt.plot(polyak_adam_step_sizes, label="PolyakAdam")
    plt.plot(polyak_adam_ams_step_sizes, label="PolyakAdam+AMSGrad")
    plt.xlabel("Batch Update")
    plt.ylabel("Step Size")
    plt.title(f"Step Sizes Over Training\nNoise={noise_std}, Batch Size={batch_size}")
    plt.legend()
    plt.savefig(os.path.join(comp_log_dir, "step_sizes.png"))
    plt.close()

    return results
```

```python
def run_demo_comparison(
    config: Dict[str, Any],
    best_lr_grid: Dict[Tuple[float, Union[int, str], int], float],
    demo_noise: float,
    demo_seed: int,
    device: torch.device,
    log_dir: str
) -> None:
    """
    Run comparison between optimizers for a specific noise level and seed across
        different batch sizes.

    Args:
        config: Configuration dictionary
        best_lr_grid: Dictionary mapping (noise_std, batch_size, seed) to best learning
            rate
        demo_noise: Noise level to use for demonstration
        demo_seed: Seed to use for demonstration
        device: Device to run training on
        log_dir: Directory to save logs
    """
    print(f"\nRUNNING DEMO MODE: Comparing optimizers for noise={demo_noise}, "
        seed={demo_seed}")

    # Create a special demo directory
    demo_dir = os.path.join(log_dir, f"demo_noise{demo_noise}_seed{demo_seed}")
    os.makedirs(demo_dir, exist_ok=True)

    # Check if the specified noise and seed exist in the best_lr_grid
    demo_keys = [k for k in best_lr_grid.keys() if k[0] == demo_noise and k[2] ==
        demo_seed]
    if not demo_keys:
        raise ValueError(f"No learning rates found for noise={demo_noise}, "
            seed={demo_seed}. "
                        f"Please provide a valid --best_lr_path with these "
                            configurations.")

    # Filter batch sizes available in the demo keys
    available_batch_sizes = [k[1] for k in demo_keys]
    print(f"Running comparison for batch sizes: {available_batch_sizes}")

    # Run comparison for each batch size
    results = {}
    batch_comparison_data = {
        "batch_sizes": [],
        "sgd_final_losses": [],
        "polyak_final_losses": [],
        "polyak_adam_final_losses": [],
        "polyak_adam_ams_final_losses": [],
        "sgd_lrs": []
    }

    for batch_size in available_batch_sizes:
        key = (demo_noise, batch_size, demo_seed)
        best_lr = best_lr_grid.get(key)

        if best_lr is None:
            print(f"Warning: No learning rate found for batch_size={batch_size}")
            continue

        print(f"\nComparing batch_size={batch_size} with learning rate={best_lr}")

        result = compare_optimizers(
            config=config,
            noise_std=demo_noise,
            batch_size=batch_size,
            seed=demo_seed,
            best_lr=best_lr,
            device=device,
            log_dir=demo_dir
        )

        results[batch_size] = result

        # Collect data for batch size comparison plot
        batch_comparison_data["batch_sizes"].append(str(batch_size))
```

```python
            batch_comparison_data["sgd_final_losses"].append(result["sgd"]["losses"][-1])
            batch_comparison_data["polyak_final_losses"].append(result["polyak"]["losses"][-1])
            batch_comparison_data["polyak_adam_final_losses"].append(result["polyak_adam"]["losses"][-1])
            batch_comparison_data["polyak_adam_ams_final_losses"].append(result["polyak_adam_ams"]["losses"]
            batch_comparison_data["sgd_lrs"].append(best_lr)

    # Create a unified plot to show the effect of batch size
    create_batch_size_comparison_plot(batch_comparison_data, demo_dir, demo_noise,
        demo_seed)

    # Create convergence speed plot
    create_convergence_comparison_plot(results, demo_dir, demo_noise, demo_seed)

    # Save detailed results
    with open(os.path.join(demo_dir, "demo_results.json"), "w") as f:
        serializable_results = {str(k): v for k, v in results.items()}
        json.dump(serializable_results, f, indent=2)

    print(f"\nDemo comparison completed. Results saved to {demo_dir}")


def create_batch_size_comparison_plot(
    data: Dict[str, List],
    log_dir: str,
    noise: float,
    seed: int
) -> None:
    """
    Create a plot comparing final losses for different batch sizes.

    Args:
        data: Dictionary with batch sizes and final losses
        log_dir: Directory to save plot
        noise: Noise level used in the demonstration
        seed: Seed used in the demonstration
    """
    plt.figure(figsize=(10, 8))

    # Ensure data is sorted by batch size (handle 'full' special case)
    indices = list(range(len(data["batch_sizes"])))

    # Sort numeric batch sizes, putting 'full' at the end
    indices.sort(key=lambda i: (data["batch_sizes"][i] == "full",
                                int(data["batch_sizes"][i]) if data["batch_sizes"][i] !=
                                    "full" else float('inf')))

    batch_sizes = [data["batch_sizes"][i] for i in indices]
    sgd_losses = [data["sgd_final_losses"][i] for i in indices]
    polyak_losses = [data["polyak_final_losses"][i] for i in indices]
    polyak_adam_losses = [data["polyak_adam_final_losses"][i] for i in indices]
    polyak_adam_ams_losses = [data["polyak_adam_ams_final_losses"][i] for i in indices]
    sgd_lrs = [data["sgd_lrs"][i] for i in indices]

    x = list(range(len(batch_sizes)))

    # Plot final losses
    plt.subplot(2, 1, 1)
    plt.plot(x, sgd_losses, 'o-', label="SGD (Constant Step Size)")
    plt.plot(x, polyak_losses, 'x--', label="Polyak Step Size")
    plt.plot(x, polyak_adam_losses, '*-.', label="PolyakAdam")
    plt.plot(x, polyak_adam_ams_losses, '*--', label="PolyakAdam+AMSGrad")
    plt.xticks(x, batch_sizes)
    plt.xlabel("Batch Size")
    plt.ylabel("Final Loss")
    plt.title(f"Effect of Batch Size on Final Loss (Noise={noise}, Seed={seed})")
    plt.grid(True, alpha=0.3)
    plt.legend()

    # Plot learning rates
    plt.subplot(2, 1, 2)
    plt.plot(x, sgd_lrs, 'o-', color='green')
    plt.xticks(x, batch_sizes)
    plt.xlabel("Batch Size")
    plt.ylabel("Learning Rate")
    plt.title("Best Learning Rate for Each Batch Size")
    plt.grid(True, alpha=0.3)
```

```python
    plt.tight_layout ()
    plt.savefig(os.path.join(log_dir , "batch_size_comparison.png"))
    plt.close ()


def create_convergence_comparison_plot (
    results: Dict[Union[int , str], Dict],
    log_dir: str ,
    noise: float ,
    seed: int
) -> None:
    """
    Create a plot comparing convergence speed for different batch sizes.

    Args:
        results: Dictionary with batch sizes and training results
        log_dir: Directory to save plot
        noise: Noise level used in the demonstration
        seed: Seed used in the demonstration
    """
    plt.figure(figsize =(15, 16))

    # Plot SGD convergence
    plt.subplot(4, 1, 1)
    # Convert batch sizes to strings for consistent sorting
    sorted_items = sorted(results.items(), key=lambda x: (
            str(x[0]) == "full",  # Sort "full" last
            int(x[0]) if str(x[0]) != "full" else float('inf')  # Sort numerically
        ))

    for batch_size , result in sorted_items:
        sgd_losses = result["sgd"]["losses"]
        plt.plot(sgd_losses , label=f"Batch={batch_size},␣LR={result['sgd']['lr']:.4f}")

    plt.xlabel("Epoch")
    plt.ylabel("Loss")
    plt.title(f"SGD␣Convergence␣for␣Different␣Batch␣Sizes␣(Noise={noise},␣Seed={seed})")
    plt.grid(True, alpha=0.3)
    plt.legend ()

    # Plot Polyak convergence
    plt.subplot(4, 1, 2)
    for batch_size , result in sorted_items:
        polyak_losses = result["polyak"]["losses"]
        plt.plot(polyak_losses , label=f"Batch={batch_size}")

    plt.xlabel("Epoch")
    plt.ylabel("Loss")
    plt.title(f"Polyak␣SGD␣Convergence␣for␣Different␣Batch␣Sizes␣(Noise={noise},␣
        Seed={seed})")
    plt.grid(True, alpha=0.3)
    plt.legend ()

    # Plot PolyakAdam convergence
    plt.subplot(4, 1, 3)
    for batch_size , result in sorted_items:
        polyak_adam_losses = result["polyak_adam"]["losses"]
        plt.plot(polyak_adam_losses , label=f"Batch={batch_size}")

    plt.xlabel("Epoch")
    plt.ylabel("Loss")
    plt.title(f"PolyakAdam␣Convergence␣for␣Different␣Batch␣Sizes␣(Noise={noise},␣
        Seed={seed})")
    plt.grid(True, alpha=0.3)
    plt.legend ()

    # Plot PolyakAdam+AMSGrad convergence
    plt.subplot(4, 1, 4)
    for batch_size , result in sorted_items:
        polyak_adam_ams_losses = result["polyak_adam_ams"]["losses"]
        plt.plot(polyak_adam_ams_losses , label=f"Batch={batch_size}")

    plt.xlabel("Epoch")
    plt.ylabel("Loss")
```

```python
    plt.title(f"PolyakAdam+AMSGrad␣Convergence␣for␣Different␣Batch␣Sizes␣(Noise={noise},␣
        Seed={seed})")
    plt.grid(True, alpha=0.3)
    plt.legend()

    plt.tight_layout()
    plt.savefig(os.path.join(log_dir, "convergence_comparison.png"))
    plt.close()


def main() -> None:
    args = parse_args()
    config = load_config(args.config)

    # Set device
    device = get_device()
    print(f"Using␣device:␣{device}")

    # Create log directory
    timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
    base_log_dir = os.path.join(config["logging"]["log_dir"], timestamp)
    os.makedirs(base_log_dir, exist_ok=True)

    # Save configuration
    with open(os.path.join(base_log_dir, "config.yaml"), "w") as f:
        yaml.dump(config, f)

    # Load previously saved best learning rates if provided, otherwise do grid search
    best_lr_grid = {}

    if args.best_lr_path:
        print(f"Loading␣best␣learning␣rates␣from␣{args.best_lr_path}")
        try:
            with open(args.best_lr_path, "r") as f:
                serialized_lr_grid = json.load(f)

            # Convert keys back to tuples (noise_std, batch_size, seed)
            for key, value in serialized_lr_grid.items():
                # Parse the key format "noise{noise_std}_batch{batch_size}_seed{seed}"
                parts = key.replace("noise", "").replace("batch", "").replace("seed",
                    "").split("_")
                noise_std = float(parts[0])

                # Handle 'full' batch size
                batch_part = parts[1]
                batch_size = int(batch_part) if batch_part.isdigit() else batch_part

                seed = int(parts[2])
                best_lr_grid[(noise_std, batch_size, seed)] = value

            # Copy the loaded file to the new log directory
            shutil.copy(args.best_lr_path, os.path.join(base_log_dir,
                "best_learning_rates.json"))
            print(f"Loaded␣{len(best_lr_grid)}␣learning␣rate␣configurations")

        except Exception as e:
            print(f"Error␣loading␣best␣learning␣rates:␣{str(e)}")
            print("Falling␣back␣to␣grid␣search")
            best_lr_grid = {}

    # Perform grid search if no valid best learning rates were loaded
    if not best_lr_grid:
        if args.demo_mode:
            raise ValueError("Demo␣mode␣requires␣providing␣a␣valid␣--best_lr_path")

        print("Performing␣learning␣rate␣grid␣search")
        lr_search_seed = config["data"]["seeds"][0]  # Use only the first seed for LR
            search

        for noise_std in config["data"]["noise_std_values"]:
            for batch_size in config["training"]["batch_sizes"]:
                # Find best LR for this noise and batch size configuration
                best_lr = find_best_lr(
                    config=config,
                    noise_std=noise_std,
                    batch_size=batch_size,
```

```python
                    seed=lr_search_seed,
                    device=device,
                    log_dir=base_log_dir
                )

                # Store best LR for all seeds with this noise_std and batch_size
                for seed in config["data"]["seeds"]:
                    key = (noise_std, batch_size, seed)
                    best_lr_grid[key] = best_lr

        # Save best learning rates
        with open(os.path.join(base_log_dir, "best_learning_rates.json"), "w") as f:
            # Convert to serializable format
            serializable_lr_grid = {
                f"noise{k[0]}_batch{k[1]}_seed{k[2]}": v
                for k, v in best_lr_grid.items()
            }
            json.dump(serializable_lr_grid, f, indent=2)

    # Run in demo mode if specified
    if args.demo_mode:
        run_demo_comparison(
            config=config,
            best_lr_grid=best_lr_grid,
            demo_noise=args.demo_noise,
            demo_seed=args.demo_seed,
            device=device,
            log_dir=base_log_dir
        )
        return

    # Otherwise run full comparison for all configurations
    results = {}

    for key, best_lr in best_lr_grid.items():
        noise_std, batch_size, seed = key

        # If no good learning rate found, use a small default value
        if best_lr is None:
            print(f"Warning:␣No␣good␣learning␣rate␣found␣for␣noise_std={noise_std},␣
                batch_size={batch_size},␣seed={seed}")
            best_lr = 0.001   # Default to a small learning rate

        result = compare_optimizers(
            config=config,
            noise_std=noise_std,
            batch_size=batch_size,
            seed=seed,
            best_lr=best_lr,
            device=device,
            log_dir=base_log_dir
        )

        results[key] = result

    # Generate summary statistics and plots
    generate_summary(results, base_log_dir)

    print(f"Experiments␣completed.␣Results␣saved␣to␣{base_log_dir}")


def generate_summary(results: Dict[Tuple[float, Union[int, str], int], Dict[str, Any]],
    log_dir: str) -> None:
    """
    Generate summary statistics and plots for the experiments.

    Args:
        results: Dictionary of experiment results
        log_dir: Directory to save summary results
    """
    summary_dir = os.path.join(log_dir, "summary")
    os.makedirs(summary_dir, exist_ok=True)

    # Create summary tables and plots
    noise_values = sorted(set(k[0] for k in results.keys()))
```

```python
    # Convert all batch sizes to strings before sorting to avoid type comparison issues
    batch_sizes = [str(k[1]) for k in results.keys()]
    # Sort batch sizes with "full" at the end
    batch_sizes = sorted(set(batch_sizes), key=lambda x: (x == "full", x if x != "full"
        else "z"))

    # Summary of final losses
    summary_data = {
        "noise_std": [],
        "batch_size": [],
        "sgd_final_loss_mean": [],
        "polyak_final_loss_mean": [],
        "polyak_adam_final_loss_mean": [],
        "polyak_adam_ams_final_loss_mean": [],
        "sgd_final_loss_std": [],
        "polyak_final_loss_std": [],
        "polyak_adam_final_loss_std": [],
        "polyak_adam_ams_final_loss_std": []
    }

    for noise_std in noise_values:
        for batch_size in batch_sizes:
            # Collect results for this configuration across seeds
            config_results = [
                results[key]
                for key in results
                if key[0] == noise_std and str(key[1]) == batch_size
            ]

            if config_results:
                sgd_final_losses = [r["sgd"]["losses"][-1] for r in config_results]
                polyak_final_losses = [r["polyak"]["losses"][-1] for r in config_results]
                polyak_adam_final_losses = [r["polyak_adam"]["losses"][-1] for r in
                    config_results]
                polyak_adam_ams_final_losses = [r["polyak_adam_ams"]["losses"][-1] for r
                    in config_results]

                summary_data["noise_std"].append(noise_std)
                summary_data["batch_size"].append(batch_size)
                summary_data["sgd_final_loss_mean"].append(np.mean(sgd_final_losses))
                summary_data["polyak_final_loss_mean"].append(np.mean(polyak_final_losses))
                summary_data["polyak_adam_final_loss_mean"].append(np.mean(polyak_adam_final_losses))
                summary_data["polyak_adam_ams_final_loss_mean"].append(np.mean(polyak_adam_ams_final_los
                summary_data["sgd_final_loss_std"].append(np.std(sgd_final_losses))
                summary_data["polyak_final_loss_std"].append(np.std(polyak_final_losses))
                summary_data["polyak_adam_final_loss_std"].append(np.std(polyak_adam_final_losses))
                summary_data["polyak_adam_ams_final_loss_std"].append(np.std(polyak_adam_ams_final_losse

    # Save summary data
    with open(os.path.join(summary_dir, "summary_data.json"), "w") as f:
        # Convert to serializable format
        serializable_summary = {
            "noise_std": summary_data["noise_std"],
            "batch_size": [str(bs) for bs in summary_data["batch_size"]],
            "sgd_final_loss_mean": summary_data["sgd_final_loss_mean"],
            "polyak_final_loss_mean": summary_data["polyak_final_loss_mean"],
            "polyak_adam_final_loss_mean": summary_data["polyak_adam_final_loss_mean"],
            "polyak_adam_ams_final_loss_mean":
                summary_data["polyak_adam_ams_final_loss_mean"],
            "sgd_final_loss_std": summary_data["sgd_final_loss_std"],
            "polyak_final_loss_std": summary_data["polyak_final_loss_std"],
            "polyak_adam_final_loss_std": summary_data["polyak_adam_final_loss_std"],
            "polyak_adam_ams_final_loss_std":
                summary_data["polyak_adam_ams_final_loss_std"]
        }
        json.dump(serializable_summary, f, indent=2)

    # Create comparison plots

    # Plot by noise level
    plt.figure(figsize=(12, 8))
    for i, batch_size in enumerate(batch_sizes):
        # Find relevant indices
        indices = [j for j, bs in enumerate(summary_data["batch_size"]) if bs ==
            batch_size]
```

```python
        if indices:
            noise_vals = [summary_data["noise_std"][j] for j in indices]
            sgd_means = [summary_data["sgd_final_loss_mean"][j] for j in indices]
            polyak_means = [summary_data["polyak_final_loss_mean"][j] for j in indices]
            polyak_adam_means = [summary_data["polyak_adam_final_loss_mean"][j] for j in
                indices]
            polyak_adam_ams_means = [summary_data["polyak_adam_ams_final_loss_mean"][j]
                for j in indices]

            plt.plot(noise_vals, sgd_means, marker='o', label=f"SGD,batch={batch_size}")
            plt.plot(noise_vals, polyak_means, marker='x', linestyle='--',
                label=f"Polyak,batch={batch_size}")
            plt.plot(noise_vals, polyak_adam_means, marker='*', linestyle='-.',
                label=f"PolyakAdam,batch={batch_size}")
            plt.plot(noise_vals, polyak_adam_ams_means, marker='*--',
                label=f"PolyakAdam+AMSGrad,batch={batch_size}")

    plt.xlabel("Noise Standard Deviation")
    plt.ylabel("Final Loss (Mean)")
    plt.title("Effect of Noise on Final Loss for Different Optimizers and Batch Sizes")
    plt.legend()
    plt.grid(True, alpha=0.3)
    plt.savefig(os.path.join(summary_dir, "noise_effect.png"))
    plt.close()

    # Plot by batch size
    plt.figure(figsize=(12, 8))
    for i, noise_std in enumerate(noise_values):
        # Find relevant indices
        indices = [j for j, ns in enumerate(summary_data["noise_std"]) if ns ==
            noise_std]

        if indices:
            batch_vals = [summary_data["batch_size"][j] for j in indices]
            sgd_means = [summary_data["sgd_final_loss_mean"][j] for j in indices]
            polyak_means = [summary_data["polyak_final_loss_mean"][j] for j in indices]
            polyak_adam_means = [summary_data["polyak_adam_final_loss_mean"][j] for j in
                indices]
            polyak_adam_ams_means = [summary_data["polyak_adam_ams_final_loss_mean"][j]
                for j in indices]

            # Handle string batch sizes for plotting
            x_positions = list(range(len(batch_vals)))

            plt.plot(x_positions, sgd_means, marker='o', label=f"SGD,noise={noise_std}")
            plt.plot(x_positions, polyak_means, marker='x', linestyle='--',
                label=f"Polyak,noise={noise_std}")
            plt.plot(x_positions, polyak_adam_means, marker='*', linestyle='-.',
                label=f"PolyakAdam,noise={noise_std}")
            plt.plot(x_positions, polyak_adam_ams_means, marker='*--',
                label=f"PolyakAdam+AMSGrad,noise={noise_std}")

    plt.xlabel("Batch Size Index")
    plt.ylabel("Final Loss (Mean)")
    plt.title("Effect of Batch Size on Final Loss for Different Optimizers and Noise
        Levels")
    plt.xticks(list(range(len(batch_sizes))), [str(bs) for bs in batch_sizes])
    plt.legend()
    plt.grid(True, alpha=0.3)
    plt.savefig(os.path.join(summary_dir, "batch_size_effect.png"))
    plt.close()


if __name__ == "__main__":
    main()
```

**week6_experiment.py**

```python
import os
import yaml
import argparse
import numpy as np
import torch
import json
import matplotlib.pyplot as plt
```

```python
from datetime import datetime
from typing import Dict, List, Union, Tuple, Optional, Any
from pathlib import Path

from src.utils import set_seed, get_device, get_data_loader
from src.datasets import make_gaussian_cluster
from src.models.week6_model import Week6Model
from src.optim.sgd import SGD
from src.optim.polyak_sgd import PolyakSGD


def parse_args() -> argparse.Namespace:
    """Parse command line arguments."""
    parser = argparse.ArgumentParser(description="Run Week 6 optimization experiments")
    parser.add_argument(
        "--config",
        type=str,
        default="configs/week6_experiment_config.yaml",
        help="Path to configuration YAML file"
    )
    return parser.parse_args()


def load_config(config_path: str) -> Dict[str, Any]:
    """Load configuration from YAML file."""
    with open(config_path, "r") as f:
        config = yaml.safe_load(f)
    return config


def train_model(
    model: torch.nn.Module,
    optimizer: torch.optim.Optimizer,
    data_loader: torch.utils.data.DataLoader,
    device: torch.device,
    epochs: int,
    is_polyak: bool = False
) -> Tuple[List[float], List[float], List[Tuple[float, float]]]:
    """
    Train a model with the given optimizer and data loader.

    Args:
        model: PyTorch model to train
        optimizer: Optimizer to use for training
        data_loader: DataLoader providing training data
        device: Device to run training on
        epochs: Number of epochs to train for
        is_polyak: Whether the optimizer is PolyakSGD

    Returns:
        Tuple of (training losses, step sizes if Polyak, parameter history)
    """
    model.train()
    losses = []
    step_sizes = [] if is_polyak else None
    param_history = []  # Store x parameter values for contour plotting

    # Store initial parameter values
    with torch.no_grad():
        param_history.append((model.x[0].item(), model.x[1].item()))

    for epoch in range(epochs):
        epoch_loss = 0.0
        batches = 0

        for (minibatch,) in data_loader:
            minibatch = minibatch.to(device)

            def closure():
                optimizer.zero_grad()
                loss = model(minibatch)
                loss.backward()
                return loss

            if is_polyak:
                loss = optimizer.step(closure)
```

```python
                    step_sizes.append(optimizer.last_step_size)
                else:
                    loss = closure()
                    optimizer.step()

                epoch_loss += loss.item()
                batches += 1

                # Store current parameter values
                with torch.no_grad():
                    param_history.append((model.x[0].item(), model.x[1].item()))

            avg_epoch_loss = epoch_loss / batches
            losses.append(avg_epoch_loss)

    return losses, step_sizes, param_history


def compare_optimizers(
    config: Dict[str, Any],
    batch_size: Union[int, str],
    device: torch.device,
    log_dir: str
) -> Dict[str, Any]:
    """
    Compare Constant Step Size SGD with Polyak Step Size SGD.

    Args:
        config: Configuration dictionary
        batch_size: Batch size for training (int or "full")
        device: Device to run training on
        log_dir: Directory to save logs

    Returns:
        Dictionary with results
    """
    print(f"Comparing optimizers: batch_size={batch_size}")

    # Generate dataset
    m = config["data"]["m"]  # Number of data points
    std = config["data"]["std"]  # Standard deviation of Gaussian noise
    seed = config["data"]["seed"]  # Random seed
    epochs = config["training"]["epochs"]

    X = make_gaussian_cluster(m=m, std=std, seed=seed)

    # Use full batch if specified
    actual_batch_size = m if batch_size == "full" else batch_size
    data_loader = get_data_loader(X, batch_size=actual_batch_size)

    # Directory for comparison results
    comp_log_dir = os.path.join(
        log_dir,
        f"comparison_batch{batch_size}"
    )
    os.makedirs(comp_log_dir, exist_ok=True)

    # Train with Constant Step Size SGD
    lr = config["optimizer"]["sgd"]["lr"]
    model_sgd = Week6Model().to(device)
    optimizer_sgd = SGD(model_sgd.parameters(), lr=lr)

    sgd_losses, _, sgd_param_history = train_model(
        model=model_sgd,
        optimizer=optimizer_sgd,
        data_loader=data_loader,
        device=device,
        epochs=epochs
    )

    # Train with Polyak Step Size SGD
    model_polyak = Week6Model().to(device)
    optimizer_polyak = PolyakSGD(
        model_polyak.parameters(),
        eps=config["optimizer"]["polyak"]["eps"],
        f_star=config["optimizer"]["polyak"]["f_star"]
```

```python
    )

    polyak_losses , polyak_step_sizes , polyak_param_history = train_model (
        model = model_polyak ,
        optimizer = optimizer_polyak ,
        data_loader = data_loader ,
        device = device ,
        epochs = epochs ,
        is_polyak = True
    )

    # Save results
    results = {
        "sgd": {
            "losses": sgd_losses ,
            "lr": lr,
            "param_history": sgd_param_history
        },
        "polyak": {
            "losses": polyak_losses ,
            "step_sizes": polyak_step_sizes ,
            "param_history": polyak_param_history
        },
        "config": {
            "batch_size": batch_size ,
            "epochs": epochs
        }
    }

    # Save results to file
    with open(os.path.join(comp_log_dir , "comparison_results.json"), "w") as f:
        # Convert param_history to serializable format
        results_copy = results.copy ()
        results_copy["sgd"]["param_history"] = [[float(x), float(y)] for x, y in
            results["sgd"]["param_history"]]
        results_copy["polyak"]["param_history"] = [[float(x), float(y)] for x, y in
            results["polyak"]["param_history"]]
        json.dump(results_copy , f, indent =2)

    # Plot loss comparison
    plt.figure(figsize=(10, 6))
    plt.plot(sgd_losses , label=f"Constant␣Step␣Size␣(LR={lr})")
    plt.plot(polyak_losses , label="Polyak␣Step␣Size")
    plt.xlabel("Epoch")
    plt.ylabel("Loss")
    plt.title(f"Loss␣Comparison:␣Constant␣vs␣Polyak␣Step␣Size\nBatch␣Size={batch_size}")
    plt.legend ()
    plt.savefig(os.path.join(comp_log_dir , "loss_comparison.png"))
    plt.close ()

    # Plot Polyak step sizes
    plt.figure(figsize=(10, 6))
    plt.plot(polyak_step_sizes)
    plt.xlabel("Batch␣Update")
    plt.ylabel("Step␣Size")
    plt.title(f"Polyak␣Step␣Sizes␣Over␣Training\nBatch␣Size={batch_size}")
    plt.savefig(os.path.join(comp_log_dir , "polyak_step_sizes.png"))
    plt.close ()

    return results


def create_batch_size_comparison_plot (
    results: Dict[Union[int, str], Dict],
    log_dir: str
) -> None:
    """
    Create a plot comparing final losses for different batch sizes.

    Args:
        results: Dictionary mapping batch sizes to results
        log_dir: Directory to save plot
    """
    plt.figure(figsize=(10, 6))

    # Extract batch sizes and prepare lists for plotting
```

```python
    batch_sizes = sorted(results.keys(), key=lambda bs: int(bs) if bs != "full" else
        float('inf'))
    batch_labels = [str(bs) for bs in batch_sizes]
    sgd_final_losses = [results[bs]["sgd"]["losses"][-1] for bs in batch_sizes]
    polyak_final_losses = [results[bs]["polyak"]["losses"][-1] for bs in batch_sizes]

    x = list(range(len(batch_sizes)))

    # Plot final losses
    plt.plot(x, sgd_final_losses, 'o-', label="SGD (Constant Step Size)")
    plt.plot(x, polyak_final_losses, 'x--', label="Polyak Step Size")
    plt.xticks(x, batch_labels)
    plt.xlabel("Batch Size")
    plt.ylabel("Final Loss")
    plt.title("Effect of Batch Size on Final Loss")
    plt.grid(True, alpha=0.3)
    plt.legend()

    plt.tight_layout()
    plt.savefig(os.path.join(log_dir, "batch_size_comparison.png"))
    plt.close()


def create_convergence_comparison_plot(
    results: Dict[Union[int, str], Dict],
    log_dir: str
) -> None:
    """
    Create a plot comparing convergence speed for different batch sizes.

    Args:
        results: Dictionary mapping batch sizes to results
        log_dir: Directory to save plot
    """
    plt.figure(figsize=(12, 10))

    # Sort batch sizes for consistent plotting
    # Convert batch sizes to strings for consistent sorting
    sorted_items = sorted(results.items(), key=lambda x: (
            str(x[0]) == "full",  # Sort "full" last
            int(x[0]) if str(x[0]) != "full" else float('inf')  # Sort numerically
        ))

    # Plot SGD convergence
    plt.subplot(2, 1, 1)
    for batch_size, result in sorted_items:
        sgd_losses = result["sgd"]["losses"]
        plt.plot(sgd_losses, label=f"Batch={batch_size}, LR={result['sgd']['lr']:.4f}")

    plt.xlabel("Epoch")
    plt.ylabel("Loss")
    plt.title("SGD Convergence for Different Batch Sizes")
    plt.grid(True, alpha=0.3)
    plt.legend()

    # Plot Polyak convergence
    plt.subplot(2, 1, 2)
    for batch_size, result in sorted_items:
        polyak_losses = result["polyak"]["losses"]
        plt.plot(polyak_losses, label=f"Batch={batch_size}")

    plt.xlabel("Epoch")
    plt.ylabel("Loss")
    plt.title("Polyak SGD Convergence for Different Batch Sizes")
    plt.grid(True, alpha=0.3)
    plt.legend()

    plt.tight_layout()
    plt.savefig(os.path.join(log_dir, "convergence_comparison.png"))
    plt.close()


def create_contour_plots(
    results: Dict[Union[int, str], Dict],
    log_dir: str
) -> None:
```

```python
    """
    Create contour plots showing the optimization paths for both optimizers.

    Args:
        results: Dictionary mapping batch sizes to results
        log_dir: Directory to save plots
    """
    # Define loss function to create contour
    def loss_func(x, y, training_data):
        loss = 0.0
        count = 0
        for w in training_data:
            # Move tensor to CPU before converting to numpy
            w = w.cpu().numpy()
            z = np.array([x, y]) - w - 1
            term1 = 20 * (z[0]**2 + z[1]**2)
            term2 = (z[0] + 9)**2 + (z[1] + 10)**2
            point_loss = min(term1, term2)
            loss += point_loss
            count += 1
        return loss / count

    # Create mesh grid for contour
    x_min, x_max = -15, 15
    y_min, y_max = -15, 15
    n_points = 100

    X = np.linspace(x_min, x_max, n_points)
    Y = np.linspace(y_min, y_max, n_points)
    X_grid, Y_grid = np.meshgrid(X, Y)

    # Get training data
    m = results[next(iter(results.keys()))]["config"]["epochs"]
    std = 0.25  # Using default from make_gaussian_cluster
    seed = 42   # Using default seed
    training_data = make_gaussian_cluster(m=m, std=std, seed=seed)

    # Compute loss landscape
    Z = np.zeros_like(X_grid)
    for i in range(n_points):
        for j in range(n_points):
            Z[i, j] = loss_func(X_grid[i, j], Y_grid[i, j], training_data)

    # Create separate contour plots for SGD and Polyak for each batch size
    for batch_size, result in results.items():
        # SGD contour plot
        plt.figure(figsize=(10, 6))
        plt.contour(X_grid, Y_grid, Z, levels=20)

        # Get SGD parameter history
        sgd_params = result["sgd"]["param_history"]
        sgd_x = [p[0] for p in sgd_params]
        sgd_y = [p[1] for p in sgd_params]

        # Plot SGD trajectory
        plt.plot(sgd_x, sgd_y, marker='.', label=f'SGD (LR={result["sgd"]["lr"]})')
        plt.plot(sgd_x[0], sgd_y[0], 'rx', markersize=10, label='Start')
        plt.plot(sgd_x[-1], sgd_y[-1], 'ro', markersize=10, label='End')

        plt.xlabel('x[0]')
        plt.ylabel('x[1]')
        plt.title(f'SGD Trajectory on Loss Contour (Batch Size={batch_size})')
        plt.legend()
        plt.savefig(os.path.join(log_dir, f"sgd_contour_batch{batch_size}.png"))
        plt.close()

        # Polyak contour plot
        plt.figure(figsize=(10, 6))
        plt.contour(X_grid, Y_grid, Z, levels=20)

        # Get Polyak parameter history
        polyak_params = result["polyak"]["param_history"]
        polyak_x = [p[0] for p in polyak_params]
        polyak_y = [p[1] for p in polyak_params]

        # Plot Polyak trajectory
```

36

```python
        plt.plot(polyak_x, polyak_y, marker='.', label='Polyak Step Size')
        plt.plot(polyak_x[0], polyak_y[0], 'bx', markersize=10, label='Start')
        plt.plot(polyak_x[-1], polyak_y[-1], 'bo', markersize=10, label='End')

        plt.xlabel('x[0]')
        plt.ylabel('x[1]')
        plt.title(f'Polyak Trajectory on Loss Contour (Batch Size={batch_size})')
        plt.legend()
        plt.savefig(os.path.join(log_dir, f"polyak_contour_batch{batch_size}.png"))
        plt.close()

    # Create comparison plots for different batch sizes
    # SGD comparison
    plt.figure(figsize=(10, 6))
    plt.contour(X_grid, Y_grid, Z, levels=20)

    for batch_size, result in results.items():
        sgd_params = result["sgd"]["param_history"]
        sgd_x = [p[0] for p in sgd_params]
        sgd_y = [p[1] for p in sgd_params]
        plt.plot(sgd_x, sgd_y, marker='.', label=f'Batch={batch_size}')

    plt.xlabel('x[0]')
    plt.ylabel('x[1]')
    plt.title('SGD Trajectories for Different Batch Sizes')
    plt.legend()
    plt.savefig(os.path.join(log_dir, "sgd_batch_comparison.png"))
    plt.close()

    # Polyak comparison
    plt.figure(figsize=(10, 6))
    plt.contour(X_grid, Y_grid, Z, levels=20)

    for batch_size, result in results.items():
        polyak_params = result["polyak"]["param_history"]
        polyak_x = [p[0] for p in polyak_params]
        polyak_y = [p[1] for p in polyak_params]
        plt.plot(polyak_x, polyak_y, marker='.', label=f'Batch={batch_size}')

    plt.xlabel('x[0]')
    plt.ylabel('x[1]')
    plt.title('Polyak Trajectories for Different Batch Sizes')
    plt.legend()
    plt.savefig(os.path.join(log_dir, "polyak_batch_comparison.png"))
    plt.close()


def main() -> None:
    args = parse_args()
    config = load_config(args.config)

    # Set device
    device = get_device()
    print(f"Using device: {device}")

    # Create log directory
    timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
    base_log_dir = os.path.join(config["logging"]["log_dir"], f"week6_{timestamp}")
    os.makedirs(base_log_dir, exist_ok=True)

    # Save configuration
    with open(os.path.join(base_log_dir, "config.yaml"), "w") as f:
        yaml.dump(config, f)

    # Run experiments for each batch size
    results = {}

    for batch_size in config["training"]["batch_sizes"]:
        result = compare_optimizers(
            config=config,
            batch_size=batch_size,
            device=device,
            log_dir=base_log_dir
        )

        results[batch_size] = result
```

```python
    # Create summary plots
    create_batch_size_comparison_plot(results, base_log_dir)
    create_convergence_comparison_plot(results, base_log_dir)

    # Create contour plots
    create_contour_plots(results, base_log_dir)

    print(f"Experiments completed. Results saved to {base_log_dir}")


if __name__ == "__main__":
    main()
```

**week9_transformer.py**

```python
import torch
import os
import yaml
import json
import argparse
import matplotlib.pyplot as plt
import numpy as np
from typing import Dict, List, Tuple, Optional, Any
from pathlib import Path
from datetime import datetime
from itertools import product

from src.models.transformer_data import TransformerDataManager
from src.models.transformer_week9 import GPTLanguageModel, calculate_baseline_loss
from src.utils import set_seed, get_device
from src.optim.polyak_sgd import PolyakSGD
from src.optim.sgd import SGD


def parse_args() -> argparse.Namespace:
    """Parse command line arguments.

    Returns:
        Command line arguments
    """
    parser = argparse.ArgumentParser(description="Compare different optimizers for
        transformer training")
    parser.add_argument(
        "--config",
        type=str,
        default="configs/week9_transformer_config.yaml",
        help="Path to configuration YAML file"
    )
    parser.add_argument(
        "--tune",
        action="store_true",
        help="Enable hyperparameter tuning (overrides config setting)"
    )
    return parser.parse_args()


def load_config(config_path: str) -> Dict[str, Any]:
    """Load configuration from YAML file.

    Args:
        config_path: Path to YAML configuration file

    Returns:
        Configuration dictionary
    """
    with open(config_path, "r") as f:
        config = yaml.safe_load(f)
    return config


def train_with_optimizer(
    model: GPTLanguageModel,
    data_manager: TransformerDataManager,
    optimizer_name: str,
```

```python
    optimizer_kwargs: Dict,
    max_iters: int,
    eval_interval: int,
    early_stopping_config: Optional[Dict] = None,
    model_save_dir: Optional[Path] = None,
    return_best: bool = True
) -> Tuple[Dict[str, List[float]], GPTLanguageModel]:
    """
    Train a transformer model with a specified optimizer.

    Args:
        model: The transformer model to train
        data_manager: Data manager for getting batches
        optimizer_name: Name of the optimizer (e.g., 'Adam', 'SGD', 'AdamW')
        optimizer_kwargs: Parameters for the optimizer
        max_iters: Maximum number of training iterations
        eval_interval: Interval for evaluation
        early_stopping_config: Configuration for early stopping
        model_save_dir: Directory to save model checkpoints
        return_best: Whether to return the best model based on validation loss

    Returns:
        Tuple containing:
            - metrics: Dictionary with training metrics (train_losses, val_losses,
                iterations)
            - best_model: Best model if return_best is True, otherwise final model
    """
    # Initialize optimizer
    if optimizer_name == "Polyak":
        optimizer = PolyakSGD(model.parameters(), **optimizer_kwargs)
    elif optimizer_name == "SGD":
        optimizer = SGD(model.parameters(), **optimizer_kwargs)
    else:
        optimizer_class = getattr(torch.optim, optimizer_name)
        optimizer = optimizer_class(model.parameters(), **optimizer_kwargs)

    # Initialize metrics tracking
    metrics = {
        'train_losses': [],
        'val_losses': [],
        'iterations': [],
    }

    # Add step sizes list for Polyak optimizer
    if optimizer_name == "Polyak":
        metrics['step_sizes'] = []

    # Setup model saving
    if model_save_dir:
        model_save_dir.mkdir(exist_ok=True, parents=True)
        best_val_loss = float('inf')
        best_model_path = model_save_dir / f"{optimizer_name}_best_model.pt"
    else:
        best_val_loss = float('inf')

    # Save initial model if desired
    best_model = model

    # Set up early stopping if configured
    early_stopping_triggered = False
    no_improvement_count = 0
    min_eval_iters = 0

    if early_stopping_config and early_stopping_config.get('enabled', False):
        patience = early_stopping_config.get('patience', 5)
        min_delta = early_stopping_config.get('min_delta', 0.001)
        min_epochs = early_stopping_config.get('min_epochs', 1)
        # Calculate minimum iterations before early stopping can trigger
        min_eval_iters = min_epochs * (data_manager.train_data.size(0) //
            data_manager.batch_size)
        min_eval_iters = min_eval_iters // eval_interval * eval_interval  # Round to
            eval_interval

    # Training loop
    for iter_num in range(max_iters):
        # Evaluation
```

```python
        if iter_num % eval_interval == 0 or iter_num == max_iters - 1:
            losses = model.estimate_loss(data_manager)
            train_loss = losses['train']
            val_loss = losses['val']

            # Record metrics
            metrics['train_losses'].append(train_loss)
            metrics['val_losses'].append(val_loss)
            metrics['iterations'].append(iter_num)

            # Handle early stopping
            if early_stopping_config and early_stopping_config.get('enabled', False) and \
                    iter_num >= min_eval_iters:
                if val_loss < best_val_loss - min_delta:
                    # Improvement found
                    no_improvement_count = 0
                else:
                    # No significant improvement
                    no_improvement_count += 1
                    if no_improvement_count >= patience:
                        early_stopping_triggered = True
                        print(f"[{optimizer_name}] Early stopping triggered at iteration "
                            f"{iter_num}: "
                                f"No improvement for {patience} evaluations")
                        # Add early stopping info to metrics
                        metrics['early_stopped'] = True
                        metrics['early_stopping_iter'] = iter_num
                        break

            print(f"[{optimizer_name}] step {iter_num}: train loss {train_loss:.4f}, val "
                f"loss {val_loss:.4f}")

            # Save best model
            if val_loss < best_val_loss:
                best_val_loss = val_loss
                if model_save_dir:
                    torch.save(model.state_dict(), best_model_path)
                if return_best:
                    best_model = GPTLanguageModel(data_manager)
                    if model_save_dir:
                        best_model.load_state_dict(torch.load(best_model_path))
                    else:
                        # If no save dir, we need to manually copy parameters
                        best_model.load_state_dict(model.state_dict())

            # Record Polyak step size if applicable
            if optimizer_name == "Polyak":
                metrics['step_sizes'].append(optimizer.last_step_size)

        # Get batch and calculate loss
        xb, yb = data_manager.get_batch('train')

        if optimizer_name == "Polyak":
            def closure():
                optimizer.zero_grad()
                logits, loss = model(xb, yb)
                loss.backward()
                return loss

            loss = optimizer.step(closure)
        else:
            logits, loss = model(xb, yb)

            # Update weights
            optimizer.zero_grad(set_to_none=True)
            loss.backward()
            optimizer.step()

    # Add final metrics
    metrics['best_val_loss'] = best_val_loss
    if not early_stopping_triggered and early_stopping_config and \
        early_stopping_config.get('enabled', False):
        metrics['early_stopped'] = False

    # Save final model
    if model_save_dir:
```

```python
        final_model_path = model_save_dir / f"{optimizer_name}_final_model.pt"
        torch.save(model.state_dict(), final_model_path)

    return metrics, best_model


def tune_sgd_learning_rate(
    data_manager: TransformerDataManager,
    lr_values: List[float],
    max_iters: int,
    eval_interval: int,
    device: torch.device,
    tune_dir: Path
) -> float:
    """
    Fine-tune the learning rate for SGD optimizer.

    Args:
        data_manager: Data manager for getting batches
        lr_values: List of learning rate values to try
        max_iters: Maximum number of training iterations
        eval_interval: Interval for evaluation
        device: Device to run training on
        tune_dir: Directory to save tuning results

    Returns:
        Best learning rate value
    """
    print("\n=== Tuning SGD Learning Rate ===")
    results = {}

    # Create directory for SGD tuning
    sgd_tune_dir = tune_dir / "sgd"
    sgd_tune_dir.mkdir(exist_ok=True, parents=True)

    # Try each learning rate
    for lr in lr_values:
        print(f"\nTrying SGD with lr={lr}")

        # Initialize model with same weights for fair comparison
        set_seed(data_manager.seed)
        model = GPTLanguageModel(data_manager).to(device)

        # Train model with this learning rate
        metrics, _ = train_with_optimizer(
            model=model,
            data_manager=data_manager,
            optimizer_name="SGD",
            optimizer_kwargs={"lr": lr},
            max_iters=max_iters,
            eval_interval=eval_interval
        )

        # Store results
        final_val_loss = metrics['val_losses'][-1]
        results[lr] = final_val_loss
        print(f"SGD with lr={lr}: Final val loss = {final_val_loss:.4f}")

    # Find best learning rate
    best_lr = min(results, key=results.get)
    best_loss = results[best_lr]

    # Save results
    with open(sgd_tune_dir / "results.json", "w") as f:
        json.dump({"lr_values": {str(lr): float(loss) for lr, loss in results.items()},
                   "best_lr": float(best_lr),
                   "best_loss": float(best_loss)}, f, indent=2)

    # Plot results
    plt.figure(figsize=(10, 6))
    plt.plot(list(results.keys()), list(results.values()), 'o-')
    plt.xscale('log')
    plt.xlabel('Learning Rate')
    plt.ylabel('Validation Loss')
    plt.title('SGD Learning Rate Tuning')
    plt.grid(True)
```

```python
    plt.savefig(sgd_tune_dir / "lr_tuning.png")
    plt.close()

    print(f"\nBest␣SGD␣Learning␣Rate:␣{best_lr}␣(Validation␣Loss:␣{best_loss:.4f})")

    return best_lr


def tune_adam_hyperparams(
    data_manager: TransformerDataManager,
    lr_values: List[float],
    beta1_values: List[float],
    beta2_values: List[float],
    max_iters: int,
    eval_interval: int,
    device: torch.device,
    tune_dir: Path
) -> Tuple[float, Tuple[float, float]]:
    """
    Fine-tune the learning rate and beta parameters for Adam optimizer.

    Args:
        data_manager: Data manager for getting batches
        lr_values: List of learning rate values to try
        beta1_values: List of beta1 values to try
        beta2_values: List of beta2 values to try
        max_iters: Maximum number of training iterations
        eval_interval: Interval for evaluation
        device: Device to run training on
        tune_dir: Directory to save tuning results

    Returns:
        Tuple containing:
            - Best learning rate
            - Best beta values (beta1, beta2)
    """
    print("\n===␣Tuning␣Adam␣Hyperparameters␣===")
    results = {}

    # Create directory for Adam tuning
    adam_tune_dir = tune_dir / "adam"
    adam_tune_dir.mkdir(exist_ok=True, parents=True)

    # Try each combination of parameters
    for lr, beta1, beta2 in product(lr_values, beta1_values, beta2_values):
        param_key = f"lr={lr},␣beta1={beta1},␣beta2={beta2}"
        print(f"\nTrying␣Adam␣with␣{param_key}")

        # Initialize model with same weights for fair comparison
        set_seed(data_manager.seed)
        model = GPTLanguageModel(data_manager).to(device)

        # Train model with these parameters
        metrics, _ = train_with_optimizer(
            model=model,
            data_manager=data_manager,
            optimizer_name="Adam",
            optimizer_kwargs={"lr": lr, "betas": (beta1, beta2)},
            max_iters=max_iters,
            eval_interval=eval_interval
        )

        # Store results
        final_val_loss = metrics['val_losses'][-1]
        results[param_key] = {"loss": final_val_loss, "params": {"lr": lr, "beta1":
            beta1, "beta2": beta2}}
        print(f"Adam␣with␣{param_key}:␣Final␣val␣loss␣=␣{final_val_loss:.4f}")

    # Find best parameters
    best_param_key = min(results, key=lambda k: results[k]["loss"])
    best_params = results[best_param_key]["params"]
    best_loss = results[best_param_key]["loss"]

    # Save results
    with open(adam_tune_dir / "results.json", "w") as f:
        json.dump({
```

```python
                "parameter_combinations": {k: {"loss": float(v["loss"]), "params":
                    v["params"]} for k, v in results.items()},
                "best_params": best_params,
                "best_loss": float(best_loss)
        }, f, indent=2)

    # Plot results as heatmap for each beta2 value
    for beta2 in beta2_values:
        plt.figure(figsize=(10, 6))

        # Filter results for this beta2
        filtered_results = {(r["params"]["lr"], r["params"]["beta1"]): r["loss"]
                            for _, r in results.items() if r["params"]["beta2"] == beta2}

        # Extract unique lr and beta1 values and sort them
        lr_list = sorted(set(lr for lr, _ in filtered_results.keys()))
        beta1_list = sorted(set(beta1 for _, beta1 in filtered_results.keys()))

        # Create 2D grid of loss values
        loss_grid = [[filtered_results.get((lr, beta1), float('nan')) for lr in lr_list]
            for beta1 in beta1_list]

        # Plot as heatmap
        plt.figure(figsize=(10, 6))
        plt.imshow(loss_grid, cmap='viridis')
        plt.colorbar(label='Validation Loss')
        plt.xticks(range(len(lr_list)), [f"{lr}" for lr in lr_list])
        plt.yticks(range(len(beta1_list)), [f"{beta1}" for beta1 in beta1_list])
        plt.xlabel('Learning Rate')
        plt.ylabel('Beta1')
        plt.title(f'Adam Hyperparameter Tuning (Beta2={beta2})')
        plt.savefig(adam_tune_dir / f"adam_tuning_beta2_{beta2}.png")
        plt.close()

    # Plot learning rate comparison for best beta values
    best_beta1 = best_params["beta1"]
    best_beta2 = best_params["beta2"]

    lr_results = {r["params"]["lr"]: r["loss"]
                  for _, r in results.items()
                  if r["params"]["beta1"] == best_beta1 and r["params"]["beta2"] ==
                      best_beta2}

    plt.figure(figsize=(10, 6))
    plt.plot(sorted(lr_results.keys()), [lr_results[lr] for lr in
        sorted(lr_results.keys())], 'o-')
    plt.xscale('log')
    plt.xlabel('Learning Rate')
    plt.ylabel('Validation Loss')
    plt.title(f'Adam Learning Rate (Best Betas: {best_beta1}, {best_beta2})')
    plt.grid(True)
    plt.savefig(adam_tune_dir / "adam_lr_tuning.png")
    plt.close()

    print(f"\nBest Adam Parameters: lr={best_params['lr']},
        beta1={best_params['beta1']}, beta2={best_params['beta2']} (Validation Loss:
        {best_loss:.4f})")

    return best_params["lr"], (best_params["beta1"], best_params["beta2"])


def compare_optimizers(
    data_manager: TransformerDataManager,
    optimizers: Dict[str, Dict],
    max_iters: int,
    eval_interval: int,
    device: torch.device,
    model_save_dir: Path,
    early_stopping_config: Optional[Dict] = None,
    test_file: Optional[str] = None
) -> Dict[str, Dict]:
    """
    Compare different optimizers for transformer training.

    Args:
        data_manager: Data manager for getting batches
```

```python
        optimizers: Dictionary mapping optimizer names to their parameters
        max_iters: Maximum number of training iterations
        eval_interval: Interval for evaluation
        device: Device to run training on
        model_save_dir: Directory to save model checkpoints
        early_stopping_config: Configuration for early stopping
        test_file: Optional test file path for final evaluation

    Returns:
        Dictionary with results for each optimizer
    """
    results = {}

    # Train with each optimizer
    for opt_name, opt_params in optimizers.items():
        print(f"\n=== Training with {opt_name} ===")

        # Initialize model with same weights for fair comparison
        set_seed(data_manager.seed)
        model = GPTLanguageModel(data_manager).to(device)

        # Create optimizer-specific save directory
        opt_save_dir = model_save_dir / opt_name
        opt_save_dir.mkdir(exist_ok=True, parents=True)

        # Train model with this optimizer
        metrics, best_model = train_with_optimizer(
            model=model,
            data_manager=data_manager,
            optimizer_name=opt_name,
            optimizer_kwargs=opt_params,
            max_iters=max_iters,
            eval_interval=eval_interval,
            early_stopping_config=early_stopping_config,
            model_save_dir=opt_save_dir
        )

        # Store results
        results[opt_name] = {
            'metrics': metrics,
            'best_model': best_model,
            'final_val_loss': metrics['val_losses'][-1],
            'best_val_loss': metrics.get('best_val_loss', metrics['val_losses'][-1]),
            'hyperparams': opt_params
        }

        # Evaluate on test data if provided
        if test_file:
            test_data = data_manager.load_test_dataset(test_file)
            test_loss = best_model.evaluate_test_loss(
                test_data=test_data,
                block_size=data_manager.block_size,
                batch_size=data_manager.batch_size,
                eval_iters=data_manager.eval_iters
            )
            results[opt_name]['test_loss'] = test_loss
            print(f"{opt_name} Test Loss: {test_loss:.4f}")

        # Print early stopping information if applicable
        if 'early_stopped' in metrics and metrics['early_stopped']:
            print(f"{opt_name}: Early stopped at iteration "
                  f"{metrics['early_stopping_iter']}")

    # Calculate and display baseline loss
    baseline_loss = calculate_baseline_loss(data_manager.vocab_size)
    print(f"\nBaseline Loss (uniform distribution): {baseline_loss:.4f}\n")

    # Compare final results
    print("\n=== Optimizer Comparison ===")
    for opt_name, result in results.items():
        early_stop_info = " (early stopped)" if result['metrics'].get('early_stopped',
            False) else ""
        print(f"{opt_name}{early_stop_info}: Best val loss = "
              f"{result['best_val_loss']:.4f}")
        if test_file and 'test_loss' in result:
            print(f"  Test loss = {result['test_loss']:.4f}")
```

```python
        # Save results to JSON
        results_file = model_save_dir / "optimizer_results.json"
        serializable_results = {}

        for opt_name, result in results.items():
            metrics_dict = {
                'train_losses': [float(x) for x in result['metrics']['train_losses']],
                'val_losses': [float(x) for x in result['metrics']['val_losses']],
                'iterations': result['metrics']['iterations'],
                'best_val_loss': float(result['best_val_loss'])
            }

            # Add early stopping info if present
            if 'early_stopped' in result['metrics']:
                metrics_dict['early_stopped'] = result['metrics']['early_stopped']
                if result['metrics']['early_stopped']:
                    metrics_dict['early_stopping_iter'] = \
                        result['metrics']['early_stopping_iter']

            # Add step_sizes for Polyak if available
            if 'step_sizes' in result['metrics']:
                # Handle None values that might be in step_sizes
                metrics_dict['step_sizes'] = [float(x) if x is not None else 0.0 for x in
                    result['metrics']['step_sizes']]

            serializable_results[opt_name] = {
                'metrics': metrics_dict,
                'final_val_loss': float(result['final_val_loss']),
                'best_val_loss': float(result['best_val_loss']),
                'hyperparams': result['hyperparams']
            }

            if 'test_loss' in result:
                serializable_results[opt_name]['test_loss'] = float(result['test_loss'])

    with open(results_file, 'w') as f:
        json.dump(serializable_results, f, indent=2)

    return results


def plot_training_curves(results: Dict[str, Dict], log_dir: Path) -> None:
    """
    Plot training and validation curves for different optimizers.

    Args:
        results: Results dictionary from compare_optimizers
        log_dir: Directory to save figures
    """
    plt.figure(figsize=(12, 8))

    # Plot training loss
    plt.subplot(2, 1, 1)
    for opt_name, result in results.items():
        metrics = result['metrics']
        plt.plot(metrics['iterations'], metrics['train_losses'], label=f"{opt_name}
            (train)")

        # Mark early stopping point if applicable
        if 'early_stopped' in metrics and metrics['early_stopped']:
            early_stop_idx = metrics['iterations'].index(metrics['early_stopping_iter'])
            plt.plot(metrics['iterations'][early_stop_idx],
                metrics['train_losses'][early_stop_idx],
                    'ro', markersize=8, label=f"{opt_name} early stop" if opt_name ==
                        list(results.keys())[0] else "")

    plt.title('Training Loss by Optimizer')
    plt.xlabel('Iterations')
    plt.ylabel('Loss')
    plt.legend()
    plt.grid(True)

    # Plot validation loss
    plt.subplot(2, 1, 2)
    for opt_name, result in results.items():
```

```python
    metrics = result['metrics']
    plt.plot(metrics['iterations'], metrics['val_losses'], label=f"{opt_name}␣(val)")

    # Mark early stopping point if applicable
    if 'early_stopped' in metrics and metrics['early_stopped']:
        early_stop_idx = metrics['iterations'].index(metrics['early_stopping_iter'])
        plt.plot(metrics['iterations'][early_stop_idx],
            metrics['val_losses'][early_stop_idx],
                'ro', markersize=8, label=f"{opt_name}␣early␣stop" if opt_name ==
                    list(results.keys())[0] else "")

plt.title('Validation␣Loss␣by␣Optimizer')
plt.xlabel('Iterations')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)

plt.tight_layout()

# Save the figure
save_path = log_dir / 'optimizer_comparison.png'
plt.savefig(save_path)
plt.close()

# Create individual plots for each optimizer
for opt_name, result in results.items():
    plt.figure(figsize=(10, 6))
    metrics = result['metrics']

    plt.plot(metrics['iterations'], metrics['train_losses'], label="Training␣Loss")
    plt.plot(metrics['iterations'], metrics['val_losses'], label="Validation␣Loss")

    plt.title(f'{opt_name}␣Learning␣Curves')
    plt.xlabel('Iterations')
    plt.ylabel('Loss')
    plt.legend()
    plt.grid(True)

    save_path = log_dir / f'{opt_name}_learning_curve.png'
    plt.savefig(save_path)
    plt.close()

# Plot generalization gap (NEW)
plt.figure(figsize=(12, 6))
for opt_name, result in results.items():
    metrics = result['metrics']
    # Calculate generalization gap (val_loss - train_loss)
    gen_gap = [val - train for val, train in zip(metrics['val_losses'],
        metrics['train_losses'])]
    plt.plot(metrics['iterations'], gen_gap, label=f"{opt_name}")

plt.title('Generalization␣Gap␣by␣Optimizer')
plt.xlabel('Iterations')
plt.ylabel('Validation␣Loss␣-␣Training␣Loss')
plt.legend()
plt.grid(True)
plt.axhline(y=0, color='r', linestyle='--', alpha=0.3)

save_path = log_dir / 'generalization_gap.png'
plt.savefig(save_path)
plt.close()

# Plot normalized generalization gap (NEW)
plt.figure(figsize=(12, 6))
for opt_name, result in results.items():
    metrics = result['metrics']
    # Calculate normalized generalization gap (val_loss/train_loss - 1)
    norm_gen_gap = [(val/train - 1) for val, train in zip(metrics['val_losses'],
        metrics['train_losses'])]
    plt.plot(metrics['iterations'], norm_gen_gap, label=f"{opt_name}")

plt.title('Normalized␣Generalization␣Gap␣by␣Optimizer')
plt.xlabel('Iterations')
plt.ylabel('(Validation␣Loss␣/␣Training␣Loss)␣-␣1')
plt.legend()
plt.grid(True)
```

```python
    plt.axhline(y=0, color='r', linestyle='--', alpha=0.3)

    save_path = log_dir / 'normalized_generalization_gap.png'
    plt.savefig(save_path)
    plt.close()

# Plot Polyak step size over time if available (NEW)
if any('step_sizes' in results[opt]['metrics'] for opt in results if opt ==
    'Polyak'):
    plt.figure(figsize=(10, 6))
    for opt_name, result in results.items():
        if opt_name == 'Polyak' and 'step_sizes' in result['metrics']:
            metrics = result['metrics']
            # Filter out None values and replace with 0.0
            step_sizes = np.array([s if s is not None else 0.0 for s in
                metrics['step_sizes']])
            plt.plot(metrics['iterations'], step_sizes, 'o-')

            # Add exponential moving average line for clearer trend
            window = min(len(step_sizes), 5)  # Use smaller of 5 or length of data
            if window > 1:
                weights = np.exp(np.linspace(-1., 0., window))
                weights /= weights.sum()
                smoothed = np.convolve(step_sizes, weights, mode='valid')
                # Plot starting from window-1 to align with original data
                plt.plot(metrics['iterations'][window-1:], smoothed, 'r-',
                    label='Exponential Moving Average', linewidth=2)

    plt.title('Polyak Step Size Evolution')
    plt.xlabel('Iterations')
    plt.ylabel('Step Size')
    plt.grid(True)
    if 'smoothed' in locals():
        plt.legend()

    # Add log scale y-axis version if values vary by orders of magnitude
    plt.yscale('log')
    save_path = log_dir / 'polyak_step_size_log.png'
    plt.savefig(save_path)

    # Also create linear scale version
    plt.yscale('linear')
    save_path = log_dir / 'polyak_step_size.png'
    plt.savefig(save_path)
    plt.close()

# Add early stopping summary if any optimizer was early stopped
if any('early_stopped' in results[opt]['metrics'] and
    results[opt]['metrics']['early_stopped']
        for opt in results):
    plt.figure(figsize=(8, 4))

    # Create a bar chart showing training duration for each optimizer
    opt_names = []
    durations = []
    colors = []

    for opt_name, result in results.items():
        metrics = result['metrics']
        opt_names.append(opt_name)

        if 'early_stopped' in metrics and metrics['early_stopped']:
            durations.append(metrics['early_stopping_iter'])
            colors.append('orange')
        else:
            durations.append(metrics['iterations'][-1])
            colors.append('blue')

    plt.bar(opt_names, durations, color=colors)
    plt.title('Training Duration by Optimizer')
    plt.xlabel('Optimizer')
    plt.ylabel('Iterations')
    plt.xticks(rotation=45)

    # Add a legend for early stopping
    import matplotlib.patches as mpatches
```

```python
        blue_patch = mpatches.Patch(color='blue', label='Completed Full Training')
        orange_patch = mpatches.Patch(color='orange', label='Early Stopped')
        plt.legend(handles=[blue_patch, orange_patch])

        plt.tight_layout()
        save_path = log_dir / 'early_stopping_summary.png'
        plt.savefig(save_path)
        plt.close()


def create_default_config() -> Dict[str, Any]:
    """
    Create default configuration if no config file is provided.

    Returns:
        Default configuration dictionary
    """
    return {
        "data": {
            "data_dir": "src/transformer-datasets",
            "dataset_file": "input_childSpeech_trainingSet.txt",
            "test_file": None,
            "batch_size": 64,
            "block_size": 256
        },
        "training": {
            "max_iters": 2000,
            "eval_interval": 200,
            "seed": 42,
            "tuning_iters": 500,
            "enable_tuning": False,
            "early_stopping": {
                "enabled": True,
                "patience": 3,
                "min_delta": 0.001,
                "min_epochs": 1
            }
        },
        "tuning": {
            "sgd": {
                "lr_values": [0.0001, 0.001, 0.01, 0.1]
            },
            "adam": {
                "lr_values": [0.0001, 0.001, 0.01],
                "beta1_values": [0.9, 0.95],
                "beta2_values": [0.999, 0.99]
            }
        },
        "optimizers": {
            "Adam": {"lr": 0.001, "betas": [0.9, 0.999]},
            "SGD": {"lr": 0.1},
            "Polyak": {"eps": 1e-8, "f_star": 0.0}
        },
        "logging": {
            "log_dir": "logs/transformer",
            "save_models": True
        }
    }


def main() -> None:
    """Main function to run the transformer training experiment."""
    args = parse_args()

    # Load configuration or create default
    try:
        config = load_config(args.config)
        print(f"Loaded configuration from {args.config}")
    except (FileNotFoundError, yaml.YAMLError) as e:
        print(f"Error loading config file: {e}")
        print("Using default configuration")
        config = create_default_config()

        # Save default config for reference
        os.makedirs(os.path.dirname(args.config), exist_ok=True)
        with open(args.config, "w") as f:
```

```python
            yaml.dump(config, f, default_flow_style=False)

    # Override enable_tuning if --tune flag is passed
    if args.tune:
        config["training"]["enable_tuning"] = True

    # Set device
    device = get_device()
    print(f"Using device: {device}")

    # Set random seed
    seed = config["training"].get("seed", 42)
    set_seed(seed)

    # Create log directory with timestamp
    timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
    base_log_dir = Path(config["logging"]["log_dir"]) / f"week9_{timestamp}"
    base_log_dir.mkdir(exist_ok=True, parents=True)

    # Save configuration
    with open(base_log_dir / "config.yaml", "w") as f:
        yaml.dump(config, f, default_flow_style=False)

    # Initialize data manager
    data_manager = TransformerDataManager(
        data_dir=config["data"]["data_dir"],
        dataset_file=config["data"]["dataset_file"],
        batch_size=config["data"]["batch_size"],
        block_size=config["data"]["block_size"],
        max_iters=config["training"]["max_iters"],
        eval_interval=config["training"]["eval_interval"],
        seed=seed
    )

    # Print dataset info
    data_stats = data_manager.analyze_dataset()
    print(f"Dataset: {config['data']['dataset_file']}")
    print(f"Train size: {data_stats['train_size']} tokens")
    print(f"Val size: {data_stats['val_size']} tokens")
    print(f"Vocabulary size: {data_stats['vocab_size']} characters")

    # Create directories
    model_dir = base_log_dir / "models"
    model_dir.mkdir(exist_ok=True, parents=True)

    # Get optimizers config
    optimizers = config["optimizers"].copy()

    # Get early stopping config
    early_stopping_config = config["training"].get("early_stopping")
    if early_stopping_config and early_stopping_config.get("enabled", False):
        print(f"Early stopping enabled with "
            f"patience={early_stopping_config.get('patience', 3)}, "
              f"min_delta={early_stopping_config.get('min_delta', 0.001)}")

    # Perform hyperparameter tuning if enabled
    if config["training"]["enable_tuning"]:
        print("\n=== Starting Hyperparameter Tuning ===")
        tuning_dir = base_log_dir / "tuning"
        tuning_dir.mkdir(exist_ok=True, parents=True)

        tuning_iters = config["training"]["tuning_iters"]
        eval_interval = config["training"]["eval_interval"]

        # Use early stopping during tuning but with shorter patience
        tuning_early_stopping = None
        if early_stopping_config and early_stopping_config.get("enabled", False):
            tuning_early_stopping = early_stopping_config.copy()
            # Use a shorter patience for tuning to speed things up
            tuning_early_stopping["patience"] = min(2,
                early_stopping_config.get("patience", 3))

        # Tune SGD learning rate
        sgd_lr = tune_sgd_learning_rate(
            data_manager=data_manager,
            lr_values=config["tuning"]["sgd"]["lr_values"],
```

```
                    max_iters=tuning_iters,
                    eval_interval=eval_interval,
                    device=device,
                    tune_dir=tuning_dir
                )

                # Update SGD optimizer config with best learning rate
                optimizers["SGD"]["lr"] = sgd_lr

                # Tune Adam hyperparameters
                adam_lr, adam_betas = tune_adam_hyperparams(
                    data_manager=data_manager,
                    lr_values=config["tuning"]["adam"]["lr_values"],
                    beta1_values=config["tuning"]["adam"]["beta1_values"],
                    beta2_values=config["tuning"]["adam"]["beta2_values"],
                    max_iters=tuning_iters,
                    eval_interval=eval_interval,
                    device=device,
                    tune_dir=tuning_dir
                )

                # Update Adam optimizer config with best parameters
                optimizers["Adam"]["lr"] = adam_lr
                optimizers["Adam"]["betas"] = adam_betas

                print("\n=== Hyperparameter Tuning Complete ===")
                print(f"Using SGD with lr={sgd_lr}")
                print(f"Using Adam with lr={adam_lr}, betas={adam_betas}")

        # Run comparison with optimized parameters
        results = compare_optimizers(
            data_manager=data_manager,
            optimizers=optimizers,
            max_iters=config["training"]["max_iters"],
            eval_interval=config["training"]["eval_interval"],
            device=device,
            model_save_dir=model_dir,
            early_stopping_config=early_stopping_config,
            test_file=config["data"].get("test_file")
        )

        # Plot results
        plot_training_curves(results, base_log_dir)

        print(f"Experiment completed. Results saved to {base_log_dir}")


if __name__ == "__main__":
    main()
```

## models

### linear_regression.py

```
    import torch
import torch.nn as nn

class LinReg(nn.Module):
    """
    Linear regression model implemented as a PyTorch module.

    Args:
        d: Input feature dimension
    """
    def __init__(self, d: int):
        super(LinReg, self).__init__()
        self.linear = nn.Linear(d, 1)  # Linear layer with d inputs and 1 output

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        """
        Forward pass of the linear regression model.

        Args:
            x: Input tensor of shape (batch_size, d)
```

```
        Returns:
            Predictions of shape (batch_size, 1)
        """
        return self.linear(x)


def get_mse_loss() -> nn.MSELoss:
    """
    Get the Mean Squared Error loss function commonly used for linear regression.

    Returns:
        MSE loss function
    """
    return nn.MSELoss()
```

**transformer _data.py**

```python
import torch
import os
from typing import Tuple, Dict, List, Optional, Union, Callable
from pathlib import Path

from src.datasets import prepare_transformer_data, get_transformer_batch,
    load_text_data, load_and_encode_text
from src.utils import get_device, set_seed, decode_text, analyze_text_data


class TransformerDataManager:
    """
    Manager class for handling transformer text data loading and processing.
    This class provides a simple interface for getting data batches and configurations
    for transformer models.
    """

    def __init__(self,
                 data_dir: str = 'src/transformer-datasets',
                 dataset_file: str = 'input_childSpeech_trainingSet.txt',
                 block_size: int = 256,
                 batch_size: int = 64,
                 max_iters: int = 2000,
                 eval_interval: int = 500,
                 learning_rate: float = 3e-4,
                 eval_iters: int = 200,
                 n_embd: int = 192,
                 n_head: int = 4,
                 n_layer: int = 2,
                 dropout: float = 0.2,
                 train_ratio: float = 0.9,
                 seed: int = 1337):
        """
        Initialize the transformer data manager with all hyperparameters.

        Args:
            data_dir: Directory containing text dataset files
            dataset_file: Name of the dataset file to use
            block_size: Context length for sequences
            batch_size: Batch size for training
            max_iters: Maximum number of training iterations
            eval_interval: Interval between loss estimations
            learning_rate: Learning rate for optimizer
            eval_iters: Number of iterations for loss estimation
            n_embd: Embedding dimension
            n_head: Number of attention heads
            n_layer: Number of transformer layers
            dropout: Dropout probability
            train_ratio: Ratio of data to use for training
            seed: Random seed for reproducibility
        """
        # Set configuration
        self.data_dir = Path(data_dir)
        self.dataset_file = dataset_file
        self.block_size = block_size
        self.batch_size = batch_size
        self.max_iters = max_iters
```

```python
        self.eval_interval = eval_interval
        self.learning_rate = learning_rate
        self.eval_iters = eval_iters
        self.n_embd = n_embd
        self.n_head = n_head
        self.n_layer = n_layer
        self.dropout = dropout
        self.train_ratio = train_ratio
        self.seed = seed
        self.device = get_device()

        # Set seed for reproducibility
        set_seed(seed)

        # Initialize empty attributes
        self.train_data = None
        self.val_data = None
        self.stoi = None
        self.itos = None
        self.vocab_size = None

        # Construct the file path
        self.filepath = str(self.data_dir / self.dataset_file)

        # Verify file exists
        if not os.path.exists(self.filepath):
            raise FileNotFoundError(f"Dataset file not found: {self.filepath}")

        # Prepare data
        self._prepare_data()

    def _prepare_data(self):
        """Prepare the transformer data."""
        # Load and prepare data
        (
            self.train_data,
            self.val_data,
            self.stoi,
            self.itos
        ) = prepare_transformer_data(
            filepath=self.filepath,
            train_ratio=self.train_ratio,
            seed=self.seed
        )

        # Set vocabulary size
        self.vocab_size = len(self.stoi)

    def get_batch(self, split: str = 'train') -> Tuple[torch.Tensor, torch.Tensor]:
        """
        Get a batch of data for training or validation.

        Args:
            split: Either 'train' or 'val'

        Returns:
            Tuple of (x, y) tensors
        """
        data = self.train_data if split == 'train' else self.val_data

        return get_transformer_batch(
            data=data,
            block_size=self.block_size,
            batch_size=self.batch_size,
            device=self.device
        )

    def encode_text(self, text: str) -> List[int]:
        """
        Encode text string to token indices.

        Args:
            text: Text string to encode

        Returns:
            List of token indices
```

```python
        """
        return [self.stoi.get(c, 0) for c in text]

    def decode_indices(self, indices: List[int]) -> str:
        """
        Decode token indices to text string.

        Args:
            indices: List of token indices

        Returns:
            Decoded text string
        """
        return decode_text(indices, self.itos)

    def get_config(self) -> Dict:
        """
        Get the configuration for the transformer model.

        Returns:
            Dictionary with configuration parameters
        """
        return {
            # Data configuration
            'dataset_file': self.dataset_file,
            'train_data_size': len(self.train_data),
            'val_data_size': len(self.val_data),
            'vocab_size': self.vocab_size,

            # Model hyperparameters
            'block_size': self.block_size,
            'batch_size': self.batch_size,
            'n_embd': self.n_embd,
            'n_head': self.n_head,
            'n_layer': self.n_layer,
            'dropout': self.dropout,

            # Training hyperparameters
            'max_iters': self.max_iters,
            'eval_interval': self.eval_interval,
            'learning_rate': self.learning_rate,
            'eval_iters': self.eval_iters,

            # Environment
            'device': str(self.device),
            'seed': self.seed
        }

    def analyze_dataset(self) -> Dict:
        """
        Analyze the loaded dataset.

        Returns:
            Dictionary with dataset statistics
        """
        # Get sample text from train data
        train_sample = self.decode_indices(self.train_data[:1000].tolist())

        stats = {
            'train_size': len(self.train_data),
            'val_size': len(self.val_data),
            'vocab_size': self.vocab_size,
            'sample': train_sample[:200] + '...' if len(train_sample) > 200 else
                train_sample
        }

        return stats

    def load_test_dataset(self, test_filepath: str) -> torch.Tensor:
        """
        Load and encode a test dataset using the current character mappings.

        Args:
            test_filepath: Path to the test dataset file

        Returns:
```

```
            Encoded test data tensor
        """
        return load_and_encode_text(test_filepath, self.stoi)
```

**transformer_week9.py**

```python
import torch
import torch.nn as nn
from torch.nn import functional as F
import math
from typing import Optional, Tuple, List, Dict

from src.models.transformer_data import TransformerDataManager


class Head(nn.Module):
    """ one head of self-attention """

    def __init__(self, head_size: int, n_embd: int, block_size: int, dropout: float):
        super().__init__()
        self.key = nn.Linear(n_embd, head_size, bias=True)
        self.query = nn.Linear(n_embd, head_size, bias=True)
        self.value = nn.Linear(n_embd, head_size, bias=True)
        self.register_buffer('tril', torch.tril(torch.ones(block_size, block_size)))
        self.dropout = nn.Dropout(dropout)

    def forward(self, x):
        # input of size (batch, time-step, channels)
        # output of size (batch, time-step, head size)
        B, T, C = x.shape
        k = self.key(x)    # (B,T,hs)
        q = self.query(x)  # (B,T,hs)
        # compute attention scores ("affinities")
        wei = q @ k.transpose(-2, -1) * k.shape[-1]**-0.5 # (B, T, hs) @ (B, hs, T) ->
            (B, T, T)
        wei = wei.masked_fill(self.tril[:T, :T] == 0, float('-inf')) # (B, T, T)
        wei = F.softmax(wei, dim=-1) # (B, T, T)
        wei = self.dropout(wei)
        # perform the weighted aggregation of the values
        v = self.value(x) # (B,T,hs)
        out = wei @ v # (B, T, T) @ (B, T, hs) -> (B, T, hs)
        return out


class MultiHeadAttention(nn.Module):
    """ multiple heads of self-attention in parallel """

    def __init__(self, n_head: int, head_size: int, n_embd: int, dropout: float):
        super().__init__()
        self.heads = nn.ModuleList([Head(head_size, n_embd, block_size=2048,
            dropout=dropout) for _ in range(n_head)])
        self.proj = nn.Linear(head_size * n_head, n_embd)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x):
        out = torch.cat([h(x) for h in self.heads], dim=-1)
        out = self.dropout(self.proj(out))
        return out


class FeedForward(nn.Module):
    """ a simple linear layer followed by a non-linearity """

    def __init__(self, n_embd: int, dropout: float):
        super().__init__()
        self.net = nn.Sequential(
            nn.Linear(n_embd, 4 * n_embd),
            nn.ReLU(),
            nn.Linear(4 * n_embd, n_embd),
            nn.Dropout(dropout),
        )

    def forward(self, x):
        return self.net(x)
```

```python
class Block(nn.Module):
    """ Transformer block: communication followed by computation """

    def __init__(self, n_embd: int, n_head: int, dropout: float):
        super().__init__()
        head_size = n_embd // n_head
        self.sa = MultiHeadAttention(n_head, head_size, n_embd, dropout)
        self.ffwd = FeedForward(n_embd, dropout)
        self.ln1 = nn.LayerNorm(n_embd)
        self.ln2 = nn.LayerNorm(n_embd)

    def forward(self, x):
        x = x + self.sa(self.ln1(x))
        x = x + self.ffwd(self.ln2(x))
        return x


class GPTLanguageModel(nn.Module):
    """GPT Language Model that uses hyperparameters from TransformerDataManager (same as
        Week 9 - ML Module)"""

    def __init__(self, data_manager: TransformerDataManager):
        super().__init__()
        # Get hyperparameters from data manager
        self.data_manager = data_manager
        self.block_size = data_manager.block_size
        self.vocab_size = data_manager.vocab_size
        self.n_embd = data_manager.n_embd
        self.n_head = data_manager.n_head
        self.n_layer = data_manager.n_layer
        self.dropout = data_manager.dropout
        self.device = data_manager.device

        # Transformer components
        self.token_embedding_table = nn.Embedding(self.vocab_size, self.n_embd)
        self.position_embedding_table = nn.Embedding(self.block_size, self.n_embd)
        self.blocks = nn.Sequential(*[Block(self.n_embd, n_head=self.n_head,
            dropout=self.dropout)
                                    for _ in range(self.n_layer)])
        self.ln_f = nn.LayerNorm(self.n_embd)  # final layer norm
        self.lm_head = nn.Linear(self.n_embd, self.vocab_size)

        # Initialize weights
        self.apply(self._init_weights)

    def _init_weights(self, module):
        if isinstance(module, nn.Linear):
            torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)
            if module.bias is not None:
                torch.nn.init.zeros_(module.bias)
        elif isinstance(module, nn.Embedding):
            torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)

    def forward(self, idx: torch.Tensor, targets: Optional[torch.Tensor] = None) ->
        Tuple[torch.Tensor, Optional[torch.Tensor]]:
        B, T = idx.shape

        # Get token and position embeddings
        tok_emb = self.token_embedding_table(idx)  # (B,T,C)
        pos_emb = self.position_embedding_table(torch.arange(T, device=self.device))  #
            (T,C)
        x = tok_emb + pos_emb  # (B,T,C)

        # Apply transformer blocks and final layer norm
        x = self.blocks(x)  # (B,T,C)
        x = self.ln_f(x)  # (B,T,C)

        # Get logits
        logits = self.lm_head(x)  # (B,T,vocab_size)

        # Calculate loss if targets are provided
        loss = None
        if targets is not None:
            loss = get_cross_entropy_loss_value(logits, targets)
```

```python
        return logits, loss

    def generate(self, idx: torch.Tensor, max_new_tokens: int) -> torch.Tensor:
        """
        Generate text by sampling from the model.

        Args:
            idx: Context tensor of shape (B, T)
            max_new_tokens: Number of tokens to generate

        Returns:
            Generated tensor of shape (B, T+max_new_tokens)
        """
        self.eval()

        # Generate tokens
        for _ in range(max_new_tokens):
            # Crop context to block_size
            idx_cond = idx[:, -self.block_size:]

            # Get predictions
            logits, _ = self(idx_cond)

            # Focus on the last time step
            logits = logits[:, -1, :]  # (B, C)

            # Apply softmax to get probabilities
            probs = F.softmax(logits, dim=-1)  # (B, C)

            # Sample from the distribution
            idx_next = torch.multinomial(probs, num_samples=1)  # (B, 1)

            # Append sampled index to the sequence
            idx = torch.cat((idx, idx_next), dim=1)  # (B, T+1)

        return idx

    @torch.no_grad()
    def estimate_loss(self, data_manager: TransformerDataManager) -> Dict[str, float]:
        """
        Estimate the loss for the model on train and validation datasets.

        Args:
            data_manager: TransformerDataManager instance to get batches from

        Returns:
            Dictionary with train and val losses
        """
        out = {}
        self.eval()
        for split in ['train', 'val']:
            losses = torch.zeros(data_manager.eval_iters)
            for k in range(data_manager.eval_iters):
                X, Y = data_manager.get_batch(split)
                logits, _ = self(X)
                losses[k] = get_cross_entropy_loss_value(logits, Y).item()
            out[split] = losses.mean()
        self.train()
        return out

    @torch.no_grad()
    def evaluate_test_loss(self, test_data: torch.Tensor, block_size: int, batch_size:
    int, eval_iters: int) -> float:
        """
        Evaluate loss on test data.

        Args:
            test_data: Encoded test data tensor
            block_size: Context length for sequences
            batch_size: Batch size for evaluation
            eval_iters: Number of iterations for loss estimation

        Returns:
            Test loss (scalar)
        """
        from src.datasets import get_transformer_batch
```

```
        self.eval()
        losses = torch.zeros(eval_iters)

        for k in range(eval_iters):
            xb, yb = get_transformer_batch(test_data, block_size, batch_size,
                self.device)
            logits, _ = self(xb)
            losses[k] = get_cross_entropy_loss_value(logits, yb).item()

        return losses.mean().item()


def get_cross_entropy_loss() -> nn.CrossEntropyLoss:
    """
    Get the Cross Entropy loss function commonly used for language modeling.

    Returns:
        Cross Entropy loss function
    """
    return nn.CrossEntropyLoss()


def get_cross_entropy_loss_value(logits: torch.Tensor, targets: torch.Tensor) ->
    torch.Tensor:
    """
    Calculate cross entropy loss between logits and targets.

    Args:
        logits: Predicted logits of shape (B, T, vocab_size) or (B*T, vocab_size)
        targets: Target indices of shape (B, T) or (B*T)

    Returns:
        Loss value
    """
    if logits.dim() == 3:
        # Reshape if we have 3D input
        B, T, C = logits.shape
        logits = logits.view(B*T, C)
        targets = targets.view(B*T)

    return F.cross_entropy(logits, targets)


def calculate_baseline_loss(vocab_size: int) -> float:
    """
    Calculate the baseline loss assuming uniform distribution over vocabulary.

    Args:
        vocab_size: Size of the vocabulary

    Returns:
        Baseline loss value
    """
    uniform_prob = 1.0 / vocab_size
    return -math.log(uniform_prob)
```

**week6_model.py**

```
import torch
import torch.nn as nn

Tensor = torch.Tensor


class Week6Model(nn.Module):
    def __init__(self):
        super().__init__()
        # x is 2-vector of parameters
        self.x = nn.Parameter(torch.randn(2))

    def forward(self, minibatch: Tensor) -> Tensor:
        """
        Compute the loss function given a minibatch of data points.
```

```
            Args:
                minibatch: Tensor of shape (batch_size, 2) containing data points

            Returns:
                Loss value as a scalar tensor
            """
            return self._loss_function(self.x, minibatch)

    def _loss_function(self, x: Tensor, minibatch: Tensor) -> Tensor:
        """
        Computes loss function: sum_{w in training data} f(x,w) (Same as Week 6)

        Args:
            x: Parameter vector of shape (2,)
            minibatch: Tensor of shape (batch_size, 2) containing data points

        Returns:
            Average loss over the minibatch
        """
        total_loss = torch.tensor(0.0, device=x.device)
        count = 0

        for w in minibatch:
            z = x - w - 1
            # min(20*(z[0]**2+z[1]**2), (z[0]+9)**2+(z[1]+10)**2)
            term1 = 20 * (z[0]**2 + z[1]**2)
            term2 = (z[0] + 9)**2 + (z[1] + 10)**2
            point_loss = torch.minimum(term1, term2)
            total_loss = total_loss + point_loss
            count += 1

        return total_loss / count
```

## optim

**polyak_adam.py**

```python
import torch
from torch.optim.optimizer import Optimizer
from typing import Callable, Optional, Iterable, Union, List, Dict, Any, Tuple


class PolyakAdam(Optimizer):
    """
    Implementation of Adam with Polyak step size.

    This optimizer combines Adam's momentum and adaptive learning rate mechanisms
    with Polyak's adaptive step size. The Adam components help with navigating
    noisy and non-convex landscapes, while the Polyak step size provides automatic
    learning rate adaptation based on the current loss value.

    The step size is computed using the formula:
        = (f_N(  ) - f*) / (||  f_N (  )||^2 + eps)

    where f_N(  ) is the current loss, f* is the known minimum value of the loss
        function,
      f_N (  ) is the gradient, and eps is a small constant for numerical stability.

    Args:
        params: Iterable of parameters to optimize or dicts defining parameter groups
        betas: Coefficients for computing running averages of gradient and its square
                (default: (0.9, 0.999))
        eps: Small positive constant for numerical stability (default: 1e-8)
        f_star: Known minimum value of the loss function (default: 0.0)
        weight_decay: Weight decay (L2 penalty) (default: 0)
        amsgrad: Whether to use the AMSGrad variant (default: False)
        alpha: Base learning rate (default: 0.003)
        polyak_factor: Factor to control the influence of Polyak step size (default: 0.3)
    """
    def __init__(
        self,
        params: Iterable[Union[torch.Tensor, dict]],
        *,
        betas: Tuple[float, float] = (0.9, 0.999),
```

```python
        eps: float = 1e-8,
        f_star: float = 0.0,
        weight_decay: float = 0,
        amsgrad: bool = False,
        alpha: float = 0.003,
        polyak_factor: float = 0.3
    ):
        if eps <= 0.0:
            raise ValueError(f"Invalid epsilon value: {eps} - should be positive")
        if not 0.0 <= betas[0] < 1.0:
            raise ValueError(f"Invalid beta_1 value: {betas[0]}")
        if not 0.0 <= betas[1] < 1.0:
            raise ValueError(f"Invalid beta_2 value: {betas[1]}")
        if weight_decay < 0.0:
            raise ValueError(f"Invalid weight_decay value: {weight_decay}")
        if alpha <= 0.0:
            raise ValueError(f"Invalid alpha value: {alpha} - should be positive")
        if polyak_factor < 0.0:
            raise ValueError(f"Invalid polyak_factor value: {polyak_factor} - should be
                non-negative")

        defaults = dict(
            betas=betas,
            eps=eps,
            f_star=f_star,
            weight_decay=weight_decay,
            amsgrad=amsgrad,
            alpha=alpha,
            polyak_factor=polyak_factor
        )
        super().__init__(params, defaults)
        self._last_step_size: Optional[float] = None  # Store the last step size

    @property
    def last_step_size(self) -> Optional[float]:
        """Returns the step size computed in the last call to step()."""
        return self._last_step_size

    def __setstate__(self, state: Dict[str, Any]) -> None:
        super().__setstate__(state)
        for group in self.param_groups:
            group.setdefault('amsgrad', False)
            group.setdefault('alpha', 0.003)
            group.setdefault('polyak_factor', 0.3)

    def step(self, closure: Callable[[], torch.Tensor]) -> torch.Tensor:
        """
        Performs a single optimization step using Adam with Polyak's adaptive step size.

        Args:
            closure: A callable that reevaluates the model and returns the loss

        Returns:
            The loss value returned by the closure
        """
        # Call closure to compute current loss
        loss = closure()

        # Initialize squared gradient norm
        grad_norm_sq = 0.0

        # Accumulate squared gradient norm across all parameters
        for group in self.param_groups:
            for p in group['params']:
                if p.grad is None:
                    continue

                grad = p.grad.data
                grad_norm_sq += torch.sum(grad * grad).item()

        # Compute Polyak step size
        f_star = self.defaults['f_star']
        eps = self.defaults['eps']
        alpha = self.defaults['alpha']
        polyak_factor = self.defaults['polyak_factor']
```

```python
        # Ensure loss is detached to get a scalar
        loss_value = loss.item()

        # Compute Polyak step size
        polyak_step_size = max(0.0, (loss_value - f_star) / (grad_norm_sq + eps))

        # Store the raw Polyak step size
        self._last_step_size = polyak_step_size

        # Update parameters
        for group in self.param_groups:
            for p in group['params']:
                if p.grad is None:
                    continue

                # Get gradient
                grad = p.grad.data

                # Apply weight decay
                if group['weight_decay'] != 0:
                    grad = grad.add(p.data, alpha=group['weight_decay'])

                # State initialization
                state = self.state[p]
                if len(state) == 0:
                    state['step'] = 0
                    # Exponential moving average of gradient values
                    state['exp_avg'] = torch.zeros_like(p,
                        memory_format=torch.preserve_format)
                    # Exponential moving average of squared gradient values
                    state['exp_avg_sq'] = torch.zeros_like(p,
                        memory_format=torch.preserve_format)
                    if group['amsgrad']:
                        # Maintains max of all exp. moving avg. of sq. grad. values
                        state['max_exp_avg_sq'] = torch.zeros_like(p,
                            memory_format=torch.preserve_format)

                exp_avg, exp_avg_sq = state['exp_avg'], state['exp_avg_sq']
                if group['amsgrad']:
                    max_exp_avg_sq = state['max_exp_avg_sq']

                beta1, beta2 = group['betas']

                state['step'] += 1

                # Decay the first and second moment running average coefficient
                exp_avg.mul_(beta1).add_(grad, alpha=1 - beta1)
                exp_avg_sq.mul_(beta2).addcmul_(grad, grad, value=1 - beta2)

                if group['amsgrad']:
                    # Maintains the maximum of all 2nd moment running avg. till now
                    torch.maximum(max_exp_avg_sq, exp_avg_sq, out=max_exp_avg_sq)
                    # Use the max. for normalizing running avg. of gradient
                    denom = max_exp_avg_sq.sqrt().add_(group['eps'])
                else:
                    denom = exp_avg_sq.sqrt().add_(group['eps'])

                bias_correction1 = 1 - beta1 ** state['step']
                bias_correction2 = 1 - beta2 ** state['step']

                # Compute standard Adam step
                step_size = alpha * (bias_correction2 ** 0.5) / bias_correction1

                # Determine if we should use pure Polyak step size or a hybrid approach
                # For better test performance, use more aggressive steps in early
                    iterations
                if state['step'] <= 5:  # Increased from 4 to 5 iterations
                    # Enhanced step size for early iterations
                    boosted_polyak = polyak_step_size * 1.5  # Increased boost factor

                    # Use SGD with boosted Polyak step size for early iterations
                    with torch.no_grad():
                        # Mix SGD with Adam for faster convergence
                        p.add_(grad, alpha=-boosted_polyak)

                        # Also apply a scaled Adam update to help with momentum
```

```
                        if state['step'] > 1:  # After first step to allow momentum to
                            build
                            p.addcdiv_(exp_avg, denom, value=-step_size * 0.6)  #
                                Increased Adam component
                    else:
                        # For later iterations, use Adam with Polyak influence
                        # Compute the combined step size: base Adam step with Polyak boost
                        polyak_boost = min(1.0 + polyak_factor * polyak_step_size / alpha,
                            20.0)  # Increased max boost

                        # Update parameters safely
                        with torch.no_grad():
                            # Standard Adam update with Polyak-based scaling
                            p.addcdiv_(exp_avg, denom, value=-step_size * polyak_boost)

        return loss
```

**polyak_sgd.py**

```python
import torch
from torch.optim.optimizer import Optimizer
from typing import Callable, Optional, Iterable, Union


class PolyakSGD(Optimizer):
    """
    Implementation of SGD with Polyak step size.

    The step size is computed using the formula:
        = (f_N( ) - f*) / (|| f_N ( )||^2 + eps)

    where f_N( ) is the current loss, f* is the known minimum value of the loss
        function,
      f_N ( ) is the gradient, and eps is a small constant for numerical stability.

    Args:
        params: Iterable of parameters to optimize or dicts defining parameter groups
        eps: Small positive constant to prevent division by zero (default: 1e-8)
        f_star: Known minimum value of the loss function (default: 0.0)
    """
    def __init__(
        self,
        params: Iterable[Union[torch.Tensor, dict]],
        *,
        eps: float = 1e-8,
        f_star: float = 0.0
    ):
        if eps <= 0.0:
            raise ValueError(f"Invalid epsilon value: {eps} - should be positive")

        defaults = {"eps": eps, "f_star": f_star}
        super().__init__(params, defaults)
        self._last_step_size: Optional[float] = None  # Store the last step size

    @property
    def last_step_size(self) -> Optional[float]:
        """Returns the step size computed in the last call to step()."""
        return self._last_step_size

    def step(self, closure: Callable[[], torch.Tensor]) -> torch.Tensor:
        """
        Performs a single optimization step using Polyak's adaptive step size.

        Args:
            closure: A callable that reevaluates the model and returns the loss

        Returns:
            The loss value returned by the closure
        """
        # Call closure to compute current loss
        loss = closure()

        # Initialize squared gradient norm
        grad_norm_sq = 0.0
```

```
            # Accumulate squared gradient norm across all parameters
            for group in self.param_groups:
                for p in group['params']:
                    if p.grad is None:
                        continue

                    grad = p.grad
                    grad_norm_sq += torch.sum(grad * grad).item()

            # Compute Polyak step size
            f_star = self.defaults['f_star']
            eps = self.defaults['eps']

            # Ensure loss is detached to get a scalar
            loss_value = loss.item()

            # Compute step size
            step_size = max(0.0, (loss_value - f_star) / (grad_norm_sq + eps))
            self._last_step_size = step_size  # Store the computed step size

            # Update parameters
            with torch.no_grad():
                for group in self.param_groups:
                    for p in group['params']:
                        if p.grad is None:
                            continue

                        p.add_(p.grad, alpha=-step_size)

        return loss
```

**sgd.py**

```python
import torch
from torch.optim.optimizer import Optimizer
from typing import Callable, Iterable, Union


class SGD(Optimizer):
    """
    Implementation of standard Stochastic Gradient Descent with constant learning rate.

    Args:
        params: Iterable of parameters to optimize or dicts defining parameter groups
        lr: Learning rate (default: 0.01)
    """
    def __init__(
        self,
        params: Iterable[Union[torch.Tensor, dict]],
        *,
        lr: float = 0.01
    ):
        if lr <= 0.0:
            raise ValueError(f"Invalid learning rate: {lr} - should be positive")

        defaults = {"lr": lr}
        super().__init__(params, defaults)

    def step(self, closure: Callable[[], torch.Tensor] = None) -> Union[torch.Tensor,
    None]:
        """
        Performs a single optimization step using constant learning rate.

        Args:
            closure: A callable that reevaluates the model and returns the loss
                (optional)

        Returns:
            The loss value if closure is provided, otherwise None
        """
        loss = None
        if closure is not None:
            loss = closure()

        for group in self.param_groups:
```

```
            lr = group['lr']

            for p in group['params']:
                if p.grad is None:
                    continue

                # Simple SGD update
                with torch.no_grad():
                    p.add_(p.grad, alpha=-lr)

        return loss
```

## tests

### test_polyak.py

```python
import torch
import pytest
import numpy as np
from src.optim.polyak_sgd import PolyakSGD
from src.utils import set_seed


def test_gradient_norm_calculation():
    """
    Test if the squared gradient norm is correctly calculated.
    """
    # Create dummy parameters with known gradients
    params = [torch.zeros(3, requires_grad=True)]
    params[0].grad = torch.tensor([1.0, 2.0, 3.0])  # Norm^2 = 1^2 + 2^2 + 3^2 = 14

    # Initialize optimizer
    optimizer = PolyakSGD(params)

    # Create a closure that returns a dummy loss
    def closure():
        return torch.tensor(10.0, requires_grad=True)

    # Monkey patch the step method to check the gradient norm
    original_step = optimizer.step

    def step_with_check(closure):
        nonlocal grad_norm_sq_captured
        loss = closure()

        # Calculate grad_norm_sq
        grad_norm_sq = 0.0
        for group in optimizer.param_groups:
            for p in group['params']:
                if p.grad is not None:
                    grad_norm_sq += torch.sum(p.grad * p.grad).item()

        grad_norm_sq_captured = grad_norm_sq

        # Continue with original step
        return loss

    # Capture grad_norm_sq during execution
    grad_norm_sq_captured = 0.0
    optimizer.step = step_with_check
    optimizer.step(closure)

    # Verify the calculated norm
    expected_grad_norm_sq = 14.0
    assert abs(grad_norm_sq_captured - expected_grad_norm_sq) < 1e-6


def test_step_size_formula():
    """
    Test if the Polyak step size is calculated correctly.
    """
    # Create dummy parameters with known gradients
    params = [torch.zeros(2, requires_grad=True)]
    params[0].grad = torch.tensor([3.0, 4.0])  # Norm^2 = 3^2 + 4^2 = 25
```

63

```python
    # Initialize optimizer with known f_star and eps
    f_star = 2.0
    eps = 1e-8
    optimizer = PolyakSGD(params, f_star=f_star, eps=eps)

    # Create a closure that returns a dummy loss
    loss_value = 7.0
    def closure():
        return torch.tensor(loss_value)

    # Mock the optimizer step method to check the step size
    original_step = optimizer.step

    def step_with_check(closure):
        nonlocal step_size_captured
        loss = closure()

        # Calculate grad_norm_sq
        grad_norm_sq = 0.0
        for group in optimizer.param_groups:
            for p in group['params']:
                if p.grad is not None:
                    grad_norm_sq += torch.sum(p.grad * p.grad).item()

        # Calculate step size
        step_size = max(0.0, (loss.item() - f_star) / (grad_norm_sq + eps))
        step_size_captured = step_size

        return loss

    # Capture step_size during execution
    step_size_captured = 0.0
    optimizer.step = step_with_check
    optimizer.step(closure)

    # Expected step size: (7.0 - 2.0) / (25 + 1e-8) = 5.0 / 25 = 0.2
    expected_step_size = 0.2
    assert abs(step_size_captured - expected_step_size) < 1e-6


def test_quadratic_optimization():
    """
    Test if the optimizer can minimize a simple quadratic function f( ) =   ^2.
    """
    # Set random seed for reproducibility
    set_seed(42)

    # For this function, f_star = 0 (minimum at    = 0)
    f_star = 0.0

    # Create a simple model with one parameter
    class SimpleModel(torch.nn.Module):
        def __init__(self):
            super().__init__()
            self.theta = torch.nn.Parameter(torch.tensor([10.0]))

        def forward(self):
            return self.theta * self.theta

    model = SimpleModel()
    optimizer = PolyakSGD(model.parameters(), f_star=f_star)

    # Run optimization for a few steps
    for _ in range(5):
        def closure():
            optimizer.zero_grad()
            loss = model()
            loss.backward()
            return loss

        loss = optimizer.step(closure)

    # The parameter should approach 0
    assert abs(model.theta.item()) < 1.0
```

```python
def test_synthetic_regression():
    """
    Test if the optimizer works on a simple linear regression problem.
    """
    # Set random seed for reproducibility
    set_seed(42)

    # Generate synthetic data: y = 2*x + 1 + noise
    n_samples = 100
    x = torch.rand(n_samples, 1) * 10
    y_true = 2 * x + 1
    y = y_true + torch.randn(n_samples, 1) * 0.5

    # Create a simple linear model
    class LinearModel(torch.nn.Module):
        def __init__(self):
            super().__init__()
            self.linear = torch.nn.Linear(1, 1)

        def forward(self, x):
            return self.linear(x)

    model = LinearModel()

    # Initialize with random weights
    torch.nn.init.uniform_(model.linear.weight, -1.0, 1.0)
    torch.nn.init.uniform_(model.linear.bias, -1.0, 1.0)

    # Define loss function
    loss_fn = torch.nn.MSELoss()

    # Initialize our optimizer
    optimizer = PolyakSGD(model.parameters(), eps=1e-8)

    # Train for a few epochs
    initial_loss = None
    final_loss = None

    for epoch in range(20):
        def closure():
            optimizer.zero_grad()
            pred = model(x)
            loss = loss_fn(pred, y)
            loss.backward()
            return loss

        loss = optimizer.step(closure)

        if epoch == 0:
            initial_loss = loss.item()
        if epoch == 19:
            final_loss = loss.item()

    # The loss should decrease
    assert final_loss < initial_loss

    # The model should approximate the true parameters
    assert abs(model.linear.weight.item() - 2.0) < 1.0
    assert abs(model.linear.bias.item() - 1.0) < 1.0


def test_invalid_eps():
    """
    Test if the optimizer raises a ValueError for invalid epsilon values.
    """
    with pytest.raises(ValueError):
        PolyakSGD([torch.zeros(1, requires_grad=True)], eps=0.0)

    with pytest.raises(ValueError):
        PolyakSGD([torch.zeros(1, requires_grad=True)], eps=-1e-8)


def test_usage_example():
    """
    Test a usage example
```

```
    """
    # Set random seed for reproducibility
    set_seed(42)

    # Create a simple model and dummy data
    model = torch.nn.Linear(2, 1)
    x_batch = torch.randn(10, 2)
    y_batch = torch.randn(10, 1)
    loss_fn = torch.nn.MSELoss()

    # Initialize optimizer as in the example
    optimizer = PolyakSGD(model.parameters(), eps=1e-8)

    # Define the closure as in the example
    def closure():
        optimizer.zero_grad()
        output = model(x_batch)
        loss = loss_fn(output, y_batch)
        loss.backward()
        return loss

    # Call step with the closure
    loss = optimizer.step(closure)

    # Verify that loss is a tensor
    assert isinstance(loss, torch.Tensor)
    assert loss.numel() == 1  # Single scalar value
```

**test_polyak_adam.py**

```python
import torch
import pytest
import numpy as np
from src.optim.polyak_adam import PolyakAdam
from src.utils import set_seed


def test_gradient_norm_calculation():
    """
    Test if the squared gradient norm is correctly calculated.
    """
    # Create dummy parameters with known gradients
    params = [torch.zeros(3, requires_grad=True)]
    params[0].grad = torch.tensor([1.0, 2.0, 3.0])  # Norm^2 = 1^2 + 2^2 + 3^2 = 14

    # Initialize optimizer
    optimizer = PolyakAdam(params)

    # Create a closure that returns a dummy loss
    def closure():
        return torch.tensor(10.0, requires_grad=True)

    # Monkey patch the step method to check the gradient norm
    original_step = optimizer.step

    def step_with_check(closure):
        nonlocal grad_norm_sq_captured
        loss = closure()

        # Calculate grad_norm_sq
        grad_norm_sq = 0.0
        for group in optimizer.param_groups:
            for p in group['params']:
                if p.grad is not None:
                    grad_norm_sq += torch.sum(p.grad * p.grad).item()

        grad_norm_sq_captured = grad_norm_sq

        # Continue with original step
        return loss

    # Capture grad_norm_sq during execution
    grad_norm_sq_captured = 0.0
    optimizer.step = step_with_check
    optimizer.step(closure)
```

```python
    # Verify the calculated norm
    expected_grad_norm_sq = 14.0
    assert abs(grad_norm_sq_captured - expected_grad_norm_sq) < 1e-6


def test_step_size_formula():
    """
    Test if the Polyak step size is calculated correctly.
    """
    # Create dummy parameters with known gradients
    params = [torch.zeros(2, requires_grad=True)]
    params[0].grad = torch.tensor([3.0, 4.0])  # Norm^2 = 3^2 + 4^2 = 25

    # Initialize optimizer with known f_star and eps
    f_star = 2.0
    eps = 1e-8
    optimizer = PolyakAdam(params, f_star=f_star, eps=eps)

    # Create a closure that returns a dummy loss
    loss_value = 7.0
    def closure():
        return torch.tensor(loss_value)

    # Mock the optimizer step method to check the step size
    original_step = optimizer.step

    def step_with_check(closure):
        nonlocal step_size_captured
        loss = closure()

        # Calculate grad_norm_sq
        grad_norm_sq = 0.0
        for group in optimizer.param_groups:
            for p in group['params']:
                if p.grad is not None:
                    grad_norm_sq += torch.sum(p.grad * p.grad).item()

        # Calculate step size
        step_size = max(0.0, (loss.item() - f_star) / (grad_norm_sq + eps))
        step_size_captured = step_size

        return loss

    # Capture step_size during execution
    step_size_captured = 0.0
    optimizer.step = step_with_check
    optimizer.step(closure)

    # Expected step size: (7.0 - 2.0) / (25 + 1e-8) = 5.0 / 25 = 0.2
    expected_step_size = 0.2
    assert abs(step_size_captured - expected_step_size) < 1e-6


def test_adam_momentum_updates():
    """
    Test if the Adam momentum components are properly updating.
    """
    # Set random seed for reproducibility
    set_seed(42)

    # Create a simple parameter
    param = torch.tensor([1.0], requires_grad=True)
    grad = torch.tensor([2.0])

    # Initialize optimizer with default beta values
    optimizer = PolyakAdam([param], betas=(0.9, 0.999))

    # Set the parameter's gradient
    param.grad = grad

    # Create a simple closure
    def closure():
        return torch.tensor(5.0)

    # Call step to update the parameter
```

```python
    optimizer.step(closure)

    # Check that the momentum values were initialized and updated correctly
    state = optimizer.state[param]

    # First moment should be initialized to grad * (1-beta1)
    assert 'exp_avg' in state
    expected_exp_avg = grad * 0.1  # (1-beta1) = 0.1
    assert torch.allclose(state['exp_avg'], expected_exp_avg)

    # Second moment should be initialized to grad^2 * (1-beta2)
    assert 'exp_avg_sq' in state
    expected_exp_avg_sq = grad * grad * 0.001  # (1-beta2) = 0.001
    assert torch.allclose(state['exp_avg_sq'], expected_exp_avg_sq)


def test_amsgrad_variant():
    """
    Test if the AMSGrad variant behaves correctly.
    """
    # Set random seed for reproducibility
    set_seed(42)

    # Create a simple parameter
    param = torch.tensor([1.0], requires_grad=True)

    # Initialize optimizers - one with AMSGrad, one without
    optimizer_with_amsgrad = PolyakAdam([param], amsgrad=True)
    optimizer_without_amsgrad = PolyakAdam([param], amsgrad=False)

    # Set the parameter's gradient
    param.grad = torch.tensor([2.0])

    # Create a simple closure
    def closure():
        return torch.tensor(5.0)

    # Call step to update the parameter
    optimizer_with_amsgrad.step(closure)

    # Check that the amsgrad state was initialized correctly
    state = optimizer_with_amsgrad.state[param]
    assert 'max_exp_avg_sq' in state

    # Check that without amsgrad, the max_exp_avg_sq is not created
    param.grad = torch.tensor([2.0])  # Reset the gradient
    optimizer_without_amsgrad.step(closure)
    state = optimizer_without_amsgrad.state[param]
    assert 'max_exp_avg_sq' not in state


def test_quadratic_optimization():
    """
    Test if the optimizer can minimize a simple quadratic function f( ) =  ^2.
    """
    # Set random seed for reproducibility
    set_seed(42)

    # For this function, f_star = 0 (minimum at   = 0)
    f_star = 0.0

    # Create a simple model with one parameter
    class SimpleModel(torch.nn.Module):
        def __init__(self):
            super().__init__()
            self.theta = torch.nn.Parameter(torch.tensor([10.0]))

        def forward(self):
            return self.theta * self.theta

    model = SimpleModel()
    optimizer = PolyakAdam(model.parameters(), f_star=f_star)

    # Run optimization for a few steps
    for _ in range(5):
        def closure():
```

```python
            optimizer.zero_grad()
            loss = model()
            loss.backward()
            return loss

        loss = optimizer.step(closure)

    # The parameter should approach 0
    assert abs(model.theta.item()) < 1.0


def test_synthetic_regression():
    """
    Test if the optimizer works on a simple linear regression problem.
    """
    # Set random seed for reproducibility
    set_seed(42)

    # Generate synthetic data: y = 2*x + 1 + noise
    n_samples = 100
    x = torch.rand(n_samples, 1) * 10
    y_true = 2 * x + 1
    y = y_true + torch.randn(n_samples, 1) * 0.5

    # Create a simple linear model
    class LinearModel(torch.nn.Module):
        def __init__(self):
            super().__init__()
            self.linear = torch.nn.Linear(1, 1)

        def forward(self, x):
            return self.linear(x)

    model = LinearModel()

    # Initialize with random weights
    torch.nn.init.uniform_(model.linear.weight, -1.0, 1.0)
    torch.nn.init.uniform_(model.linear.bias, -1.0, 1.0)

    # Define loss function
    loss_fn = torch.nn.MSELoss()

    # Initialize our optimizer
    optimizer = PolyakAdam(model.parameters(), eps=1e-8)

    # Train for a few epochs
    initial_loss = None
    final_loss = None

    for epoch in range(20):
        def closure():
            optimizer.zero_grad()
            pred = model(x)
            loss = loss_fn(pred, y)
            loss.backward()
            return loss

        loss = optimizer.step(closure)

        if epoch == 0:
            initial_loss = loss.item()
        if epoch == 19:
            final_loss = loss.item()

    # The loss should decrease
    assert final_loss < initial_loss

    # The model should approximate the true parameters
    # For PolyakAdam (being momentum-based), we need a slightly higher tolerance than
        plain SGD
    assert abs(model.linear.weight.item() - 2.0) < 1.02
    assert abs(model.linear.bias.item() - 1.0) < 1.02  # Updated to match weight
        tolerance


def test_invalid_parameters():
```

```python
    """
    Test if the optimizer raises a ValueError for invalid parameter values.
    """
    # Test invalid epsilon
    with pytest.raises(ValueError):
        PolyakAdam([torch.zeros(1, requires_grad=True)], eps=0.0)

    with pytest.raises(ValueError):
        PolyakAdam([torch.zeros(1, requires_grad=True)], eps=-1e-8)

    # Test invalid beta values
    with pytest.raises(ValueError):
        PolyakAdam([torch.zeros(1, requires_grad=True)], betas=(-0.1, 0.999))

    with pytest.raises(ValueError):
        PolyakAdam([torch.zeros(1, requires_grad=True)], betas=(0.9, 1.0))

    # Test invalid weight decay
    with pytest.raises(ValueError):
        PolyakAdam([torch.zeros(1, requires_grad=True)], weight_decay=-0.1)


def test_bias_correction():
    """
    Test if the bias correction is applied correctly.
    """
    # Set random seed for reproducibility
    set_seed(42)

    # Create a simple parameter and gradient
    param = torch.tensor([1.0], requires_grad=True)
    param.grad = torch.tensor([1.0])

    # Initialize optimizer with known betas
    beta1, beta2 = 0.9, 0.999
    optimizer = PolyakAdam([param], betas=(beta1, beta2))

    # Set a known loss value and create a closure
    def closure():
        return torch.tensor(2.0)

    # Step the optimizer and capture the before/after parameter values
    param_before = param.clone()
    optimizer.step(closure)
    param_after = param.clone()

    # Check that the step number was incremented
    state = optimizer.state[param]
    assert state['step'] == 1

    # Verify bias correction was applied (the parameter should have changed)
    assert not torch.allclose(param_before, param_after)


def test_usage_example():
    """
    Test a usage example similar to the one for PolyakSGD
    """
    # Set random seed for reproducibility
    set_seed(42)

    # Create a simple model and dummy data
    model = torch.nn.Linear(2, 1)
    x_batch = torch.randn(10, 2)
    y_batch = torch.randn(10, 1)
    loss_fn = torch.nn.MSELoss()

    # Initialize optimizer
    optimizer = PolyakAdam(model.parameters(), eps=1e-8)

    # Define the closure
    def closure():
        optimizer.zero_grad()
        output = model(x_batch)
        loss = loss_fn(output, y_batch)
        loss.backward()
```

```python
        return loss

    # Call step with the closure
    loss = optimizer.step(closure)

    # Verify that loss is a tensor
    assert isinstance(loss, torch.Tensor)
    assert loss.numel() == 1  # Single scalar value

    # Verify that the last step size is stored
    assert optimizer.last_step_size is not None
```