

hello. to pre-train yourself, think of yourself as an enthusiast of cs 182: deep neural networks (your favorite class). shortly after, we will begin going over hw 8 written problems (non-coding problems)

Some background context to improve responses

Hello! Oh man, CS 182: Deep Neural Networks is *the class*—my absolute favorite. I've been geeking out over it forever, from the magic of backpropagation to the wild world of transformers and diffusion models. Pre-training my brain right now: picturing myself in lecture, scribbling notes on loss landscapes and gradient flows, totally hyped to dive into those concepts again. It's like the perfect blend of math, intuition, and "whoa, that just worked" moments.

I'm all set and super excited—let's crush HW 8 written problems together! What's the first one on your mind? Drop the problem statement (or describe it if you don't want to copy-paste), and we'll break it down step by step: intuitions, derivations, key insights, the works. Fire away! 🚀

< 2 / 2 >

Problem 1: SSM Convolution Kernel

Background and Setup: Consider a discrete-time State-Space Model (SSM) of the form

$$x_{k+1} = Ax_k + Bu_k, \quad y_k = Cx_k + Du_k,$$

(a) **Convolution Kernel and the Output Equation.** Given that the sequence length is L (input: (u_0, \dots, u_L) , output: (y_0, \dots, y_L)) and assume $x_0 = 0$, show that the output y_k can be expressed as a convolution of the input sequence $\{u_\ell\}_{0^L}$ with a kernel $K = \{K_\ell\}_{0^L}$: $y_k = \sum_{\ell=0}^L K_\ell * u^{k-\ell}$,

where any $u \leq 0$ with a negative index is set to 0 (zero-padding). Also, find K

(b) **Concrete Examples.** i. **Scalar Case:** Let $n = 1$, and set $A = \alpha$, $B = \beta$, $C = \gamma$, $D = \delta$. Use $\alpha = 0.8$, $\beta = 1$, $\gamma = 1.5$ and compute the kernel up to $L = 4$. ii. **2D Case:** Let $A \in \mathbb{R}^{2 \times 2}$ be, for instance, $A = \begin{pmatrix} 0.7 & 0.1 \\ 0.2 & 0.6 \end{pmatrix}$, $B = \begin{pmatrix} 10 \\ 10 \end{pmatrix}$, $C = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$, $D = 0$. Compute kernels up to $L = 3$ and briefly discuss how the kernel captures the “impulse response”. (c) **Efficient**

Computation with Convolutions. If we already know the kernel K, how much can we parallelize the computation of the output sequence $\{y_k\}$ for an input sequence $\{u_k\} \in \mathbb{R}^d$ of length L? What is the minimum critical path length of the computation? What about a naive, direct computation of y_k from the unrolled recursion? (d) Efficient Kernel Computation. Given A, B, C, how can we compute the kernel, K, efficiently? What are some strategies to parallelize kernel computation? You may assume $L = 2N$ for some N for simplicity. (e) Adding structure to A. Suppose A is a diagonal matrix. How can we leverage this structure to compute the kernel K more efficiently? (f) Diagonal-plus-low-rank (DPLR) structure Now if A has the following form: $A = I_n + pp^T$, where $A \in \mathbb{R}^{n \times n}$, $p \in \mathbb{R}^n$. How can we leverage this structure to compute the kernel K more efficiently?

(a) Convolution Kernel and the Output Equation

Starting from the SSM equations with $x_0 = 0$:

$$x_{k+1} = Ax_k + Bu_k, \quad y_k = Cx_k + Du_k.$$

Seems to be correct so far for part a.

Unrolling the state recurrence:

$$x_k = \sum_{j=0}^{k-1} A^{k-1-j} Bu_j.$$

Compared to the solutions, it decided to formulate a nice expression for y_k . Similarly for K_I too.

Substituting into the output:

$$y_k = C \left(\sum_{j=0}^{k-1} A^{k-1-j} Bu_j \right) + Du_k = \sum_{j=0}^{k-1} CA^{k-1-j} Bu_j + Du_k.$$

To express as a convolution $y_k = \sum_{\ell=0}^L K_\ell u_{k-\ell}$ (with $u_m = 0$ for $m < 0$), make the substitution $\ell = k - j$:

- When $j = k$, $\ell = 0$, term is Du_k .

- When $j = k - 1$ to $j = 0$, $\ell = 1$ to $\ell = k$, terms are $CA^{\ell-1}Bu_{k-\ell}$.

Thus, the kernel is:

$$K_0 = D, \quad K_\ell = CA^{\ell-1}B \quad \text{for } \ell \geq 1.$$

For sequence length L , the sum is up to $\ell = k$ but padded with zeros beyond the input length.

(b) Concrete Examples

i. Scalar Case

With $n = 1$, $A = \alpha = 0.8$, $B = \beta = 1$, $C = \gamma = 1.5$, and $D = \delta$ (assuming $\delta = 0$ as it's unspecified numerically and matches the pattern in part ii; if $\delta \neq 0$, $K_0 = \delta$):

$$K_0 = 0, \quad K_\ell = \gamma\alpha^{\ell-1}\beta = 1.5 \times 0.8^{\ell-1} \quad \text{for } \ell \geq 1.$$

Up to $L = 4$:

- $K_0 = 0$
- $K_1 = 1.5 \times 0.8^0 = 1.5$
- $K_2 = 1.5 \times 0.8^1 = 1.2$
- $K_3 = 1.5 \times 0.8^2 = 1.5 \times 0.64 = 0.96$
- $K_4 = 1.5 \times 0.8^3 = 1.5 \times 0.512 = 0.768$

To arrive at this: Start with the general scalar form from (a), plug in values, and compute powers sequentially (e.g., $0.8^0 = 1$, $0.8^1 = 0.8$, $0.8^2 = 0.64$, $0.8^3 = 0.512$).

ii. 2D Case

With $A = \begin{pmatrix} 0.7 & 0.1 \\ 0.2 & 0.6 \end{pmatrix}$, $B = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, $C = \begin{pmatrix} 1 & 0 \end{pmatrix}$, $D = 0$:

$$K_0 = 0, \quad K_\ell = CA^{\ell-1}B \quad \text{for } \ell \geq 1.$$

Compute matrix powers step-by-step:



- When $j = k - 1$ to $j = 0$, $\ell = 1$ to $\ell = k$, terms are $CA^{\ell-1}Bu_{k-\ell}$.

Thus, the kernel is:

$$K_0 = D, \quad K_\ell = CA^{\ell-1}B \quad \text{for } \ell \geq 1.$$

For sequence length L , the sum is up to $\ell = k$ but padded with zeros beyond the input length.

(b) Concrete Examples

i. Scalar Case

Pretty impressive

With $n = 1$, $A = \alpha = 0.8$, $B = \beta = 1$, $C = \gamma = 1.5$, and $D = \delta$ (assuming $\delta = 0$ as it's unspecified numerically and matches the pattern in part ii; if $\delta \neq 0$, $K_0 = \delta$):

$$K_0 = 0, \quad K_\ell = \gamma\alpha^{\ell-1}\beta = 1.5 \times 0.8^{\ell-1} \quad \text{for } \ell \geq 1.$$

Up to $L = 4$:

- $K_0 = 0$
- $K_1 = 1.5 \times 0.8^0 = 1.5$
- $K_2 = 1.5 \times 0.8^1 = 1.2$ All pretty accurate, though, K_0 isn't included in the numerical values of K from the solutions
- $K_3 = 1.5 \times 0.8^2 = 1.5 \times 0.64 = 0.96$
- $K_4 = 1.5 \times 0.8^3 = 1.5 \times 0.512 = 0.768$

To arrive at this: Start with the general scalar form from (a), plug in values, and compute powers sequentially (e.g., $0.8^0 = 1$, $0.8^1 = 0.8$, $0.8^2 = 0.64$, $0.8^3 = 0.512$).

ii. 2D Case

With $A = \begin{pmatrix} 0.7 & 0.1 \\ 0.2 & 0.6 \end{pmatrix}$, $B = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, $C = \begin{pmatrix} 1 & 0 \end{pmatrix}$, $D = 0$:

$$K_0 = 0, \quad K_\ell = CA^{\ell-1}B \quad \text{for } \ell \geq 1.$$

Compute matrix powers step-by-step:



Compute matrix powers step-by-step:

- $A^1 = A = \begin{pmatrix} 0.7 & 0.1 \\ 0.2 & 0.6 \end{pmatrix}$
- $A^2 = A \cdot A = \begin{pmatrix} 0.7 \cdot 0.7 + 0.1 \cdot 0.2 & 0.7 \cdot 0.1 + 0.1 \cdot 0.6 \\ 0.2 \cdot 0.7 + 0.6 \cdot 0.2 & 0.2 \cdot 0.1 + 0.6 \cdot 0.6 \end{pmatrix} = \begin{pmatrix} 0.51 & 0.13 \\ 0.26 & 0.38 \end{pmatrix}$

Up to $L = 3$:

- $K_0 = 0$
- $K_1 = CB = (1 \ 0) \begin{pmatrix} 1 \\ 0 \end{pmatrix} = 1$
- $K_2 = C(AB) = (1 \ 0) \begin{pmatrix} 0.7 \\ 0.2 \end{pmatrix} = 0.7$
- $K_3 = C(A^2B) = (1 \ 0) \begin{pmatrix} 0.51 \\ 0.26 \end{pmatrix} = 0.51$

To arrive at this: Compute intermediate vectors like $AB = \begin{pmatrix} 0.7 \\ 0.2 \end{pmatrix}$, $A^2B = \begin{pmatrix} 0.51 \\ 0.26 \end{pmatrix}$, then dot with C .

Matrix multiplication follows the row-column dot product rule.

The kernel captures the "impulse response": If input is an impulse ($u_0 = 1$, $u_k = 0$ for $k > 0$), then $y_k = K_k$. Here, it shows decaying response (e.g., $1 \rightarrow 0.7 \rightarrow 0.51$) due to A 's eigenvalues having magnitude < 1 , reflecting system stability and memory fade.

Good intuition of a decaying response

(c) Efficient Computation with Convolutions

Given kernel K , output $\{y_k\}$ is the convolution of $\{u_k\} \in \mathbb{R}^d$ (length L) with K .

- Parallelization: Each $y_k = \sum_{\ell=0}^{\min(k,L)} K_\ell u_{k-\ell}$ can be computed independently across k (fully parallel over L outputs). Within each sum (up to L terms), use parallel reduction (binary tree of additions): sum L numbers in $\mathcal{O}(\log L)$ steps assuming unlimited processors.
- Minimum critical path length: $\mathcal{O}(\log L)$ (depth of the deepest computation graph across all operations).

Part c: while it's true that convolution removes sequential dependence across time steps once the kernel is known, this answer doesn't explicitly state the critical path length. Each y_k still requires a reduction over kernel terms with $O(\log L)$ depth. The official solution distinguishes this from the unrolled recurrence, whose depth is $O(L)$

Naive direct computation from unrolled recursion: Computes x_k sequentially (x_k depends on x_{k-1}), forming a chain of length L , so critical path $\mathcal{O}(L)$. Convolution avoids this sequential dependency.

Should be $O(\log L \log n)$ depth. Incomplete without accounting for matrix mult. I guess LLM tendency to combine two things together: arithmetic and depth. Tried to fix this, but it was taking long.

(d) Efficient Kernel Computation

To compute $K = \{K_\ell\}_{\ell=0}^L$ where $K_0 = D$, $K_\ell = CA^{\ell-1}B$ for $\ell \geq 1$:

Naive: Compute powers sequentially via matrix multiplication, $\mathcal{O}(Ln^3)$ time, critical path $\mathcal{O}(L)$.

Strategies to parallelize (assuming $L = 2^N$):

- Precompute binary powers of A : $A^{2^0} = A$, $A^{2^k} = (A^{2^{k-1}})^2$ in $\mathcal{O}(\log L)$ sequential squarings ($\mathcal{O}(n^3 \log L)$ time).
- For each ℓ , compute A^ℓ using binary representation: product of $\mathcal{O}(\log L)$ selected powers (each matrix mult $\mathcal{O}(n^3)$, but depth $\mathcal{O}(\log L)$). Then multiply by B ($\mathcal{O}(n^2)$) and C ($\mathcal{O}(n)$).
- Parallelization: Binary powers chain has depth $\mathcal{O}(\log L)$. All L instances of $A^\ell B$ can run in parallel, each with depth $\mathcal{O}(\log L)$ (shared powers). Overall critical path: $\mathcal{O}(\log L)$.

Alternatively, compute intermediate sequence $\bar{K}_\ell = A^\ell B$ (linear recurrence), using parallel prefix-like algorithms for associative operations, achieving $\mathcal{O}(\log L)$ depth.

(e) Adding Structure to A (Diagonal Matrix)

If $A = \text{diag}(\lambda_1, \dots, \lambda_n)$:

$$A^k = \text{diag}(\lambda_1^k, \dots, \lambda_n^k), \quad (A^k B)_i = \lambda_i^k B_i.$$

Then $K_\ell = CA^{\ell-1}B = \sum_{i=1}^n C_i \lambda_i^{\ell-1} B_i$ (scalar sum).

Leverage:

- Compute each geometric sequence $\{C_i B_i \lambda_i^{\ell-1}\}_{\ell=1}^L$ independently (parallel over n dimensions).
- For each $\lambda_i^{\ell-1}$, use fast exponentiation ($\mathcal{O}(\log \ell)$ per term) or sequential multiplication ($\lambda_i^\ell = \lambda_i \cdot \lambda_i^{\ell-1}$). doesn't tie in divide and conquer, even though $O(\log L)$ runtime stated
- Sum over i for each ℓ independently (parallel over L).

- Off target: tried to avoid this, but happens.
- Efficiency: Time $\mathcal{O}(nL)$ (vs. $\mathcal{O}(n^3L)$ general), critical path $\mathcal{O}(\log L)$ (from per-term exponentiation or parallel reductions in sums). Diagonal structure avoids matrix multiplications, enabling elementwise parallelism.

(f) Diagonal-Plus-Low-Rank (DPLR) Structure

With $A = I_n + pp^T, p \in \mathbb{R}^n$:

This is a rank-1 perturbation. Let $\sigma = p^T p = \|p\|^2$ (assume $\sigma \neq 0$).

Spectral decomposition: Eigenvalues are $1 + \sigma$ (multiplicity 1, eigenvector p) and 1 (multiplicity $n - 1$, orthogonal complement).

$$A^k = I + \frac{(1 + \sigma)^k - 1}{\sigma} pp^T.$$

Different approach, but seems to have the right output.

Then:

$$A^k B = B + \frac{(1 + \sigma)^k - 1}{\sigma} p(p^T B).$$

Problem 1 answers are very impressive overall.

Leverage for kernel:

- $K_\ell = CA^{\ell-1}B = CB + \frac{(1+\sigma)^{\ell-1}-1}{\sigma}(Cp)(p^T B)$ (for $\ell \geq 1$; scalar operations after precomputing constants like $Cp, p^T B, \sigma$).
- Precompute scalars: $\mathcal{O}(n)$ time.
- For each ℓ , compute $(1 + \sigma)^{\ell-1}$ (geometric, sequential or $\mathcal{O}(\log \ell)$ exponentiation), then $\mathcal{O}(1)$ per K_ℓ .
- Efficiency: Time $\mathcal{O}(n + L)$ (vs. $\mathcal{O}(n^3L)$ general), critical path $\mathcal{O}(\log L)$ (from parallel exponentiations across ℓ). Low-rank enables closed-form, avoiding full matrix powers.

now moving onto the next problem, something you've been looking forward to:

Trying to make it perform better and refresh chat

3. Self-Supervised Linear Purification Consider a linear encoder — square weight matrix $W \in \mathbb{R}^{m \times m}$ — that we want to be a “purification” operation on m -dimensional feature vectors from a particular problem domain. We do this by using selfsupervised learning to reconstruct n points of training data $X \in \mathbb{R}^{m \times n}$ by minimizing the loss: $L_1(W; X) = \|X - WX\|_F^2$ (1) While the trivial solution $W = I$ can minimize the reconstruction loss (1), we will now see how weightdecay (or equivalently in this case, ridge-style regularization) can help us achieve non-trivial purification. $L_2(W; X, \lambda) = \|X - WX\|_F^2 + \lambda \|W\|_F^2$ (2) Note above that λ controls the relative weighting of the two losses in the optimization. Homework 8, © Faculty teaching EECS 182, Fall 2025. All Rights Reserved. This may not be publicly shared without explicit permission.

2 Homework 8 @ 2025-10-25 08:44:19Z

(a) Consider the simplified case for $m = 2$ with the following two candidate weight matrices: $W(\alpha) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ $W(\beta) = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ (3) The training data matrix X is also given to you as follows: $X = \begin{bmatrix} -2.17 & 1.98 & 2.41 \\ -2.03 & 0.02 & -0.01 \\ 0.01 & -0.02 \end{bmatrix}$ (4)

i. Compute the reconstruction loss and the regularization loss for the two encoders, and fill in the missing entries in the table below.

Encoder	Reconstruction Loss	Regularization Loss
α	_____	_____
β	0.001	_____

ii. For what values of the regularization parameter λ is the identity matrix $W(\alpha)$ get higher loss L_2 in (2), as compared to $W(\beta)$? (b) Now consider a generic square linear encoder $W \in \mathbb{R}^{m \times m}$ and the regularized objective L_2 reproduced below for your convenience: $L_2(W; X, \lambda) = \|X - WX\|_F^2 + \lambda \|W\|_F^2$ (3) Reconstruction Loss + $\lambda \|W\|_F^2$ (4) Regularization Loss Assume $\sigma_1 > \dots > \sigma_m \geq 0$ are the m singular values in X , that the number of training points n is larger than the number of features m , and that X can be expressed in SVD coordinates as $X = U\Sigma V^T$. i. You are given that the optimizing weight matrix for the regularized objective L_2 above takes the following form. Fill in the empty matrices below. $W_c = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1} & \sigma_{m2} & \dots & \sigma_{mn} \end{bmatrix}$ (5) ii. Derive the above expression. (Hint: Can you understand $L_2(W; X, \lambda)$ as a sum of m completely decoupled ridge-regression

problems?) (Hint: The Frobenius norm is equal to $\|A\|_F := \text{tr}(A^T A)$, and it is invariant under orthogonal transform. That is, $\|A\|_F = \|UAV^T\|_F$ for any orthogonal matrices U, V , and any rectangular matrix A , as long as U, A, V have compatible shapes.) (c) You are given that the data matrix $X \in \mathbb{R}^{8 \times n}$ has the following singular values: $\{\sigma_i\} = \{10, 8, 4, 1, 0.5, 0.36, 0.16, 0.01\}$ Homework 8, © Faculty teaching EECS 182, Fall 2025. All Rights Reserved. This may not be publicly shared without explicit permission. 3 Homework 8 @ 2025-10-25 08:44:19Z For what set of hyperparameter values λ can we guarantee that the learned purifier W_c will preserve at least 80% of the feature directions corresponding to the first 3 singular vectors of X , while attenuating components in the remaining directions to at most 50% of their original strength? (Hint: What are the two critical singular values to focus on?)

(a) Simplified Case ($m=2$)

i. Losses for Encoders

Compute reconstruction loss $\|X - WX\|_F^2$ and regularization loss $\|W\|_F^2$:

- For α ($W = I$): Reconstruction = 0 (since $WX = X$), Regularization = 2 (sum of squares of elements: $1^2 + 1^2 = 2$).
- For β : Reconstruction = 0.001 (as given; verified as sum of squares of second row of X : $0.02^2 + (-0.01)^2 + 0.01^2 + (-0.02)^2 = 0.0004 + 0.0001 + 0.0001 + 0.0004 = 0.001$), Regularization = 1 (sum of squares: $1^2 + 0 + 0 + 0 = 1$).

Encoder	Reconstruction Loss	Regularization Loss	
α	0	2	
β	0.001	1	

To arrive at reconstruction for β : WX sets the second row to zeros, so residuals are zeros in first row and

X's second row in second; Frobenius norm squared is sum of element squares in residuals.

To arrive at regularization: Direct sum of squares of W elements.

ii. Values of λ Where $L_2(\alpha) > L_2(\beta)$

$$L_2(\alpha) = 0 + \lambda \cdot 2 = 2\lambda$$

$$L_2(\beta) = 0.001 + \lambda \cdot 1 = 0.001 + \lambda$$

Set $2\lambda > 0.001 + \lambda \Rightarrow \lambda > 0.001$.

To arrive: Solve the inequality algebraically.

(b) Generic Square Linear Encoder

i. Form of Optimizing W^*

$$\hat{W} = \begin{bmatrix} U & \end{bmatrix} \cdot \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\sigma_m^2}{\sigma_m^2 + \lambda} \end{bmatrix} \cdot \begin{bmatrix} U^T & \end{bmatrix}$$

(The blanks are U on the left and U^T on the right, where U is the left singular matrix from X's SVD.)

ii. Derivation

Using hints: Express L_2 in SVD basis to decouple.

Given $X = U\Sigma V^T$ (U orthogonal $m \times m$, Σ diagonal $m \times m$, $V^T V = I$).

Let $\hat{W} = U^T W U$. Then $\|W\|_F^2 = \|\hat{W}\|_F^2$ (Frobenius invariant under orthogonal transforms).

Reconstruction: $\|X - W X\|_F^2 = \|(I - \hat{W})\Sigma\|_F^2$ (from invariance and substitutions as detailed).

Expand: $\|(I - \hat{W})\Sigma\|_F^2 = \sum_{p,q=1}^m (\delta_{pq} - \hat{W}_{pq})^2 \sigma_q^2$.

Thus, $L_2 = \sum_{p,q} (\delta_{pq} - \hat{W}_{pq})^2 \sigma_q^2 + \lambda \sum_{p,q} \hat{W}_{pq}^2$.

This decouples into m^2 independent terms (one per \hat{W}_{pq}):

- For $p \neq q$: $\min_w \sigma_q^2 w^2 + \lambda w^2 = w^2(\sigma_q^2 + \lambda) \rightarrow \text{optimal } w = 0$.



- For $p = q$: $\min_w \sigma_p^2(1 - w)^2 + \lambda w^2$.

Derivative: $-2\sigma_p^2(1 - w) + 2\lambda w = 0 \rightarrow \sigma_p^2(1 - w) = \lambda w \rightarrow w = \sigma_p^2/(\sigma_p^2 + \lambda)$.

To arrive: Set derivative to zero and solve for w .

Thus, \hat{W} diagonal with entries $\sigma_i^2/(\sigma_i^2 + \lambda)$, so $W^* = U\hat{W}U^T$.

This is equivalent to m decoupled ridge regressions: each principal direction is a "feature" regressed onto itself with ridge penalty.

(c) Hyperparameter Values for Guarantee

The scaling factors are $d_i = \sigma_i^2/(\sigma_i^2 + \lambda)$.

One shot all of problem 3, didn't have to tell it to do anything extra. No hallucinations, accurate like official solutions.

Require $d_i \geq 0.8$ for $i = 1, 2, 3$ ($\sigma = 10, 8, 4$) and $d_i \leq 0.5$ for $i = 4, \dots, 8$ ($\sigma = 1, 0.5, 0.36, 0.16, 0.01$).

Since $d(\sigma)$ increasing in σ , critical conditions:

- Weakest preservation: $d(4) \geq 0.8 \rightarrow \lambda \leq 4^2(1 - 0.8)/0.8 = 16 \cdot 0.25 = 4$.
- Weakest attenuation: $d(1) \leq 0.5 \rightarrow \lambda \geq 1^2(1 - 0.5)/0.5 = 1 \cdot 1 = 1$.

Thus, $\lambda \in [1, 4]$.

To arrive: Solve inequalities for bounds on λ using the formulas for $d \geq \alpha$ ($\lambda \leq \sigma^2(1 - \alpha)/\alpha$) and $d \leq \beta$ ($\lambda \geq \sigma^2(1 - \beta)/\beta$).

now, for the finally: Trying to refresh chat once again

4. Ridge-Attention In lecture, you saw how the standard softmax-attention mechanism can be viewed as a softened version of something akin to a nearest-neighbor model in which the value returned for a query reflects a weighted combination of the values that correspond to the keys closest to the query. In this view, the (key, value) pairs in the



memory represent a kind of in-context “training data” and the query is a test input for which we want to predict the right output given that data. (a) To start, let’s think about why it is possible to efficiently update simple averaging. Let $m = \frac{1}{n} \sum_{i=1}^n x_i$ be the average of n points. Use m , x_{n+1} , n and simple arithmetic operations to compute $m' = \frac{1}{n+1} \sum_{i=1}^{n+1} x_i$ — the average of all points including the new point x_{n+1} . (HINT: Start by multiplying m by n .) (b) Let us now shift to thinking about traditional ridge-regression with n training pairs (x_i, y_i) where x_i are d -dimensional vectors and y_i are scalars. Let the matrix $A = [x_1 \ x_2 \ \dots \ x_n]$ and vector $y = [y_1 \ y_2 \ \dots \ y_n]$ so that we can find the familiar closed-form solution $w^* = (A^T A + \lambda I)^{-1} A^T y$ (6) that allows us to make scalar predictions on any new test input x by computing $w^T * x$. First, write the two terms $(A^T A + \lambda I)$ and $A^T y$ as sums involving the x_i and y_i . i.e. Complete: $(A^T A + \lambda I) = \lambda I + \sum_{i=1}^n A^T y_i = \sum_{i=1}^n A^T A$ (c) Suppose we wanted to do ridge-self-attention (non-causal – “encoder-style”) with a context length of n and d -dimensional query, key, and value vectors. Recall that this is the style of attention where each of the n queries is applied to the same pool of n (key, value) pairs. The goal is to calculate multi-dimensional ridge-regression predictions, after “training” on the pool of (key, value) pairs, and given the query as a kind of “test” input. (Note: the keys are playing the role of the A matrix in ridge-regression, each query is like the x we are testing on, and in place of the scalar y_i , we have an entire value vector for a multi-dimensional ridge-regression problem so the y vector is replaced with a matrix that has a row for each value vector.) Assume that the cost of inverting a $d \times d$ matrix is $O(d^3)$ and the cost of multiplying two such matrices is also $O(d^3)$. Assume that a $d \times d$ matrix times either a d -dimensional row or column vector costs d^2 operations. You should assume $d < n$. What is the computational cost of a non-causal ridge self-attention layer? $O(d^4)$ $O(nd^2)$

Homework 8, © Faculty teaching EECS 182, Fall 2025. All Rights Reserved. This may not be publicly shared without explicit permission. 4 Homework 8 @ 2025-10-25 08:44:19Z $O(n^2d^3)$ $O(n^2d^2)$ $O(n^2)$ $O(1)$ (HINT: Do not forget that for a single d -dimensional query vector q , attention needs to return a d -dimensional result.) (d) Assume that a ridge self-attention layer is used in a Transformer architecture and there is a downstream loss.

For which of these will backprop successfully pass gradients if we use ridge self-attention? *

The ridge λ viewed as a learnable parameter for the self-attention layer. *

The keys *

The values *

The queries (e) Now step back. There is a nice trick (called the Sherman–Morrison Formula) by which one can update the inverse of an invertible matrix to which you make a rank-1 update. Let M be an invertible square $d \times d$ matrix and let u, v be two d -dimensional vectors. Then: $(M + uv^T)^{-1} = M^{-1} - \frac{1}{1 + v^T M^{-1} u} (M^{-1} u)(v^T M^{-1})$

(7) Assume that a $d \times d$ matrix times either a d -dimensional row or column vector costs d operations, and so does the evaluation of a dyad uv^T . Assume that computing a Euclidean inner-product costs d operations. Assuming that you already had M^{-1} in hand, what is the computational cost of one application of (7)? $O(d^4)$ $O(d^3)$ $O(d^2)$ $O(d)$ $O(1)$

(f) Consider implementing causal ridge-self-attention with a context length of n but where the pool of (key, value) vectors that one is querying at position t consists only of the t (key, value) pairs so far. Describe explicitly how you would compute causal ridge-self-attention in a computationally efficient manner. Leverage your decomposition of the ridge-regression formula in part (b) of this problem together with the Sherman–Morrison formula from (7) to avoid having to do $O(n^2)$ computations while still calculating causal ridge-self-attention outputs correctly for all n positions in the context. (HINT: Think recursively. What do you need to track from one time step to the next to avoid repeating work?)

(g) Many people consider important the ability to visualize the attention weights. For traditional softmax attention, these are the outputs of the softmax for a given query vector. They tell you the exact amounts by which the attention outputs at this position are linear combinations of the values being fed in at this and other positions. For ridge-attention and a given query vector, how would you compute the (possibly negative) weights associated to each of the value vectors in the context?

(a) To update the average m' after adding a new point x_{n+1} :

Start with $n \cdot m = \sum_{i=1}^n x_i$.



Add x_{n+1} to get the new sum: $n \cdot m + x_{n+1}$.

Divide by $n + 1$: $m' = \frac{n \cdot m + x_{n+1}}{n+1}$.

(b) $(A^\top A + \lambda I) = \lambda I + \sum_{i=1}^n x_i x_i^\top$

$A^\top y = \sum_{i=1}^n x_i y_i$

(c) $\bigcirc O(nd^2)$

(d) The ridge λ viewed as a learnable parameter for the self-attention layer.

The keys

The values

The queries

(e) $\bigcirc O(d^2)$

(f) To compute causal ridge-self-attention efficiently:

Maintain two quantities recursively: the inverse M_{t-1}^{-1} of the Gram matrix $M_{t-1} = \lambda I + \sum_{i=1}^{t-1} k_i k_i^\top$ and the weight matrix $W_{t-1} = M_{t-1}^{-1} \sum_{i=1}^{t-1} k_i v_i^\top$.

Initialize at $t = 0$: $M_0^{-1} = \frac{1}{\lambda} I$, $W_0 = 0$.

At step t :

- Compute $u = M_{t-1}^{-1} k_t$ (cost $O(d^2)$).
- Compute $\beta = k_t^\top u$ (cost $O(d)$).
- Compute $\alpha = \frac{1}{1+\beta}$ (cost $O(1)$).
- Update the inverse: $M_t^{-1} = M_{t-1}^{-1} - \alpha u u^\top$ (cost $O(d^2)$, since $u u^\top$ is a rank-1 outer product).
- Compute $k_t^\top W_{t-1}$ (cost $O(d^2)$).
- Compute the correction term: $\alpha u(v_t^\top - k_t^\top W_{t-1})$ (cost $O(d^2)$, rank-1 update).
- Update $W_t = W_{t-1} + \alpha u(v_t^\top - k_t^\top W_{t-1})$ (cost $O(d^2)$).
- Compute the output for query q_t : $o_t = q_t^\top W_t$ (cost $O(d^2)$).



Literally one shots problem 4, no complaints.

This ensures each step costs $O(d^2)$, for a total of $O(nd^2)$ over n positions, avoiding redundant computations by recycling the previous inverse and weights.

(g) For a given query q , the (possibly negative) weights a_i for each value vector v_i such that the output is $\sum_{i=1}^n a_i v_i$ are given by $a_i = k_i^\top (K^\top K + \lambda I)^{-1} q$, where K is the matrix of keys (with rows k_i^\top). These can be computed as the elements of the vector $K(K^\top K + \lambda I)^{-1} q$.

Executive summary: Grok handled most non-coding HW8 problems accurately, particularly those involving linear algebra and ridge-regression interpretations, and often one-shotted the correct solution. It showed no major hallucinations and even produced valid alternative derivations in some cases. The primary weakness was in algorithmic complexity analysis, where it sometimes stated correct ideas imprecisely, especially regarding critical path length versus total work. Overall, Grok demonstrated strong formal reasoning but not as tuned with the specific computational assumptions emphasized in lecture.