Weight Matrix Evolution during MiniAlexNet training

Presenter: Doyoon Kim, Enoch Yiu

UC Berkeley

2nd May, 2025

Background

According to Martin&Mahoney, Regularization \approx Generalization So that if in order for the deep learning model to perform well on unseen data(test data), it should be regularized well. To prevent being overfitted to training data, the weight matrices of the model should be 'simple'.



There are two kinds of regularization: explicit vs. implicit

- explicit regularization:
 - Dropout: randomly zeroes a fraction of activations during training
 - Weight norm constraints, e.g.
 - Lasso (L₁ penalty):

$$\min_{w} \frac{1}{N} \sum_{i=1}^{N} L(f(x_i; w), y_i) + \lambda \|w\|_1$$

3/39

• implicit regularization:

Neural network's training dynamics drive weight matrices toward a "simple" state

How can we "check" that the matrix is getting simpler?

Marchenko–Pastur (MP) Law: For a random matrix $W \in {}^{p \times n}$ with iid entries of variance σ^2 , the eigenvalue density of $\frac{1}{n}WW^T$ converges to

$$\rho_{\rm MP}(\lambda) = \frac{1}{2\pi\sigma^2 c \,\lambda} \sqrt{(\lambda_+ - \lambda)(\lambda - \lambda_-)}, \quad \lambda_{\pm} = \sigma^2 (1 \pm \sqrt{c})^2, \quad c = \frac{p}{n}.$$

The interval $[\lambda_{-}, \lambda_{+}]$ is the *noise bulk*; eigenvalues outside are *spikes* (signal).



This growth in spike count signals that the weight matrix develops a simpler, more low-rank "signal" component.

5/39

- If the weight matrix follows MP law, can we reconstruct the matrix with a few spikes(factors)?
- Is it possible to shrink Eigenvectors of covariance matrices using JSE to improve regularization?
- If it's possible, where should we shrink it?

Used the same architecture used in Martin&Marhoney, which is miniAlexnet, simplified version of alexnet. It has 2 CNN layers, MaxPooling layer with 3 Fully Connected layers, which doesn't include explicit regularization.



We used CIFAR10 dataset to train minialexnet. The data consists of 60000 32x32 colour images in 10 classes, with 50000 training images and 10000 test images.

airplane	Service and	X	-	X	*	+	ø,	-4-	-	-
automobile	æ				-	The			-	*
bird	S.	ſ	1			-	1	Y	2	4
cat			4	de			2	Å.	No.	2
deer	6	48	\mathbf{X}	RA		Y	Ŷ	1	đ.	<u>\$</u>
dog	1×1	C.	-	9.	1		9	V?	A	No.
frog	.2	1	-		? ?			5		See
horse	-	T.	P	2	1	KT	1	24	6	1
ship	-		dirin.	-	M		2	127	1	
truck		No.		Ş.					-	da

Q. If the weight matrix follows MP law, can we reconstruct the matrix with a few spikes(factors)?

A. We can approximate it, but not perfect. For FC1 (W_1) , FC2 (W_2) , and FC3 (W_3)

$$W_{i} = U\Sigma V^{T} = U_{p \times k} \Sigma_{k \times k} V_{n \times k}^{T} + Z,$$
$$\hat{W}_{i}(k) = U_{p,k} \Sigma_{k,k} V_{n,k}^{T} = \sum_{i=1}^{k} \sigma_{i} u_{i} v_{i}^{T},$$
$$\operatorname{Acc}(\hat{W}_{i}(k)) = \frac{\# \text{correct predictions}}{\# \text{test samples}}$$

Fixed FC2, FC3 matrices and varied k in FC1 weight matrix(4096x384). As we train the model, the number of factors k needed to approximate full rank decreased.



Similarly for FC2 layer(384x192),

As we train the model, the number of factors **k** needed to approximate full rank decreased.





Q. Is it possible to shrink Eigenvectors of covariance matrices using JSE to improve regularization?

A. No. If we shrink it toward the mean.

Shrink each eigenvectors towards its mean

$$W = U \Sigma V^T \quad \longrightarrow \quad H = \frac{1}{\sqrt{n}} U_{:,1:k} \Sigma_{1:k,1:k}$$
$$M = \frac{1}{p} \mathbf{1} \mathbf{1}^T H, \quad R = H - M$$
$$\nu^2 = \frac{(RR^T)}{n_+ - k}, \quad J = R^T R, \quad C = I_k - \nu^2 J^{-1}$$
$$\overline{H_{JS}} = HC + M (I_k - C) \implies \quad H_{JS} = U_{JS} \Sigma_{JS} V_{JS}^T$$

James–Stein Shrinkage for SVD Reconstruction

(a) Vector-only shrinkage:

$$\widehat{W}_{\text{vec}}(k) = U_{\text{JS}} \Sigma_{1:k,1:k} V_{:,1:k}^T$$

(b) Full shrinkage of singular values:

$$\widehat{W}_{\text{full}}(k) = U_{\text{JS}} \Sigma_{\text{JS}} V_{:,1:k}^T$$

where

$$S^{2} = \frac{1}{n}\Sigma^{2}, \quad \Psi^{2} = I - \nu^{2} S^{-2} = I - \nu^{2} n \Sigma^{-2}, \quad \Phi = S^{2} \Psi^{2} = \frac{1}{n} (\Sigma^{2} - n\nu^{2}I)$$

$$\Sigma_{\rm JS} = (n \Phi)^{\frac{1}{2}} = (\Sigma^2 - n \nu^2 I)^{\frac{1}{2}} = {\rm diag}(\sqrt{\sigma_i^2 - n \nu^2})$$

13/39

James–Stein Shrinkage for SVD Reconstruction

No dramatic increase in performance no matter what k is.



James–Stein Shrinkage for SVD Reconstruction

Zooming in, there are some factor numbers that JS shrinkage works better, but more like coincidence.



We need better shrinkage target than grand mean point. For that, we checked the behavior of leading eigenvector upon epoch. Let $x_i \in {}^p$ be the leading eigenvector at epoch i, $||x_i|| = 1$. Define the projection basis matrix A by:

$$A = \begin{cases} \begin{bmatrix} b & \frac{t}{\|t\|} \end{bmatrix} \in^{p \times 2}, & \text{for 1D GPCA (tangent line);} \\ \begin{bmatrix} b & t_1 & t_2 \end{bmatrix} \in^{p \times 3}, & \text{for 2D GPCA (tangent plane).} \end{cases}$$

Then in both cases:

$$P = A (A^T A)^{-1} A^T, \quad \hat{x}_i = \frac{P x_i}{\|P x_i\|} \in S^{p-1}.$$

Residual Sum of Squares (RSS) = $\sum_{i=1}^{L} d^2(x_i, \hat{x}_i) \quad (x_i, \hat{x}_i \in S^{p-1}),$

where $d(\cdot, \cdot)$ is the geodesic distance on the sphere.

Mixed Variance =
$$\underbrace{\sum_{i=1}^{L} d^2(x_i, \hat{x}_i)}_{\text{RSS}} + \underbrace{\sum_{i=1}^{L} d^2(\hat{x}_i, \mu)}_{\text{variance on } S^{p-1}}$$

where μ is the Fréchet mean of $\{x_i\}$, i.e. the point on S^{p-1} that

$$\mu = \arg \min_{p \in S^{p-1}} \sum_{i=1}^{L} d^2(x_i, p).$$

Fitting Score =
$$R^2 = 1 - \frac{RSS}{Mixed Variance} \in [0, 1]$$
,
measuring the proportion of total dispersion explained by the
projection.

2D GPCA fitting example

With proper learning rate(when weight doesn't converge within 100 epochs), leading eigen vector tend to follow geodesics



Figure 1: FC1 results

2D GPCA fitting example



Figure 2: FC2 results

21/39

Choosing start of trajectory

The trajectory of leading eigenvector stabilizes after few epochs. The first few vectors have high residuals.

Truncating those leads to lower RSS and higher fitting score.



However, the optimal "start" of the trajectory depends on your hyperparameters.

To choose it automatically, monitor each epoch's weight-matrix spectrum using Marchenko–Pastur law.

The first epoch at which a clear outlier (spike) appears in the spectrum is then taken as the beginning of your trajectory.

23/39

Choosing start of trajectory





2D GPCA: RSS=3.2347, FitScore=0.8651





Trained model with 20 different seeds and projected 3-100 epoch points into S1, S2 sphere. Recall: RSS of Random Walk was about 200, Fitting score was about 0.55



Multiple seed experiments



Multiple seed experiments



Multiple seed experiments



28/39

There were some cases that epochs follows three different geodesics. After 18, 48 epochs, the projected eigenvector moves abruptly. And has big resiual magnitude at that point.

29/39

Multiple geodesic during training





1-18





19-48







30 / <u>39</u>

This might be due to the leading eigenvalue crossing.

That the leading eigenvalue, vector changes.

This phenomena sometimes happened when the learning rate(step size) is big.



Multiple geodesic during training



To understand how eigenvector trajectories react under perturbations, we compare:

• Data shift at epoch 50: Swap to a disjoint set of classes mid-training.

For first 50 epochs, train the model using the airplane, automobile, bird, cat and deer images.

After 50 epochs, train the model with dog, frog, horse, ship and truck images.

Will the leading eigenvector's trajectory follow the geodesic still?

Varying training data after 50 epochs



Varying training data after 50 epochs

Direction changes during 51 53 epochs, but after that, it goes to the similar direction as before.

Still, leading eigenvectors are on the geodesic line.



Varying training data after 50 epochs

Since the residual near epoch 50 is rather small, we can say that even if we change training data, it still follows the geodesic line.

So whatever the data is, the weight matrix tend to follow geodesic line.



Since we found that the eigenvectors tend to follow geodesic when we train it, we can choose shrinkage point upon the geodesic.

- Mahoney, Michael, and Charles Martin. "Traditional and heavy tailed self regularization in neural network models." International Conference on Machine Learning. PMLR, 2019.
- Goldberg, Lisa R., and Alec N. Kercheval. "James–Stein for the leading eigenvector." Proceedings of the National Academy of Sciences 120.2 (2023): e2207046120.

39/39