

Supplementary Material: Data-Efficient Augmentation for Training Neural Networks

A Proof of Main Results

A.1 Proof for Lemma 4.1

Proof. Let $\delta_i := \tilde{\sigma}_i - \sigma_i$, where $\mathbb{P}(\delta_i < 0) = p_i$. Assuming uniform probability between $-\|\mathbf{E}\|$ to 0, and between 0 to $\|\mathbf{E}\|$, we have pdf $\rho_i(x)$ for δ_i :

$$\rho_i(x) = \begin{cases} \frac{p_i}{\|\mathbf{E}\|}, & \text{if } -\|\mathbf{E}\| \leq x < 0 \\ \frac{1-p_i}{\|\mathbf{E}\|}, & 0 \leq x \leq \|\mathbf{E}\| \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

Taking expectation,

$$\mathbb{E}(\tilde{\sigma}_i - \sigma_i) = \mathbb{E}(\delta_i) = \int_{-\infty}^{\infty} x \rho_i(x) dx \quad (2)$$

$$= \int_{-\|\mathbf{E}\|}^0 x \frac{p_i}{\|\mathbf{E}\|} dx + \int_0^{\|\mathbf{E}\|} x \frac{1-p_i}{\|\mathbf{E}\|} dx \quad (3)$$

$$= -\frac{\|\mathbf{E}\| p_i}{2} + \frac{(1-p_i)\|\mathbf{E}\|}{2} \quad (4)$$

$$= \frac{(1-2p_i)\|\mathbf{E}\|}{2} \quad (5)$$

We also have

$$\mathbb{E}(\delta_i^2) = \int_{-\infty}^{\infty} x^2 \rho_i(x) dx \quad (6)$$

$$= \int_{-\|\mathbf{E}\|}^0 x^2 \frac{p_i}{\|\mathbf{E}\|} dx + \int_0^{\|\mathbf{E}\|} x^2 \frac{1-p_i}{\|\mathbf{E}\|} dx \quad (7)$$

$$= \frac{\|\mathbf{E}\|^2 p_i}{3} + \frac{(1-p_i)\|\mathbf{E}\|^2}{3} \quad (8)$$

$$= \frac{\|\mathbf{E}\|^2}{3} \quad (9)$$

Thus, we have

$$\mathbb{E}(\tilde{\lambda}_i) = \mathbb{E}(\tilde{\sigma}_i^2) \quad (10)$$

$$= \mathbb{E}((\sigma_i + \delta_i)^2) \quad (11)$$

$$= \mathbb{E}(\sigma_i^2 + 2\sigma_i\delta_i + \delta_i^2) \quad (12)$$

$$= \sigma_i^2 + 2\sigma_i\mathbb{E}[\delta_i] + \mathbb{E}[\delta_i^2] \quad (13)$$

$$= \sigma_i^2 + 2\sigma_i \frac{(1-2p_i)\|\mathbf{E}\|}{2} + \frac{\|\mathbf{E}\|^2}{3} \quad (14)$$

$$= \sigma_i^2 + \sigma_i(1-2p_i)\|\mathbf{E}\| + \frac{\|\mathbf{E}\|^2}{3} \quad (15)$$

□

A.2 Proof of Corollary 4.2

Under the assumptions of Theorem 5.1 of [1], i.e. where the minimum eigenvalue of the NTK is $\lambda_{\min}(\mathcal{J}\mathcal{J}^T) \geq \lambda_0$ for a constant $\lambda_0 > 0$, and training data \mathbf{X} of size n sampled i.i.d. from distribution D and 1-Lipschitz loss \mathcal{L} , we have that with probability $\delta/3$, training the over-parameterized neural network with gradient descent for $t \geq \Omega\left(\frac{1}{n\lambda_0} \log \frac{n}{\delta}\right)$ iterations results in the following population loss \mathcal{L}_D (generalization error)

$$\mathcal{L}_D(\mathbf{W}^t, \mathbf{X}) \leq \sqrt{\frac{2\mathbf{y}^T(\mathcal{J}\mathcal{J}^T)^{-1}\mathbf{y}}{n}} + \mathcal{O}\left(\frac{\log \frac{n}{\lambda_0\delta}}{n}\right), \quad (16)$$

with high probability of at least $1 - \delta$ over random initialization and training samples.

Hence, using $\lambda_{\min}, \sigma_{\min}$ to denote minimum eigen and singular value respectively of the NTK corresponding to full data, we get

$$\mathcal{L}_{D_{train}}(\mathbf{W}^t, \mathbf{X}_{train}) \leq \sqrt{\frac{2\frac{1}{\lambda_{\min}}\|y\|^2}{n}} + \mathcal{O}\left(\log \frac{1}{\delta}\right) \quad (17)$$

$$\leq \sqrt{\frac{2}{\sigma_{\min}^2}} + \mathcal{O}\left(\log \frac{1}{\delta}\right). \quad (18)$$

For augmented dataset \mathbf{X}_{aug} , we have $\tilde{\sigma}_i \leq \sigma_i + \sqrt{n}L\epsilon_0$, hence the improvement in the generalization error is at most

$$\mathcal{L}_{D_{aug}}(\mathbf{W}^t, \mathbf{X}_{aug}) \leq \sqrt{\frac{2}{(\sigma_{\min} + \sqrt{n}L\epsilon_0)^2}} + \mathcal{O}\left(\log \frac{1}{\delta}\right). \quad (19)$$

Combining these two results, we obtain Corollary 4.2.

A.3 Proof of Lemma 5.1

Proof.

$$\|\mathcal{J}^T(\mathbf{W}^t, \mathbf{X}_{aug})\mathbf{r} - \text{diag}(\boldsymbol{\rho}^t)\mathcal{J}^t(\mathbf{W}^t, \mathbf{X}_{S_{aug}})\mathbf{r}_S\| \quad (20)$$

$$= \|(\mathcal{J}^T(\mathbf{W}^t, \mathbf{X}) + \mathbf{E})\mathbf{r} - (\text{diag}(\boldsymbol{\rho}^t)\mathcal{J}^t(\mathbf{W}^t, \mathbf{X}_S) + \mathbf{E}_S)\mathbf{r}_S\| \quad (21)$$

$$\leq \|(\mathcal{J}^T(\mathbf{W}^t, \mathbf{X})\mathbf{r} - \text{diag}(\boldsymbol{\rho}^t)\mathcal{J}^t(\mathbf{W}^t, \mathbf{X}_S)\mathbf{r}_S) + \mathbf{E}\mathbf{r} - \mathbf{E}_S\mathbf{r}_S\| \quad (22)$$

$$\leq \|(\mathcal{J}^T(\mathbf{W}^t, \mathbf{X})\mathbf{r} - \text{diag}(\boldsymbol{\rho}^t)\mathcal{J}^t(\mathbf{W}^t, \mathbf{X}_S)\mathbf{r}_S)\| + \|\mathbf{E}\mathbf{r}\| + \|\mathbf{E}_S\mathbf{r}_S\| \quad (23)$$

Applying definition of coresets, we obtain

$$\|(\mathcal{J}^T(\mathbf{W}^t, \mathbf{X})\mathbf{r} - \text{diag}(\boldsymbol{\rho}^t)\mathcal{J}^t(\mathbf{W}^t, \mathbf{X}_S)\mathbf{r}_S)\| + \|\mathbf{E}\mathbf{r}\| + \|\mathbf{E}_S\mathbf{r}_S\| \quad (24)$$

$$\leq \xi + \|\mathbf{E}\mathbf{r}\| + \|\mathbf{E}_S\mathbf{r}_S\| \quad (25)$$

$$\leq \xi + 2n^{\frac{3}{2}}L\epsilon_0 \quad (26)$$

$$\leq \xi + 2\sqrt{L} \quad (27)$$

□

A.4 Proof of Theorem 5.2

Proof. In this proof, as shorthand notation, we use \mathbf{X}_f and \mathbf{X}_{train} interchangeably. We further use \mathbf{X}_c to represent the coreset selected from the full data, and $\mathbf{X}_{c_{aug}}$ to represent the augmented coreset.

By Theorem 1 of [2], under the α -PL assumption for \mathcal{L} and interpolation assumption (i.e. for every sequence $\mathbf{W}^1, \mathbf{W}^2, \dots$ such that $\lim_{t \rightarrow \infty} \mathcal{L}(\mathbf{W}^t, \mathbf{X}) = 0$, we have that the loss for each data point $\lim_{t \rightarrow \infty} \mathcal{L}(\mathbf{W}^t, \mathbf{x}_i) = 0$), the convergence of SGD with constant step size is given by

$$\mathbb{E}[\|\nabla \mathcal{L}(\mathbf{W}^t, \mathbf{X}_{f+c_{aug}})\|^2] \leq \left(1 - \frac{\alpha\eta}{2}\right)^t \mathcal{L}(\mathbf{W}^0, \mathbf{X}_{f+c_{aug}}) \quad (28)$$

$$\leq \frac{1}{\alpha} \left(1 - \frac{\alpha\eta}{2}\right)^t \|\nabla \mathcal{L}(\mathbf{W}^0, \mathbf{X}_{f+c_{aug}})\|^2 \quad (29)$$

Using Jensen's inequality, we have

$$E[\text{krL}(W^0; X_{f+C_{\text{aug}}})k] \tag{30}$$

$$\frac{1}{q} \frac{E[\text{krL}(W^t; X_{f+C_{\text{aug}}})k^2]}{E[\text{krL}(W^t; X_{f+C_{\text{aug}}})k]} \tag{31}$$

$$\frac{1}{p} = \frac{1}{2} \frac{1}{2} \text{krL}(W^0; X_{f+C_{\text{aug}}})k \tag{32}$$

$$\frac{1}{p} = \frac{1}{2} \frac{1}{2} \text{krL}(W^0; X_f)k + \text{krL}(W^0; X_{C_{\text{aug}}})k \tag{33}$$

$$\frac{1}{p} = \frac{1}{2} \frac{1}{2} G_0 + k(J(W^0; X_c) + E)(y - f(W^0; X_c +))k \tag{34}$$

$$\frac{1}{p} = \frac{1}{2} \frac{1}{2} \tag{35}$$

$$G_0 + k(J(W^0; X_c) + E)^T (y - f(W^0; X_c) - r_{xf}(W^0; X_c)^T - O(\tau))k \tag{36}$$

$$= \frac{1}{p} = \frac{1}{2} \frac{1}{2} (G_0 + \text{krL}(W^0; X_c) - (J(W^0; X_c)^T (r_{xf}(W^0; X_c)^T + O(\tau))) + \tag{37}$$

$$E(y - f(W^0; X_c +))k) \tag{38}$$

$$\frac{1}{p} = \frac{1}{2} \frac{1}{2} (G_0 + \text{krL}(W^0; X_c) - (J(W^0; X_c)^T (r_{xf}(W^0; X_c)^T + O(\tau)))k + \tag{39}$$

$$p \sqrt{2kE}k) \tag{40}$$

$$\frac{1}{p} = \frac{1}{2} \frac{1}{2} (G_0 + \text{krL}(W^0; X_c)k + \max L^p \bar{n}_0 + \max O(n_0^2)) + p \sqrt{2nL_0} \tag{41}$$

$$= \frac{1}{p} = \frac{1}{2} \frac{1}{2} (G_0 + \text{krL}(W^0; X_f)k + \max L^p \bar{n}_0 + \max O(n_0^2)) + p \sqrt{2nL_0} \tag{42}$$

$$\frac{1}{p} = \frac{1}{2} \frac{1}{2} (2G_0 + \max L^p \bar{n}_0 + \max O(n_0^2)) + p \sqrt{2nL_0} \tag{43}$$

□

A.5 Finding Subsets

Let S be a subset of training data points. Furthermore, assume that there is a mapping $\rho: S \rightarrow V$ that for every W assigns every data point $i \in V$ to its closest element $j \in S$, i.e. $j = \rho_{w;S}(i) = \arg \max_{j \in S} s_{ij}(W)$, where $s_{ij}(W) = C - kJ^T(W^t; x_i)r_i - J^T(W^t; x_j)r_jk$ is the similarity between gradients of f and j , and $C = \max_{ij} s_{ij}(W)$ is a constant. Consider a matrix $G_{w;S} \in \mathbb{R}^{n \times m}$, in which every row contains gradient of $w(i)$, i.e., $[G_{w;S}]_i = J^T(W^t; x_{\rho_{w;S}(i)})r_{\rho_{w;S}(i)}$. The Frobenius norm of the matrix $G_{w;S}$ provides an upper-bound on the error of the weighted subset in capturing the alignment of the residuals of the full training data with the Jacobian matrix. Formally,

$$kJ^T(W^t; X_{\text{train}})r_{\text{train}} - \sum_{S^t} J^T(W^t; [X_{\text{train}}]_{:S^t})r_{S^t}k \leq kG_{w;S}k_F; \tag{44}$$

where the weight vector $s^t \in \mathbb{R}^{|S^t|}$ contains the number of elements that are mapped to every element $j \in S$ by mapping $\rho_{w;S}$, i.e. $s_j = \sum_{i \in V} 1[\rho_{w;S}(i) = j]$. Hence, the set of training points that closely estimate the projection of the residuals of the full training data on the Jacobian spectrum can be obtained by finding a subset S that minimizes the Frobenius norm of matrix $G_{w;S}$.

B Additional Theoretical Results

B.1 Convergence analysis for training on augmented full data

Theorem B.1. Gradient descent with learning rate η applied to a neural network with constant NTK and Lipschitz constant L , and data points \mathcal{D}_{aug} augmented with m additive perturbations bounded by σ results in the following training dynamics:

$$\mathbb{E}[\|k_{\text{y}} - f(X_{\text{aug}}; W^t)\|_2^2] = \frac{1}{X^n} \prod_{i=1}^n \left(1 - \frac{\eta^2}{\lambda_i^2} + \frac{\eta^2}{\lambda_i} (1 - 2p_i) \frac{\eta L^2}{3} \right)^{2t} \left((u_i y)^2 + 2n \frac{p}{2} \frac{\eta L^2}{3} \right) \quad (45)$$

where \mathbb{E} with ηL^2 $\frac{p}{2} \frac{\eta L^2}{3}$ is the perturbation to the Jacobian, and $p := P(\lambda_i < 0)$ is the probability that λ_i decreases as a result of data augmentation.

B.2 Proof of Theorem B.1

Using Jensen's inequality, we have

$$\mathbb{E}[\|k_{\text{y}} - f(X_{\text{aug}}; W^t)\|_2^2] \leq \frac{2 \sum_{i=1}^n \mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2]}{X^n} \quad (46)$$

$$= \frac{2 \sum_{i=1}^n \mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2]}{X^n} \quad (47)$$

$$\mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2] = \frac{\mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2]}{\mathbb{E}[(1 - \tilde{\gamma}_i)^{2t} (u_i^T y)^2]} \quad (48)$$

$$\mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2] = \frac{\mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2]}{\mathbb{E}[(1 - \tilde{\gamma}_i)^{2t} ((u_i y)^2 + 2n \frac{p}{2} \frac{\eta L^2}{3})]} \quad (49)$$

$$\mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2] = \frac{\mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2]}{\mathbb{E}[(1 - \tilde{\gamma}_i)^{2t} ((u_i y)^2 + 2n \frac{p}{2} \frac{\eta L^2}{3})]} \quad (50)$$

$$= \frac{\mathbb{E}[\|k_{\text{y}} - f(X_i; W^t)\|_2^2]}{X^n \prod_{i=1}^n \left(1 - \frac{\eta^2}{\lambda_i^2} + \frac{\eta^2}{\lambda_i} (1 - 2p_i) \frac{\eta L^2}{3} \right)^{2t} \left((u_i y)^2 + 2n \frac{p}{2} \frac{\eta L^2}{3} \right)} \quad (51)$$

B.3 Convergence analysis for training on the coreset and its augmentation

Theorem B.2. Let L_i be μ -smooth, L be L -smooth and satisfy the PL condition, that is for $\gamma > 0$, $\text{kr}(L(W; X)) \geq \gamma \|L(W; X) - L(W; X^*)\|_2^2$ for all weights W . Let ϵ upper-bound the normed difference in gradients between the weighted coreset and full dataset. Assume that the network $f(W; X)$ is Lipschitz in W , X with Lipschitz constant L and L' respectively, and $\epsilon = \max\{L, L'\} \sigma$. Let G_0 the gradient over the full dataset at initialization, λ_{max} the maximum Jacobian singular value at initialization. Choosing perturbation bound $\frac{\epsilon}{\lambda_{\text{max}}}$ where λ_{max} is the maximum singular value of the coreset Jacobian and the size of the original dataset, running SGD on the coreset and its augmentation using constant step size $\frac{\epsilon}{\lambda_{\text{max}}}$, we get the following convergence bound:

$$\mathbb{E}[\|k_{\text{y}} - f(X_{\text{C} + \text{C}_{\text{aug}}}; W^t)\|_2^2] \leq \frac{1}{2} \left(1 - \frac{\gamma}{2} \right)^{\frac{t}{2}} (2G_0 + 2\epsilon) + O\left(\frac{L}{\lambda_{\text{max}}}\right); \quad (52)$$

where $X_{\text{C} + \text{C}_{\text{aug}}}$ represents the dataset containing the (weighted) coreset and its augmentation.

Proof. As in the proof for Theorem 5.2, we begin with the following inequality

$$E[|kL(W^t; X_{c+c_{aug}})k^2] \leq \frac{1}{2} L(W^0; X_{c+c_{aug}}) \quad (53)$$

$$\leq \frac{1}{2} |kL(W^0; X_{c+c_{aug}})k^2 \quad (54)$$

Thus, we can write

$$E[|kL(W^0; X_{c+c_{aug}})k] \leq \frac{E[|kL(W^t; X_{c+c_{aug}})k^2]}{2} \quad (55)$$

$$\leq \frac{1}{2} |kL(W^0; X_{c+c_{aug}})k \quad (56)$$

$$\leq \frac{1}{2} |kL(W^0; X_c)k + |kL(W^0; X_{c_{aug}})k \quad (57)$$

$$\leq \frac{1}{2} G_0 + |k(J(W^0; X_c) + E)(f(W^0; X_{c+c_{aug}}))k \quad (58)$$

$$\leq \frac{1}{2} G_0 + |k(J(W^0; X_c) + E)(f(W^0; X_{c+c_{aug}}))k \quad (59)$$

The rest of the proof is similar to that of Theorem 5.2. \square

B.4 Lemma for eigenvalues of coreset

The following Lemma characterizes the sum of eigenvalues of the NTK associated with the coreset.

Lemma B.3. Let ϵ be an upper bound of the normed difference in gradient of the weighted coreset and the original dataset, i.e. for full data X and its corresponding coreset X_s with weights s , and respective residuals r_s , we have the bound $|kL(W^t; X)r^t - sJ^T(W^t; X_s)r_s^t| \leq \epsilon$. Let λ_i be the eigenvalues of the NTK associated with the coreset. Then we have that

$$\sum_{i=1}^V \lambda_i \leq \frac{|kL(W^t; X)r^t - sJ^T(W^t; X_s)r_s^t|}{|kr_s^t|}.$$

Proof. Let singular values of coreset Jacobian be λ_i . Let $J^T(W^t; X)r^t = sJ^T(W^t; X_s)r_s^t + \delta$ where $|\delta| \leq \epsilon$.

Taking Frobenius norm, we get

$$|kL(W^t; X)r^t - sJ^T(W^t; X_s)r_s^t| \leq \epsilon \quad (60)$$

$$\implies |kL(W^t; X)r^t - sJ^T(W^t; X_s)r_s^t| \leq \epsilon \quad (61)$$

$$\implies |kL(W^t; X)r^t - sJ^T(W^t; X_s)r_s^t| \leq \epsilon \quad (62)$$

$$\implies \sum_{i=1}^V \lambda_i \leq \frac{|kL(W^t; X)r^t - sJ^T(W^t; X_s)r_s^t|}{|kr_s^t|} \quad (63)$$

$$\implies \sum_{i=1}^V \lambda_i \leq \frac{|kL(W^t; X)r^t - sJ^T(W^t; X_s)r_s^t|}{|kr_s^t|} \quad (64)$$

$$\implies \sum_{i=1}^V \lambda_i \leq \frac{|kL(W^t; X)r^t - sJ^T(W^t; X_s)r_s^t|}{|kr_s^t|} \quad \text{by reverse triangle inequality} \quad (65)$$

\square

We can make the following observations: For overparameterized networks, with bounded activation functions and labels, e.g. softmax and one-hot encoding, the norm of the residual vector is bounded,

and the gradient norm is likely to be much larger than residual, especially when dimension of gradient is large. In this case, the Jacobian matrix associated with small weighted coresets found by solving Eq. (9), have large singular values.

B.5 Augmentation as Linear Transformation: Linear Model Analysis

We introduce a simplified linear model to extend our theoretical analysis to augmentations modelled as linear transformation matrices applied to the original training data. These augmentations are also originally studied by [7]. In this section, we specially study the effect of these augmentations using a linear model when applied to coresets.

Lemma B.4 (Augmented coreset gradient bounds: Linear) Let f be a simple linear model with weights $W \in \mathbb{R}^{d \times C}$ where $f(W; x_i) = W^T x_i$, trained on mean squared loss function. Let $F \in \mathbb{R}^{d \times d}$ be a common linear augmentation matrix with norm $\|F\|_F = k$ with augmentation x_i^{aug} given by $F x_i$. Let coreset be of size v and full dataset be of size n . Further assume that the predicted label of x_i and its augmentation $F x_i$ are sufficiently close, i.e. there exists δ_i such that $\|W^T(F x_i) - W^T x_i\| \leq \delta_i$. Let ρ upper-bound the normed difference in gradients between the weighted coreset and full dataset. Then, the normed difference between the gradient of the augmented full data and augmented coreset is given by

$$\sum_{i=1}^v \rho L(W; x_i^{aug}) - \sum_{j=1}^k s_j \rho L(W; x_{s_j}^{aug}) \leq k \|F\|_F \left(\sum_{j=1}^P \delta_j \right)$$

for some (small) constant

Proof. By our assumption, we can begin with,

$$\sum_{i=1}^v \rho L(W; x_i) - \sum_{j=1}^k s_j \rho L(W; x_{s_j}) \tag{66}$$

Furthermore, by [6], we know that sum of the coreset weights is given by $\sum_{j=1}^P s_j = n$.

Hence,

$$\sum_{i=1}^v \rho L(W; x_i^{aug}) - \sum_{j=1}^k s_j \rho L(W; x_{s_j}^{aug}) \tag{67}$$

$$= \sum_{i=1}^v \rho (J(W; x_i^{aug}))^T [W^T(F x_i) - y_i] - \sum_{j=1}^k s_j \rho (J(W; x_{s_j}^{aug}))^T [W^T(F x_{s_j}) - y_{s_j}] \tag{68}$$

$$= \sum_{i=1}^v \rho F x_i [W^T(F x_i) - y_i] - \sum_{j=1}^k s_j \rho F x_{s_j} [W^T(F x_{s_j}) - y_{s_j}] \tag{69}$$

$$= k \sum_{i=1}^v \rho x_i (W^T x_i - y_i) - F \sum_{j=1}^k s_j \rho x_{s_j} (W^T x_{s_j} + z_j - y_{s_j}) \tag{70}$$

$$= k \sum_{i=1}^v \rho L(W; x_i) - F \sum_{j=1}^k s_j \rho L(W; x_{s_j}) - F \sum_{j=1}^k s_j \rho x_{s_j} z_{s_j} \tag{71}$$

$$= k \|F\|_F \sum_{i=1}^v \rho L(W; x_i) - \sum_{j=1}^k s_j \rho L(W; x_{s_j}) + k \sum_{j=1}^k s_j \rho x_{s_j} z_{s_j} \tag{72}$$

$$= k \|F\|_F + \rho \|F\|_F n \tag{73}$$

$$= k \|F\|_F \left(\sum_{j=1}^P \delta_j \right) \tag{74}$$

□

Corollary B.5. In the simplified linear case above, the difference in gradients of the full training data with its augmentations $\| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{f+aug}) - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{c+caug}) \|$ and gradients of the coreset with its augmentations $\| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{c+caug}) - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{c+caug}) \|$ can be bounded by

$$\| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{f+aug}) - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{c+caug}) \| \leq (kFk+1) + \frac{p}{d} kFkn!$$

Proof. Applying Eq. (66) and Lemma B.4, we obtain

$$\| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{f+aug}) - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{c+caug}) \| \leq \tag{75}$$

$$= \| (\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_f) + \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{aug})) - (\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_c) + \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{caug})) \| \tag{76}$$

$$= \| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_f) - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_c) \| + \| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{aug}) - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{caug}) \| \tag{77}$$

$$\leq \| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_f) - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_c) \| + \| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{aug}) - \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}_{caug}) \| \tag{78}$$

$$\leq (kFk+1) + \frac{p}{d} kFkn! \tag{79}$$

$$= (kFk+1) + \frac{p}{d} kFkn! \tag{80}$$

□

Theorem B.6 (Convergence of linear model). Let \mathcal{L} be a linear model with weight \mathbf{W} and augmentation be represented by the common linear transformation \mathbf{A} . Let \mathcal{L} be μ -smooth and satisfy the μ -PL condition, that is for $\epsilon > 0$, $\| \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{X}) \|^2 \geq \mu \mathcal{L}(\mathbf{W}; \mathbf{X})$ for all weights \mathbf{W} . Let ϵ upper-bound the normed difference in gradients between the weighted coreset and full dataset and bound $\mathbf{W}^T (\mathbf{F} \mathbf{x}_i) = \mathbf{W}^T \mathbf{x}_i + z_i$, $\| z_i \| \leq \epsilon$. Let \mathbf{G}_0^0 be the gradient over the full dataset and its augmentations at initialization. Then, running SGD on the coreset with its augmentation using constant step size $\frac{\epsilon}{2}$, we get the following convergence bound:

$$\mathbb{E}[\| \nabla_{\mathbf{W}^t} \mathcal{L}(\mathbf{W}^t; \mathbf{X}_{c+caug}) \|^2] \leq \frac{1}{2} \left(1 - \frac{\epsilon}{2} \right)^{\frac{t}{2}} \mathbf{G}_0^0 + (kFk+1) + \frac{p}{d} kFkn!$$

Proof. From [2], we have

$$\mathbb{E}[\| \nabla_{\mathbf{W}^t} \mathcal{L}(\mathbf{W}^t; \mathbf{X}_{c+caug}) \|^2] \leq \frac{1}{2} \left(1 - \frac{\epsilon}{2} \right)^{\frac{t}{2}} \mathcal{L}(\mathbf{W}^0; \mathbf{X}_{c+caug}) \tag{81}$$

$$\leq \frac{1}{2} \left(1 - \frac{\epsilon}{2} \right)^{\frac{t}{2}} \| \nabla_{\mathbf{W}^0} \mathcal{L}(\mathbf{W}^0; \mathbf{X}_{c+caug}) \|^2 \tag{82}$$

$$\tag{83}$$

Using Jensen's inequality, we have

$$\mathbb{E}[\| \nabla_{\mathbf{W}^t} \mathcal{L}(\mathbf{W}^t; \mathbf{X}_{c+caug}) \|^2] \tag{84}$$

$$\leq \frac{1}{2} \left(1 - \frac{\epsilon}{2} \right)^{\frac{t}{2}} \mathbb{E}[\| \nabla_{\mathbf{W}^0} \mathcal{L}(\mathbf{W}^0; \mathbf{X}_{c+caug}) \|^2] \tag{85}$$

$$\leq \frac{1}{2} \left(1 - \frac{\epsilon}{2} \right)^{\frac{t}{2}} \| \nabla_{\mathbf{W}^0} \mathcal{L}(\mathbf{W}^0; \mathbf{X}_{c+caug}) \|^2 \tag{86}$$

$$\leq \frac{1}{2} \left(1 - \frac{\epsilon}{2} \right)^{\frac{t}{2}} \mathbf{G}_0^0 + (kFk+1) + \frac{p}{d} kFkn! \tag{87}$$

where the last inequality follows from applying Corollary B.5. □

C Singular spectrum analysis

C.1 Experiment details

We generate singular spectrum plots for both MNIST and CIFAR10 datasets in Figure 1. Due to the computational infeasibility of computing the network Jacobian for the full datasets in deep network settings, we instead construct and use a reduced version of these datasets by uniformly select 900 images from the first 3 classes. For our experiments on MNIST, we pretrain a MLP model with 1

hidden layer for 15 epochs. For our experiments on CIFAR10, we pretrain a ResNet20 model for 15 epochs. We then compute the singular spectrums for augmented and non-augmented data based on these pretrained networks.

Since it is difficult to perform a one-to-one matching of singular values produced from augmented and non-augmented datasets, we instead bin our singular values into 30 separate and uniformly distributed bins each containing the same number of singular values. To measure perturbation to singular values resulted from augmentation, we compute the mean difference between each bin. On the other hand, to measure perturbation to singular vectors, we compute mean subspace angle between the singular subspace spanned by singular vectors in each bin.

C.2 Real-world strong augmentations

We study the effects of real-world, unbounded augmentations on the singular spectrum of the network Jacobian. In particular, in addition to the plots in the main paper, we show the effect of strong augmentations through (1) random rotation (up to 30°) and AutoAugment [3] for MNIST and (2) random horizontal flips/random crops and AutoAugment for CIFAR10. The policies implemented by AutoAugment include translations, shearing, as well as contrast and brightness transforms. We study the effects of these augmentations on the singular spectrum in Figure 1. Despite these augmentations being unbounded transformations, we note that the results of our theory still holds. In particular, it can be observed that data augmentation increases smaller singular values relatively more with a higher probability. On the other hand, data augmentation affects the prominent singular vectors of the Jacobian to a smaller extent, and preserves the prominent directions. As such, our argument empirically extends to real-world, unbounded label-invariant transformations characteristic of strong augmentations.

D Experiment Setup and Additional Experiments

D.1 Experiment setup

For all experiments, we train using SGD with 0.9 momentum and learning rate decay. For experiments on CIFAR10 and variants/ResNet20, we train for 200 epochs, for Caltech256 (ImageNet pretrained)/ResNet18, we trained for 40 epochs starting at learning rate of 0.1 and batch size 64. We also report results for Caltech256 without ImageNet pretraining in Sec. D.8, where we train for 400 epochs to ensure convergence with a starting learning rate of 0.1 and batch size 64. For experiments on ImageNet/ResNet50 and TinyImageNet/ResNet50, we use the standard 90 epoch learning schedule starting at learning rate of 0.1 and batch size 64.

Data and augmentation. We apply our method to training ResNet20 and Wide-ResNet-28-10 on CIFAR10, and ResNet32 on CIFAR10-IMB (Long-Tailed CIFAR10 with Imbalance factor of 100 following [5]) and SVHN datasets. We train Caltech256 on ImageNet-pretrained ResNet18, and include experiments with random initialization in Appendix D. TinyImageNet and ImageNet are trained on with ResNet50. We use [3] for CIFAR10/SVHN, and AutoAugment [3] for Caltech256, TinyImageNet, and ImageNet as the strong augmentation method. Note that we append strong augmentations rather than apply them in-place, which we show to be more effective in Appendix D. All results are averaged over 5 runs using an Nvidia A40 GPU.

D.2 Experiment setup

For all experiments, we train using SGD with 0.9 momentum and learning rate decay. We also set weight decay as For experiments on CIFAR10 and variants/ResNet20, we train for 200 epochs, for Caltech256 (ImageNet pretrained)/ResNet18, we trained for 40 epochs starting at learning rate of 0.1 and batch size 64. We also report results for Caltech256 without ImageNet pretraining in Sec. D.8, where we train for 400 epochs to ensure convergence with a starting learning rate of 0.1 and batch size 64. For experiments on ImageNet/ResNet50 and TinyImageNet/ResNet50, we use the standard 90 epoch learning schedule starting at learning rate of 0.1 and batch size 64.

Data and augmentation. We apply our method to training ResNet20 and Wide-ResNet-28-10 on CIFAR10, and ResNet32 on CIFAR10-IMB (Long-Tailed CIFAR10 with Imbalance factor of 100 following [5]) and SVHN datasets. We train Caltech256 on ImageNet-pretrained ResNet18, and

(a) MNIST $\sigma_0 = 8$ - Values (b) MNIST $\sigma_0 = 8$ - Vectors (c) CIFAR10 $\sigma_0 = 8$ - Values (d) CIFAR10 $\sigma_0 = 8$ - Vectors

(e) MNIST $\sigma_0 = 16$ - Values (f) MNIST $\sigma_0 = 16$ - Vectors (g) CIFAR10 $\sigma_0 = 16$ - Values (h) CIFAR10 $\sigma_0 = 16$ - Vectors

(i) MNIST Rotate - Values (j) MNIST Rotate - Vectors (k) CIFAR10 Flip-Crop - Values (l) CIFAR10 Flip-Crop - Vectors

(m) MNIST AutoAugment - Values (n) MNIST AutoAugment - Vectors (o) CIFAR10 AutoAugment - Values (p) CIFAR10 AutoAugment - Vectors

(q) MNIST Rotate + AutoAugment - Values (r) MNIST Rotate + AutoAugment - Vectors (s) CIFAR10 Flip + Crop + AutoAugment - Values (t) CIFAR10 Flip + Crop + AutoAugment - Vectors

Figure 1: Difference in mean singular values (Cols 1 & 3) between augmented and non-augmented data and mean angular difference (Cols 2 & 4) between subspaces spanned by singular vectors for augmented and non-augmented data.

include experiments with random initialization in Appendix D. TinyImageNet and ImageNet are trained on with ResNet50. We use `AutoAugment` for CIFAR10/SVHN, and `AutoAugment` for Caltech256, TinyImageNet, and ImageNet as the strong augmentation method. Note that we append strong augmentations rather than apply them in-place, which we show to be more effective in Appendix D. All results are averaged over 5 runs using an Nvidia A40 GPU.

D.3 Full Results for Table 4

This section contains full experiment results including standard deviations and the full augmentation benchmark for Table 4. Augmenting coresets of size 10% achieves 51% of the improvement obtained from augmentation of the full data and further enjoys a 6x speedup in training time on CIFAR10.

This speedup becomes more significant when using strong augmentation techniques which are mostly computationally demanding, especially when applied to the entire dataset.

Table 1: Supplementary table for Table 4 - Test accuracy on CIFAR10 + ResNet20, SVHN + ResNet32, CIFAR10-Imbalanced + ResNet32 including standard deviation errors and full dataset augmentation accuracy.

Method	Size	CIFAR10		CIFAR10-IMB		SVHN	
None	0%	89:46	0:17%	87:08	0:50%	95:676	0:108%
Random	5%	90:34	0:18%	88:48	0:25%	95:760	0:084%
	10%	91:07	0:13%	89:52	0:15%	96:187	0:112%
	30%	92:11	0:12%	91:11	0:18%	96:569	0:073%
Max-Loss	5%	90:79	0:19%	88:77	0:35%	96:165	0:108%
	10%	91:39	0:08%	89:22	0:48%	96:370	0:076%
	30%	92:43	0:07%	91:11	0:25%	96:735	0:068%
Coreset	5%	90:87	0:05%	89:10	0:41%	96:121	0:055%
	10%	91:54	0:19%	89:75	0:52%	96:354	0:091%
	30%	92:49	0:15%	91:12	0:26%	96:791	0:051%
All	100%	93:50	0:25%	92:48	0:34%	97:068	0:030%

D.4 Supplementary results for Table 1

Table 2: Supplementary results for Tab. 1. Training ResNet20 (R20) and WideResnet-28-10 (W2810) on CIFAR10 (C10) using small subsets, and ResNet18 (R18) on Caltech256 (Cal).

Model/Dataset	Subset	Random				Ours			
		Weak Aug.		Strong Aug.		Weak Aug.		Strong Aug.	
C10/R20	0.1% (5)	31:7	3:2	33:5	2:7	29:6	3:8	37:8	4:5
	0.2% (10)	35:9	2:1	42:7	3:9	33:6	3:2	45:1	2:3
	0.5% (25)	51:1	2:3	58:7	1:3	55:8	3:1	63:9	2:1
	1% (50)	66:2	1:0	74:4	0:8	65:9	4:0	74:7	1:1
C10/W2810	1% (50)	61:3	2:4	57:7	0:8	59:9	2:4	62:1	3:1
Cal/R18	5% (3)	24:8	1:5	41:5	0:5	33:8	1:7	52:7	1:2
	10% (6)	49:5	0:6	61:8	0:8	55:7	0:3	65:4	0:3
	20% (12)	66:6	0:2	72:5	0:1	67:5	0:3	73:1	0:1
	30% (18)	72:0	0:1	75:7	0:2	71:9	0:2	76:3	0:2
	40% (24)	74:6	0:3	77:6	0:4	74:2	0:4	77:7	0:5
	50% (30)	76:1	0:5	78:5	0:3	76:1	0:1	78:9	0:2

D.5 Training dynamics vs generalization

Figure 2 demonstrates the relationship between training loss and validation accuracy resulted from data augmentation. While training loss of augmented datasets do not decrease as quickly as non-augmented datasets, generalization performance (shown by val. acc.) improves.

Figure 2: Training loss vs validation accuracy of CIFAR10/ResNet20 using AutoAugment.

D.6 Augmentations applied through appending vs in-place

Our experiments on Caltech256/ResNet18/AutoAugment (R=5) show that even for cheaper strong augmentation methods (AutoAugment), while in-place augmentation may decrease the performance, appending Random (R) and Coresets (C) augmentations (Append) outperforms in-place augmentation of 2x data points (In-place 2x) for various subset sizes.

Table 3: Caltech256/AutoAugment in-place vs. appending for Caltech256.

	No Aug.	In-place	In-place (2x)	Append
C5%	33.8%	26.4%	48.2%	52.7%
R5%	24.8%	17.4%	40.2%	41.5%
C10%	55.7%	48.2%	62.8%	65.4%
R10%	50.6%	40.2%	62.0%	61.8%
C30%	71.9%	68.8%	74.9%	76.3%
R30%	72.0%	68.7%	75.1%	75.7%

D.7 Speed-up measurements

We measure the improvement in training time in the case of training on full data and augmenting subsets of various sizes. While our method yields similar or slightly lower speed-up to the max-loss policy and random approach respectively, our resulting accuracy outperforms these two approaches. We show this in Fig. D.7. For example, for SVHN/Resnet32 using 10% coresets, we sacrifice 41% of the speed-up to obtain an additional 2.18% of the gain in accuracy from full data augmentation when compared to a random subset of the same size. We show the speed-up obtained for our method and various subset sizes in Tab. 4, and provide wall-clock times for our method in Tab. D.7.

Table 4: Speedup on CIFAR10 + ResNet20 (C10/R20), SVHN + ResNet32 (SVHN/R32).

Dataset	Full Aug. 100%	Ours						Max loss. 30%	Random. 30%
		5%	10%	15%	20%	25%	30%		
C10 / R20	1x	7:93x	6:31x	4:46x	4:27x	3:41x	3:43x	3:48x	4:03x
SVHN / R32	1x	5:35x	3:93x	3:40x	2:80x	2:49x	2:18x	2:21x	2:43x

(a) CIFAR10/ResNet20 (b) SVHN/Resnet32

Figure 3: Speedup/Accuracy of augmented 10% coresets compared to original max-loss policy for (a) ResNet20 trained on CIFAR10 and (b) ResNet32 trained on SVHN.

Table 5: Wall-clock times to find various sized coresets from all classes of Caltech256 and TinyImageNet at 1 epoch. Note, training ResNet20/CIFAR10 with [7] takes 14.4 hrs. In practice, coresets can be found in parallel (threads) from different classes, and selection happens every 15 epochs. Hence, the numbers divide by R.

Caltech256			TinyImageNet		
10%	30%	50%	10%	30%	50%
10.50s	10.52s	10.53s	7.85s	8.09s	8.17s

D.8 End to end training on Caltech256

As Caltech256 contains many classes and higher resolution images, training on smaller subset without pretraining has a low accuracy. Thus, many works (e.g. Achille et al., 2020) netune from ImageNet pretrained initialization. However, we show that our results still hold even when training from scratch. We demonstrate our results in Tab. 6, where we train Caltech256 on ResNet50 without pretraining for 400 epochs, and with $R = 40$, where our method consistently outperforms random subsets for multiple subset sizes (5%, 10%, 30%, 50%).

Table 6: Caltech256 (w/o pretraining) /ResNet50, 400 epochs, $R=40$

Random				Ours			
5%	10%	30%	50%	5%	10%	30%	50%
17.26	35.38	58.2	64.67	20.58	38.20	60.30	65.17

D.9 Training on full data and augmenting small subsets re-selected every epoch

We apply our proposed method to select a new subset for augmentation every epoch (i.e. using $R = 1$) and compare our results with other approaches using accuracy and percentage of data not selected (NS). We see that while the max-loss policy selects a small fraction of data points over and over and random uniformly selects all the data points, our approach successfully finds the smallest subset of data points that are the most crucial for data augmentation. Hence, it can achieve a superior accuracy than max-loss policy, while augmenting only slightly more examples. This confirms the data-efficiency of our approach. This is especially evident when using coresets of size 0.2%. Furthermore, despite the random baseline using a significantly larger percentage of data, it is outperformed by our approach in both data-efficiency and accuracy. We emphasize that results in this table is different from that of Table 1, as default augmentations on the full training data are performed once every $R = 1$ epochs instead of every $R = 20$ epochs. Since selecting subsets at every epoch can be computationally expensive, we only perform these experiments on small coresets and hence still enjoy good speedups compared to full data augmentation. This shows that our approach is still effective at very small subset sizes, hence can be computationally efficient even when subsets are re-selected every epoch.

Table 7: Training on full data and selecting a new subset for augmentation every epoch)

Subset	Random		Max-loss Policy				Ours					
	Acc	NS (%)	Acc	NS (%)	Acc	NS (%)						
0%	91:96	0:12	91:96	0:12	91:96	0:12						
0.2%	92:22	0:22	67:03	0:04	91:94	0:12	86:70	0:15	92:26	0:13	79:19	1:10
0.5%	92:06	0:17	36:70	0:18	92:20	0:13	76:80	0:31	92:27	0:08	63:23	0:35

D.10 Additional visualizations for training on coresets and its augmentations - Measuring training dynamics over time

We include additional visualizations in Figure 4 for training on coresets and its augmentations as supplementary plots to Figure 7(c) and Table 1. We plot metrics obtained during each point (epoch)

(a) CIFAR10 - 0.1%

(b) CIFAR10 - 0.2%

(c) CIFAR10 - 0.5%

(d) CIFAR10 - 1%

(e) CIFAR10 - 5%

Figure 4: Supplementary plots for Figure 7(c): Training on coreset and its augmentation compared to random baseline, measured using test accuracy against percentage of data used on CIFAR10 dataset across various subset sizes. Accuracy and percentage of data used are measured at every epoch and averaged over 5 runs.

of the training process based on percentage of data selected/used and test accuracy achieved. All metrics are averaged over 5 runs and obtained using avg . These plots demonstrate that coreset augmentation approaches outperform random augmentation baselines throughout the training process. Furthermore, they show that augmentation of coresets result in a larger increase in test accuracy compared to augmentation of randomly selected training examples, especially for small subset sizes.

Figure 5: Intersection between max-loss and coresets in the top points selected aggregated across the entire training process. Here, we show the increasing overlap between max-loss and coreset points as N grows.

Figure 6: Qualitative evaluation of coreset and max-loss points.

D.11 Intersection of max-loss policy and coresets

Figure 7(a) depicts the increase in intersection between max-loss subsets and coresets over time. In addition, we also aggregate 10% subsets selected every $\tau = 20$ epochs using both approaches over the entire training process to compute intersection between the selected data points. Our plots

