# RDM-DC: Poisoning Resilient Dataset Condensation with Robust Distribution Matching (Supplementary Material)

**Tianhang Zheng**[1]

**Baochun Li**[1]

[1]Department of Electrical and Computer Engineering, University of Toronto

## 1 OMITTED PROOF

**Lemma 1.1** *Assuming that $\mathcal{D}$ and $\mathcal{B}$ have bounded covariance matrices $\Sigma_{\mathcal{D}}, \Sigma_{\mathcal{B}} \leq \sigma^2 I$, and their means have an apparent difference,* i.e., $\|\mu_{\mathcal{D}} - \mu_{\mathcal{B}}\|_2^2 \geq \frac{\alpha\sigma^2}{\epsilon}$ *where $\alpha > \frac{2665}{576}$, then if we drop all the representations that satisfies $|\langle r - \mu_{\mathcal{P}}, v\rangle| \geq t$ with a certain $t$, then we can reduce the scale of the poisoned deviation from $O(\epsilon\sqrt{d_r})$ to $\Theta(\epsilon^2\sqrt{d_r})$.*

To prove the above lemma, we need the help of Chebyshev's inequality, which is introduced in the following.

**Lemma 1.2 (Chebyshev's inequality)** *Given a scalar random variable $X$, if $\mathbb{E}[X] = \mu$ and $Var[X] = \sigma^2$, then*

$$\mathbb{P}(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2} \tag{1}$$

Given Chebyshev's inequality, we have the following corollary, which will be used in the proof of Lemma 1.1.

**Corollary 1.1** *Given a multi-dimensional variable $X$, if $\mathbb{E}[X] = \mu$ and $Cov[X] \leq \sigma^2 I$, then for any unit vector $u$, we have*

$$\mathbb{P}(|\langle X - \mu, u\rangle| > t) \leq \frac{\sigma^2}{t^2} \tag{2}$$

**Proof** [Proof of Corollary 1.1]

Considering $\langle X, u\rangle$ as a scalar random variable, we have $\mathbb{E}[\langle X, u\rangle] = \langle \mu, u\rangle$ and,

$$Var[\langle X, u\rangle] = u^T Cov[X]u \leq \sigma^2. \tag{3}$$

With Chebyshev's inequality, we know that

$$\mathbb{P}(|\langle X, u\rangle - \langle \mu, u\rangle| \geq t) \leq \frac{Var[\langle X, u\rangle]}{t^2} \leq \frac{\sigma^2}{t^2} \tag{4}$$

Beyond Corollary 1.1, we also need to use the following lemma and corollary in the proof of Lemma 1.1.

**Lemma 1.3** *Given two distributions $P$ and $Q$ with mean $\mu_P$ and $\mu_Q$ and covariance matrices $\Sigma_P, \Sigma_Q \leq \sigma^2 I$, if $\|\mu_P - \mu_Q\|_2^2 \geq \frac{\alpha\sigma^2}{\epsilon}$, then $\langle v, \mu_P - \mu_Q\rangle^2 \geq \frac{\alpha\sigma^2 - \sigma^2/(1-\epsilon)}{\epsilon}$ where $v$ is the first eigenvector of the covariance matrix of $(1-\epsilon)P + \epsilon Q$.*

**Proof** [Proof of Lemma 1.3] The mean of the mixture $(1-\epsilon)P + \epsilon Q$ is $(1-\epsilon)\mu_P + \epsilon\mu_Q$, which is denoted by $\mu_M$. We denote $\mu_P - \mu_Q$ by $\delta$. The covariance matrix of $(1-\epsilon)P + \epsilon Q$ can be expressed as

$$\mathbb{E}_{X \sim (1-\epsilon)P+\epsilon Q}[(X - \mu_M)(X - \mu_M)^T]$$
$$= (1-\epsilon)\mathbb{E}_{X \sim P}[(X - \mu_M)(X - \mu_M)^T]$$
$$+ \epsilon\mathbb{E}_{X \sim Q}[(X - \mu_M)(X - \mu_M)^T] \tag{5}$$

Since we have

$$\mathbb{E}_{X \sim P}[(X - \mu_M)(X - \mu_M)^T]$$
$$= \mathbb{E}_{X \sim P}[(X - \mu_P + \epsilon\delta)(X - \mu_P + \epsilon\delta)^T]$$
$$= \Sigma_P + \epsilon^2\delta\delta^T$$
$$\mathbb{E}_{X \sim Q}[(X - \mu_M)(X - \mu_M)^T]$$
$$= \mathbb{E}_{X \sim Q}[(X - \mu_Q - (1-\epsilon)\delta)^T)(X - \mu_Q - (1-\epsilon)\delta)^T]$$
$$= \Sigma_Q + (1-\epsilon)^2\delta\delta^T ,$$

we have a lower bound on the covariance matrix of the mixture $(1-\epsilon)P + \epsilon Q$ as

$$\Sigma_M = \mathbb{E}_{X \sim (1-\epsilon)P+\epsilon Q}[(X - \mu_M)(X - \mu_M)^T]$$
$$= (1-\epsilon)\Sigma_P + \epsilon\Sigma_Q + \epsilon(1-\epsilon)\delta\delta^T \geq \epsilon(1-\epsilon)\delta\delta^T. \tag{6}$$

Suppose that $v$ is the first eigenvector of $\Sigma_M$ and $u = \frac{\delta}{\|\delta\|_2}$, we then have

$$v^T \Sigma_M v \geq u^T \Sigma_M u \geq \epsilon(1-\epsilon)u^T \delta\delta^T u = \epsilon(1-\epsilon)\|\delta\|_2^2. \tag{7}$$

Since $\Sigma_P, \Sigma_Q \leq \sigma^2 I$, we also have

$$v^T \Sigma_M v = (1-\epsilon)v^T \Sigma_P v + \epsilon v^T \Sigma_Q v + \epsilon(1-\epsilon)v^T \delta\delta^T v$$
$$\leq \sigma^2 + \epsilon(1-\epsilon)\langle v, \delta \rangle^2 \tag{8}$$

Thus, we have

$$\langle v, \delta \rangle^2 \geq \frac{v^T \Sigma_M v - \sigma^2}{\epsilon(1-\epsilon)} \geq \|\delta\|_2^2 - \frac{\sigma^2}{\epsilon(1-\epsilon)} \tag{9}$$

Given the assumption that $\|\delta\|_2^2 \geq \frac{\alpha\sigma^2}{\epsilon}$,

$$\langle v, \delta \rangle^2 \geq \frac{\alpha\sigma^2 - \sigma^2/(1-\epsilon)}{\epsilon} \tag{10}$$

∎

Based on Lemma 1.3, we have the following corollary.

**Corollary 1.2** *Given the definitions and conditions in Lemma 1.3, if $\epsilon \leq \frac{1}{10}$ and $\alpha > \frac{2665}{576}$, then we have $(1-2\epsilon)|\langle \delta, v \rangle| > \frac{3\sigma}{2\sqrt{\epsilon}}$.*

**Proof** Given Lemma 1.3, we have

$$(1-2\epsilon)|\langle \delta, v \rangle| \geq (1-2\epsilon)\sqrt{\alpha - \frac{1}{1-\epsilon}}\frac{\sigma}{\sqrt{\epsilon}} \tag{11}$$

Since $1-2\epsilon$ and $-\frac{1}{1-\epsilon}$ are decreasing functions w.r.t. $\epsilon$, they achieve the minimum at $\epsilon = \frac{1}{10}$. Thus, we have

$$(1-2\epsilon)|\langle \delta, v \rangle| \geq \frac{4}{5}\sqrt{\alpha - \frac{10}{9}}\frac{\sigma}{\sqrt{\epsilon}}. \tag{12}$$

So if $\alpha > \frac{2665}{576}$, we have $(1-2\epsilon)|\langle \delta, v \rangle| > \frac{3\sigma}{2\sqrt{\epsilon}}$. ∎

**Proof** [Proof of Lemma 1.1] The mean of the poisoned representation distribution $\mathcal{P}$ is $\mu_\mathcal{P} = (1-\epsilon)\mu_\mathcal{D} + \epsilon\mu_\mathcal{B}$. Let $\delta = \mu_\mathcal{B} - \mu_\mathcal{D}$ and $t = |\epsilon\langle \delta, v \rangle| + \frac{\sigma}{\sqrt{\epsilon}}$. We denote the covariance matrix of $\mathcal{P}$ by $\Sigma_\mathcal{P}$ and its first eigenvector by $v$.

For the original representation distribution, we have

$$\mathbb{P}_{r \sim \mathcal{D}}[|\langle r - \mu_\mathcal{P}, v \rangle| > t]$$
$$= \mathbb{P}_{r \sim \mathcal{D}}[|\langle r - \mu_\mathcal{D}, v \rangle - \epsilon\langle \delta, v \rangle| > t] \ \textcircled{1}$$
$$\leq \mathbb{P}_{r \sim \mathcal{D}}[|\langle r - \mu_\mathcal{D}, v \rangle| > \frac{\sigma}{\sqrt{\epsilon}}] \ \textcircled{2}$$
$$\leq \epsilon \ \textcircled{3} \tag{13}$$

$\textcircled{1}$ is because $\mu_\mathcal{P} = \mu_\mathcal{D} + \epsilon\delta$. $\textcircled{2}$ is because if $|\langle r - \mu_\mathcal{D}, v \rangle - \epsilon\langle \delta, v \rangle| > t$, then either $\langle r - \mu_\mathcal{D} \rangle > t + \epsilon\langle \delta, v \rangle > \frac{\sigma}{\sqrt{\epsilon}}$ or $\langle r - \mu_\mathcal{D} \rangle < -t + \epsilon\langle \delta, v \rangle < -\frac{\sigma}{\sqrt{\epsilon}}$ holds true. Thus, we have $|\langle r - \mu_\mathcal{D} \rangle| > \frac{\sigma}{\sqrt{\epsilon}}$, and $\{r, |\langle r - \mu_\mathcal{D}, v \rangle - \epsilon\langle \delta, v \rangle| > t\} \subseteq \{r, |\langle r - \mu_\mathcal{D}, v \rangle| > \frac{\sigma}{\sqrt{\epsilon}}\}$. Therefore, $\textcircled{2}$ holds true. $\textcircled{3}$ is because of Corollary 1.1.

For the poisoned distribution, we have

$$\mathbb{P}_{r \sim \mathcal{B}}[|\langle r - \mu_\mathcal{P}, v \rangle| < t]$$
$$= \mathbb{P}_{r \sim \mathcal{B}}[|\langle r - \mu_\mathcal{B}, v \rangle + (1-\epsilon)\langle \delta, v \rangle| < t] \ \textcircled{1}$$
$$\leq \mathbb{P}_{r \sim \mathcal{B}}[|\langle r - \mu_\mathcal{B}, v \rangle| > (1-2\epsilon)|\langle \delta, v \rangle| - \frac{\sigma}{\sqrt{\epsilon}}] \ \textcircled{2}$$
$$\leq \mathbb{P}_{r \sim \mathcal{B}}[|\langle r - \mu_\mathcal{B}, v \rangle| > \frac{\sigma}{2\sqrt{\epsilon}}] \leq 4\epsilon \ \textcircled{3} \tag{14}$$

$\textcircled{1}$ is because $\mu_\mathcal{P} = \mu_\mathcal{B} - (1-\epsilon)\delta$. In the following, we prove $\textcircled{2}$: Given $|\langle r - \mu_\mathcal{B}, v \rangle + (1-\epsilon)\langle \delta, v \rangle| < t$, we have $-t - (1-\epsilon)\langle \delta, v \rangle < \langle r - \mu_\mathcal{B}, v \rangle < t - (1-\epsilon)\langle \delta, v \rangle$. Since $t = |\epsilon\langle \delta, v \rangle| + \frac{\sigma}{\sqrt{\epsilon}}$, $-|\epsilon\langle \delta, v \rangle| - \frac{\sigma}{\sqrt{\epsilon}} - (1-\epsilon)\langle \delta, v \rangle < \langle r - \mu_\mathcal{B}, v \rangle < |\epsilon\langle \delta, v \rangle| + \frac{\sigma}{\sqrt{\epsilon}} - (1-\epsilon)\langle \delta, v \rangle$.

Then, we consider two cases: If $\langle \delta, v \rangle \geq 0$, we have $\langle r - \mu_\mathcal{B}, v \rangle < \frac{\sigma}{\sqrt{\epsilon}} - (1-2\epsilon)|\langle \delta, v \rangle|$. Given Corollary 1.2, we have $|\langle r - \mu_\mathcal{B}, v \rangle| > (1-2\epsilon)|\langle \delta, v \rangle| - \frac{\sigma}{\sqrt{\epsilon}}$. If $\langle \delta, v \rangle < 0$, we have $\langle r - \mu_\mathcal{B}, v \rangle > (1-2\epsilon)|\langle \delta, v \rangle| - \frac{\sigma}{\sqrt{\epsilon}}$. Given Corollary 1.2, we also have $|\langle r - \mu_\mathcal{B}, v \rangle| > (1-2\epsilon)|\langle \delta, v \rangle| - \frac{\sigma}{\sqrt{\epsilon}}$. Therefore, $\textcircled{2}$ holds true. $\textcircled{3}$ is because of Corollary 1.2.

Suppose after filtering out the data that satisfies $|\langle r - \mu_\mathcal{P}, v \rangle| \geq t$, the remaining deviation caused by $\mathcal{B}$ is expected to be

$$|\epsilon\mathbb{E}_{r \sim \mathcal{B}, |\langle r - \mu_\mathcal{P}, v \rangle| < t}[r]| < \epsilon t \mathbb{P}_{r \sim \mathcal{B}}[|\langle r - \mu_\mathcal{P}, v \rangle| < t]$$
$$\leq 4\epsilon^2 t = 4\epsilon^2(|\epsilon\langle \delta, v \rangle| + \frac{\sigma}{\sqrt{\epsilon}}) \tag{15}$$

Since $\frac{\sigma}{\sqrt{\epsilon}} \leq \frac{2}{3}|\epsilon\langle \delta, v \rangle|$ according to Corollary 1.2, we have

$$|\epsilon\mathbb{E}_{r \sim \mathcal{B}, |\langle r - \mu_\mathcal{P}, v \rangle| < t}[r]| \leq \frac{20}{3}\epsilon^3|\langle \delta, v \rangle| \leq \frac{20}{3}\epsilon^3\|\delta\|_2 \tag{16}$$

Since $\epsilon \leq \frac{1}{10}$ and $\|\delta\|_2 \sim \Theta(\sqrt{d_r})$, we have

$$|\epsilon\mathbb{E}_{r \sim \mathcal{B}, |\langle r - \mu_\mathcal{P}, v \rangle| < t}[r]| \sim \Theta(\epsilon^2 \sqrt{d_r}). \tag{17}$$

∎