

1 A Appendix of Proofs

2 A.1 Proof of Thm.3.2

3 **Theorem 3.2.** *By choosing KL divergence $D_{KL}(Q||Q_0) = \int Q \log \frac{Q}{Q_0} dx$, optimizing CL-DRO (cf. Eqn. 3) is equivalent to optimizing CL (InfoNCE, cf. Eqn. 1):*

$$\begin{aligned} \mathcal{L}_{CL-DRO}^{KL} &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \max_{Q \in \mathcal{Q}} \{ \mathbb{E}_Q[f_\theta] - \alpha [D_{KL}(Q||Q_0) - \eta] + \eta_1 (\mathbb{E}_{Q_0}[\frac{Q}{Q_0}] - 1) \} \\ &= -\mathbb{E}_{P_0} \left[\alpha^* \log \frac{e^{f_\theta/\alpha^*}}{\mathbb{E}_{Q_0}[e^{f_\theta/\alpha^*}]} \right] + Constant = \alpha^* \mathcal{L}_{InfoNCE} + Constant \end{aligned} \quad (4)$$

5 where α, η_1 represent the Lagrange multipliers, and the optimal α^* finally serves as the temperature
6 τ in CL.

7 *Proof.* To complete the proof, we start with giving some important notations and theorem.

8 **Definition A.1** (ϕ -divergence [13]). For any convex function ϕ with $\phi(1) = 0$, the ϕ -divergence
9 between Q and Q_0 is:

$$D_\phi(Q||Q_0) := \mathbb{E}_{Q_0}[\phi(dQ/dQ_0)] \quad (13)$$

10 where $D_\phi(Q||Q_0) = \infty$ if P is not absolutely continuous with respect to Q_0 . Specially, when
11 $\phi(x) = x \log x - x + 1$, ϕ -divergence degenerates to the well-known KL divergence.

12 **Definition A.2** (Convex conjugate [9]). We consider a pair (A, B) of topological vector spaces and a
13 bilinear form $\langle \cdot, \cdot \rangle \rightarrow \mathbb{R}$ such that $(A, B, \langle \cdot, \cdot \rangle)$ form a dual pair. For a convex function $f : \mathbb{R} \rightarrow \mathbb{R}$,
14 $dom f := \{x \in \mathbb{R} : f(x) < \infty\}$ is the effective domain of f . The convex conjugate, also known as
15 the Legendre-Fenchel transform, of $f : A \rightarrow \mathbb{R}$ is the function $f^* : B \rightarrow \mathbb{R}$ defined as

$$f^*(b) = \sup_a \{ab - f(a)\}, \quad b \in B \quad (14)$$

16 **Theorem A.3** (Interchange of minimization and integration [2]). *Let (Ω, \mathcal{F}) be a measurable space
17 equipped with σ -algebra \mathcal{F} , $L^p(\Omega, \mathcal{F}, P)$ be the linear space of measurable real valued functions
18 $f : \Omega \rightarrow \mathbb{R}$ with $\|f\|_p < \infty$, and let $\mathcal{X} := L^p(\Omega, \mathcal{F}, P)$, $p \in [1, +\infty]$. Let $g : \mathbb{R} \times \Omega \rightarrow \mathbb{R}$ be a
19 normal integrand, and define on \mathcal{X} . Then,*

$$\min_{x \in \mathcal{X}} \int_{\Omega} g(x(\omega), \omega) dP(\omega) = \int_{\Omega} \min_{s \in \mathbb{R}} g(s, \omega) dP(\omega) \quad (15)$$

20 To ease the derivation, we denote the likelihood ratio $L(x, y) = Q(x, y)/Q_0(x, y)$. Note that the
21 ϕ -divergence between Q and Q_0 is constrained, and thus $L(\cdot)$ is fine definition. For brevity, we
22 usually short $L(x, y)$ as L . And in terms of definition A.1 of ϕ -divergence, the expression of CL-DRO
23 becomes:

$$\mathcal{L}_{CL-DRO}^\phi = -\mathbb{E}_{P_0}[f_\theta] + \max_L \mathbb{E}_{Q_0}[f_\theta L] \quad s.t. \quad \mathbb{E}_{Q_0}[\phi(L)] \leq \eta \quad (16)$$

24 Note that $\mathbb{E}_{Q_0}[f_\theta L]$ and $\mathbb{E}_{Q_0}[\phi(L)]$ are both convex in L . We use the Lagrangian function solver:

$$\begin{aligned} \mathcal{L}_{CL-DRO}^\phi &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \max_L \{ \mathbb{E}_{Q_0}[f_\theta L] - \alpha [\mathbb{E}_{Q_0}[\phi(L)] - \eta] + \eta_1 (\mathbb{E}_{Q_0}[L] - 1) \} \\ &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha \eta - \eta_1 + \alpha \max_L \left\{ \mathbb{E}_{Q_0} \left[\frac{f_\theta + \eta_1}{\alpha} L - \phi(L) \right] \right\} \right\} \\ &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha \eta - \eta_1 + \alpha \mathbb{E}_{Q_0} \left[\max_L \left\{ \frac{f_\theta + \eta_1}{\alpha} L - \phi(L) \right\} \right] \right\} \\ &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha \eta - \eta_1 + \alpha \mathbb{E}_{Q_0} \left[\phi^* \left(\frac{f_\theta + \eta_1}{\alpha} \right) \right] \right\} \end{aligned} \quad (17)$$

25 The first equality holds due to the strong duality [3]. The second equality is a re-arrangement for
26 optimizing L . The third equation follows by the Thm. A.3. The last equality is established based

27 on the definition of convex conjugate A.2. When we choose KL-divergence, we have $\phi_{KL}(x) =$
 28 $x \log x - x + 1$. It can be deduced that $\phi_{KL}^*(x) = e^x - 1$. Then, we have:

$$\begin{aligned}
 \mathcal{L}_{\text{CL-DRO}}^{KL} &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha \eta - \eta_1 + \alpha \mathbb{E}_{Q_0} \left[\phi^* \left(\frac{f_\theta + \eta_1}{\alpha} \right) \right] \right\} \\
 &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha \eta - \eta_1 + \alpha \mathbb{E}_{Q_0} \left[e^{\frac{f_\theta + \eta_1}{\alpha}} - 1 \right] \right\} \\
 &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0} \left\{ \alpha \eta + \alpha \log \mathbb{E}_{Q_0} \left[e^{\frac{f_\theta}{\alpha}} \right] \right\} \\
 &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0} \left\{ \alpha \eta + \alpha \log \mathbb{E}_{Q_0} \left[e^{\frac{f_\theta}{\alpha}} \right] \right\} \\
 &= -\mathbb{E}_{P_0} \left[\alpha^* \log \frac{e^{f_\theta/\alpha^*}}{\mathbb{E}_{Q_0} [e^{f_\theta/\alpha^*}]} \right] + \alpha \eta \\
 &= \alpha^* \mathcal{L}_{\text{InfoNCE}} + \text{Constant}
 \end{aligned} \tag{18}$$

29 Here the α^* represents the optimal value of $\min_{\alpha \geq 0} \left\{ \alpha \eta + \alpha \log \mathbb{E}_{Q_0} \left[e^{\frac{f_\theta}{\alpha}} \right] \right\}$. □

30 A.2 Proof of Thm.3.3

31 **Theorem 3.3.** [Generalization Bound] Let $\widehat{\mathcal{L}}_{\text{InfoNCE}}$ be an estimation of InfoNCE with N negative
 32 samples. Then if Q^{ideal} satisfied $D_{KL}(Q^{\text{ideal}}||Q_0) \leq \eta$, we have that with probability at least $1 - \rho$:

$$\mathcal{L}_{\text{unbiased}} \leq \tau \widehat{\mathcal{L}}_{\text{InfoNCE}} + \mathcal{B}(\rho, N, \tau) \tag{5}$$

33 where $\mathcal{B}(\rho, N, \tau) = \frac{1}{N-1+\exp(\frac{1}{\tau})} \sqrt{\frac{N \exp(\frac{2}{\tau}) \log(\frac{1}{\rho})}{2}}$.

34 Here we simply disregard the constant term present in Eqn. 4 as it does not impact optimization, and
 35 omit the error from the positive instances.

36 *Proof.* Before detailing the proof process, we first introduce a pertinent theorem:

37 **Theorem A.4** (McDiarmid's inequality [10]). Let X_1, \dots, X_n be independent random variables,
 38 where X_i has range \mathcal{X} . Let $f : \mathcal{X}_1 \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$ be any function with the (c_1, \dots, c_n) -
 39 bounded difference property: for every $i = 1, \dots, n$ and every $(x_1, \dots, x_n), (x'_1, \dots, x'_n) \in$
 40 $\mathcal{X}_1 \times \dots \times \mathcal{X}_n$ that differ only in the i -th coordinate ($x_j = x'_j$ for all $j \neq i$), we have
 41 $|f(x_1, \dots, x_n) - f(x'_1, \dots, x'_n)| \leq c_i$. For any $\epsilon > 0$,

$$\mathbb{P}(f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)] \geq \epsilon) \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2}\right) \tag{19}$$

42 Now we delve into the proof. As Q^{ideal} satisfies $D_{KL}(Q||Q_0) \leq \eta$, we can bound $\mathcal{L}_{\text{unbiased}}$ with:

$$\begin{aligned}
 \mathcal{L}_{\text{unbiased}} &= -\mathbb{E}_{P_0}[f_\theta] + \mathbb{E}_{Q^{\text{ideal}}}[f_\theta] \\
 &\leq -\mathbb{E}_{P_0}[f_\theta] + \max_{D_{KL}(Q||Q_0) \leq \eta} \mathbb{E}_Q[f_\theta] \\
 &= \mathcal{L}_{\text{CL-DRO}}^{KL}
 \end{aligned} \tag{20}$$

43 where Q^{ideal}, Q^* denotes the ideal negative distribution and the worst-case distribution in CL-DRO.
 44 From the Thm.3.2, we have the equivalence between InfoNCE and CL-DRO. Thus here we choose CL-
 45 DRO for analyses. Suppose we have N negative samples, and for any pair of samples $(x_i, y_i), (x_j, y_j)$,
 46 we have the following bound:

$$|Q^*(x_i, y_i) f_\theta(x_i, y_i) - Q^*(x_j, y_j) f_\theta(x_j, y_j)| \leq \sup_{(x, y) \sim Q_0} |Q^*(x, y) f_\theta(x, y)| \leq \frac{\exp(\frac{1}{\tau})}{N-1+\exp(\frac{1}{\tau})} \tag{21}$$

47 where the first inequality holds as $Q^*(x, y) f_\theta(x, y) > 0$. The second inequality holds based on
 48 the expression of $Q^* = Q_0 \frac{\exp[f_\theta/\tau]}{\mathbb{E}_{Q_0} \exp[f_\theta/\tau]}$ (refer to Appendix A.6). Suppose $f_\theta \in [M_1, M_2]$, the

49 upper bound of $\sup_{(x,y) \sim Q_0} |Q^*(x,y)f_\theta(x,y)|$ arrives if $f_\theta(x,y) = M_2$ for the sample (x,y) and
 50 $f_\theta(x,y) = M_1$ for others. We have $\sup_{(x,y) \sim Q_0} |Q^*(x,y)f_\theta(x,y)| \leq \frac{M_2 \exp((M_2 - M_1)/\tau)}{N - 1 + \exp((M_2 - M_1)/\tau)}$. In
 51 this work, for brevity, here we simply consider $M_1 = 0, M_2 = 1$ for analyses. It shares the common
 52 properties with the general interval $[M_1, M_2]$.

53 By using McDiarmid's inequality in Thm A.4, for any ϵ , we have:

$$\begin{aligned} & \mathbb{P}[(\mathcal{L}_{\text{CL-DRO}}^{KL} - \tau \widehat{\mathcal{L}}_{\text{InfoNCE}}) \geq \epsilon] \\ & \leq \exp\left(\frac{-2\epsilon^2(N-1 + \exp(\frac{1}{\tau}))^2}{N \exp(\frac{2}{\tau})}\right) \end{aligned} \quad (22)$$

54 Let

$$\rho = \exp\left(\frac{-2\epsilon^2(N-1 + \exp(\frac{1}{\tau}))^2}{N \exp(\frac{2}{\tau})}\right) \quad (23)$$

55 we get:

$$\epsilon = \frac{1}{N-1 + \exp(\frac{1}{\tau})} \sqrt{\frac{N \exp(\frac{2}{\tau}) \log(\frac{1}{\rho})}{2}} \quad (24)$$

56 Thus, for $\forall \rho \in (0, 1)$, we conclude that with probability at least $1 - \rho$.

$$\mathcal{L}_{\text{unbiased}} \leq \widehat{\mathcal{L}}_{\text{InfoNCE}} + \frac{1}{N-1 + \exp(\frac{1}{\tau})} \sqrt{\frac{N \exp(\frac{2}{\tau}) \log(\frac{1}{\rho})}{2}} \quad (25)$$

57 □

58 A.3 Proof of Coro.3.4

59 **Corollary 3.4.** *[The optimal α - Lemma 5 of [6]] The value of the optimal α (i.e., τ) can be*
 60 *approximated as follow:*

$$\tau \approx \sqrt{\mathbb{V}_{Q_0}[f_\theta]/2\eta}. \quad (6)$$

61 where $\mathbb{V}_{Q_0}[f_\theta]$ denotes the variance of f_θ under the distribution Q_0 .

62 *Proof.* While Corollary 3.4 has already been proven in [6], we present a brief outline of the proof here
 63 for the sake of completeness and to ensure that our article is self-contained. To verify the relationship
 64 between τ and η , we could utilize the approximate expression of InfoNCE (cf. Eqn. 29) and focus on
 65 the first order conditions for τ . In detail, we have:

$$-\mathbb{E}_{P_0}[f_\theta] + \inf_{\alpha \geq 0} \left\{ \mathbb{E}_{Q_0}[f_\theta] - \frac{1}{2\alpha} \frac{1}{\phi^{(2)}(1)} \mathbb{V}_{Q_0}[f_\theta] - \alpha\eta \right\}$$

66 To find the optimal value of α (or equivalently, τ), we differentiate the above equation and set it to 0.
 67 This yields a fixed-point equation

$$\tau = \sqrt{\frac{\mathbb{V}_{Q_0}[f_\theta]}{2\eta}}$$

68 The corollary gets proved. □

69 A.4 Proof of Thm.3.5

70 **Theorem 3.5.** *Given any ϕ -divergence, the corresponding CL-DRO objective could be approximated*
 71 *as a mean-variance objective:*

$$\mathcal{L}_{\text{CL-DRO}}^\phi(f_\theta) \approx -\mathbb{E}_{P_0}[f_\theta] + (\mathbb{E}_{Q_0}[f_\theta] + \frac{1}{2\tau} \frac{1}{\phi^{(2)}(1)} \cdot \mathbb{V}_{Q_0}[f_\theta]) \quad (7)$$

72 where $\phi^{(2)}(1)$ denotes the the second derivative value of $\phi(\cdot)$ at point 1, and $\mathbb{V}_{Q_0}[f_\theta]$ denotes the
 73 variance of f under the distribution Q_0 .

74 Specially, if we consider KL divergence, the approximation transforms:

$$\mathcal{L}_{\text{CL-DRO}}^{KL}(f_\theta) \approx -\mathbb{E}_{P_0}[f_\theta] + (\mathbb{E}_{Q_0}[f_\theta] + \frac{1}{2\tau} \mathbb{V}_{Q_0}[f_\theta]) \quad (8)$$

75 *Proof.* We start with introducing a useful lemma.

76 **Lemma A.5** (Lemma A.2 of [8]). *Suppose that $\phi : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ is a closed, convex function*
 77 *such that $\phi(z) \geq \phi(1) = 0$ for all z , is two times continuously differentiable around $z = 1$, and*
 78 *$\phi(1) > 0$, Then*

$$\begin{aligned}\phi^*(\zeta) &= \max_z \{z\zeta - \phi(z)\} \\ &= \zeta + \frac{1}{2!} \left[\frac{1}{\phi''(1)} \right] \zeta^2 + o(\zeta^2)\end{aligned}\tag{26}$$

79 Note that most of the ϕ -divergences [13] (e.g., KL divergence, Cressie-Read divergence, Burg entropy,
 80 J-divergence, χ^2 -distance, modified χ^2 -distance, and Hellinger distance) satisfy the smoothness
 81 conditions. When $n = 2$, $\phi^*[\zeta] \approx \zeta + \frac{1}{2} \left[\frac{1}{\phi^{(2)}(1)} \right] \zeta^2$. Substituting this back to Eqn.17 we have:

$$\begin{aligned}\mathcal{L}_{\text{CL-DRO}}^\phi &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha\eta - \eta_1 + \alpha \mathbb{E}_{Q_0} \left[\phi^* \left(\frac{f_\theta + \eta_1}{\alpha} \right) \right] \right\} \\ &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha\eta - \eta_1 + \alpha \mathbb{E}_{Q_0} \left[\frac{f_\theta + \eta_1}{\alpha} + \frac{1}{2} \frac{1}{\phi^{(2)}(1)} \left(\frac{f_\theta + \eta_1}{\alpha} \right)^2 \right] \right\} \\ &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha\eta - \eta_1 + \mathbb{E}_{Q_0}[f_\theta + \eta_1] + \frac{1}{2} \frac{1}{\phi^{(2)}(1)\alpha} (f_\theta + \eta_1)^2 \right\} \\ &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\alpha \geq 0, \eta_1} \left\{ \alpha\eta + \mathbb{E}_{Q_0}[f_\theta] + \frac{1}{2} \frac{1}{\phi^{(2)}(1)\alpha} (f_\theta + \eta_1)^2 \right\}\end{aligned}\tag{27}$$

82 If we differentiate it w.r.t. η_1 :

$$\frac{\partial \left\{ \alpha\eta + \mathbb{E}_{Q_0} \left[f_\theta + \frac{1}{2} \frac{1}{\phi^{(2)}(1)\alpha} (f_\theta + \eta_1)^2 \right] \right\}}{\partial \eta_1} = 0\tag{28}$$

83 we have $\eta_1^* = -\mathbb{E}_{Q_0}[f_\theta]$, and the objective transforms into:

$$\begin{aligned}& -\mathbb{E}_{P_0}[f_\theta] + \inf_{\alpha \geq 0, \eta_1} \left\{ \alpha\eta + \mathbb{E}_{Q_0} \left[f_\theta + \frac{1}{2} \frac{1}{\phi^{(2)}(1)\alpha} (f_\theta + \eta_1)^2 \right] \right\} \\ &= -\mathbb{E}_{P_0}[f_\theta] + \inf_{\alpha \geq 0} \left\{ \mathbb{E}_{Q_0} \left[f_\theta + \frac{1}{2\alpha} \frac{1}{\phi^{(2)}(1)} (f_\theta - \mathbb{E}_{Q_0}[f_\theta])^2 \right] - \alpha\eta \right\} \\ &= -\mathbb{E}_{P_0}[f_\theta] + \inf_{\alpha \geq 0} \left\{ \mathbb{E}_{Q_0}[f_\theta] + \frac{1}{2\alpha} \frac{1}{\phi^{(2)}(1)} \mathbb{V}_{Q_0}[f_\theta] - \alpha\eta \right\}\end{aligned}\tag{29}$$

84 Choosing KL-divergence, we have $\phi^{(2)}(1) = 1$. Substituting $\alpha^*(\tau)$ into Eqn. 29 and ignoring the
 85 constant $\alpha\eta$:

$$-\mathbb{E}_{P_0}[f_\theta] + \mathbb{E}_{Q_0}[f_\theta] + \frac{1}{2\tau} \mathbb{V}_{Q_0}[f_\theta]$$

86 Then Theorem gets proved. □

87 A.5 Proof of Thm.4.2

88 **Theorem 4.2.** *For distributions P, Q such that $P \ll Q$, let \mathcal{F} be a set of bounded measurable*
 89 *functions. Let CL-DRO draw positive and negative instances from P and Q , marked as $\mathcal{L}_{\text{CL-DRO}}^\phi(P, Q)$.*
 90 *Then the CL-DRO objective is the tight variational estimation of ϕ -divergence. In fact, we have:*

$$D_\phi(P||Q) = \sup_{f \in \mathcal{F}} -\mathcal{L}_{\text{CL-DRO}}^\phi(P, Q) = \sup_{f \in \mathcal{F}} \mathbb{E}_P[f] - \min_{\lambda \in \mathbb{R}} \{ \lambda + \mathbb{E}_Q[\phi^*(f - \lambda)] \}\tag{10}$$

91 Here, the choice of ϕ in CL-DRO corresponds to the probability measures in $D_\phi(P||Q)$.

92 *Proof.* Regarding this theorem, our proof primarily relies on the variational representation of ϕ -
 93 divergence and optimized certainty equivalent (OCE) risk. Towards this end, we start to introduce the
 94 basic concepts:

95 **Definition A.6** (OCE [2]). Let X be a random variable and let u be a convex, lower-semicontinuous
 96 function satisfies $u(0) = 0, u^*(1) = 0$, then optimized certainty equivalent (OCE) risk $\rho(X)$ is
 97 defines as:

$$\rho(X) = \inf_{\lambda \in \mathbb{R}} \{\lambda + \mathbb{E}[u(f - \lambda)]\} \quad (30)$$

98 OCE is a type of risk measure that is widely used by both practitioners and academics [1, 2]. With
 99 duality theory, its various properties have been inspiring in our study of DRO.

Definition A.7 (Variational formulation).

$$D_\phi(P||Q) := \sup_{f \in \mathcal{F}} \{\mathbb{E}_P[f] - \mathbb{E}_Q[\phi^*(f)]\} \quad (31)$$

100 where the supremum is taken over all bounded real-valued measurable functions \mathcal{F} defined on \mathcal{X} .

101 Note that in order to keep consistent with the definition of CL-DRO, we transform Eqn.10 to :

$$D_\phi(P||Q) = \sup_{f \in \mathcal{F}} -\mathcal{L}_{\text{CL-DRO}}^\phi(P, Q) = \sup_{f \in \mathcal{F}} \mathbb{E}_P[f] - \inf_{\eta_1 \in \mathbb{R}} \{-\eta_1 + \mathbb{E}_Q[\phi^*(f + \eta_1)]\} \quad (32)$$

102 Our proof for this theorem primarily relies on utilizing OCE risk as a bridge and can be divided into
 103 two distinct steps:

104 Step 1: $\sup_{f \in \mathcal{F}} -\mathcal{L}_{\text{CL-DRO}}^\phi(P, Q) = \sup_{f \in \mathcal{F}} \mathbb{E}_P[f] - \inf_{\eta_1 \in \mathbb{R}} \{-\eta_1 + \mathbb{E}_Q[\phi^*(f + \eta_1)]\}$.

105 Step 2: $D_\phi(P||Q) = \sup_{f \in \mathcal{F}} \mathbb{E}_P[f] - \inf_{\eta_1 \in \mathbb{R}} \{-\eta_1 + \mathbb{E}_Q[\phi^*(f + \eta_1)]\}$

106 1. **We show that** $\sup_{f \in \mathcal{F}} -\mathcal{L}_{\text{CL-DRO}}^\phi(P, Q) = \sup_{f \in \mathcal{F}} \mathbb{E}_P[f] - \inf_{\eta_1 \in \mathbb{R}} \{-\eta_1 + \mathbb{E}_Q[\phi^*(f + \eta_1)]\}$.

$$\begin{aligned} -\mathcal{L}_{\text{CL-DRO}}^\phi &= \mathbb{E}_P[f] - \inf_{\alpha \geq 0, \eta_1} \sup_L \{\mathbb{E}_Q[fL] - \alpha[\mathbb{E}_Q[\phi(L)] - \eta] + \eta_1(\mathbb{E}_Q[L] - 1)\} \\ &= \mathbb{E}_P[f] - \inf_{\alpha \geq 0, \eta_1} \{\alpha\eta - \eta_1 + \alpha\mathbb{E}_Q[\phi^*(\frac{f + \eta_1}{\alpha})]\} \\ &= \mathbb{E}_P[f] - \inf_{\eta_1} \{\alpha^*\eta - \eta_1 + \alpha^*\mathbb{E}_Q[\phi^*(\frac{f + \eta_1}{\alpha^*})]\} \\ &= \mathbb{E}_P[f] - \inf_{\eta_1} \{-\eta_1 + \alpha^*\mathbb{E}_Q[\phi^*(\frac{f + \eta_1}{\alpha^*})] + \text{Constant}\} \end{aligned} \quad (33)$$

107 When $\alpha^* = 1$, step 1 gets proved.

108 2. **We show that** $D_\phi(P||Q) = \sup_{f \in \mathcal{F}} \mathbb{E}_P[f] - \inf_{\eta_1 \in \mathbb{R}} \{-\eta_1 + \mathbb{E}_Q[\phi^*(f + \eta_1)]\}$.

109 Firstly, we transform $\mathbb{E}_P[f]$ to $\mathbb{E}_Q[f \frac{dP}{dQ}]$ as:

$$\begin{aligned} &\sup_{f \in \mathcal{F}} \mathbb{E}_P[f] - \inf_{\eta_1} \{-\eta_1 + \mathbb{E}_Q[\phi^*(f + \eta_1)]\} \\ &= \sup_{f \in \mathcal{F}} \mathbb{E}_Q[f \frac{dP}{dQ}] - \inf_{\eta_1} \{-\eta_1 + \mathbb{E}_Q[\phi^*(f + \eta_1)]\} \end{aligned} \quad (34)$$

110 Let $f + \eta_1 = Y$, we have:

$$\begin{aligned} &\sup_{f \in \mathcal{F}} \mathbb{E}_Q[f \frac{dP}{dQ}] - \inf_{\eta_1} \{-\eta_1 + \mathbb{E}_Q[\phi^*(f + \eta_1)]\} \\ &= \sup_{Y \in \mathcal{F}} \inf_{\eta_1} \mathbb{E}_Q[(Y - \eta_1) \frac{dP}{dQ}] - \{-\eta_1 + \mathbb{E}_Q[\phi^*(Y)]\} \\ &= \sup_{Y \in \mathcal{F}} \mathbb{E}_Q[Y \frac{dP}{dQ} - \phi^*(Y)] + \inf_{\eta_1} \eta_1 \mathbb{E}_Q[1 - \frac{dP}{dQ}] \\ &= \sup_{Y \in \mathcal{F}} \mathbb{E}_Q[Y \frac{dP}{dQ} - \phi^*(Y)] + 0 \end{aligned} \quad (35)$$

111 The first equality follows from replacing $f + \theta_1$ with Y . The second equality is a re-arrangement
 112 for optimizing η_1 . The third equation holds as $\mathbb{E}_Q[1 - \frac{dP}{dQ}] = 0$.

113 Applying Thm. A.3, the last supremum reduces to:

$$\begin{aligned}
 & \sup_{Y \in \mathcal{F}} \mathbb{E}_Q[Y \frac{dP}{dQ} - \phi^*(Y)] \\
 &= \mathbb{E}_Q[\sup_{Y \in \mathcal{F}} \{Y \frac{dP}{dQ} - \phi^*(Y)\}] \\
 &= \mathbb{E}_Q[\phi^{**}(\frac{dP}{dQ})] \\
 &= \mathbb{E}_Q[\phi(\frac{dP}{dQ})] \\
 &= D_\phi(P||Q)
 \end{aligned} \tag{36}$$

114 where the last equality follows from the fact that $\phi^{**} = \phi$. This concludes the proof.

115 □

116 A.6 Proof of Q^*

117 *Proof.* From theorem3.2, CL-DRO can be rewritten as:

$$\begin{aligned}
 \mathcal{L}_{\text{CL-DRO}}^\phi &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\eta_1} \{ \alpha^* \eta - \eta_1 + \alpha^* \mathbb{E}_{Q_0}[\max_L \{ \frac{f_\theta + \eta_1}{\alpha^*} L - \phi(L) \}] \} \\
 &= -\mathbb{E}_{P_0}[f_\theta] + \min_{\eta_1} \{ \alpha^* \eta - \eta_1 + \alpha^* \mathbb{E}_{Q_0}[\phi^*(\frac{f_\theta + \eta_1}{\alpha^*})] \}
 \end{aligned} \tag{37}$$

118 For the inner optimization, we can draw the optimal L for $\max_L \{ \frac{f_\theta + \eta_1}{\alpha^*} L - \phi(L) \}$ as:

$$L = e^{\frac{f_\theta + \eta_1}{\alpha^*}} \tag{38}$$

119 For the outer optimization, we can draw the optimal η_1 for $\min_{\eta_1} \{ \alpha^* \eta - \eta_1 + \alpha^* \mathbb{E}_{Q_0}[e^{\frac{f_\theta + \eta_1}{\alpha^*}} - 1] \}$
 120 as

$$\eta_1 = -\alpha^* \log \mathbb{E}_{Q_0} e^{\frac{f_\theta}{\alpha^*}} \tag{39}$$

121 Then we plug Eqn. 39 into Eqn. 38.

$$L = \frac{e^{\frac{f_\theta}{\alpha^*}}}{\mathbb{E}_{Q_0}[e^{\frac{f_\theta}{\alpha^*}}]} \tag{40}$$

122 Based on the definition of L , we can derive the expression for Q^* :

$$Q^* = \frac{e^{\frac{f_\theta}{\alpha^*}}}{\mathbb{E}_{Q_0}[e^{\frac{f_\theta}{\alpha^*}}]} Q_0 \tag{41}$$

123 □

124 B Experiments

125 Figure 5 shows PyTorch-style pseudocode for the standard objective, the adjusted InfoNCE objective.
 126 The proposed adjusted reweighting loss is very simple to implement, requiring only two extra lines of
 127 code compared to the standard objective.

```

1 # pos      : exp of inner products for positive examples
2 # neg      : exp of inner products for negative examples
3 # N        : number of negative examples
4 # t        : temperature scaling
5 # mu       : center position
6 # sigma    : height scale
7
8 #InfoNCE
9 standard_loss = -log(pos.sum() / (pos.sum() + neg.sum()))
10
11 #ADNCE
12 weight=1/(sigma * sqrt(2*pi)) * exp( -0.5 * ((neg-mu)/sigma)**2 )
13 weight=weight/weight.mean()
14 Adjusted_loss = -log(pos.sum() / (pos.sum() + (neg * weight.detach() ).sum()
    )

```

Figure 5: Pseudocode for our proposed adjusted InfoNCE objective, as well as the original NCE contrastive objective. The implementation of our adjusted reweighting method only requires two additional lines of code compared to the standard objective.

Table 6: hyperparameters setting on each datasets.

DATASETS	CIFAR10	STL10	CIFAR100
BEST τ	{ 0.1, 0.2, 0.3, 0.4 , 0.5, 0.6 }	{ 0.1, 0.2 , 0.3, 0.4, 0.5, 0.6 }	{ 0.1, 0.2, 0.3 , 0.4, 0.5, 0.6 }
μ	{ 0.5, 0.6, 0.7 , 0.8, 0.9 }	{ 0.5, 0.6, 0.7, 0.8 , 0.9 }	0.5 , 0.6, 0.7, 0.8, 0.9 }
σ	{ 0.5 , 1.0 }	{0.5, 1.0 }	{0.5, 1.0 }

128 B.1 Visual Representation

129 **Model.** For contrastive learning on images, we adopt SimCLR [4] as our baseline and follow the same
130 experimental setup as [5]. Specifically, we use the ResNet-50 network as the backbone. To ensure a
131 fair comparison, we set the embedded dimension to 2048 (the representation used in linear readout)
132 and project it into a 128-dimensional space (the actual embedding used for contrastive learning).
133 Regarding the temperature parameter τ , we use the default value τ_0 of 0.5 in most researches, and we
134 also perform grid search on τ varying from 0.1 to 1.0 at an interval of 0.1, denoted by τ^* . The best
135 parameters for each dataset is reported in Table 6. Note that $\{\cdot\}$ indicates the range of hyperparameters
136 that we tune and the numbers in **bold** are the final settings. For α -CL, we follow the setting of [14],
137 where $p = 4$ and $\tau = 0.5$. We use the Adam optimizer with a learning rate of 0.001 and weight decay
138 of $1e - 6$. All models are trained for 400 epochs.

139 **Noisy experiments in Sec.3.4.** To investigate the relationship between the temperature parameter τ
140 (or η) and the noise ratio, we follow the approach outlined in [5] and utilize the class information
141 of each image to select negative samples as a combination of true negative samples and false
142 negative samples. Specifically, $r_{ratio} = 0$ indicates all negative samples are true negative samples,
143 $r_{ratio} = 0.5$ suggests 50% of true positive samples existing in negative samples, $r_{ratio} = 1$ means
144 uniform sampling.

145 **Variance analysis in Sec.3.4.** To verify the mean-variance objective of InfoNCE, we adopt the
146 approach outlined in [16] and record the negative prediction scores for 256 samples (assuming a batch
147 size of 256) in each minibatch. Specifically, we randomly select samples from a batch to calculate
148 the statistics and visualize them. (1) For positive samples, we calculate cosine similarity by taking
149 the inner product after normalization, and retain the mean value of the 256 positive scores as ‘*pos*
150 *mean*’. (2) For negative samples, we average the means and variances of 256 negative samples to
151 show the statistical characteristics of these N negative samples ‘(*mean neg*; *var neg*)’. We record this
152 data at each training step to track score distribution throughout the training process.

153 B.2 Sentence Representation

154 For the sentence contrastive learning, we adopt the approach outlined in [7] and evaluate our
155 method on 7 popular STS datasets: STS tasks from 2012-2016, STS-B and SICK-R. We utilize
156 the SentEval toolkit to obtain all 7 datasets. Each dataset includes sentence pairs which are rated

Table 7: hyperparameters setting on sentence CL. Note that $\{\cdot\}$ indicates the range of hyperparameters that we tune and the numbers in **bold** are the final settings.

DATASETS	SIMCSE-BERT _{BASE}	SIMCSE-ROBERTA _{BASE}
BEST τ	{ 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07 , 0.08, 0.09, 0.10, 0.15, 0.20 }	{ 0.01, 0.02, 0.03, 0.04, 0.05, 0.06 , 0.07, 0.08, 0.09, 0.10, 0.15, 0.20 }
μ	{ 0.3, 0.4 , 0.5, 0.6, 0.7, 0.8, 0.9 }	{ 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1.5, 2.0 , 2.5, 3.0 }
σ	{ 0.5, 1.0 }	{ 0.5, 1.0 }

157 on a scale of 0 to 5, indicating the degree of semantic similarity. To validate the effective of our
 158 proposed method, we utilize several methods as baselines: average GloVe embeddings, BERT-
 159 flow, BERT-whitening, CT-BERT and SimCSE. The best parameters for each dataset is reported
 160 in Table 7. To ensure fairness, we employed the official code, which can be accessed at <https://github.com/princeton-nlp/SimCSE>.
 161

162 B.3 Graph Representation

163 For the graph contrastive learning experiments on TU-Dataset [12], we adopted the same experimental
 164 setup as outlined in [15]. The dataset statistics can be found in Tab.8. To ensure fairness, we
 165 employed the official code, which can be accessed at https://github.com/Shen-Lab/GraphCL/tree/master/unsupervised_TU. We made only modifications to the script by incorporating our
 166 ADNCE method and conducting experiments on the hyper-parameter $\mu \in \{0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$
 167 and $\sigma = 1$ on most datasets. Each parameter was repeated from scratch five times, and the best
 168 parameter was selected by evaluating on the validation dataset. The best parameters for each dataset
 169 is reported in Table 9.
 170

171 We summarize the statistics of TU-datasets [12] for unsupervised learning in Table 8. Tab. 10
 demonstrates the consistent superiority of our proposed ADNCE approach.

Table 8: Statistics for unsupervised learning TU-datasets.

DATASETS	CATEGORY	GRAPHS#	AVG. N#	AVG. DEGREE
NCII	BIOCHEMICAL MOLECULES	4,110	29.87	1.08
PROTEINS	BIOCHEMICAL MOLECULES	1,113	39.06	1.86
DD	BIOCHEMICAL MOLECULES	1,178	284.32	715.66
MUTAG	BIOCHEMICAL MOLECULES	188	17.93	19.79
COLLAB	SOCIAL NETWORKS	5,000	74.49	32.99
RDT-B	SOCIAL NETWORKS	2,000	429.63	1.15
RDT-M	SOCIAL NETWORKS	2,000	429.63	497.75
IMDB-B	SOCIAL NETWORKS	1,000	19.77	96.53

172

Table 9: hyperparameters setting on graph CL. Note that $\{\cdot\}$ indicates the range of hyperparameters that we tune and the numbers in **bold** are the final settings.

DATASETS	BEST τ	μ	σ
NCII	{ 0.05 , 0.10, 0.15, 0.20, 0.25 }	{ 0.2, 0.3, 0.4, 0.5, 0.6, 0.7 , 0.8, 0.9 }	{ 0.5, 1.0 }
PROTEINS	{ 0.05 , 0.10, 0.15, 0.20, 0.25 }	{ 0.5, 1.0, 1.5 , 2.0 }	{ 0.5, 1.0 }
DD	{ 0.05, 0.10, 0.15, 0.20 , 0.25 }	{ 0.2 , 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9 }	{ 0.5, 1.0 }
MUTAG	{ 0.05, 0.10, 0.15 , 0.20, 0.25 }	{ 0.2, 0.3, 0.4, 0.5, 0.6, 0.7 , 0.8, 0.9 }	{ 0.5, 1.0 }
COLLAB	{ 0.05, 0.10 , 0.15, 0.20, 0.25 }	{ 0.2 , 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9 }	{ 0.5, 1.0 }
RDT-B	{ 0.05, 0.10, 0.15 , 0.20, 0.25 }	{ 0.2, 0.3 , 0.4, 0.5, 0.6, 0.7, 0.8, 0.9 }	{ 0.5, 1.0 }
RDT-M	{ 0.05, 0.10, 0.15 , 0.20, 0.25 }	{ 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8 , 0.9 }	{ 0.5, 1.0 }
IMDB-B	{ 0.10, 0.20, 0.30, 0.40, 0.50 }	{ 0.2, 0.3, 0.4, 0.5, 0.6, 0.7 , 0.8, 0.9 }	{ 0.5, 1.0 }

Table 10: Unsupervised representation learning classification accuracy (%) on TU datasets. The compared numbers are from except AD-GCL, whose statistics are reproduced on our platform. **Bold** indicates the best performance while underline indicates the second best on each dataset.

DATASET	NCI1	PROTEINS	DD	MUTAG	COLLAB	RDT-B	RDT-M5K	IMDB-B	AVG.
NO PRE-TRAIN	65.40±0.17	72.73±0.51	75.67±0.29	87.39±1.09	65.29±0.16	76.86 ±0.25	48.48±0.28	69.37±0.37	70.15
INFOGRAPH	76.20±1.06	74.44±0.31	72.85±1.78	<u>89.01±1.13</u>	70.05±1.13	82.50±1.42	53.46±1.03	73.03±0.87	74.02
GRAPHCL	77.87±0.41	74.39±0.45	78.62±0.40	86.80±1.34	71.36±1.15	89.53±0.84	55.99±0.28	71.14±0.44	75.71
AD-GCL	73.91±0.77	73.28±0.46	75.79±0.87	88.74±1.85	<u>72.02±0.56</u>	90.07±0.85	54.33±0.32	70.21±0.68	74.79
RGCL	78.14±1.08	75.03±0.43	78.86±0.48	87.66±1.01	<u>70.92±0.65</u>	90.34±0.58	56.38±0.40	71.85±0.84	76.15
ADNCE	79.30±0.67	75.10±0.25	79.23±0.59	89.04±1.30	72.26±1.10	91.39±0.31	<u>56.01±0.35</u>	<u>71.58±0.72</u>	76.74

References

- 173
- 174 [1] Aharon Ben-Tal and Marc Teboulle. Expected utility, penalty functions, and duality in stochastic
175 nonlinear programming. *Management Science*, 32(11):1445–1466, 1986.
- 176 [2] Aharon Ben-Tal and Marc Teboulle. An old-new concept of convex risk measures: the optimized
177 certainty equivalent. *Mathematical Finance*, 17(3):449–476, 2007.
- 178 [3] Stephen P Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press,
179 2004.
- 180 [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework
181 for contrastive learning of visual representations. In *International conference on machine*
182 *learning*, pages 1597–1607. PMLR, 2020.
- 183 [5] Ching-Yao Chuang, Joshua Robinson, Yen-Chen Lin, Antonio Torralba, and Stefanie Jegelka.
184 Debaised contrastive learning. *Advances in neural information processing systems*, 33:8765–
185 8775, 2020.
- 186 [6] Louis Faury, Ugo Tanielian, Elvis Dohmatob, Elena Smirnova, and Flavian Vasile. Distribu-
187 tionally robust counterfactual risk minimization. In *Proceedings of the AAAI Conference on*
188 *Artificial Intelligence*, volume 34, pages 3850–3857, 2020.
- 189 [7] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence
190 embeddings. *arXiv preprint arXiv:2104.08821*, 2021.
- 191 [8] Jun-ya Gotoh, Michael Jong Kim, and Andrew EB Lim. Robust empirical optimization is
192 almost the same as mean–variance optimization. *Operations research letters*, 46(4):448–452,
193 2018.
- 194 [9] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of convex analysis*. Springer
195 Science & Business Media, 2004.
- 196 [10] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*.
197 2018.
- 198 [11] Christopher Morris, Nils M Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion
199 Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. *arXiv*
200 *preprint arXiv:2007.08663*, 2020.
- 201 [12] Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion
202 Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. *CoRR*,
203 abs/2007.08663, 2020.
- 204 [13] XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence func-
205 tionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information*
206 *Theory*, 56(11):5847–5861, 2010.
- 207 [14] Yuandong Tian. Understanding deep contrastive learning via coordinate-wise optimization. In
208 *Advances in Neural Information Processing Systems*, 2022.
- 209 [15] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen.
210 Graph contrastive learning with augmentations. *Advances in Neural Information Processing*
211 *Systems*, 33:5812–5823, 2020.

- 212 [16] Rui Zhu, Bingchen Zhao, Jingen Liu, Zhenglong Sun, and Chang Wen Chen. Improving
213 contrastive learning by visualizing feature transformation. In *Proceedings of the IEEE/CVF*
214 *International Conference on Computer Vision*, pages 10306–10315, 2021.