
Label Noise: Ignorance Is Bliss [6 pages]

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We establish a new theoretical framework for learning under multi-class, instance-
2 dependent label noise. At the heart of our framework is the concept of *relative*
3 *signal strength* (RSS), which is a point-wise measure of noisiness. We use relative
4 signal strength to establish matching upper and lower bounds for excess risk. Our
5 theoretical findings reveal a surprising result: the extremely simple *Noise Ignorant*
6 *Empirical Risk Minimization (NI-ERM)* principle, which conducts empirical risk
7 minimization as if no label noise exists, is minimax optimal. Finally, we translate
8 these theoretical insights into practice: by using NI-ERM to fit a linear classifier
9 on top of a frozen foundation model, we achieve state-of-the-art performance on
10 the CIFAR-N data challenge.

11 1 Introduction

12 The problem of classification with label noise can be stated in terms of variables (X, Y, \tilde{Y}) , where
13 X is the feature vector, $Y \in \{1, \dots, K\}$ is the true label associated to X , and $\tilde{Y} \in \{1, \dots, K\}$ as
14 a noisy version of Y . The learner has access to i.i.d. realizations of (X, \tilde{Y}) , and the objective is to
15 learn a classifier that optimizes the risk associated with (X, Y) .

16 In recent years, there has been a surge of interest in the challenging setting of instance (i.e., feature)
17 dependent label noise, in which \tilde{Y} can depend on both Y and X . While several algorithms have been
18 developed, there remains relatively little theory regarding algorithm performance and the fundamental
19 limits of this learning paradigm.

20 This work develops a theoretical framework for learning under multi-class, instance-dependent label
21 noise. Our framework hinges on the concept of *relative signal strength*, which is a point-wise measure
22 of “noisiness” in a label noise problem. Using relative signal strength, we establish matching upper
23 and lower bounds for excess risk. We further identify distributional assumptions that ensure that the
24 lower bound tends to zero as the sample size n grows, implying that consistent learning is possible.

25 Surprisingly, our theoretical findings reveal that *Noise Ignorant Empirical Risk Minimization (NI-*
26 *ERM)*, which conducts empirical risk minimization as if no label noise exists, is minimax optimal. To
27 translate this insight into practice, we use NI-ERM to fit a linear classifier on top of a self-supervised
28 feature extractor, achieving state-of-the-art performance on the CIFAR-N data challenge.

29 2 Literature review

30 Theory and algorithms for classification with label noise are often based on different probabilistic
31 models. These can be categorized depending on how \tilde{Y} depends on Y and X . The simplest model is
32 symmetric noise, where the distribution of \tilde{Y} is independent of Y and X [Angluin and Laird, 1988].
33 In this case, the probability that $\tilde{Y} = k$ is the same for all $k \neq Y$, regardless of Y and X . In this

34 setting, it is easy to show that minimizing the noisy excess risk (associated to the 0/1 loss) implies
 35 minimizing the clean excess risk, a property known as *immunity*. When immunity holds, there is no
 36 need to modify the learning algorithm on account of noisy labels. In other words, the learner may be
 37 *ignorant* of the label noise and still learn consistently.

38 A more general model is classification with label dependent noise, in which the distribution of \tilde{Y}
 39 depends on Y , but not X . Many practical algorithms have been developed over the years, based
 40 on principles including data re-weighting [Liu and Tao, 2015], robust training [Han et al., 2018,
 41 Liu et al., 2020, Foret et al., 2021] and data cleaning [Brodley and Friedl, 1999, Northcutt et al.,
 42 2021]. Consistent learning algorithms still exist, such as those based on loss correction [Natarajan
 43 et al., 2013, Patrini et al., 2017, Van Rooyen and Williamson, 2018, Liu and Guo, 2020, Zhang et al.,
 44 2022]. However, the primary limitation of these methods is their reliance on the knowledge of noise
 45 transition probabilities, which are generally not identifiable [Zhang et al., 2021b].

46 In the most general setting, that of instance dependent label noise, the distribution of \tilde{Y} depends on
 47 both Y and X . While algorithms are emerging [Cheng et al., 2021, Zhu et al., 2021, Wang et al., 2022,
 48 Yang et al., 2023], the theoretical developments have primarily focused on the binary setting. Scott
 49 [2019] establishes immunity for a Neyman-Pearson-like performance criterion under a *posterior drift*
 50 model, discussed in more detail below. Cannings et al. [2020] establish an upper bound for excess
 51 risk under the strong assumption that the optimal classifiers for the clean and noisy distributions are
 52 the same. Closest to our work, Im and Grigas [2023] derive excess risk upper and lower bounds, and
 53 reach a similar conclusion, that noise-ignorant ERM attains the lower bound. Our results, based on
 54 the new concept of relative signal strength, provide a more refined analysis.

55 Additional connections between our contributions and prior work are made throughout the paper.

56 3 Problem statement and relative signal strength

57 *Notation.* \mathcal{X} denotes a feature space and $\mathcal{Y} = \{1, 2, \dots, K\}$ a label space, with $K \in \mathbb{N}$. The
 58 K -simplex is $\Delta^K := \{p \in \mathbb{R}^K : \forall i, p_i \geq 0, \sum p_i = 1\}$. A $K \times K$ matrix is *row stochastic* if all
 59 of its rows are in Δ^K . Denote the i -th element of a vector v as $[v]_i$, and the (i, j) -th element of a
 60 matrix M as $[M]_{i,j}$.

61 3.1 Learning from label noise

62 In conventional multiclass classification, we observe training data $(X_1, Y_1), \dots, (X_n, Y_n)$ drawn
 63 i.i.d. from a joint distribution P_{XY} . The marginal distribution of X is denoted by P_X , and the *class*
 64 *posterior* probabilities $P_{Y|X=x}$ are captured by a K -simplex-valued vector $\boldsymbol{\eta} : \mathcal{X} \rightarrow \Delta^K$, where
 65 the j -th component of the vector is $[\boldsymbol{\eta}(x)]_j = \mathbb{P}(Y = j | X = x)$. A classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ maps an
 66 instance x to a class $f(x) \in \mathcal{Y}$. Denote the risk of a classifier f with respect to distribution P_{XY} as
 67 $R(f) = \mathbb{E}_{(X,Y) \sim P_{XY}} [\mathbb{1}_{\{f(X) \neq Y\}}]$. The Bayes optimal classifier for P_{XY} is $f^*(x) \in \arg \max \boldsymbol{\eta}(x)$.
 68 The Bayes risk, which is the minimum achievable risk, is then denoted as $R^* = R(f^*) = \inf_f R(f)$.

69 We consider the setting where, instead of the true class label Y , a noisy label \tilde{Y} is observed.
 70 The training data $(X_1, \tilde{Y}_1), \dots, (X_n, \tilde{Y}_n)$ can be viewed as an i.i.d. sample drawn from a “noisy”
 71 distribution $P_{X\tilde{Y}}$. We define $P_{\tilde{Y}|X=x}$, $\tilde{\boldsymbol{\eta}}$, \tilde{R} and \tilde{f}^* analogously to the “clean” distribution P_{XY} .

72 We view (X, Y, \tilde{Y}) as a jointly distributed triple. Given $(X, Y) \sim P_{XY}$, the distribution of the noisy
 73 label \tilde{Y} is determined by the matrix-valued function $\mathbf{E} : \mathcal{X} \rightarrow \{\mathbf{M} \in \mathbb{R}^{K \times K} : \mathbf{M} \text{ is row stochastic}\}$,
 74 which is called the *noise transition matrix*, and whose (i, j) -th element is

$$[\mathbf{E}(x)]_{i,j} = \mathbb{P}(\tilde{Y} = j | Y = i, X = x).$$

75 Note that the noisy and clean class posteriors satisfy $\tilde{\boldsymbol{\eta}}(x) = \mathbf{E}(x)^\top \boldsymbol{\eta}(x)$, with $^\top$ denoting transpose.

76 The goal of learning from label noise is to find a classifier that is able to minimize the “clean test
 77 error,” that is, the risk R defined w.r.t. P_{XY} , even though the learner’s access is limited to corrupted
 78 training data $(X_i, \tilde{Y}_i) \stackrel{\text{i.i.d.}}{\sim} P_{X\tilde{Y}}$. In essence, label noise can be conceptualized as a type of domain
 79 adaptation problem, where P_{XY} is the source domain, $P_{X\tilde{Y}}$ is the target domain, and the source

80 and target are linked by “posterior drift”, which means that the source and target have the same
 81 X -marginal, but the “posteriors” (distribution of label given X) may change [Scott, 2019, Cai and
 82 Wei, 2021, Maity et al., 2023]. Thus, a label noise problem is captured by a triple $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})$.

83 3.2 Relative signal strength

84 To study label noise, we introduce the concept of *relative signal strength* (RSS). This is a pointwise
 85 measure of how much “signal” (certainty about the label) is contained in the noisy distribution relative
 86 to the clean distribution. Previous work [Cannings et al., 2020, Cai and Wei, 2021] has examined a
 87 related concept within the context of binary classification, under the restriction that clean and noisy
 88 Bayes classifiers are identical. Our definition incorporates multi-class classification and relaxes the
 89 requirement that the clean and noisy Bayes classifiers agree.

90 **Definition 1 (Relative Signal Strength)** For any class probability vectors $\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}$, define the relative
 91 signal strength (RSS) at $x \in \mathcal{X}$ as

$$\mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = \min_{j \in \mathcal{Y}} \frac{\max_i [\tilde{\boldsymbol{\eta}}(x)]_i - [\tilde{\boldsymbol{\eta}}(x)]_j}{\max_i [\boldsymbol{\eta}(x)]_i - [\boldsymbol{\eta}(x)]_j}, \quad (1)$$

92 where $0/0 := +\infty$. Furthermore, for $\kappa \in [0, \infty)$, denote the set of points whose RSS exceeds κ as

$$\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = \{x \in \mathcal{X} : \mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) > \kappa\}.$$

93 $\mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})$ is a point-wise measure of how much “signal” the noisy posterior contains about the
 94 clean posterior. To gain some intuition, first notice that if the noisy Bayes classifier predicts a different
 95 class than the clean Bayes classifier, the RSS is 0 by taking $j = \arg \max \tilde{\boldsymbol{\eta}}$ (assuming for simplicity
 96 that the $\arg \max$ is a singleton set). Now suppose the clean and noisy Bayes classifiers *do* make the
 97 same prediction at x , say i^* , and consider a fixed j . If

$$\frac{[\tilde{\boldsymbol{\eta}}(x)]_{i^*} - [\tilde{\boldsymbol{\eta}}(x)]_j}{[\boldsymbol{\eta}(x)]_{i^*} - [\boldsymbol{\eta}(x)]_j}$$

98 is small, it means that the clean Bayes classifier is relatively certain that j is not the correct clean
 99 label, while the noisy Bayes classifier is less certain that j is not the correct noisy label. Taking the
 100 minimum over j gives the relative signal strength at x . As we formalize in the next section, a large
 101 RSS at x ensures that a small (pointwise) noisy excess risk at x implies a small (pointwise) clean
 102 excess risk. To gain more intuition, consider the following examples.

103 **Example 1** When $\boldsymbol{\eta}(x) = [0 \ 1 \ 0]^\top$ and $\tilde{\boldsymbol{\eta}}(x) = [0.3 \ 0.6 \ 0.1]^\top$,

$$\mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = \min_{j \in \mathcal{Y}} \frac{\max_i [\tilde{\boldsymbol{\eta}}(x)]_i - [\tilde{\boldsymbol{\eta}}(x)]_j}{\max_i [\boldsymbol{\eta}(x)]_i - [\boldsymbol{\eta}(x)]_j} = \frac{[\tilde{\boldsymbol{\eta}}(x)]_2 - [\tilde{\boldsymbol{\eta}}(x)]_1}{[\boldsymbol{\eta}(x)]_2 - [\boldsymbol{\eta}(x)]_1} = \frac{0.6 - 0.3}{1 - 0} = 0.3.$$

104 Here, the clean Bayes classifier is absolutely certain about its prediction, while the noisy Bayes
 105 classifier is much less certain.

106 **Example 2** When $\boldsymbol{\eta}(x) = [0 \ 1 \ 0]^\top$ and $\tilde{\boldsymbol{\eta}}(x) = [0 \ 0 \ 1]^\top$,

$$\mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = \min_{j \in \mathcal{Y}} \frac{\max_i [\tilde{\boldsymbol{\eta}}(x)]_i - [\tilde{\boldsymbol{\eta}}(x)]_j}{\max_i [\boldsymbol{\eta}(x)]_i - [\boldsymbol{\eta}(x)]_j} = \frac{[\tilde{\boldsymbol{\eta}}(x)]_3 - [\tilde{\boldsymbol{\eta}}(x)]_3}{[\boldsymbol{\eta}(x)]_2 - [\boldsymbol{\eta}(x)]_3} = \frac{1 - 1}{1 - 0} = 0.$$

107 The zero signal strength results from $\tilde{\boldsymbol{\eta}}$ and $\boldsymbol{\eta}$ leading to different predictions.

108 **Example 3 (Comparison to KL divergence)** When $\boldsymbol{\eta}(x) = [0.05 \ 0.7 \ 0.25]^\top$, and $\tilde{\boldsymbol{\eta}}^{(1)}(x) =$
 109 $[0.25 \ 0.7 \ 0.05]^\top$, $\tilde{\boldsymbol{\eta}}^{(2)}(x) = [0.1 \ 0.6 \ 0.3]^\top$,

$$\frac{1}{\mathcal{D}_{\text{KL}}(\boldsymbol{\eta} \parallel \tilde{\boldsymbol{\eta}}^{(1)})} < \frac{1}{\mathcal{D}_{\text{KL}}(\boldsymbol{\eta} \parallel \tilde{\boldsymbol{\eta}}^{(2)})} \quad \text{while} \quad \mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}^{(1)}) > \mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}^{(2)}).$$

110 Here, $\tilde{\boldsymbol{\eta}}^{(2)}$ is “closer” to $\boldsymbol{\eta}$ in terms of KL divergence, but $\tilde{\boldsymbol{\eta}}^{(1)}$ provides more information in terms of
 111 predicting the $\arg \max$ of $\boldsymbol{\eta}$. There is no conflict: KL divergence considers the similarity between
 112 two (whole) distributions, while the task of classification only focuses on predicting the $\arg \max$.

113 This also illustrates why our notion of RSS is better suited for the label noise problem than other
 114 general-purpose distance measures between distributions.

115 A desirable learning scenario would be if $\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = \mathcal{X}$ for some large κ , indicating that the signal
 116 strength is big across the entire space. Unfortunately, this ideal situation is generally not achievable.
 117 To understand this limitation, we begin by making a mild assumption about the label noise:

$$\left| \arg \max \boldsymbol{\eta}(x) \right| \leq \left| \arg \max \tilde{\boldsymbol{\eta}}(x) \right|, \quad \text{almost surely } (P_X) \quad (\text{C})$$

118 where $|\cdot|$ denotes set cardinality, that is, that $\arg \max \tilde{\boldsymbol{\eta}}(x)$ contains at least as many entries as
 119 $\arg \max \boldsymbol{\eta}$. In other words, when there are ties for the clean Bayes prediction, the noisy Bayes
 120 prediction has at least as many ties. This assumption is reasonable because label noise typically
 121 introduces ambiguity by “confusing” the class probabilities, resulting in more ties.

122 **Proposition 1** *If (C) holds, then $\mathcal{A}_0(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = \{x \in \mathcal{X} : \arg \max \tilde{\boldsymbol{\eta}}(x) = \arg \max \boldsymbol{\eta}(x)\}$.*

123 This proposition, which is proved in Appendix A.1.1, indicates that \mathcal{A}_0 , the region with positive
 124 RSS, precisely corresponds to the area where the true and noisy Bayes classifiers agree. Accordingly,
 125 $\mathcal{X} \setminus \mathcal{A}_0$, the zero signal region, is the region where Bayes decision rules differ. The “region of strong
 126 signal,” \mathcal{A}_κ , is a subset of \mathcal{A}_0 . Since the clean and noisy Bayes classifiers will typically disagree for
 127 at least some x , $\mathcal{A}_0 \neq \mathcal{X}$ in general. We note that the (somewhat strong) assumption that $\mathcal{A}_0 = \mathcal{X}$
 128 has been made in prior studies [Cannings et al., 2020, Cai and Wei, 2021].

129 **Posterior Drift Model Class.** Now putting definitions together, we consider the posterior drift model
 130 Π defined over the triple $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})$. Let $\epsilon \in [0, 1]$, $\kappa \in (0, +\infty)$, and define

$$\Pi(\epsilon, \kappa) := \left\{ (P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) : P_X(\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})) \geq 1 - \epsilon \right\}.$$

131 This is a set of triples (label noise problems) such that \mathcal{A}_κ , the region with RSS at least κ , covers at
 132 least $1 - \epsilon$ of the probability mass. In the next section, we will demonstrate that the performance
 133 within \mathcal{A}_κ can be guaranteed, whereas learning outside the region \mathcal{A}_κ is provably challenging.

134 4 Upper and lower bounds

135 In this section, we establish both upper and lower bounds for excess risk under multi-class instance-
 136 dependent label noise.

137 4.1 Minimax lower bound

138 Our first theorem reveals a fundamental limit: no classifier trained using noisy data can surpass the
 139 constraints imposed by relative signal strength in a minimax sense. To state the theorem, we employ
 140 the following notation and terminology. Denote the noisy training data by $Z^n = \{(X_i, \tilde{Y}_i)\}_{i=1}^n$ *i.i.d.*
 141 $P_{X\tilde{Y}}$. A *learning rule* \hat{f} is an algorithm that takes Z^n and outputs a classifier. The risk $R(\hat{f})$ of a
 142 learning rule is a random variable, where the randomness is due to the draw Z^n .

143 **Theorem 1 (Minimax Lower Bound)** *Let $\epsilon \in [0, 1]$, $\kappa > 0$. Then*

$$\inf_{\hat{f}} \sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa)} \mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] \geq \frac{K-1}{K} \epsilon + \Omega \left(\frac{1}{\kappa} \sqrt{\frac{1}{n}} \right),$$

144 *where the inf is over all learning rules.*

145 *Proof Sketch.* The idea is to pick a jointly distributed triple (X, Y, \tilde{Y}) such that: Y is independent of
 146 \tilde{Y} for $x \in \mathcal{X} \setminus \mathcal{A}_\kappa$ and Y and \tilde{Y} are supported on the same classes for $x \in \mathcal{A}_\kappa$. Consider the two
 147 regions separately: the excess risk on $\mathcal{X} \setminus \mathcal{A}_\kappa$ is no less than $\frac{K-1}{K} \epsilon$, while the excess risk on \mathcal{A}_κ
 148 decreases with sample size n , and the rate is affected by the RSS κ . See Appendix A.1.2. ■

149 The proof offers insights into how label noise impacts the learning process: if the signal is low,
 150 learning is hopeless, and the learner incurs an irreducible error of $(1 - 1/K)\epsilon$; in the high signal
 151 region, the signal strength κ determines the rate of convergence. These aspects determine fundamental
 152 limits that no classifier trained only on noisy data can overcome without additional assumptions.

153 4.2 Upper bound

154 This subsection establishes an upper bound for NI-ERM, the empirical risk minimizer trained on
155 noisy data. This result implies that NI-ERM is minimax optimal, a potentially surprising result given
156 that NI-ERM is arguably the simplest approach one might consider.

157 **Theorem 2 (Excess Risk Upper Bound of NI-ERM)** *Let $\epsilon \in [0, 1], \kappa > 0$. Consider any*
158 *$(P_X, \eta, \tilde{\eta}) \in \Pi(\epsilon, \kappa)$, assume function class \mathcal{F} has Natarajan dimension V , and the noisy Bayes*
159 *classifier \tilde{f}^* belongs to \mathcal{F} . Let $\hat{f} \in \mathcal{F}$ be the ERM trained on $Z^n = \{(X_i, \tilde{Y}_i)\}_{i=1}^n$. Then for*
160 *$n > 2V$,*

$$\mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] \leq \epsilon + \mathcal{O} \left(\frac{1}{\kappa} \sqrt{\frac{V}{n}} \right).$$

161 Natarajan dimension is a multiclass analogue of VC dimension. The upper bound (Theorem 2)
162 matches the lower bound (which remains unchanged under the setting of Theorem 2 – see Appendix
163 A.1.2 for a refined statement) in terms of both the irreducible error and the rate of convergence. This
164 result is surprising as it indicates that the simplest possible approach, which ignores the presence of
165 noise, is optimal. No learning rule could perform significantly better in this minimax sense.

166 5 Practical implication

167 The modern practice of machine learning often involves training a deep neural network. In complex
168 tasks involving noisy labels, the naïve NI-ERM is often outperformed by state-of-the-art methods by
169 a significant extent [Li et al., 2020, Xiao et al., 2023]. This is consistent with the finding that directly
170 training a large neural network on noisy data frequently leads to overfitting [Zhang et al., 2021a].

171 Yet this is not grounds for abandoning NI-ERM altogether as a practical strategy. Instead of using
172 NI-ERM for end-to-end training of a deep neural network, we instead propose the following simple,
173 two-step procedure, termed ‘feature extraction + NI-ERM’.

- 174 1. Perform feature extraction using any method (e.g., transfer learning or self-supervised learning)
175 that does not require labels.
- 176 2. Learn a simple classifier (e.g., a linear classifier) on top of these extracted features, using the
177 noisily labelled data, in a noise-ignorant way.

178 We acknowledge the that practical idea of ignoring label noise is not new [Ghosh and Lan, 2021],
179 but the full power of this approach has not been previously recognized. For example, prior work
180 that has suggested ignoring the label noise usually augments this approach with additional heuristics
181 [Zheltonozhskii et al., 2022, Xue et al., 2022].

182 Remarkably, this two-step approach attains extremely strong performance. We conducted experiments
183 on the CIFAR image data under two scenarios: synthetic label flipping and realistic human label
184 errors [Wei et al., 2022], as shown in Figure 1. We examine three different feature extractors: the
185 DINOv2 foundation model [Oquab et al., 2023], ResNet-50 features extracted from training on
186 ImageNet [He et al., 2016], and self-supervised ResNet-50 using contrastive loss [Chen et al., 2020].
187 We also compared to a simple linear model trained on the raw pixel intensities, and a ResNet-50
188 trained end-to-end. We observed that ResNet-50 exhibits a “tempered overfitting” behavior, consistent
189 with previous findings [Zhang et al., 2021a, Mallinar et al., 2022]. The linear model demonstrates
190 robustness to noise, but suffers from significant approximation error.

191 Conversely, the FE+NI-ERM approach enjoys the best of both worlds. Regardless of how the feature
192 extraction is carried out, the resulting models exhibit robustness to label noise, while the overall
193 accuracy depends entirely on the quality of the extracted features. This is illustrated in Figure 1,
194 where the flatness of the accuracy curves as noise increases indicates the robustness, and the intercept
195 at zero label noise is a measure of the feature quality. Importantly, this property holds true even under
196 realistic label noise of CIFAR-N [Wei et al., 2022]. In fact, we find that using the DINOv2 [Oquab
197 et al., 2023] extracted features in our FE+NI-ERM approach yields state of the art results on the
198 CIFAR-10N and CIFAR-100N benchmarks, across the noise levels, as shown in Table 1.

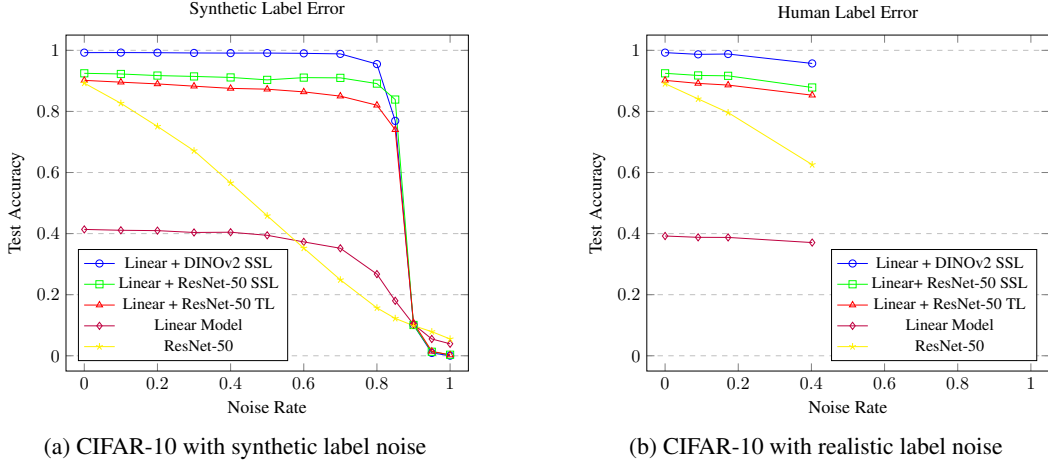


Figure 1: A linear model trained on features obtained from either transfer learning (pretrained ResNet-50 on ImageNet [He et al., 2016]), self-supervised learning (ResNet-50 trained on CIFAR-10 images with contrastive loss [Chen et al., 2020]), or a pretrained self-supervised foundation model DINOv2 [Oquab et al., 2023] significantly boosts the performance of the original linear model. In contrast, directly training a ResNet-50 leads to overfitting.

199 We reiterate that the only hyperparameters of our model are the hyperparameters of the linear classifier,
 200 which are tuned automatically using standard cross-validation on the noisy labels. This in contrast to
 201 the implementations of many methods on the CIFAR-N leaderboard (<http://noisylabels.com/>),
 202 where the hyperparameters are hard-coded. Altogether, this dominant performance, along with
 203 the simplicity of the approach and the lack of any untunable hyperparameters, suggest that the
 204 FE+NI-ERM is very powerful, and indicates a need for further investigation of its properties.

Table 1: Performance comparison with CIFAR-N leaderboard (<http://noisylabels.com/>) in terms of testing accuracy. “Aggre”, “Rand1”, ..., “Noisy” denote various types of human label noise. We compare with four methods that covers the top three performance for all noise categories: ProMix [Xiao et al., 2023], ILL [Chen et al., 2023], PLS [Albert et al., 2023] and DivideMix [Li et al., 2020]. Our approach, a Noise Ignorant linear model trained on features extracted by the self-supervised foundation model DINOv2 [Oquab et al., 2023] achieves new state-of-the-art results, highlighted in bold. We employed Python’s sklearn logistic regression and cross-validation functions without data augmentation; the results are deterministic and directly reproducible.

| Leaderboard | CIFAR-10N | | | | | CIFAR-100N |
|-------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|
| | Aggre | Rand1 | Rand2 | Rand3 | Worst | Noisy |
| ProMix | 97.65 ± 0.19 | 97.39 ± 0.16 | 97.55 ± 0.12 | 97.52 ± 0.09 | 96.34 ± 0.23 | 73.79 ± 0.28 |
| ILL | 96.40 ± 0.03 | 96.06 ± 0.07 | 95.98 ± 0.12 | 96.10 ± 0.05 | 93.55 ± 0.14 | 68.07 ± 0.33 |
| PLS | 96.09 ± 0.09 | 95.86 ± 0.26 | 95.96 ± 0.16 | 96.10 ± 0.07 | 93.78 ± 0.30 | 73.25 ± 0.12 |
| DivideMix | 95.01 ± 0.71 | 95.16 ± 0.19 | 95.23 ± 0.07 | 95.21 ± 0.14 | 92.56 ± 0.42 | 71.13 ± 0.48 |
| FE + NI-ERM | 98.69 ± 0.00 | 98.80 ± 0.00 | 98.65 ± 0.00 | 98.67 ± 0.00 | 95.71 ± 0.00 | 83.17 ± 0.00 |

205 6 Conclusions

206 This work presents a rigorous theory for learning under multi-class, instance-dependent label noise.
 207 We establish matching upper and lower bounds for excess risk. Our theory reveals the minimax
 208 optimality of Noise Ignorant Empirical Risk Minimizer (NI-ERM). To make this theory practical, we
 209 provide a simple modification leveraging a feature extractor with NI-ERM, demonstrating significant
 210 performance enhancements. A limitation of this work is that our methodology warrants more
 211 extensive experimental evaluation.

212 References

- 213 Paul Albert, Eric Arazo, Tarun Krishna, Noel E O'Connor, and Kevin McGuinness. Is your noise
214 correction noisy? pls: Robustness to label noise with two stage detection. In *Proceedings of the*
215 *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 118–127, 2023.
- 216 Dana Angluin and Philip Laird. Learning from noisy examples. *Machine learning*, 2:343–370, 1988.
- 217 Carla E Brodley and Mark A Friedl. Identifying mislabeled training data. *Journal of artificial*
218 *intelligence research*, 11:131–167, 1999.
- 219 T Tony Cai and Hongji Wei. Transfer learning for nonparametric classification: Minimax rate and
220 adaptive classifier. *The Annals of Statistics*, 49(1):100–128, 2021.
- 221 Timothy I Cannings, Yingying Fan, and Richard J Samworth. Classification with imperfect training
222 labels. *Biometrika*, 107(2):311–330, 2020.
- 223 Hao Chen, Ankit Shah, Jindong Wang, Ran Tao, Yidong Wang, Xing Xie, Masashi Sugiyama, Rita
224 Singh, and Bhiksha Raj. Imprecise label learning: A unified framework for learning with various
225 imprecise label configurations. *arXiv preprint arXiv:2305.12715*, 2023.
- 226 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
227 contrastive learning of visual representations. In *International conference on machine learning*,
228 pages 1597–1607. PMLR, 2020.
- 229 Hao Cheng, Zhaowei Zhu, Xingyu Li, Yifei Gong, Xing Sun, and Yang Liu. Learning with instance-
230 dependent label noise: A sample sieve approach. In *International Conference on Learning*
231 *Representations*, 2021.
- 232 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
233 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
234 pages 248–255. IEEE, 2009.
- 235 Luc Devroye, László Györfi, and Gabor Lugosi. *A Probabilistic Theory of Pattern Recognition*,
236 volume 31. Springer, 1996.
- 237 Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization
238 for efficiently improving generalization. In *International Conference on Learning Representations*,
239 2021. URL <https://openreview.net/forum?id=6Tm1mposlRM>.
- 240 Aritra Ghosh and Andrew Lan. Contrastive learning improves model robustness under label noise. In
241 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
242 2703–2708, 2021.
- 243 Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi
244 Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels.
245 *Advances in neural information processing systems*, 31, 2018.
- 246 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
247 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
248 pages 770–778, 2016.
- 249 Hyunki Im and Paul Grigas. Binary classification with instance and label dependent label noise.
250 *arXiv preprint arXiv:2306.03402*, 2023.
- 251 Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as semi-
252 supervised learning. In *International Conference on Learning Representations*, 2020.
- 253 Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-learning
254 regularization prevents memorization of noisy labels. *Advances in neural information processing*
255 *systems*, 33:20331–20342, 2020.
- 256 Tongliang Liu and Dacheng Tao. Classification with noisy labels by importance reweighting. *IEEE*
257 *Transactions on pattern analysis and machine intelligence*, 38(3):447–461, 2015.

- 258 Yang Liu and Hongyi Guo. Peer loss functions: Learning from noisy labels without knowing noise
259 rates. In *International conference on machine learning*, pages 6226–6236. PMLR, 2020.
- 260 Subha Maity, Diptavo Dutta, Jonathan Terhorst, Yuekai Sun, and Moulinath Banerjee. A linear
261 adjustment based approach to posterior drift in transfer learning. *Biometrika*, 2023. URL <https://doi.org/10.1093/biomet/asad029>.
262
- 263 Neil Mallinar, James Simon, Amirhesam Abedsoltan, Parthe Pandit, Misha Belkin, and Preetum
264 Nakkiran. Benign, tempered, or catastrophic: Toward a refined taxonomy of overfitting. *Advances
265 in Neural Information Processing Systems*, 35:1182–1195, 2022.
- 266 Balas K Natarajan. On learning sets and functions. *Machine Learning*, 4:67–97, 1989.
- 267 Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with
268 noisy labels. *Advances in neural information processing systems*, 26, 2013.
- 269 Curtis Northcutt, Lu Jiang, and Isaac Chuang. Confident learning: Estimating uncertainty in dataset
270 labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021.
- 271 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V Vo, Marc Szafraniec, Vasil Khalidov,
272 Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
273 robust visual features without supervision. *Transactions on Machine Learning Research*, 2023.
- 274 Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making
275 deep neural networks robust to label noise: A loss correction approach. In *Proceedings of the
276 IEEE conference on computer vision and pattern recognition*, pages 1944–1952, 2017.
- 277 Clayton Scott. A generalized Neyman-Pearson criterion for optimal domain adaptation. In *Algorithmic
278 Learning Theory*, pages 738–761. PMLR, 2019.
- 279 Brendan Van Rooyen and Robert C Williamson. A theory of learning with corrupted labels. *Journal
280 of Machine Learning Research*, 18(228):1–50, 2018.
- 281 V. N. Vapnik and A. Ya. Chervonenkis. On the uniform convergence of relative frequencies of
282 events to their probabilities. *Theory of Probability & Its Applications*, 16(2):264–280, 1971. doi:
283 10.1137/1116025. URL <https://doi.org/10.1137/1116025>.
- 284 Jialu Wang, Eric Xin Wang, and Yang Liu. Estimating instance-dependent label-noise transition
285 matrix using a deep neural network. In *International Conference on Machine Learning*, 2022.
- 286 Jiaheng Wei, Zhaowei Zhu, Hao Cheng, Tongliang Liu, Gang Niu, and Yang Liu. Learning with
287 noisy labels revisited: A study using real-world human annotations. In *International Conference
288 on Learning Representations*, 2022.
- 289 Ruixuan Xiao, Yiwen Dong, Haobo Wang, Lei Feng, Runze Wu, Gang Chen, and Junbo Zhao.
290 Promix: Combating label noise via maximizing clean sample utility. In Edith Elkind, editor,
291 *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-
292 23*, pages 4442–4450. International Joint Conferences on Artificial Intelligence Organization, 8
293 2023. doi: 10.24963/ijcai.2023/494. URL <https://doi.org/10.24963/ijcai.2023/494>.
294 Main Track.
- 295 Yihao Xue, Kyle Whitecross, and Baharan Mirzasoleiman. Investigating why contrastive learning
296 benefits robustness against label noise. In *International Conference on Machine Learning*, pages
297 24851–24871. PMLR, 2022.
- 298 Shuo Yang, Songhua Wu, Erkun Yang, Bo Han, Yang Liu, Min Xu, Gang Niu, and Tongliang Liu. A
299 parametrical model for instance-dependent label noise. *IEEE Transactions on Pattern Analysis
300 and Machine Intelligence*, 2023.
- 301 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep
302 learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115,
303 2021a.

- 304 Jianxin Zhang, Yutong Wang, and Clayton Scott. Learning from label proportions by learning with
305 label noise. *Advances in Neural Information Processing Systems*, 35:26933–26942, 2022.
- 306 Mingyuan Zhang, Jane Lee, and Shivani Agarwal. Learning from noisy labels with no change to the
307 training process. In *International Conference on Machine Learning*, pages 12468–12478. PMLR,
308 2021b.
- 309 Evgenii Zheltonozhskii, Chaim Baskin, Avi Mendelson, Alex M Bronstein, and Or Litany. Contrast to
310 divide: Self-supervised pre-training for learning with noisy labels. In *Proceedings of the IEEE/CVF*
311 *Winter Conference on Applications of Computer Vision*, pages 1657–1667, 2022.
- 312 Zhaowei Zhu, Tongliang Liu, and Yang Liu. A second-order approach to learning with instance-
313 dependent label noise. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
314 *recognition*, pages 10113–10123, 2021.

315 **A Appendix / supplemental material**

316 **A.1 Proofs**

317 **A.1.1 Proof of Proposition 1**

318 **Proposition** *Assume (C) holds. Then for $\kappa \geq 0$,*

$$\mathcal{A}_0(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = \left\{ x \in \mathcal{X} : \arg \max \tilde{\boldsymbol{\eta}}(x) = \arg \max \boldsymbol{\eta}(x) \right\} \text{ almost surely.}$$

319 *Proof.* Notice that

$$\mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = 0 \iff \arg \max \tilde{\boldsymbol{\eta}}(x) \not\subseteq \arg \max \boldsymbol{\eta}(x).$$

320 This is because $\mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = 0$ when the numerator is zero and the denominator is non-zero, which
321 happens when $\arg \max \tilde{\boldsymbol{\eta}}(x) \not\subseteq \arg \max \boldsymbol{\eta}(x)$. An equivalent statement of this is

$$\mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) > 0 \iff \arg \max \tilde{\boldsymbol{\eta}}(x) \subseteq \arg \max \boldsymbol{\eta}(x).$$

322 Under assumption (C), the set $\arg \max \tilde{\boldsymbol{\eta}}(x)$ cannot be a proper subset of $\arg \max \boldsymbol{\eta}(x)$, and therefore

$$\mathcal{M}(x; \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) > 0 \iff \arg \max \tilde{\boldsymbol{\eta}}(x) = \arg \max \boldsymbol{\eta}(x).$$

323 Thus,

$$\mathcal{A}_0(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) = \left\{ x \in \mathcal{X} : \arg \max \tilde{\boldsymbol{\eta}}(x) = \arg \max \boldsymbol{\eta}(x) \right\},$$

324 almost everywhere. ■

325 **A.1.2 Proof of Lower Bound: Theorem 1**

326 Now we provide a more formal statement of the minimax lower bound and its proof. We begin with
327 the scenario where the noisy distribution $P_{X\tilde{Y}}$ has zero Bayes risk as an introductory example. The
328 proof for the general case follows a similar strategy but involves more complex bounding techniques.
329 We recommend that interested readers first review the proof of the zero-error version to build a solid
330 understanding before tackling the general case.

331 Now consider a more restricted set of $\Pi(\epsilon, \kappa)$:

$$\Pi(\epsilon, \kappa, V, 0) := \left\{ (P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) : P_X \left(\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \right) \geq 1 - \epsilon, P_X \text{ supported on } V + 1 \text{ points, } \tilde{R}^* = 0 \right\}.$$

332 **Theorem (Minimax Lower Bound: when $\tilde{R}^* = 0$)** *Let $\epsilon \in [0, 1], \kappa > 0, V > 1$. For any learning*
333 *rule \hat{f} based upon $Z^n = \{(X_i, \tilde{Y}_i)\}_{i=1}^n$, and $n > \max(V - 1, 2)$,*

$$\begin{aligned} \sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa)} \mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] &\geq \sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa, V, 0)} \mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] \\ &\geq \frac{K-1}{K} \epsilon + \frac{1}{\kappa} \frac{(V-1)(1-\epsilon)}{8en} \end{aligned}$$

334 *Proof.* Consider any $V + 1$ distinct points x_0, x_1, \dots, x_V . Choose

$$P_X(x) = \begin{cases} \epsilon & x = x_0 \\ (1-\epsilon) \cdot \frac{1}{n} & x = x_1, \dots, x_{V-1} \\ (1-\epsilon) \cdot \left(1 - \frac{V-1}{n}\right) & x = x_V. \end{cases}$$

335 where $n > V - 1$, and define the clean and noisy class posteriors by

$$\text{If } x = x_0, \text{ then } \boldsymbol{\eta}(x) = \mathbf{e}_j, \tilde{\boldsymbol{\eta}}(x) = \mathbf{e}_1, \quad j \in \{1, 2, \dots, K\} \quad (2)$$

$$\text{If } x = x_t, 1 \leq t \leq V - 1, \text{ then } \boldsymbol{\eta}(x) = \begin{bmatrix} \frac{1}{2} + \frac{1}{2(\kappa+\delta)} \cdot (-1)^{b_t+1} \\ \frac{1}{2} - \frac{1}{2(\kappa+\delta)} \cdot (-1)^{b_t+1} \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \tilde{\boldsymbol{\eta}}(x) = \mathbf{e}_{b_t}, b_t \in \{1, 2\}, \delta > 0, \quad (3)$$

$$\text{If } x = x_V, \text{ then } \boldsymbol{\eta}(x) = \begin{bmatrix} \frac{1}{2} + \frac{1}{2(\kappa+\delta)} \\ \frac{1}{2} - \frac{1}{2(\kappa+\delta)} \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \tilde{\boldsymbol{\eta}}(x) = \mathbf{e}_1, \quad (4)$$

336 where \mathbf{e}_i denotes the one-hot vector whose i -th element is one.

337 The triple $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})$ is thus parameterized by $j, \mathbf{b} := [b_1 \ b_2 \ \dots \ b_{V-1}]^\top$, and δ .

338 This construction ensures $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa, V, 0)$. In particular,

$$\begin{aligned} \mathcal{A}_\kappa &\supseteq \{x_1, x_2, \dots, x_V\}, & P_X(\mathcal{A}_\kappa) &\geq 1 - \epsilon, \\ \mathcal{X} \setminus \mathcal{A}_\kappa &\subseteq \{x_0\}, & P_X(\mathcal{X} \setminus \mathcal{A}_\kappa) &\leq \epsilon, \end{aligned}$$

339 and $\tilde{R}^* = 0$ because $\tilde{\boldsymbol{\eta}}(x)$ is one-hot for all x .

340 For classifier any f , by definition, its risk equals to

$$\begin{aligned} R(f) &= \mathbb{E}_{X,Y} [\mathbb{1}_{f(X) \neq Y}] \\ &= \mathbb{E}_X \mathbb{E}_{Y|X} [\mathbb{1}_{f(X) \neq Y}] \\ &= \mathbb{E}_X \mathbb{E}_{Y|X} [1 - \mathbb{1}_{f(X)=Y}] \\ &= \mathbb{E}_X [1 - [\boldsymbol{\eta}(X)]_{f(X)}] \\ &= \int_{\mathcal{X}} (1 - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x), \end{aligned}$$

341 Under our construction of P_X , $R(f)$ can be decomposed into two parts

$$R(f) = \underbrace{\int_{\{x_0\}} (1 - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x)}_{:=R_0(f)} + \underbrace{\int_{\{x_1, \dots, x_V\}} (1 - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x)}_{:=R_V(f)},$$

342 so does the excess risk

$$R(f) - R(f^*) = (R_0(f) - R_0(f^*)) + (R_V(f) - R_V(f^*)).$$

343 Recall that in our construction, $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})$ is parameterized by j, \mathbf{b} , and δ . Therefore

$$\begin{aligned} \sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa, V, 0)} \mathbb{E}_{Z^n} [R(\hat{f}) - R(f^*)] &\geq \sup_{j, \mathbf{b}, \delta} \mathbb{E}_{Z^n} [R(\hat{f}) - R(f^*)] \\ &= \sup_{j, \mathbf{b}, \delta} \left\{ \mathbb{E}_{Z^n} [R_0(\hat{f}) - R_0(f^*)] \right. \\ &\quad \left. + \mathbb{E}_{Z^n} [R_V(\hat{f}) - R_V(f^*)] \right\} \\ &= \sup_j \mathbb{E}_{Z^n} [R_0(\hat{f}) - R_0(f^*)] \\ &\quad + \sup_{\mathbf{b}, \delta} \mathbb{E}_{Z^n} [R_V(\hat{f}) - R_V(f^*)] \end{aligned}$$

344 where the last equality holds because region $\{x_0\}$ only depends on j , while region $\{x_1, \dots, x_V\}$
 345 only depends on \mathbf{b}, δ .

346 In the remaining part of the proof, we will examine

$$\sup_j \mathbb{E}_{Z^n} \left[R_0(\hat{f}) - R_0(f^*) \right] \quad (5)$$

347 and

$$\sup_{\mathbf{b}, \delta} \mathbb{E}_{Z^n} \left[R_V(\hat{f}) - R_V(f^*) \right] \quad (6)$$

348 separately.

349 Let's start with the first term (5), which acts over the “low signal strength” region $\{x_0\}$. Since $\boldsymbol{\eta}$ is
 350 one-hot on $\{x_0\}$, its Bayes risk over that is zero

$$\begin{aligned} \sup_j \mathbb{E}_{Z^n} \left[R_0(\hat{f}) - R_0(f^*) \right] &= \sup_j \mathbb{E}_{Z^n} \left[R_0(\hat{f}) \right] \\ &= \sup_j \mathbb{E}_{Z^n} \left[\int_{\{x_0\}} \mathbb{1}_{\hat{f}(x) \neq j} dP_X(x) \right]. \end{aligned}$$

351 To deal with \sup_j , we use a technique called “the probabilistic method”: replace j with a random
 352 variable $J \sim \text{Uniform}\{1, 2, \dots, K\}$:

$$\begin{aligned} \sup_j \mathbb{E}_{Z^n} \left[\int_{\{x_0\}} \mathbb{1}_{\hat{f}(x) \neq j} dP_X(x) \right] &\geq \mathbb{E}_{J, Z^n} \left[\int_{\{x_0\}} \mathbb{1}_{\hat{f}(x) \neq J} dP_X(x) \right] \\ &= \mathbb{E}_{Z^n} \left[\mathbb{E}_{J|Z^n} \left[\int_{\{x_0\}} \mathbb{1}_{\hat{f}(x) \neq J} dP_X(x) \right] \right]. \end{aligned}$$

353 Again, notice that J is an independent draw. Even if the point x_0 is observed in Z^n , the associated
 354 noisy label $\tilde{Y} = 1$ does not give any information about the clean label $Y = J$. Thus

$$\begin{aligned} \mathbb{E}_{Z^n} \left[\mathbb{E}_{J|Z^n} \left[\int_{\{x_0\}} \mathbb{1}_{\hat{f}(x) \neq J} dP_X(x) \right] \right] &= \mathbb{E}_{Z^n} \left[\mathbb{E}_J \left[\int_{\{x_0\}} \mathbb{1}_{\hat{f}(x) \neq J} dP_X(x) \right] \right] \\ &= \mathbb{E}_{Z^n} \left[\int_{\{x_0\}} \mathbb{E}_J \left[\mathbb{1}_{\hat{f}(x) \neq J} \right] dP_X(x) \right] \\ &= \mathbb{E}_{Z^n} \left[\int_{\{x_0\}} \left(1 - \frac{1}{K} \right) dP_X(x) \right] \\ &= \left(1 - \frac{1}{K} \right) \epsilon. \end{aligned}$$

355 Now we have the minimax lower bound for the first part (5):

$$\sup_j \mathbb{E}_{Z^n} \left[R_{\{x_0\}}(\hat{f}) - R_{\{x_0\}}(f^*) \right] \geq \left(1 - \frac{1}{K} \right) \epsilon.$$

356 For the second part (6), which is over $\{x_1, \dots, x_V\}$, due to the relative signal strength condition, and
 357 from our explicit construction in Eqn. (3) and (4), the excess risk w.r.t. true and noisy distribution are
 358 related by

$$\begin{aligned} R_V(f) - R_V(f^*) &= \int_{\{x_1, \dots, x_V\}} \left(\max \boldsymbol{\eta}(x) - [\boldsymbol{\eta}(x)]_{f(x)} \right) dP_X(x) \\ &= \int_{\{x_1, \dots, x_V\}} \frac{1}{\kappa + \delta} \left(\max \tilde{\boldsymbol{\eta}}(x) - [\tilde{\boldsymbol{\eta}}(x)]_{f(x)} \right) dP_X(x) \quad \because \text{by construction of } \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}} \\ &= \frac{1}{\kappa + \delta} \left(\tilde{R}_V(f) - \tilde{R}_V(f^*) \right), \end{aligned}$$

359 where $\tilde{R}_V(f) := \int_{\{x_1, \dots, x_V\}} (1 - [\tilde{\eta}(x)]_{f(x)}) dP_X(x)$. Also note that $f^*(x) = \tilde{f}^*(x)$ for $x \in$
 360 $\{x_1, \dots, x_V\}$, which is a result of our construction of $\eta, \tilde{\eta}$.

361 Then

$$\sup_{\mathbf{b}, \delta} \mathbb{E}_{Z^n} \left[R_V(\hat{f}) - R_V(f^*) \right] = \sup_{\mathbf{b}, \delta} \mathbb{E}_{Z^n} \left[\frac{1}{\kappa + \delta} \left(\tilde{R}_V(f) - \tilde{R}_V(\tilde{f}^*) \right) \right].$$

362 This allows us to reduce the label noise problem into a standard learning problem: we have an iid
 363 sample Z^n from $P_{X\tilde{Y}}$ and consider the risk evaluated on the same distribution $P_{X\tilde{Y}}$. The remainder
 364 of the proof is similar to the proof of Theorem 14.1 in Devroye et al. [1996].

365 Notice that by our construction, \tilde{Y} is a deterministic function of X . To be specific, $\tilde{Y} = \tilde{f}^*(X)$,
 366 where

$$\tilde{f}^*(x) = \begin{cases} 1 & x = x_0, \\ b_t & x = x_t, 1 \leq t \leq V-1 \\ 1 & x = x_V \end{cases}$$

367 is the noisy Bayes classifier.

368 We use the shorthand $f_{\mathbf{b}} := \tilde{f}^*$ to denote that the noisy Bayes classifier depends on \mathbf{b} , we also write
 369 the learning rule $\hat{f}(x; Z^n) := \hat{f}(x)$ to indicate the dependence of the learning rule \hat{f} on the random
 370 sample Z^n .

371 Since the noisy Bayes risk is zero,

$$\sup_{\mathbf{b}, \delta} \mathbb{E}_{Z^n} \left[\frac{1}{\kappa + \delta} \left(\tilde{R}_V(f) - \tilde{R}_V(\tilde{f}^*) \right) \right] = \sup_{\mathbf{b}, \delta} \frac{1}{\kappa + \delta} \mathbb{E}_{Z^n} \left[\tilde{R}_V(\hat{f}) \right].$$

372 Again, use the probabilistic method, replace \mathbf{b} with $\mathbf{B} \sim \text{Uniform}\{1, 2\}^{V-1}$,

$$\begin{aligned} \sup_{\mathbf{b}, \delta} \frac{1}{\kappa + \delta} \mathbb{E}_{Z^n} \left[\tilde{R}_V(\hat{f}) \right] &\geq \sup_{\delta} \frac{1}{\kappa + \delta} \mathbb{E}_{\mathbf{B}, Z^n} \left[\tilde{R}_V(\hat{f}) \right] \\ &= \sup_{\delta} \frac{1}{\kappa + \delta} \mathbb{E}_{Z^n} \left[\mathbb{E}_{\mathbf{B}|Z^n} \left[\int_{\{x_1, \dots, x_V\}} \mathbb{1}_{\hat{f}(x) \neq f_{\mathbf{B}}(x)} dP_X(x) \right] \right] \end{aligned}$$

373 Since we have $\mathbf{B} \sim \text{Uniform}\{1, 2\}^{V-1}$ and also $Z|\mathbf{B} \sim P_{X\tilde{Y}}$, then by Bayes rule (or eye-balling),
 374 we get the posterior distribution of $\mathbf{B}|Z^n$, to be specific:

$$\forall x \in \{x_1, \dots, x_V\}, \mathbb{P}(f_{\mathbf{B}}(x) = 1|Z^n) = \begin{cases} \frac{1}{2} & x \neq X_1, \dots, x \neq X_n, x \neq x_V \\ 0 \text{ or } 1 & \text{otherwise,} \end{cases}$$

375 where we overload the notation \mathbb{P} to denote conditional probability of $\mathbf{B}|Z^n$.

376 Then the optimal decision rule for predicting \mathbf{B} based on sample Z^n is:

$$f_{\mathbf{B}}^*(x; Z^n) := \begin{cases} \tilde{Y}_i & x = X_i, i \in \{1, 2, \dots, n\} \\ 1 & x = x_V \\ \text{random guess from } \{1, 2\} & x \neq X_1, \dots, x \neq X_n, x \neq x_V. \end{cases}$$

377 Therefore, roughly speaking, the error comes from the probability of $X \in \{x_1, \dots, x_V\}$ not being
 378 one of observed X_i : for any \hat{f} ,

$$\begin{aligned}
\mathbb{E}_{\mathbf{B}, Z^n} \left[\tilde{R}_V(\hat{f}) \right] &= \mathbb{E}_{Z^n} \left[\mathbb{E}_{\mathbf{B}|Z^n} \left[\int_{\{x_1, \dots, x_V\}} \mathbb{1}_{\hat{f}(x) \neq f_{\mathbf{B}}(x)} dP_X(x) \right] \right] \\
&\geq \mathbb{P}(X \in \{x_1, \dots, x_V\}, f_{\mathbf{B}}^*(X; Z^n) \neq f_{\mathbf{B}}(X)) \quad \because \text{error of } \hat{f} > \text{error of } f_{\mathbf{B}}^* \\
&= \left(1 - \frac{1}{2}\right) \mathbb{P}(X \neq X_1, \dots, X \neq X_n, X \neq x_V, X \in \{x_1, \dots, x_V\}) \\
&= \frac{1}{2} \sum_{t=1}^V \mathbb{P}(X \neq X_1, \dots, X \neq X_n, X \neq x_V, X = x_t) \\
&= \frac{1}{2} \sum_{t=1}^V \mathbb{P}(X_1 \neq x_t, \dots, X_n \neq x_t, x_V \neq x_t, X = x_t) \quad \because \text{replace all } X \text{ with } x_t \\
&= \frac{1}{2} \sum_{t=1}^{V-1} \mathbb{P}(X_1 \neq x_t, \dots, X_n \neq x_t, X = x_t) \\
&= \frac{1}{2} \sum_{t=1}^{V-1} \mathbb{P}(X_1 \neq x_t, \dots, X_n \neq x_t | X = x_t) \mathbb{P}(X = x_t) \\
&= \frac{1}{2} \sum_{t=1}^{V-1} (1 - \mathbb{P}(X = x_t))^n \mathbb{P}(X = x_t) \\
&= \frac{1}{2} (V-1) \left(1 - \frac{1-\epsilon}{n}\right)^n \left(\frac{1-\epsilon}{n}\right) \\
&= \frac{(V-1)(1-\epsilon)}{2n} \left(1 - \frac{1-\epsilon}{n}\right)^n \\
&= \frac{(V-1)(1-\epsilon)}{2n} \left(1 - \frac{1-\epsilon}{n}\right)^{1+\epsilon} \left(1 - \frac{1-\epsilon}{n}\right)^{n-1-\epsilon} \\
&\geq \frac{(V-1)(1-\epsilon)}{2n} \left(1 - \frac{1-\epsilon}{n}\right)^{1+\epsilon} e^{-1+\epsilon} \quad \because \left(1 - \frac{1-\epsilon}{n}\right)^{n-1-\epsilon} \downarrow e^{-1+\epsilon} \\
&\geq \frac{(V-1)(1-\epsilon)}{2n} \left(1 - \frac{1}{n}\right)^2 e^{-1} \quad \because \epsilon \in [0, 1] \\
&\geq \frac{(V-1)(1-\epsilon)}{2n} \frac{e^{-1}}{4} = \frac{(V-1)(1-\epsilon)}{8en}. \quad \text{take } n > 2
\end{aligned}$$

379 Now we get the minimax risk for the second part (6)

$$\begin{aligned}
\sup_{\mathbf{b}, \delta} \mathbb{E}_{Z^n} \left[R_{\mathcal{A}_\kappa}(\hat{f}) - R_{\mathcal{A}_\kappa}(f^*) \right] &\geq \sup_{\delta} \frac{1}{\kappa + \delta} \frac{(V-1)(1-\epsilon)}{8en} \\
&\geq \frac{1}{\kappa} \frac{(V-1)(1-\epsilon)}{8en} \quad \text{let } \delta \downarrow 0
\end{aligned}$$

380 Combine the two parts together, we get the final result, for $n > \max(V-1, 2)$

$$\sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa, V, 0)} \mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] \geq \frac{K-1}{K} \epsilon + \frac{1}{\kappa} \frac{(V-1)(1-\epsilon)}{8en}.$$

382 As for the general version of the lower bound, now consider the set of triples:

$$\Pi(\epsilon, \kappa, V, L) := \left\{ (P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) : P_X(\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})) \geq 1 - \epsilon, \right. \\ \left. P_X \text{ supported on } V + 1 \text{ points, } \frac{\tilde{R}_{\mathcal{A}_\kappa}(\tilde{f}^*)}{P_X(\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}))} = L \right\},$$

383 where $\tilde{R}_C(f) = \int_C (1 - [\tilde{\boldsymbol{\eta}}(x)]_{f(x)}) dP_X(x)$.

384 **Theorem (Minimax Lower Bound (General Version))** Let $\epsilon \in [0, 1], \kappa > 0, V > 1, L \in (0, 1/2)$.

385 For any learning rule \hat{f} based upon $Z^n = \{(X_i, \tilde{Y}_i)\}_{i=1}^n$, for $n \geq \frac{V-1}{2L} \max\left\{9, \frac{1}{(1-2L)^2}\right\}$

$$\sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa)} \mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] \geq \sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa, V, L)} \mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] \\ \geq \frac{K-1}{K} \epsilon + \frac{1-\epsilon}{\kappa} \sqrt{\frac{(V-1)L}{72n}} e^{-8}.$$

386 *Proof.*

387 Consider any $V + 1$ distinct points x_0, x_1, \dots, x_V . Choose

$$P_X(x) = \begin{cases} \epsilon & x = x_0 \\ (1-\epsilon) \cdot p & x = x_1, \dots, x_{V-1} \\ (1-\epsilon) \cdot (1-(V-1)p) & x = x_V. \end{cases}$$

388 This imposes the constraint $(V-1)p \leq 1$, which will be satisfied in the end. Notice the differ-
389 ence compared to the previous zero-error proof: we place probability mass p , rather than $1/n$, on
390 x_1, \dots, x_{V-1} .

391 As for the clean and noisy class probabilities, choose

$$\text{If } x = x_0, \text{ then } \boldsymbol{\eta}(x) = \mathbf{e}_j, \tilde{\boldsymbol{\eta}}(x) = \mathbf{e}_1, \quad j \in \{1, 2, \dots, k\} \quad (7)$$

$$\text{If } x = x_t, 1 \leq t \leq V-1, \text{ then } \boldsymbol{\eta}(x) = \begin{bmatrix} \frac{1}{2} + \frac{c}{\kappa+\delta} \cdot (-1)^{b_t+1} \\ \frac{1}{2} - \frac{c}{\kappa+\delta} \cdot (-1)^{b_t+1} \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \tilde{\boldsymbol{\eta}}(x) = \begin{bmatrix} \frac{1}{2} + c \cdot (-1)^{b_t+1} \\ \frac{1}{2} - c \cdot (-1)^{b_t+1} \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

$$b_t \in \{1, 2\}, \delta > 0, c \in \left(0, \frac{1}{2}\right) \quad (8)$$

$$x = x_V, \quad \boldsymbol{\eta}(x) = \begin{bmatrix} \frac{1}{2} + \frac{1}{2(\kappa+\delta)} \\ \frac{1}{2} - \frac{1}{2(\kappa+\delta)} \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \tilde{\boldsymbol{\eta}}(x) = \mathbf{e}_1, \quad (9)$$

392 where \mathbf{e}_i denotes the one-hot vector whose i -th element is one.

393 The construction for class posterior is also similar to the previous proof, except that for $x =$
394 $x_t, t \in \{1, \dots, V-1\}$, $\tilde{\boldsymbol{\eta}}$ is no longer a one-hot vector, rather has class probability separated by $2c$:

$$395 \left| [\tilde{\boldsymbol{\eta}}(x)]_1 - [\tilde{\boldsymbol{\eta}}(x)]_2 \right| = 2c.$$

396 The triple $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})$ can be parameterized by $j, \mathbf{b} := [b_1 \ b_2 \ \dots \ b_{V-1}]^\top, \delta, c$ and p .

397 Again, this construction ensures $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa)$, to be specific:

$$\mathcal{A}_\kappa \supseteq \{x_1, x_2, \dots, x_V\}, \quad P_X(\mathcal{A}_\kappa) \geq 1 - \epsilon, \\ \mathcal{X} \setminus \mathcal{A}_\kappa \subseteq \{x_0\}, \quad P_X(\mathcal{X} \setminus \mathcal{A}_\kappa) \leq \epsilon,$$

398 For any classifier f , its risk can be decomposed into two parts

$$R(f) = \underbrace{\int_{\{x_0\}} (1 - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x)}_{:=R_0(f)} + \underbrace{\int_{\{x_1, \dots, x_V\}} (1 - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x)}_{:=R_V(f)},$$

399 so does the excess risk

$$R(f) - R(f^*) = (R_0(f) - R_0(f^*)) + (R_V(f) - R_V(f^*)).$$

400 In our construction, $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})$ is parameterized by $j, \mathbf{b} := [b_1 \ b_2 \ \dots \ b_{V-1}]^\top, \delta, c$ and p , therefore

$$\sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa, V, L)} \mathbb{E}_{Z^n} [R(\hat{f}) - R(f^*)] \geq \sup_j \mathbb{E}_{Z^n} [R_0(\hat{f}) - R_0(f^*)] \quad (10)$$

$$+ \sup_{\mathbf{b}, \delta, c, p} \mathbb{E}_{Z^n} \left[\frac{1}{\kappa + \delta} (R_V(f) - R_V(\tilde{f}^*)) \right]. \quad (11)$$

401 Note that we have used the fact that

$$R_V(f) - R_V(f^*) = \frac{1}{\kappa + \delta} (R_V(f) - R_V(\tilde{f}^*)),$$

402 where $R_V(f) := \int_{\{x_1, \dots, x_V\}} (1 - [\tilde{\boldsymbol{\eta}}(x)]_{f(x)}) dP_X(x)$.

403 The first part (10) is exactly the same as the that of the zero-error proof, we have

$$\sup_j \mathbb{E}_{Z^n} [R_0(\hat{f}) - R_0(f^*)] \geq \left(1 - \frac{1}{K}\right) \epsilon.$$

404 From this point forward, the procedure is similar to the proof of Theorem 14.5 in Devroye et al.
405 [1996]. For the second part (11), the noisy Bayes classifier is still

$$\tilde{f}^*(x) = \begin{cases} j & x = x_0, \\ b_t & x = x_t, 1 \leq t \leq V \\ 1 & x = x_V. \end{cases}$$

406 We also use the shorthand $f_{\mathbf{b}} := \tilde{f}^*$ to denote that the noisy Bayes classifier depends on \mathbf{b} . The noisy
407 Bayes risk is no longer zero. In fact

$$\tilde{R}_V(\tilde{f}^*) = (V-1)(1-\epsilon)p \left(\frac{1}{2} - c \right)$$

408 Notice that our construction $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa, V, L)$, by definition

$$\frac{\tilde{R}_{\mathcal{A}_\kappa}(\tilde{f}^*)}{P_X(\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}))} = L,$$

409 therefore

$$L = \frac{\tilde{R}_{\mathcal{A}_\kappa}(\tilde{f}^*)}{P_X(\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}))} \leq \frac{\tilde{R}_V(\tilde{f}^*)}{P_X(\{x_1, \dots, x_V\})} = (V-1)p \left(\frac{1}{2} - c \right), \quad (12)$$

410 where the inequality holds from $\tilde{R}_{\mathcal{A}_\kappa}(\tilde{f}^*) = \tilde{R}_V(\tilde{f}^*)$ and $P_X(\mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})) \geq P_X(\{x_1, \dots, x_V\})$.

411 It should be noted that since $(V-1)p \leq 1$ is required, and that $c > 0$, we have $L < 1 \cdot 1/2$. This is
412 the origin of our condition $L < 1/2$ in the statement of the theorem. Naturally, the statement can be
413 adjusted to $\min(L, 1/2)$ instead.

414 For fixed \mathbf{b} , the excess risk in region $\{x_1, \dots, x_V\}$ becomes

$$\begin{aligned}\tilde{R}_V(\hat{f}) - \tilde{R}_V(\tilde{f}^*) &= \int_{\{x_1, \dots, x_V\}} 2c \mathbb{1}_{\hat{f}(x) \neq f_{\mathbf{b}}(x)} dP_X(x) \\ &\geq 2c \sum_{t=1}^{V-1} (1-\epsilon) p \mathbb{1}_{\hat{f}(x_t) \neq f_{\mathbf{b}}(x_t)},\end{aligned}$$

415 where the inequality follows from the fact that we ignore the risk on point x_V .

416 Using the probabilistic method, replace \mathbf{b} with $\mathbf{B} \sim \text{Uniform}\{1, 2\}^{V-1}$,

$$\begin{aligned}\sup_{\mathbf{b}, \delta, c, p} \mathbb{E}_{Z^n} \left[\frac{1}{\kappa + \delta} \left(\tilde{R}_V(\hat{f}) - \tilde{R}_V(\tilde{f}^*) \right) \right] &\geq \sup_{\delta, c, p} \mathbb{E}_{\mathbf{B}, Z^n} \left[\frac{1}{\kappa + \delta} \left(\tilde{R}_V(\hat{f}) - \tilde{R}_V(\tilde{f}^*) \right) \right] \\ &= \sup_{\delta, c, p} \frac{1}{\kappa + \delta} \mathbb{E}_{Z^n} \left[\mathbb{E}_{\mathbf{B}|Z^n} \left[\left(\tilde{R}_V(\hat{f}) - \tilde{R}_V(\tilde{f}^*) \right) \right] \right]\end{aligned}$$

417 Now, we need to calculate $\mathbf{B}|Z^n$, which can be calculated using Bayes rule because we have
418 $\mathbf{B} \sim \text{Uniform}\{1, 2\}^{V-1}$ and also $Z|\mathbf{B} \sim P_{X\tilde{Y}}$.

419 To be specific, for any $x \in \{x_0, x_1, \dots, x_{V-1}\}$, assume point x_t is observed k times in training
420 sample Z^n ,

$$\mathbb{P}(f_{\mathbf{B}}(x) = 1|Z^n) = \begin{cases} \frac{1}{2} & x \neq X_1, \dots, x \neq X_n, x \neq x_V \\ \mathbb{P}(B_t = 1|Y_{t_1}, \dots, Y_{t_k}) & x = x_t = X_{t_1} = \dots = X_{t_k}, 1 \leq t \leq V-1, \end{cases}$$

421 where B_t denotes the t -th element of vector \mathbf{B} .

422 Next we compute $\mathbb{P}(B_t = 1|Y_{t_1} = y_1, \dots, Y_{t_k} = y_k)$ for $y_1, \dots, y_k \in \{1, 2\}$. Denote the numbers
423 of ones and twos by $k_1 = |\{j \leq k : y_j = 1\}|$ and $k_2 = |\{j \leq k : y_j = 2\}|$. Using Bayes rule, we
424 get

$$\begin{aligned}\mathbb{P}(B_t = 1|Y_{t_1}, \dots, Y_{t_k}) &= \frac{\mathbb{P}(B_t = 1 \cap Y_{t_1}, \dots, Y_{t_k})}{\mathbb{P}(Y_{t_1}, \dots, Y_{t_k})} \\ &= \frac{\mathbb{P}(Y_{t_1}, \dots, Y_{t_k}|B_t = 1) \mathbb{P}(B_t = 1)}{\sum_{i=1}^2 \mathbb{P}(Y_{t_1}, \dots, Y_{t_k}|B_t = i) \mathbb{P}(B_t = i)} \\ &= \frac{(1/2 + c)^{k_1} (1/2 - c)^{k_2} (1/2)}{(1/2 + c)^{k_1} (1/2 - c)^{k_2} (1/2) + (1/2 + c)^{k_2} (1/2 - c)^{k_1} (1/2)}\end{aligned}$$

425 After some calculation, following Theorem 14.5 in Devroye et al. [1996], we get

$$\begin{aligned}\sup_{\mathbf{b}, \delta, c, p} \mathbb{E}_{Z^n} \left[\frac{1}{\kappa + \delta} \left(\tilde{R}_{\mathcal{A}_\kappa}(\hat{f}) - \tilde{R}_{\mathcal{A}_\kappa}(\tilde{f}^*) \right) \right] &\geq \sup_{\delta, c, p} \frac{1}{\kappa + \delta} c(V-1)(1-\epsilon) p e^{-\frac{8n(1-\epsilon)pc^2}{1-2c} - \frac{4c\sqrt{n(1-\epsilon)p}}{1-2c}} \\ &\geq \frac{1-\epsilon}{\kappa} \sup_{c, p} c(V-1) p e^{-\frac{8npc^2}{1-2c} - \frac{4c\sqrt{np}}{1-2c}} \quad \because \epsilon \geq 0, \text{ take } \delta \downarrow 0 \\ &= \frac{1-\epsilon}{\kappa} \sup_{c, p} c \frac{L}{1/2 - c} e^{-\frac{8npc^2}{1-2c} - \frac{4c\sqrt{np}}{1-2c}}\end{aligned}$$

426 Take $c = \sqrt{\frac{(V-1)}{8nL}}$, and p is automatically specified by Eqn. (12). Then

$$\sup_{\mathbf{b}, \delta, c, p} \mathbb{E}_{Z^n} \left[\frac{1}{\kappa + \delta} \left(\tilde{R}_{\mathcal{A}_\kappa}(\hat{f}) - \tilde{R}_{\mathcal{A}_\kappa}(\tilde{f}^*) \right) \right] \geq \frac{1-\epsilon}{\kappa} \sqrt{\frac{V-1}{8nL}} \frac{L}{1/2 - c} e^{-\frac{p(V-1) + \sqrt{\frac{4}{1-2c}}}{1-2c}}$$

427 Furthermore, take $n \geq \frac{9(V-1)}{2L}$, such that $c \leq 1/6$. Also use the fact that $(V-1)p \leq 1$, simplify the
428 expression

$$\begin{aligned}\sup_{\mathbf{b}, \delta, c, p} \mathbb{E}_{Z^n} \left[\frac{1}{\kappa + \delta} \left(\tilde{R}_{\mathcal{A}_\kappa}(\hat{f}) - \tilde{R}_{\mathcal{A}_\kappa}(\tilde{f}^*) \right) \right] &\geq \frac{1-\epsilon}{\kappa} \sqrt{\frac{V-1}{8nL}} \frac{L}{1/2 - 1/6} e^{-\frac{1+\sqrt{6}}{1-1/3}} \\ &\geq \frac{1-\epsilon}{\kappa} \sqrt{\frac{3(V-1)L}{8n}} e^{-6}\end{aligned}$$

429 In order to satisfy the condition $L \leq (V-1)p(1/2 - c)$ and $(V-1)p \leq 1$, plug in $c = \sqrt{\frac{V-1}{8nL}}$,
 430 we have $n \geq \frac{V-1}{2L(1-2L)^2}$.

431 Compare two parts together

$$\sup_{(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa)} \mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] \geq \frac{K-1}{K} \epsilon + \frac{1-\epsilon}{\kappa} \sqrt{\frac{3(V-1)L}{8n}} e^{-6},$$

432 for $n \geq \frac{V-1}{2L} \max \left\{ 9, \frac{1}{(1-2L)^2} \right\}$.

433 ■

434 A.1.3 Proof of Upper Bounds: Lemma 1 and Theorem 2

435 **Lemma 1 (Oracle Inequality under Feature-dependent Label Noise)** For any $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})$ and
 436 any classifier f , we have

$$R(f) - R(f^*) \leq \inf_{\kappa} \left\{ P_X(\mathcal{X} \setminus \mathcal{A}_{\kappa}(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})) + \frac{1}{\kappa} \left(\tilde{R}(f) - \tilde{R}(\tilde{f}^*) \right) \right\}.$$

437

438 *Proof.* For any $\kappa \geq 0$, the input space \mathcal{X} can be divided into two regions: $\mathcal{X} \setminus \mathcal{A}_{\kappa}$ and \mathcal{A}_{κ} .

439 For any f , its risk can be decomposed into two parts

$$\begin{aligned} R(f) &= \mathbb{E}_{X,Y} [\mathbb{1}_{f(X) \neq Y}] \\ &= \mathbb{E}_X \mathbb{E}_{Y|X} [\mathbb{1}_{f(X) \neq Y}] \\ &= \mathbb{E}_X \mathbb{E}_{Y|X} [1 - \mathbb{1}_{f(X)=Y}] \\ &= \mathbb{E}_X [1 - [\boldsymbol{\eta}(X)]_{f(X)}] \\ &= \int_{\mathcal{X}} (1 - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x) \\ &= \int_{\mathcal{X} \setminus \mathcal{A}_{\kappa}} (1 - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x) + \int_{\mathcal{A}_{\kappa}} (1 - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x) \end{aligned}$$

440 Therefore, excess risk equals to

$$R(f) - R(f^*) \underbrace{\int_{\mathcal{X} \setminus \mathcal{A}_{\kappa}} (\max \boldsymbol{\eta}(x) - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x)}_{\text{a}} + \underbrace{\int_{\mathcal{A}_{\kappa}} (\max \boldsymbol{\eta}(x) - [\boldsymbol{\eta}(x)]_{f(x)}) dP_X(x)}_{\text{b}}$$

441 Now examine the two terms separately,

$$\text{a} \leq \int_{\mathcal{X} \setminus \mathcal{A}_{\kappa}} 1 dP_X(x) = P_X(\mathcal{X} \setminus \mathcal{A}_{\kappa}(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})),$$

442 and

$$\begin{aligned} \text{b} &< \int_{\mathcal{A}_{\kappa}} \frac{1}{\kappa} (\max \tilde{\boldsymbol{\eta}}(x) - [\tilde{\boldsymbol{\eta}}(x)]_{f(x)}) dP_X(x) \quad \because \text{by definition of relative signal strength} \\ &\leq \int_{\mathcal{X}} \frac{1}{\kappa} (\max \tilde{\boldsymbol{\eta}}(x) - [\tilde{\boldsymbol{\eta}}(x)]_{f(x)}) dP_X(x) \\ &= \frac{1}{\kappa} (\tilde{R}(f) - \tilde{R}(\tilde{f}^*)) \quad \because \text{by definition of } \tilde{R}. \end{aligned}$$

443 Since this works for any $\kappa > 0$, we then have

$$R(f) - R(f^*) \leq \inf_{\kappa > 0} \left\{ P_X(\mathcal{X} \setminus \mathcal{A}_{\kappa}(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})) + \frac{1}{\kappa} (\tilde{R}(f) - \tilde{R}(\tilde{f}^*)) \right\}.$$

445 Now, to prepare for the rate of convergence proof, we first introduce the concept of shattering in
 446 multiclass and the Natarajan dimension [Natarajan, 1989], which can be viewed as a multi-class
 447 analogy of VC dimension [Vapnik and Chervonenkis, 1971].

448 **Definition 2 (Shattering (Multiclass))** Let \mathcal{H} be a class of functions from \mathcal{X} to $\mathcal{Y} = \{1, 2, \dots, K\}$.
 449 For any set containing n distinct elements $C_n = \{x_1, \dots, x_n\} \subset \mathcal{X}$, denote

$$\mathcal{H}_{C_n} = \{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\},$$

450 and therefore $|\mathcal{H}_{C_n}|$ is the number of distinct vectors of length n that can be realized by functions in
 451 \mathcal{H} .

452 The n^{th} shatter coefficient is defined as

$$S(\mathcal{H}, n) := \max_{C_n} |\mathcal{H}_{C_n}|.$$

453 We say that a set C_n is shattered by \mathcal{H} if there exists $f, g : C_n \rightarrow \mathcal{Y}$ such that for every $x \in C_n$,
 454 $f(x) \neq g(x)$, and

$$\mathcal{H}_C \supseteq \{f(x_1), g(x_1)\} \times \{f(x_2), g(x_2)\} \times \dots \times \{f(x_n), g(x_n)\}$$

455 If $\mathcal{Y} = \{1, 2\}$, this definition reduces the binary notion of shattering which says all labeling of points
 456 can be realized by some function in the hypothesis class \mathcal{H} , i.e., $\mathcal{H}_C = \{1, 2\}^{|C|}$. Contrary to the
 457 intuition that multiclass shattering means being able to label all K possible labels for each point
 458 $x \in C$: $\mathcal{H}_C = \{1, 2, \dots, K\}^{|C|}$. This definition is more like “embed the binary cube into multiclass”,
 459 where every $x \in C$ is allowed to pick from two labels rather than from K labels.

460 **Definition 3 (Natarajan Dimension)** The Natarajan dimension of \mathcal{H} , denoted $Ndim(\mathcal{H})$, is the
 461 maximal size of a shattered set $C \in \mathcal{X}$.

462 **Theorem (Excess Risk Upper Bound of NI-ERM)** Let $\epsilon \in [0, 1], \kappa \in (0, +\infty)$. Consider any
 463 $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa)$, assume function class \mathcal{F} has Natarajan dimension V , and the noisy Bayes
 464 classifier \tilde{f}^* belongs to \mathcal{F} . Let $\hat{f} \in \mathcal{F}$ be the ERM trained on $Z^n = \{(X_i, \tilde{Y}_i)\}_{i=1}^n$, then

$$\begin{aligned} \mathbb{E}_{Z^n} [R(\hat{f}) - R(f^*)] &\leq \epsilon + \frac{1}{\kappa} \cdot 8\sqrt{\frac{V \log 2n + 2V \log K + 4}{2n}} \\ &\leq \epsilon + \mathcal{O}\left(\frac{1}{\kappa} \sqrt{\frac{V}{n}}\right) \quad \text{up to log factor.} \end{aligned}$$

465 *Proof.* Following directly from Lemma 1, with $(P_X, \boldsymbol{\eta}, \tilde{\boldsymbol{\eta}}) \in \Pi(\epsilon, \kappa)$, we already have

$$\begin{aligned} R(f) - R(f^*) &\leq P_X(\mathcal{X} \setminus \mathcal{A}_\kappa(\boldsymbol{\eta}, \tilde{\boldsymbol{\eta}})) + \frac{1}{\kappa} (\tilde{R}(f) - \tilde{R}(\tilde{f}^*)) \\ &\leq \epsilon + \frac{1}{\kappa} (\tilde{R}(f) - \tilde{R}(\tilde{f}^*)). \end{aligned}$$

466 Now replace f with NI-ERM \hat{f} , the following procedure is similar to the derivation of generalization
 467 error based on VC dimension in binary classification setup, except using a “Natarajan version” of
 468 Sauer’s lemma.

469 Let’s start from excess risk bound based on Shattering coefficient.

Lemma 2

$$\mathbb{E}_{Z^n} [\tilde{R}(\hat{f}) - \tilde{R}(\tilde{f}^*)] \leq 8\sqrt{\frac{\log(8eS(\mathcal{H}, 2n))}{2n}}$$

470 *Proof Sketch.* First convert excess risk $\tilde{R}(\hat{f}) - \tilde{R}(\tilde{f}^*)$ to generalization error with the inequality
 471 ([Devroye et al., 1996, Lemma 8.2])

$$\tilde{R}(\hat{f}) - \tilde{R}(\tilde{f}^*) \leq 2 \sup_{f \in \mathcal{F}} |\tilde{R}(f) - \tilde{R}_n(f)|,$$

472 where \tilde{R}_n denotes the empirical risk evaluated on training sample Z^n .

473 Then use the “symmetrization lemma” [Vapnik and Chervonenkis, 1971, Lemma 2], for $n\epsilon^2 \geq 2$,

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}} |\tilde{R}(f) - \tilde{R}_n(f)| \geq \epsilon \right) \leq 2\mathbb{P} \left(\sup_{f \in \mathcal{F}} |\tilde{R}'_n(f) - \tilde{R}_n(f)| \geq \frac{\epsilon}{2} \right)$$

474 to convert the population risk $\tilde{R}(f)$ into $\tilde{R}'_n(f)$, the risk evaluated on the “ghost sample”, a second,
 475 independent sample of n points.

476 Finally, with union bound and Hoeffding’s inequality, we can bound the error probability with
 477 shattering coefficient on $2n$ points

$$\begin{aligned} 2\mathbb{P} \left(\sup_{f \in \mathcal{F}} |\tilde{R}'_n(f) - \tilde{R}_n(f)| \geq \frac{\epsilon}{2} \right) &\leq 2S(\mathcal{H}, 2n) \mathbb{P} \left(|\tilde{R}'_n(f) - \tilde{R}_n(f)| \geq \frac{\epsilon}{2} \right) \\ &\leq 2S(\mathcal{H}, 2n) \left\{ \mathbb{P} \left(|\tilde{R}'_n(f) - \tilde{R}(f)| \geq \frac{\epsilon}{4} \right) + \mathbb{P} \left(|\tilde{R}'_n(f) - \tilde{R}(f)| \geq \frac{\epsilon}{4} \right) \right\} \\ &\leq 2S(\mathcal{H}, 2n) \left\{ 2 \cdot 2 \cdot e^{-2n \left(\frac{\epsilon}{4}\right)^2} \right\} \\ &\leq 8S(\mathcal{H}, 2n) e^{-\frac{n\epsilon^2}{8}}, \end{aligned}$$

478 notice that the inequality also holds for $n\epsilon^2 < 2$ because in this case the upper bound is bigger than 1.

479 Therefore we get the convergence in probability bound for ERM \hat{f} ,

$$\mathbb{P} \left(\tilde{R}(\hat{f}) - \tilde{R}(\tilde{f}^*) \geq \epsilon \right) \leq 8S(\mathcal{H}, 2n) e^{-\frac{n\epsilon^2}{32}}.$$

480 Convert it to expectation using the inequality from Devroye et al. [1996, Problem 12.1], we finally get

$$\mathbb{E}_{Z^n} \left[\tilde{R}(\hat{f}) - \tilde{R}(\tilde{f}^*) \right] \leq 8 \sqrt{\frac{\log(8eS(\mathcal{H}, 2n))}{2n}}.$$

481 ■

482 The last step is to bound multiclass shattering coefficient with Natarajan dimension. We leverage the
 483 following lemma.

484 **Lemma 3 (Natarajan [1989])** *Let C and \mathcal{Y} be two finite sets and let \mathcal{H} be a set of functions from C
 485 to \mathcal{Y} .*

$$|\mathcal{H}| \leq |C|^{Ndim(\mathcal{H})} \cdot |\mathcal{Y}|^{2Ndim(\mathcal{H})}.$$

486 Replacing shattering coefficient with Natarajan dimension, we get

$$\begin{aligned} \mathbb{E}_{Z^n} \left[\tilde{R}(\hat{f}) - \tilde{R}(\tilde{f}^*) \right] &\leq 8 \sqrt{\frac{\log 8e + \log((2n)^V K^{2V})}{2n}} \\ &\leq 8 \sqrt{\frac{V \log 2n + 2V \log K + 4}{2n}} \end{aligned}$$

487 Putting things together,

$$\mathbb{E}_{Z^n} \left[R(\hat{f}) - R(f^*) \right] \leq \epsilon + \frac{1}{\kappa} \cdot 8 \sqrt{\frac{V \log 2n + 2V \log K + 4}{2n}}.$$

488 ■

489 A.2 Experimental Details

490 A.2.1 2D Gaussian with synthetic label noise

491 For 2D Gaussian mixture data, we draw from two Gaussian centered at $[1 \ 1]^\top$ and $[-1 \ -1]^\top$, with
 492 covariance matrix being identity, 200 data points from each, with label $Y = 1, 2$ respectively. To
 493 generate noisy labels, we flip every label uniformly with some probability. We use Sklearn’s logistic
 494 regression (with no ℓ_2 regularization). The experiment was conducted on AMD Ryzen 5 3600 CPU.
 495 The goal of the simulation is to experimentally verify noise immunity results in Section ???. Notice
 496 that different trial corresponds to different draw of both instances and noisy labels.

Table 2: Testing accuracy of logistic regression on gaussian mixture data with uniform label noise. “Noise rate” refers to $\mathbb{P}(\tilde{Y} \neq Y)$, the percentage of wrong labels in the training data. As theory in Section ??? predicts, when $\mathbb{P}(\tilde{Y} \neq Y)$ reach 50%, there is a sharp decrease in performance.

| Noise rate | 0 | 0.05 | 0.1 | 0.15 | 0.2 | 0.25 | 0.3 | 0.35 | 0.4 | 0.45 | 0.5 | 0.55 | 0.6 | 0.65 | 0.7 | 0.75 | 0.8 | 0.85 | 0.9 | 0.95 | 1 |
|------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|------|------|------|------|------|------|
| Trial #1 | 93.00 | 92.83 | 92.38 | 92.08 | 91.78 | 91.93 | 92.25 | 92.90 | 91.83 | 92.58 | 74.68 | 25.12 | 9.70 | 7.73 | 7.52 | 7.25 | 7.38 | 7.15 | 7.18 | 7.10 | 7.00 |
| Trial #2 | 91.73 | 91.60 | 92.05 | 91.63 | 91.78 | 91.78 | 91.68 | 91.63 | 91.55 | 91.48 | 80.40 | 21.10 | 9.93 | 8.55 | 8.38 | 8.22 | 8.20 | 8.35 | 8.33 | 8.40 | 8.28 |
| Trial #3 | 92.73 | 92.75 | 92.78 | 92.78 | 92.58 | 92.45 | 91.68 | 88.15 | 82.58 | 59.83 | 49.53 | 35.80 | 21.28 | 14.35 | 9.33 | 8.53 | 8.12 | 7.70 | 7.13 | 7.23 | 7.28 |
| Trial #4 | 91.55 | 91.58 | 91.60 | 91.63 | 91.68 | 91.60 | 91.25 | 90.98 | 89.98 | 86.38 | 60.53 | 9.95 | 8.75 | 10.00 | 10.45 | 9.08 | 9.00 | 9.53 | 9.20 | 9.03 | 8.45 |
| Trial #5 | 91.55 | 91.58 | 91.60 | 91.63 | 91.68 | 91.60 | 91.25 | 90.98 | 89.98 | 86.38 | 60.53 | 9.95 | 8.75 | 10.00 | 10.45 | 9.08 | 9.00 | 9.53 | 9.20 | 9.03 | 8.45 |
| Mean | 92.11 | 92.07 | 92.08 | 91.95 | 91.90 | 91.87 | 91.62 | 90.93 | 89.18 | 83.33 | 65.13 | 20.40 | 11.68 | 10.10 | 9.23 | 8.43 | 8.34 | 8.45 | 8.21 | 8.16 | 7.89 |
| Std | 0.70 | 0.66 | 0.51 | 0.50 | 0.38 | 0.35 | 0.41 | 1.74 | 3.79 | 13.44 | 12.35 | 10.94 | 5.39 | 2.56 | 1.29 | 0.75 | 0.68 | 1.07 | 1.03 | 0.94 | 0.70 |

497 A.2.2 MNIST with synthetic label noise

498 We flip the clean training label of MNIST (<http://yann.lecun.com/exdb/mnist/>) uniformly
 499 (to any of the wrong classes). We use a shallow neural network with two convolution layers and
 500 two fully connected layers. We train with stochastic gradient descent with learning rate 0.01 for 10
 501 epochs, batch size equals 64. We use the same hyperparamters for all tests. The experiments were
 502 conducted on a single NVIDIA GTX 1660S GPU. The goal of the simulation is to experimentally
 503 verify noise immunity results in Section ???. Here randomness corresponds to different realization of
 504 noisy labels and stochastic gradient descent.

Table 3: Testing accuracy of a shallow CNN (2 conv layers with 2 fully connected layers) on MNIST with uniform label noise. “Noise rate” refers to $\mathbb{P}(\tilde{Y} \neq Y)$, the percentage of wrong labels in the training data. As theory in Section ??? predicts, when $\mathbb{P}(\tilde{Y} \neq Y)$ reach 90%, there is a sharp decrease in performance.

| Noise rate | 0 | 0.05 | 0.1 | 0.15 | 0.2 | 0.25 | 0.3 | 0.35 | 0.4 | 0.45 | 0.5 | 0.55 | 0.6 | 0.65 | 0.7 | 0.75 | 0.8 | 0.85 | 0.9 | 0.95 | 1 |
|------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|------|------|
| Trial #1 | 98.97 | 98.89 | 98.81 | 98.46 | 98.49 | 98.16 | 98.46 | 98.07 | 97.98 | 97.57 | 97.88 | 97.84 | 97.19 | 97.10 | 96.70 | 95.02 | 89.00 | 83.72 | 11.58 | 0.17 | 0.03 |
| Trial #2 | 98.88 | 98.73 | 98.94 | 98.55 | 98.72 | 98.66 | 98.50 | 98.24 | 98.15 | 98.23 | 97.86 | 97.98 | 97.70 | 97.10 | 96.91 | 95.76 | 91.99 | 88.49 | 9.99 | 0.08 | 0.04 |
| Trial #3 | 99.00 | 99.04 | 98.86 | 98.56 | 98.69 | 98.66 | 98.51 | 98.49 | 98.37 | 98.25 | 97.39 | 97.37 | 97.18 | 96.66 | 94.88 | 92.15 | 81.48 | 6.19 | 0.14 | 0.04 | |
| Trial #4 | 99.04 | 98.86 | 98.70 | 98.76 | 98.83 | 98.65 | 98.34 | 98.42 | 98.58 | 98.47 | 98.00 | 97.41 | 97.63 | 97.09 | 96.46 | 95.94 | 93.19 | 84.78 | 8.68 | 0.19 | 0.01 |
| Trial #5 | 99.05 | 98.58 | 98.89 | 98.82 | 98.72 | 98.83 | 98.34 | 98.55 | 98.40 | 98.38 | 98.01 | 97.31 | 97.33 | 96.21 | 96.29 | 94.92 | 90.38 | 85.84 | 8.98 | 0.13 | 0.08 |
| Mean | 98.99 | 98.82 | 98.84 | 98.63 | 98.69 | 98.59 | 98.43 | 98.35 | 98.30 | 98.18 | 98.00 | 97.59 | 97.44 | 96.94 | 96.60 | 95.30 | 91.34 | 84.86 | 9.08 | 0.14 | 0.04 |
| Std | 0.07 | 0.17 | 0.09 | 0.15 | 0.12 | 0.25 | 0.08 | 0.20 | 0.23 | 0.36 | 0.16 | 0.30 | 0.21 | 0.41 | 0.24 | 0.51 | 1.65 | 2.59 | 1.98 | 0.04 | 0.03 |

505 A.2.3 CIFAR with synthetic label noise

506 We flip the clean training label of CIFAR-10 (<https://www.cs.toronto.edu/~kriz/cifar.html>) uniformly (to any of the wrong classes). To have a fair comparison between different methods,
 507 we fix the realization of noisy labels. Follow the 2-step procedure described in Section 5, we use
 508 different pre-trained neural networks as feature extractor: forward-passing the training image through
 509 the network and record the feature. Then use sklearn’s (<https://scikit-learn.org/stable/>)
 510 logistic regression function to fit the (feature, noisy label) pair in a full batch manner. We pre-
 511 specify a range of values for ℓ_2 regularization ($\{0.0001, 0.001, 0.01, 0.1, 1, 10, 100\}$) and number of
 512 iterations for lbfgs optimizer ($\{10, 20, 50, 100\}$), then do cross-validation on noisy data to pick the
 513 best hyper-parameters. We use the same range of hyper-parameters in all tests. The experiments were
 514 conducted on a single NVIDIA Tesla V100 GPU. The result is deterministic.

Table 4: Performance on CIFAR-10 with synthetic label noise. We apply linear model on top of different feature extractors: “ResNet-50 TL” refers to using a pre-trained ResNet-50 on ImageNet [Deng et al., 2009] (available in Pytorch model library) in a transfer learning fashion, “ResNet-50 SSL” refers to using a pre-trained ResNet-50 on unlabeled CIFAR data with self-supervised loss [Chen et al., 2020] (publicly downloadable weights <https://github.com/ContrastToDivide/C2D?tab=readme-ov-file>) and “DINOv2 SSL” refers to using the self-supervised foundation model DINOv2 [Oquab et al., 2023] (available at <https://github.com/facebookresearch/dinov2>) as the feature extractor. “Noise rate” refers to $\mathbb{P}(\tilde{Y} \neq Y)$, the percentage of wrong labels in the training data. As theory in Section ?? predicts, when $\mathbb{P}(\tilde{Y} \neq Y)$ reach 90%, there is a sharp decrease in performance. We employed Python’s sklearn logistic regression and cross-validation functions without data augmentation. The results are deterministic.

| Noise rate | 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.85 | 0.9 | 0.95 | 1 |
|------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|------|------|
| Linear | 41.37 | 41.09 | 40.97 | 40.37 | 40.45 | 39.44 | 37.28 | 35.20 | 26.74 | 18.00 | 10.28 | 5.50 | 3.92 |
| Linear + ResNet-50 TL | 90.17 | 89.58 | 89.01 | 88.27 | 87.55 | 87.28 | 86.40 | 85.01 | 82.03 | 74.02 | 10.82 | 1.47 | 0.26 |
| Linear + ResNet-50 SSL | 92.48 | 92.26 | 91.74 | 91.46 | 91.13 | 90.33 | 91.07 | 90.99 | 89.11 | 83.89 | 10.08 | 1.31 | 0.34 |
| Linear + DINOv2 SSL | 99.25 | 99.27 | 99.23 | 99.14 | 99.10 | 99.11 | 99.02 | 98.84 | 95.50 | 76.91 | 10.13 | 0.92 | 0.03 |

516 A.2.4 CIFAR with human label error

517 We load the noisy human labels provided by <http://noisylabels.com/>, then follow exact the
518 same procedure as above.

Table 5: Performance on CIFAR-N dataset (<http://noisylabels.com/>) in terms of testing accuracy. “Aggre”, “Rand1”, ..., “Noisy” denote various types of human label noise. We apply linear model on top of different feature extractors: “ResNet-50 TL” refers to using a pre-trained ResNet-50 on ImageNet [Deng et al., 2009] in a transfer learning fashion, “ResNet-50 SSL” refers to using a pre-trained ResNet-50 on unlabeled CIFAR data with self-supervised loss [Chen et al., 2020] and “DINOv2 SSL” refers to using the self-supervised foundation model DINOv2 [Oquab et al., 2023] as the feature extractor. We employed Python’s sklearn logistic regression and cross-validation functions without data augmentation; the results are deterministic and directly reproducible.

| Methods | CIFAR-10N | | | | | CIFAR-100N |
|------------------------|-----------|-------|-------|-------|-------|------------|
| | Aggre | Rand1 | Rand2 | Rand3 | Worst | Noisy |
| Linear | 40.73 | 40.41 | 40.31 | 40.63 | 38.43 | 16.61 |
| Linear + ResNet-50 TL | 89.18 | 88.63 | 88.61 | 88.66 | 85.32 | 62.89 |
| Linear + ResNet-50 SSL | 91.78 | 91.66 | 91.39 | 91.28 | 87.84 | 57.95 |
| Linear + DINOv2 SSL | 98.69 | 98.80 | 98.65 | 98.67 | 95.71 | 83.17 |

519 **A.2.5 Linear probing, then fine tuning (LP-FT)**

520 We study whether LP-FT works in label noise learning scenario.

Table 6: Performance on CIFAR-N dataset (<http://noisylabels.com/>) in terms of testing accuracy. “Aggre”, “Rand1”, ..., “Noisy” denote various types of human label noise. We apply linear model on top of different feature extractors: “ResNet-50 TL” refers to using a pre-trained ResNet-50 on ImageNet [Deng et al., 2009] in a transfer learning fashion, “ResNet-50 SSL” refers to using a pre-trained ResNet-50 on unlabeled CIFAR data with contrastive loss [Chen et al., 2020] and “DINOv2 SSL” refers to using the self-supervised foundation model DINOv2 [Oquab et al., 2023] as the feature extractor.

| Feature | Method | CIFAR-10N | | | | | | CIFAR-100N | |
|--------------------|--------|-----------|-------|-------|-------|-------|-------|------------|-------|
| | | Clean | Aggre | Rand1 | Rand2 | Rand3 | Worst | Clean | Noisy |
| ResNet-50 TL | LP | 90.17 | 89.18 | 88.63 | 88.61 | 88.66 | 85.32 | 71.79 | 62.89 |
| | LP-FT | 95.94 | 92.03 | 88.55 | 87.78 | 87.82 | 71.88 | 82.3 | 63.85 |
| ResNet-50 SSL | LP | 92.54 | 91.78 | 91.66 | 91.46 | 91.17 | 87.85 | 69.88 | 57.98 |
| | LP-FT | 94.11 | 89.11 | 84.49 | 83.75 | 84.15 | 65.00 | 74.41 | 54.49 |
| DINOv2 (small) SSL | LP | 96.09 | 94.8 | 94.39 | 94.42 | 94.35 | 91.14 | 83.82 | 72.46 |
| | LP-FT | 98.23 | 93.29 | 88.03 | 87.27 | 86.94 | 67.42 | 89.97 | 64.81 |

521 **A.2.6 Robust learning strategy over DINOv2 feature**

522 This section examines how different robust learning strategy works over DINOv2 feature, compared
 523 with only training with cross entropy.

| Feature | Method | CIFAR-10N | | | | | | CIFAR-100N | |
|------------|----------|-----------|-------|-------|-------|-------|-------|------------|-------|
| | | Clean | Aggre | Rand1 | Rand2 | Rand3 | Worst | Clean | Noisy |
| DINOv2 SSL | CE | 99.25 | 98.69 | 98.8 | 98.65 | 98.67 | 95.71 | 92.85 | 83.17 |
| | MAE | 99.27 | 99.04 | 99.01 | 99.09 | 99.11 | 95.55 | 90.68 | 82.55 |
| | Sigmoid | 99.26 | 98.86 | 98.91 | 98.87 | 98.96 | 96.66 | 92.82 | 82.03 |
| | CE + SAM | 99.09 | 97.66 | 98.47 | 98.53 | 98.47 | 95.47 | 89.97 | 82.85 |

525 **NeurIPS Paper Checklist**

526 **1. Claims**

527 Question: Do the main claims made in the abstract and introduction accurately reflect the
528 paper's contributions and scope?

529 Answer: [\[Yes\]](#)

530 Justification: The claims are either supported by theory statements or by reproducible
531 experiment results.

532 Guidelines:

- 533 • The answer NA means that the abstract and introduction do not include the claims
534 made in the paper.
- 535 • The abstract and/or introduction should clearly state the claims made, including the
536 contributions made in the paper and important assumptions and limitations. A No or
537 NA answer to this question will not be perceived well by the reviewers.
- 538 • The claims made should match theoretical and experimental results, and reflect how
539 much the results can be expected to generalize to other settings.
- 540 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
541 are not attained by the paper.

542 **2. Limitations**

543 Question: Does the paper discuss the limitations of the work performed by the authors?

544 Answer: [\[Yes\]](#)

545 Justification: Limitations about our practical method is described.

546 Guidelines:

- 547 • The answer NA means that the paper has no limitation while the answer No means that
548 the paper has limitations, but those are not discussed in the paper.
- 549 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 550 • The paper should point out any strong assumptions and how robust the results are to
551 violations of these assumptions (e.g., independence assumptions, noiseless settings,
552 model well-specification, asymptotic approximations only holding locally). The authors
553 should reflect on how these assumptions might be violated in practice and what the
554 implications would be.
- 555 • The authors should reflect on the scope of the claims made, e.g., if the approach was
556 only tested on a few datasets or with a few runs. In general, empirical results often
557 depend on implicit assumptions, which should be articulated.
- 558 • The authors should reflect on the factors that influence the performance of the approach.
559 For example, a facial recognition algorithm may perform poorly when image resolution
560 is low or images are taken in low lighting. Or a speech-to-text system might not be
561 used reliably to provide closed captions for online lectures because it fails to handle
562 technical jargon.
- 563 • The authors should discuss the computational efficiency of the proposed algorithms
564 and how they scale with dataset size.
- 565 • If applicable, the authors should discuss possible limitations of their approach to
566 address problems of privacy and fairness.
- 567 • While the authors might fear that complete honesty about limitations might be used by
568 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
569 limitations that aren't acknowledged in the paper. The authors should use their best
570 judgment and recognize that individual actions in favor of transparency play an impor-
571 tant role in developing norms that preserve the integrity of the community. Reviewers
572 will be specifically instructed to not penalize honesty concerning limitations.

573 **3. Theory Assumptions and Proofs**

574 Question: For each theoretical result, does the paper provide the full set of assumptions and
575 a complete (and correct) proof?

576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628

Answer: [\[Yes\]](#)

Justification: Assumptions are stated in the theorem statement. Full proofs are included in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: Important information about the experiments are in main text. Details on the experimental setup is described in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

629 **5. Open access to data and code**

630 Question: Does the paper provide open access to the data and code, with sufficient instruc-
631 tions to faithfully reproduce the main experimental results, as described in supplemental
632 material?

633 Answer: [Yes]

634 Justification: Code is provided, common benchmark dataset were used, instructions are given,
635 the result is easily reproducible.

636 Guidelines:

- 637 • The answer NA means that paper does not include experiments requiring code.
- 638 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
639 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 640 • While we encourage the release of code and data, we understand that this might not be
641 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
642 including code, unless this is central to the contribution (e.g., for a new open-source
643 benchmark).
- 644 • The instructions should contain the exact command and environment needed to run to
645 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
646 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 647 • The authors should provide instructions on data access and preparation, including how
648 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 649 • The authors should provide scripts to reproduce all experimental results for the new
650 proposed method and baselines. If only a subset of experiments are reproducible, they
651 should state which ones are omitted from the script and why.
- 652 • At submission time, to preserve anonymity, the authors should release anonymized
653 versions (if applicable).
- 654 • Providing as much information as possible in supplemental material (appended to the
655 paper) is recommended, but including URLs to data and code is permitted.

656 **6. Experimental Setting/Details**

657 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
658 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
659 results?

660 Answer: [Yes]

661 Justification: See appendix and attached code.

662 Guidelines:

- 663 • The answer NA means that the paper does not include experiments.
- 664 • The experimental setting should be presented in the core of the paper to a level of detail
665 that is necessary to appreciate the results and make sense of them.
- 666 • The full details can be provided either with the code, in appendix, or as supplemental
667 material.

668 **7. Experiment Statistical Significance**

669 Question: Does the paper report error bars suitably and correctly defined or other appropriate
670 information about the statistical significance of the experiments?

671 Answer: [Yes]

672 Justification: We have done repeated experiments for data simulation, e.g., different iid draw.
673 For the result on CIFAR-N data challenge, our result is deterministic and thus no error bar.

674 Guidelines:

- 675 • The answer NA means that the paper does not include experiments.
- 676 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
677 dence intervals, or statistical significance tests, at least for the experiments that support
678 the main claims of the paper.

- 679 • The factors of variability that the error bars are capturing should be clearly stated (for
680 example, train/test split, initialization, random drawing of some parameter, or overall
681 run with given experimental conditions).
- 682 • The method for calculating the error bars should be explained (closed form formula,
683 call to a library function, bootstrap, etc.)
- 684 • The assumptions made should be given (e.g., Normally distributed errors).
- 685 • It should be clear whether the error bar is the standard deviation or the standard error
686 of the mean.
- 687 • It is OK to report 1-sigma error bars, but one should state it. The authors should
688 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
689 of Normality of errors is not verified.
- 690 • For asymmetric distributions, the authors should be careful not to show in tables or
691 figures symmetric error bars that would yield results that are out of range (e.g. negative
692 error rates).
- 693 • If error bars are reported in tables or plots, The authors should explain in the text how
694 they were calculated and reference the corresponding figures or tables in the text.

695 8. Experiments Compute Resources

696 Question: For each experiment, does the paper provide sufficient information on the com-
697 puter resources (type of compute workers, memory, time of execution) needed to reproduce
698 the experiments?

699 Answer: [Yes]

700 Justification: See appendix.

701 Guidelines:

- 702 • The answer NA means that the paper does not include experiments.
- 703 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
704 or cloud provider, including relevant memory and storage.
- 705 • The paper should provide the amount of compute required for each of the individual
706 experimental runs as well as estimate the total compute.
- 707 • The paper should disclose whether the full research project required more compute
708 than the experiments reported in the paper (e.g., preliminary or failed experiments that
709 didn't make it into the paper).

710 9. Code Of Ethics

711 Question: Does the research conducted in the paper conform, in every respect, with the
712 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

713 Answer: [Yes]

714 Justification: The authors have read the NeurIPS Code of Ethics and confirm that this
715 research follows the code of ethics.

716 Guidelines:

- 717 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 718 • If the authors answer No, they should explain the special circumstances that require a
719 deviation from the Code of Ethics.
- 720 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
721 eration due to laws or regulations in their jurisdiction).

722 10. Broader Impacts

723 Question: Does the paper discuss both potential positive societal impacts and negative
724 societal impacts of the work performed?

725 Answer: [NA]

726 Justification: This is a theory-oriented paper towards better understanding of label noise
727 problem.

728 Guidelines:

- 729 • The answer NA means that there is no societal impact of the work performed.
- 730 • If the authors answer NA or No, they should explain why their work has no societal
- 731 impact or why the paper does not address societal impact.
- 732 • Examples of negative societal impacts include potential malicious or unintended uses
- 733 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 734 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 735 groups), privacy considerations, and security considerations.
- 736 • The conference expects that many papers will be foundational research and not tied
- 737 to particular applications, let alone deployments. However, if there is a direct path to
- 738 any negative applications, the authors should point it out. For example, it is legitimate
- 739 to point out that an improvement in the quality of generative models could be used to
- 740 generate deepfakes for disinformation. On the other hand, it is not needed to point out
- 741 that a generic algorithm for optimizing neural networks could enable people to train
- 742 models that generate Deepfakes faster.
- 743 • The authors should consider possible harms that could arise when the technology is
- 744 being used as intended and functioning correctly, harms that could arise when the
- 745 technology is being used as intended but gives incorrect results, and harms following
- 746 from (intentional or unintentional) misuse of the technology.
- 747 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 748 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 749 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 750 feedback over time, improving the efficiency and accessibility of ML).

751 11. Safeguards

752 Question: Does the paper describe safeguards that have been put in place for responsible
 753 release of data or models that have a high risk for misuse (e.g., pretrained language models,
 754 image generators, or scraped datasets)?

755 Answer: [NA]

756 Justification: The paper poses no such risks

757 Guidelines:

- 758 • The answer NA means that the paper poses no such risks.
- 759 • Released models that have a high risk for misuse or dual-use should be released with
- 760 necessary safeguards to allow for controlled use of the model, for example by requiring
- 761 that users adhere to usage guidelines or restrictions to access the model or implementing
- 762 safety filters.
- 763 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 764 should describe how they avoided releasing unsafe images.
- 765 • We recognize that providing effective safeguards is challenging, and many papers do
- 766 not require this, but we encourage authors to take this into account and make a best
- 767 faith effort.

768 12. Licenses for existing assets

769 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
 770 the paper, properly credited and are the license and terms of use explicitly mentioned and
 771 properly respected?

772 Answer: [Yes]

773 Justification: Citations and urls are included.

774 Guidelines:

- 775 • The answer NA means that the paper does not use existing assets.
- 776 • The authors should cite the original paper that produced the code package or dataset.
- 777 • The authors should state which version of the asset is used and, if possible, include a
- 778 URL.
- 779 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 780
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- 781
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 782
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 783
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 784
- 785
- 786
- 787
- 788
- 789

790 13. New Assets

791 Question: Are new assets introduced in the paper well documented and is the documentation
792 provided alongside the assets?

793 Answer: [NA]

794 Justification: The paper does not release new assets.

795 Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

804 14. Crowdsourcing and Research with Human Subjects

805 Question: For crowdsourcing experiments and research with human subjects, does the paper
806 include the full text of instructions given to participants and screenshots, if applicable, as
807 well as details about compensation (if any)?

808 Answer: [NA]

809 Justification: The paper does not involve crowdsourcing nor research with human subjects.

810 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

819 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 820 Subjects

821 Question: Does the paper describe potential risks incurred by study participants, whether
822 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
823 approvals (or an equivalent approval/review based on the requirements of your country or
824 institution) were obtained?

825 Answer: [NA]

826 Justification: The paper does not involve with this matter.

827 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

830
831
832
833
834
835
836
837

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.