

When Disagreement Meets Noise: Noise Robust Annotator Embeddings for Subjective NLP

Anonymous ACL submission

Abstract

Subjective NLP tasks such as sentiment analysis and hate speech classification often involve inherent annotator disagreement, reflecting diverse perspectives shaped by annotators’ lived experiences. Although conventional approaches resolve disagreement through majority voting or aggregation, these methods risk erasing valuable nuances and minority viewpoints. Recent embedding-based/multitask models have advanced the modeling of annotator-specific judgments, yet their robustness to annotation noise remains underexplored. In this work, we systematically investigate how state-of-the-art multi-annotator learning models perform in the presence of noisy labels and observe a significant performance degradation under such conditions. To address this, we propose Noise Robust Annotator Embedding (NRA-Embed), which integrates Robust InfoNCE (RINCE) contrastive loss to enhance models’ robustness under noisy annotation conditions. Through extensive experiments, we demonstrate that NRA-Embed effectively captures subjective disagreement while maintaining robustness to noise, achieving performance that matches or exceeds state-of-the-art methods.

1 Introduction

Collecting multiple annotations per instance is a standard practice in NLP to improve label reliability and mitigate individual annotator errors (Snow et al., 2008; Nowak and Ruger, 2010). In conventional supervised learning, disagreement among annotators is typically resolved through majority voting, averaging (Sabou et al., 2014), or expert adjudication (Waseem and Hovy, 2016) to produce a single ground-truth label. However, for *subjective NLP tasks*, such as sentiment analysis, hate speech detection, and political stance classification, a single objectively correct label often does not exist (Alm, 2011). Enforcing consensus in these

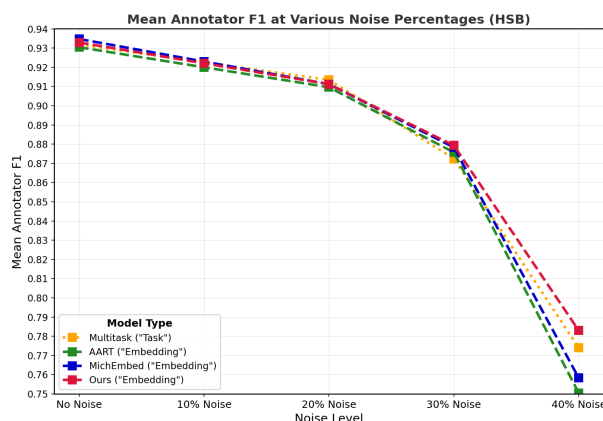


Figure 1: Comparison of Mean Annotator F1 scores for different multi-annotator modeling approaches on the HSB (Akhtar et al., 2021) Dataset under increasing noise levels. We highlight that our embedding-based method demonstrates greater robustness, maintaining higher performance compared to other approaches as noise level increases.

settings can obscure meaningful variation in human judgments and suppress socially situated perspectives (Cheplygina and Plum, 2018). At the same time, real-world annotation pipelines are inevitably affected by *annotation noise*, arising from misunderstanding, inattention, or malicious behavior. Despite this, the joint effect of subjective disagreement and annotation noise remains insufficiently understood.

Annotator disagreement in subjective tasks is not random but is often systematically shaped by annotators’ demographic backgrounds, social identities, and political beliefs. For instance, feminist and anti-racism activists have been shown to annotate hate speech differently from crowd workers (Waseem and Hovy, 2016), while political affiliations influence labeling behavior in stance detection (Luo et al., 2020). Similar effects have been observed in sentiment analysis (Dıaz et al., 2018) and toxic language annotation (Patton et al., 2019).

062	Collapsing such diverse judgments into a single	robustness to annotation noise.	114
063	label via majority voting risks marginalizing mi-	We evaluate NRA-Embed across multiple subjec-	115
064	nority viewpoints and introducing representational	tive classification tasks, noise types, and annotator	116
065	harms (Prabhakaran et al., 2021). Motivated by	conditions. Our contributions are:	117
066	these concerns, a growing body of work has sought		
067	to <i>model</i> annotator disagreement rather than elimi-	• We conduct the first systematic evaluation	118
068	nate it. Benchmark datasets such as Multi-Domain	of multi-annotator models under controlled	119
069	Agreement (MDA) (Leonardelli et al., 2021) en-	noise, revealing significant vulnerabilities in	120
070	able systematic study, while probabilistic reliability	existing approaches.	121
071	modeling (Reidsma and Carletta, 2008), large-scale		
072	disagreement learning (Uma et al., 2022), and dis-	• We propose NRA-Embed, achieving several	122
073	tributational supervision (Uma et al., 2021) provide	significant performance improvements such	123
074	principled alternatives to label aggregation.	as a 3.25% increase over AART (Mokhberian	124
075	Building on these foundations, recent methods	et al., 2023) on the HSB (Akhtar et al., 2021)	125
076	have moved beyond label aggregation to explicitly	dataset.	126
077	model annotator perspectives. These include multi-		
078	task architectures with per-annotator classification	• We benchmark across three binary and one	127
079	heads (Davani et al., 2022a; Jinadu and Ding, 2024)	multiclass datasets and two noise types (sym-	128
080	and embedding-based approaches that represent	metric and renegade), establishing best prac-	129
081	annotators as learned vectors (Mokhberian et al.,	tices for noise-robust subjective learning.	130
082	2023; Deng et al., 2023). Embedding methods		
083	using contrastive learning have achieved state-of-	2 Related Works	131
084	the-art performance on disagreement benchmarks	2.1 Annotator Disagreement	132
085	under clean supervision (Mokhberian et al., 2023).	Annotator disagreement is a recognized challenge	133
086	However, contrastive objectives like InfoNCE	in NLP. Traditional approaches treat this variability	134
087	fundamentally rely on correct positive/negative pair	as noise through majority voting or label aggre-	135
088	construction, when labels are corrupted, the con-	gation (Pavlick and Callison-Burch, 2016; Sabou	136
089	trastive signal becomes unreliable. In this work,	et al., 2014). However, many tasks admit multiple	137
090	we conduct the first systematic study of how state-	valid interpretations where no single correct label	138
091	of-the-art disagreement learning methods behave	exists (Plank, 2022; Passonneau et al., 2012; Nie	139
092	under controlled annotation noise. We show em-	et al., 2020; Jiang and Marneffe, 2022; Jinadu and	140
093	pirically in (Figure 1) that methods which account	Ding, 2024). In subjective settings like toxicity de-	141
094	for multi-annotator perspectives degrade in perfor-	tection, judgments differ due to annotators’ beliefs	142
095	mance with increasing annotation noise. These	and identities, reflecting genuine human diversity	143
096	findings indicate that existing embedding-based	rather than errors (Waseem, 2016; Al Kuwatly et al.,	144
097	disagreement models are not inherently robust to	2020; Sap et al., 2021; Pavlick and Kwiatkowski,	145
098	noisy supervision.	2019). Recent work leverages this disagreement	146
099	To address this limitation, while keeping the	for personalization (Plepi et al., 2022).	147
100	benefits of contrastive loss, we argue that robust-	2.2 Modeling Annotator Disagreement	148
101	ness must be incorporated directly into the con-	Methods for modeling disagreement have evolved	149
102	trastive learning objective itself. We therefore pro-	from treating it as noise to embracing it as sig-	150
103	pose the Noise-Robust Annotator Embedding	nal. Early approaches include Dawid–Skene for	151
104	method (NRA-Embed) , which is an annotator	estimating annotator reliability (Dawid and Skene,	152
105	embedding-based model that replaces the conven-	1979), uncertainty quantification via MC Dropout	153
106	tional InfoNCE loss with <i>Robust InfoNCE (RINCE)</i>	and Deep Ensembles (Zhou et al., 2021), and full	154
107	(Chuang et al., 2022). RINCE introduces a princi-	label distribution modeling (Meissner et al., 2021).	155
108	pled robustness mechanism that reduces sensitivity		
109	to corrupted positive and negative pairs through	Multitask approaches. Davani et al. (2022a)	156
110	a tunable robustness parameter. By integrating	proposed a multitask framework with separate clas-	157
111	RINCE into the annotator embedding framework,	sification heads per annotator, enabling individual-	158
112	NRA-Embed preserves the benefits of contrastive	ized predictions while sharing text representations.	159
113	annotator modeling while substantially improving	Jinadu and Ding (2024) extended this with explicit	160

noise tolerance. However, multitask methods suffer from under determination for sparse annotators who contribute few labels.

Embedding-based approaches. Recent work represents annotators as learned embeddings integrated early in the network. [Deng et al. \(2023\)](#) jointly learn annotator and annotation embeddings to capture fine-grained subjectivity signals. [Mokhberian et al. \(2023\)](#) introduced AART, which uses contrastive learning (InfoNCE) to align annotator embeddings with labeling patterns, demonstrating more equitable performance for sparse annotators. However, they noted that “differentiating between labeling noise and natural disagreement is a challenge that remains unaddressed,” a gap our work directly targets.

2.3 Learning from Noisy Labels

Noise correction targets errors in labels, such as random flips and annotation mistakes, to improve data quality and robustness ([Zhan et al., 2019](#)). Cross-entropy easily fits noisy labels ([Zhang et al., 2016](#)), while bootstrapping adds prediction terms to reduce this effect ([Reed et al., 2014](#)). Deep models learn clean patterns before memorizing noise ([Arazo et al., 2019](#); [Liu et al., 2020](#); [Li et al., 2020](#); [Nishi et al., 2021](#)), a property exploited by many denoising methods. Effective correction must remove misleading labels without discarding true signal ([Arazo et al., 2019](#); [Jinadu and Ding, 2024](#)).

3 Preliminary

We first set up the general multi-annotator learning problem. We then discuss two high-level network architecture approaches for modeling multi-annotator datasets. We then propose a framework for optimizing the embedding-based approaches under conditions of label noise.

3.1 Problem Definition

We consider an annotated dataset $D = \{(x_i, a_j, y_{ij})\}$, which consists of triplets formed from input text items $X = \{x_i\}_{i=1}^N$, annotators $A = \{a_j\}_{j=1}^m$, and annotations $Y = \{1, \dots, Q\}$. A pair of (i, j) can appear at most once in the dataset D , which means label y_{ij} is assigned to text item x_i by the annotator a_j . Y contains numerous missing values in most annotated datasets since each annotator labels only a subset of the instances. The problem is modeled as a classification task where

a classifier is trained to predict the label to be assigned to the text item x_i . All methods explored in this study utilize pre-trained transformer-based language models for text encoding, specifically using RoBERTa ([Liu et al., 2019](#)). For a given input text x_i , we obtain its text representation by extracting the [CLS] token embedding from the final layer of the language model, denoted as $e(x_i)$.

3.2 Task-based Multi-Annotator Learning

The most frequent approach, called *single-task*, aims to predict the aggregate label to be assigned to the text item x_i through majority voting or averaging over annotators’ labels $\{y_{ij}\}_{j=1}^M$.

Multi-task approaches are seen in several previous works ([Fornaciari et al., 2021](#); [Davani et al., 2022a](#); [Jinadu and Ding, 2024](#)), and basically are a generalization of single-task approach. They train a separate, fully connected layer for each annotator to learn the annotator-specific labeling behavior.

3.3 Embedding-based Multi-Annotator Learning

Another method is to embed annotators in a latent space and integrate this information early in the model architecture. In these approaches, a learnable matrix encodes the representations of the annotators. Given a text instance x_i and an annotator embedding, we compute the annotator-aware embedding as:

$$g(x_i, a_j) = e(x_i) \oplus f(a_j),$$

Where $e(x_i)$ is the text embedding, $f(a_j)$ is the embedding of annotator a_j , and \oplus is the fusion operation that we treat it as a simple addition in our work. This fused, annotator-aware representation is then fed to a classification model to make a prediction.

A few methods utilize this approach. For example, the approach proposed by ([Mokhberian et al., 2023](#)) directly incorporates the annotator embeddings into the text representations. We refer to this method as Annotator Aware Representations for Texts (AART) in this paper. Another technique is ([Deng et al., 2023](#)), which additionally incorporates annotation embeddings. We refer to this method as MichEmbed in this paper.

4 Methodology: Noise Robust Annotator Embedding (NRA-Embed)

Recent SOTA studies ([Mokhberian et al., 2023](#)), [Deng et al. \(2023\)](#) demonstrated that embedding-

based approaches have outstanding potential for multi-annotator classification. Additionally, [Mokhberian et al. \(2023\)](#), showed that contrastive loss enables learning from similarities and dissimilarities among annotators’ label choices, a capability absent in multitask approaches that optimize separate heads independently. However, our experiments (depicted in Figure 1) demonstrate that the performance of these methods is degraded in the presence of noisy annotations. The underlying reason is that conventional contrastive losses, such as InfoNCE ([Oord et al., 2018](#); [Chen et al., 2020](#)), often fail to learn accurate embeddings that represent annotators’ opinions in noisy environments, since inconsistent or noisy annotations can distort the learning signals and hinder the model’s ability to form coherent representations. Therefore, we need a contrastive loss that is robust to annotation noise and tunable to emphasize confident and informative annotation signals while down-weighting uncertain or potentially noisy ones. Motivated by this, we propose utilizing Robust InfoNCE (RINCE) ([Chuang et al., 2022](#)) and adapting it to our problem setting.

RINCE builds on the insight that contrastive learning with noisy representations can be interpreted as a binary classification with noisy labels over pairwise views (in image context), assigning a label of 1 if the views co-occur (joint distribution) and -1 if sampled independently (product of marginals)([Chuang et al., 2022](#)). This interpretation aligns well with our setting, where each view corresponds to an annotator’s label on a given input; we treat annotator pairs as positive (label 1) if they agree on the label of a text instance and negative (label -1) if they disagree. This positive/negative dichotomy justifies our focus on binary classification datasets, enabling clear alignment between label agreement, disagreement, and contrastive objectives.[Ghosh et al. \(2015\)](#) demonstrates that *symmetric* loss functions offer robustness to label noise in binary classification tasks. RINCE introduces a *symmetric adaptation of InfoNCE* that satisfies the symmetry condition in binary classification, and thus guarantees robustness against noisy representations. Specifically, a symmetric contrastive learning objective should have the following form([Chuang et al., 2022](#)):

$$\mathcal{L}(s) = \underbrace{\ell(s^+; 1)}_{\text{Positive Pair}} + \lambda \sum_{i=1}^K \underbrace{\ell(s_i^-; -1)}_{K \text{ Negative Pairs}} \quad (1)$$

Where the first term is the loss of the positive pair, and the second term is the sum of losses of K negative pairs. $\lambda > 0$ is a density weighting term controlling the ratio between positive (class 1) and negative (class -1) pairs.

Based on the idea of robust symmetric classification loss, the Robust InfoNCE (RINCE) loss is defined as ([Chuang et al., 2022](#)):

$$\mathcal{L}_{\text{RINCE}}^{\lambda, q}(s) = \frac{e^{q \cdot s^+}}{q} + \frac{\left(\lambda \cdot \left(e^{s^+} + \sum_{i=1}^K e^{s_i^-} \right) \right)^q}{q} \quad (2)$$

Where s^+ is the score (similarity measure like cosine similarity) for a positive (agreement) pair and s_i^- are scores for negative (disagreement) pairs. A tunable parameter $q \in (0, 1]$ is introduced to interpolate between the robustness of RINCE and the expressive power of InfoNCE. When $q = 1$, RINCE becomes a contrastive loss that fully satisfies the symmetry property in Equation (1) and offers resistance to annotation noise.

To jointly learn task performance and annotator embeddings, we pass combined embeddings $g(x_i, a_j)$ through a classification layer to predict the annotator’s label for each instance. We optimize the following composite objective function:

$$\mathcal{L} = \mathcal{L}_{\text{CE}} + \lambda_1 \sum_j \|f(a_j)\|_2^2 + \lambda_2 \sum_{j, j'} \mathcal{L}_{\text{RINCE}}(j, j') \quad (3)$$

The first term, \mathcal{L}_{CE} , is a standard cross-entropy loss used to predict annotator a_j ’s label assignment for input x_i from the joint representation $g(x_i, a_j)$. The second term applies an ℓ_2 regularization penalty to the annotator embeddings $f(a_j)$ to mitigate overfitting and promote generalization. The third term incorporates the RINCE contrastive loss to structure the annotator embedding space: label-concordant annotator pairs $(a_j, a_{j'})$ are treated as positives and encouraged to align, while label-discordant pairs are repelled. This contrastive objective enforces consistency among reliable annotators while remaining robust to annotation noise.

Our NRA-Embed framework adopts a similar

embedding-based multi-annotator architecture to AART (Mokhberian et al., 2023), with a critical modification: replacing the conventional InfoNCE objective with RINCE for contrastive learning. This design choice enables a controlled evaluation of robust contrastive learning while preserving architectural consistency and comparability with existing annotator embedding approaches. Importantly, the RINCE loss is architecture-agnostic and can be integrated into any similarity-based annotator embedding framework that relies on pairwise representation comparisons, without requiring changes to the underlying model structure or training pipeline. While our experiments focus on a specific instantiation of this framework, extending RINCE to other annotator-aware architectures and tasks remains an important direction for future work. We provide a theoretical analysis of RINCE’s robustness properties—including its gradient dynamics and the q -parameterized exploration–exploitation trade-off—in Appendix C.

5 Experimental Setup

5.1 Datasets

We use the following datasets in our evaluation.

- **The Multi-Domain Agreement Dataset (MDA)** This dataset comprises 9,814 English tweets drawn from three topical domains (the Black Lives Matter movement, the 2020 U.S. election, and the COVID-19 pandemic), each independently annotated for offensiveness by five crowdworkers via Amazon Mechanical Turk (Leonardelli et al., 2021).
- **English Perspectivist Irony Corpus (EPIC)** EPIC is a disaggregated dataset for irony detection consisting of 3,000 short social-media conversations represented as *post–reply* pairs collected from Twitter and Reddit, spanning five regional varieties of English (UK, US, Ireland, Australia, India) (Frenda et al., 2023; epi). We preprocess EPIC into one instance per annotated post–reply pair and the assigned binary label to the reply. We encode each instance’s text item as a two-segment input where the post is segment A, and the reply is segment B using the tokenizer’s native pair encoding (i.e., passing text and text_pair to the tokenizer). Each instance is annotated by 5 annotators in this dataset.

- **HS-Brexit Dataset (HSB)**

The HS-Brexit Dataset (Akhtar et al., 2021) contains Brexit-related tweets annotated for hate speech, aggressiveness, offensiveness, and stereotypes by six annotators from diverse demographic groups. The dataset captures multiple perspectives on abusive language and uses a polarization index to measure annotator disagreement. The HS-Brexit Dataset has 6 crowdworkers per sample.

5.2 Baseline Models

We compare against the following baseline methods:

- **Multitask:** We follow the approach proposed by (Davani et al., 2022b) which involves one fully-connected layer for each annotator with a shared RoBERTA model.
- **AART:** We evaluate the approach introduced by (Mokhberian et al., 2023) which utilizes an embedding for each annotator as well as a contrastive loss objective and a single fully-connected classification layer built off of a RoBERTA backbone in our evaluations.
- **MichEmbed:** We follow the approach by (Deng et al., 2023) which utilizes annotator embeddings as well as annotation embeddings and a single fully-connected classification layer built off of RoBERTA in our experiments.

5.3 Noise Injection

For MDA, EPIC, and HSB, we inject label noise by independently flipping individual annotator labels with a specified probability. Specifically, for each annotator, we assign a noise level (ranging from 0% to 40% in our experiments), and then independently flip each of that annotator’s labels with probability equal to their assigned noise level. For binary classification tasks, label flips are implemented by inverting the label (i.e., $0 \rightarrow 1$ or $1 \rightarrow 0$). This per-annotator, per-label flipping approach allows us to model annotator-specific noise patterns while maintaining the multi-annotator structure of the data. Any annotators who did not contribute to a given sample were excluded from the noise injection process and thus did not affect the training loss. We also evaluated on the multiclass Sentiment Analysis (SNT) (Díaz et al., 2018) dataset using symmetric noise injection for completeness.

Results and discussion are reported in Appendix B. All label noise is injected only in the training set. Test sets remain clean to enable fair evaluation and see how effective the various methods are at filtering out noise while keeping subjective opinions and minority voices.

5.4 Implementation Details

We implemented the classification models using the HuggingFace transformers library (version 4.39) (Wolf et al., 2020). Our experimental setup for the annotator embedding approach for subjective classification closely resembled that of Mokhberian et al. (2023) where Annotator embeddings are initialized from a standard random distribution with the same dimensionality as the model’s hidden layers. For all dataset experiments, we trained the models for 10 epochs. We used this to train our baseline and the other models, and then introduced our noise correction method. We used the pre-trained Roberta-base (Liu et al., 2019) model as the underlying architecture. Optimization was conducted using the AdamW optimizer with a learning rate of 1e-5 and a weight decay of 0.01. A linear decay scheduler with zero warm-up steps was then applied. The hyperparameters $\lambda_1 = 0.1$ and $\lambda_2 = 0.1$ were used for all experiments. The q parameter values used for each experiment are specified in the results tables. All reported results (unless specified otherwise) represent means over 5 independent runs with different random seeds, with standard deviations reported in tables.

5.5 Evaluation Metrics

Mean-Annotator F1 Score

This study is driven by the need to preserve minority annotator perspectives that are often lost when labels are naively aggregated. To that end, we evaluate the model’s performance for each annotator a_j on each test item x_i , comparing the true labels y_{ij} against the model’s predictions. We then summarize these per-annotator results via the **Mean-Annotator F1** (Mokhberian et al. (2023)), defined as the average macro-F1 score across all J annotators:

$$\text{Mean-Annotator F1} = \frac{1}{J} \sum_{j=1}^J \text{F1}(a_j),$$

where $\text{F1}(a_j)$ is the macro-F1 score computed for annotator a_j over all test items x_i .

6 Results

We present the main results in Table 1 as mean-annotator f1 scores. This metric weights all annotators equally regardless of contribution volume, ensuring sparse annotator representation. Because Majority Vote produces a single item-level prediction rather than annotator-specific outputs, we evaluate it using accuracy as well as F1, while all annotator-aware models are evaluated using Mean-Annotator F1. We see that our method performs the best or competitively in many conditions, with the clearest gains at high noise levels.

6.1 Annotator Embedding Approach + Noise Robustness Enhancements

We compare our NRA-Embed Approach against several baselines to measure gains in embedding stability under synthetic annotation noise. Table 1 presents mean Annotator-Aware F1 scores on three benchmark datasets, MDA, EPIC, and HSB, at no noise, 20% noise, 30% noise, and 40% noise. Our Noise-Robust Annotator Embedding Approach consistently improves or performs competitively to the baselines.

6.2 Impact of Parameter q on Noise-Robustness in Annotator Embedding

Table 2 examines the effect of the robustness parameter q across noise levels. The parameter q controls the trade-off between learning from hard samples (low q , closer to standard InfoNCE) and robustness to noisy labels (high q). We observe that $q=1.0$ consistently achieves the highest or near-highest numerical performance at 40% noise across all datasets. This aligns with our intuition that stronger robustness helps under severe label corruption. At moderate noise levels (20%), the pattern is less consistent, with optimal q varying by dataset. However, pairwise comparisons across q values do not reach statistical significance (Welch’s t-test, $p > 0.05$ for all comparisons, $n=5$ runs per configuration). This suggests that while $q=1.0$ offers a slight numerical advantage at high noise, performance is generally stable across q values. In practice, this is useful, as practitioners can default to $q=1.0$ for annotation settings they expect to be noisy without requiring extensive hyperparameter tuning.

Dataset	Noise Level	Majority Vote	Multitask	MichEmbed	AART	Ours
MDA	No Noise	0.6031 (0.5206)	0.5121 ± 0.006	0.7867 ± 0.006	0.7912 ± 0.005	0.7927 ± 0.005
	20% Noise	0.6007 (0.5191)	0.4961 ± 0.018	0.7679 ± 0.018	<u>0.7902 ± 0.0033</u>	0.7905 ± 0.003
	30% Noise	0.6127 (0.5253)	0.4740 ± 0.013	0.7124 ± 0.018	<u>0.7814 ± 0.011</u>	0.7839 ± 0.01
	40% Noise	0.5850 (0.5114)	0.4265 ± 0.035	0.5211 ± 0.059	<u>0.7183 ± 0.026</u>	0.7208 ± 0.03
EPIC	No Noise	0.6222(0.5306)	0.5131 ± 0.008	0.7879 ± 0.007	<u>0.7925 ± 0.005</u>	0.7935 ± 0.005
	20% Noise	0.6101(0.5236)	0.4938 ± 0.023	0.7728 ± 0.022	<u>0.7881 ± 0.0015</u>	0.7891 ± 0.003
	30% Noise	0.5996(0.5177)	0.4741 ± 0.018	0.7073 ± 0.017	<u>0.7783 ± 0.012</u>	0.7792 ± 0.008
	40% Noise	0.5767(0.5074)	0.4342 ± 0.032	0.5116 ± 0.058	<u>0.7195 ± 0.027</u>	0.7228 ± 0.028
HSB	No Noise	0.8524 (0.7486)	0.9318 ± 0.004	0.9347 ± 0.004	0.9304 ± 0.003	<u>0.9329 ± 0.004</u>
	20% Noise	0.7776 (0.6509)	0.9134 ± 0.008	0.9112 ± 0.011	0.9096 ± 0.008	<u>0.9112 ± 0.011</u>
	30% Noise	0.7129 (0.5798)	0.8724 ± 0.014	<u>0.8783 ± 0.02</u>	0.8757 ± 0.03	0.8793 ± 0.019
	40% Noise	0.6639 (0.5288)	<u>0.7740 ± 0.03</u>	0.7583 ± 0.05	0.7505 ± 0.08	0.7830 ± 0.082

Table 1: Mean annotator-level F1 across three datasets (MDA, HSB, and EPIC) under varying synthetic label–noise levels. Best is **bold**; second best is underlined. Our method performs best or second best in most conditions, especially under high-noise conditions (The q-value chosen varies depending on what provides the best results). In brackets after each F1 in Majority Vote (MV), we report the baseline accuracy computed with the same item–annotator metric: for each item x_i , we collapse annotator labels $\{y_{ij}\}_j$ to a single consensus $y_i^{\text{MV}} = \text{mode}(\{y_{ij}\}_j)$ (ties broken deterministically by the training-set class prior), predict y_i^{MV} for all annotators, and score $\frac{|\{(i,j): y_i^{\text{MV}}=y_{ij}\}|}{|\{(i,j)\}|}$.

6.3 Impact of Renegade Annotators on Model Performance

We analyze model robustness in realistic scenarios involving renegade annotators, individuals who intentionally provide malicious or random annotations. To do this, we randomly choose 10% of annotators to have very high noise, that is 70% of their annotations are perturbed. Experiments compare our proposed Noise-Robust Annotation Embedding method against the Task-Based approach. Results demonstrate that no method consistently outperforms the others under conditions of renegade annotation. This suggests that robustness to systematic adversarial annotators may require explicit detection or weighting mechanisms beyond contrastive robustness. Detailed performance metrics are provided in Table 3.

7 Discussion

Our results demonstrate several key trends that hold consistently across datasets and noise configurations, offering both theoretical and practical insight into designing models for subjective classification under noisy annotation.

RINCE often improves model robustness across noise scenarios. Across most tested noise levels, our approach led to a notable increase in

mean annotator F1 and reduced degradation under high-noise conditions. This supports our hypothesis that subjective NLP tasks require not only modeling of annotator identity but also a mechanism to counteract annotation noise.

Annotator embedding models outperform multitask learning. Our results show that models such as MichEmbed, AART, and our approach NRA-Embed, which learn annotator embeddings to modulate shared representations, outperform multitask approaches with separate prediction heads per annotator. We hypothesize that this advantage arises from parameter sharing and regularization effects, embedding-based models can exploit commonalities across annotators while still personalizing behavior, whereas multitask heads may overfit when annotation coverage is sparse or imbalanced. Additionally, embedding approaches inherently support more efficient transfer across annotators and can generalize better when annotators have limited individual data. We further analyze per-annotator performance across methods on HSB in Appendix B.

Multitask models degrade in performance with sparse annotators. The multitask model performed the worst with the SNT dataset (Appendix A). This is likely due to how many annotators there are compared to how many samples they

	No Noise			20% Noise			40% Noise		
	$q=0.5$	$q=0.75$	$q=1.0$	$q=0.5$	$q=0.75$	$q=1.0$	$q=0.5$	$q=0.75$	$q=1.0$
Rince_q									
MDA	0.7915 ± 0.005	0.7906 ± 0.005	0.7927 ± 0.004	0.7891 ± 0.003	0.7905 ± 0.003	0.7901 ± 0.004	0.7118 ± 0.043	0.7099 ± 0.037	0.7174 ± 0.027
HSB	0.9315 ± 0.003	0.9302 ± 0.005	0.9321 ± 0.004	0.9112 ± 0.011	0.9070 ± 0.009	0.9028 ± 0.016	0.7635 ± 0.096	0.7678 ± 0.100	0.7782 ± 0.110
EPIC	0.7926 ± 0.006	0.7913 ± 0.006	0.7935 ± 0.005	0.7869 ± 0.0022	0.7891 ± 0.003	0.7883 ± 0.0028	0.7207 ± 0.04	0.7167 ± 0.027	0.7228 ± 0.028

Table 2: Mean Annotator F1 Scores (top row) with standard deviations (bottom row) showing the effect of q on robustness across noise levels.

Model	MDA	EPIC	HSB
Multitask	0.498	0.642	0.806
MichEmbed	0.784	0.764	0.897
AART	0.795	0.790	0.891
NRAEmbed	0.796	0.789	0.891

Table 3: Mean Annotator F1 Scores comparing model robustness to renegade annotators (malicious/random annotation behavior). Bolded values highlight the best-performing approach across datasets (Results reported are average of 2 runs).

DS	#A	#E/#A	#CW	#S	#L
MDA	819	60	5	44k	2
EPIC	74	192	5	14.7k	2
HSB	6	952	6	5.7k	2
SNT	1481	41	4–12	60.4k	5

Table 4: Dataset statistics. #A is the number of annotators, #E/#A is the average number of examples per annotator, #CW is the number of annotators per example, #S is the number of samples, and #L is the number of labels.

annotated on average, which is very few, creating sparse annotators (see Table 4). This aligns with Mokherian et al. (2023), who showed through statistical parity analysis that multitask models exhibit the highest performance disparity for sparse annotators compared to embedding-based approaches, with more balanced and fair results. On the other hand, the multitask model performed very well on HSB which had a much smaller amount of annotators who each labeled many samples.

Binary classification aligns with contrastive objectives. Our main experiments focus on binary tasks where contrastive agree/disagree signals

match the label structure. Supplementary experiments on the multiclass Likert-scale SNT dataset (Appendix A) show that contrastive methods underperform annotation-embedding approaches, suggesting ordinal settings require class-sensitive modifications.

These findings reinforce the need to view subjective learning as a two-fold challenge: embracing disagreement while resisting noise. Annotator-aware models alone are not sufficient if they assume all disagreement is meaningful; conversely, noise-robust objectives without subjectivity modeling may conflate diverse opinions with error. Our work shows that integrating both perspectives yields the most reliable performance, and that simple but principled interventions, like swapping In-fonCE for RINCE, can offer significant gains in real-world annotation environments.

8 Conclusion

We explored the distinction between label noise and subjective disagreement in subjective learning tasks. Most prior works only consider one or the other; however, these factors are intertwined. Disagreement is core to many human-centered activities and should be accounted for when building datasets. We address this issue by separating label disagreement and label noise through our NRA-Embed approach. Our benchmarking of existing multi-annotator models provides a strong baseline for developing advanced models that can tolerate unique noise patterns. Our results suggest that embedding-based approaches are the superior methodology for training in multi-annotator cases. Furthermore, we recommend that raw labels should be released, however noisy, so that issues with label noise can be directly addressed by model.

9 Limitations

Our evaluation relies on synthetic noise injection (label flipping, random perturbations), which may not fully capture the complexity of real-world annotation errors. Future work should explore noise models that better reflect actual annotator behavior, such as systematic biases, attention lapses, or task-specific confusion patterns.

Per-annotator modeling becomes computationally expensive in datasets with many annotators. For large-scale applications, future work should investigate annotator grouping strategies or approximation methods to reduce computational overhead while preserving perspective diversity.

Our approach depends on existing datasets with multiple annotations per sample. Like any crowd-sourced dataset, these have inherent limitations in representing population diversity due to cost constraints (typically 3-5 annotators per sample), annotator recruitment biases, and geographic skew. While our method robustly models available perspectives, it cannot capture viewpoints absent from the annotation pool. Future work should explore methods to incorporate high-quality minority perspectives, such as targeted recruitment or active learning strategies that identify underrepresented groups.

10 Ethics Statement

In data annotation, capturing the full spectrum of annotator perspectives is crucial for producing fair and representative models. However, factors like annotator fatigue and shifting judgments over time can conceal the true range of opinions present in large datasets.

To address this, we propose drawing on insights from the entire annotator pool—including those who contribute less frequently—rather than focusing solely on the most active contributors. Incorporating these “sparser” judgments broadens the diversity of viewpoints the model sees, yielding predictions that are both more robust and more nuanced.

That said, this inclusive approach carries its own risks. A small, coordinated subgroup of annotators might exert undue influence, and any biases embedded within our large language model infrastructure could further distort individual annotations. Even so, we argue that the benefits of embracing a wider array of voices—enhancing both inclusivity and resilience in AI systems—far outweigh these

potential drawbacks.

In addition, we utilized the assistance of AI (LLMs) in some research, coding, and writing.

11 Acknowledgements

References

- [Multilingual-Perspectivist-NLU/EPIC \(epicorpus\) dataset card](#). Hugging Face Datasets. Accessed 2026-01-04.
- Sohail Akhtar, Valerio Basile, and Viviana Patti. 2021. Whose opinions matter? perspective-aware models to identify opinions of hate speech victims in abusive language detection. *arXiv preprint arXiv:2106.15896*.
- Hala Al Kuwatly, Maximilian Wich, and Georg Groh. 2020. Identifying and measuring annotator bias based on annotators’ demographic characteristics. In *Proceedings of the fourth workshop on online abuse and harms*, pages 184–190.
- Cecilia Ovesdotter Alm. 2011. Subjective natural language problems: Motivations, applications, characterizations, and implications. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 107–112.
- Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor, and Kevin McGuinness. 2019. Unsupervised label noise modeling and loss correction. In *International conference on machine learning*, pages 312–321. PMLR.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmLR.
- Veronika Cheplygina and Josien PW Pluim. 2018. Crowd disagreement about medical images is informative. In *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis: 7th Joint International Workshop, CVII-STENT 2018 and Third International Workshop, LABELS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*, pages 105–111. Springer.
- Ching-Yao Chuang, R Devon Hjelm, Xin Wang, Vibhav Vineet, Neel Joshi, Antonio Torralba, Stefanie Jegelka, and Yale Song. 2022. Robust contrastive learning against noisy views. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16670–16681.
- Aida Mostafazadeh Davani, Mark Díaz, and Vinodkumar Prabhakaran. 2022a. Dealing with disagreements: Looking beyond the majority vote in subjective annotations. *Transactions of the Association for Computational Linguistics*, 10:92–110.
- Aida Mostafazadeh Davani, Mark Díaz, and Vinodkumar Prabhakaran. 2022b. [Dealing with disagreements: Looking beyond the majority vote in subjective annotations](#). *Transactions of the Association for Computational Linguistics*, 10:92–110.
- Alexander Philip Dawid and Allan M Skene. 1979. Maximum likelihood estimation of observer error-rates using the em algorithm. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 28(1):20–28.

861 Maarten Sap, Swabha Swayamdipta, Laura Vianna, Xuhui
862 Zhou, Yejin Choi, and Noah A Smith. 2021. Annotators with attitudes: How annotator beliefs and identities bias toxic language detection. *arXiv preprint arXiv:2111.07997*.

866 Rion Snow, Brendan O’connor, Dan Jurafsky, and Andrew Y
867 Ng. 2008. Cheap and fast—but is it good? evaluating non-expert annotations for natural language tasks. In *Proceedings of the 2008 conference on empirical methods in natural language processing*, pages 254–263.

871 Alexandra Uma, Dina Almanea, and Massimo Poesio. 2022. Scaling and disagreements: Bias, noise, and ambiguity. *Frontiers in Artificial Intelligence*, 5:818451.

874 Alexandra N Uma, Tommaso Fornaciari, Dirk Hovy, Silviu
875 Paun, Barbara Plank, and Massimo Poesio. 2021. Learning from disagreement: A survey. *Journal of Artificial Intelligence Research*, 72:1385–1470.

878 Zeerak Waseem. 2016. Are you a racist or am i seeing things? annotator influence on hate speech detection on twitter. In *Proceedings of the first workshop on NLP and computational social science*, pages 138–142.

882 Zeerak Waseem and Dirk Hovy. 2016. Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL student research workshop*, pages 88–93.

886 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45.

893 Xueying Zhan, Yaowei Wang, Yanghui Rao, and Qing Li. 2019. Learning from multi-annotator data: A noise-aware classification framework. *ACM Transactions on Information Systems (TOIS)*, 37(2):1–28.

897 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2016. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*.

901 Xiang Zhou, Yixin Nie, and Mohit Bansal. 2021. Distributed nli: Learning to predict human opinion distributions for language reasoning. *arXiv preprint arXiv:2104.08676*.

Dataset	Noise	MV	Multitask	MichEmbed	AART	Ours
SNT	0%	0.245	0.267 \pm .007	0.517 \pm .009	0.521\pm.006	0.494
	10%	0.242	0.268 \pm .002	<u>0.492\pm.002</u>	0.500\pm.008	0.489
	20%	0.251	0.259 \pm .008	0.483\pm.012	0.475 \pm .016	0.471
	30%	0.235	0.248 \pm .013	0.446\pm.001	0.439 \pm .010	0.441
	40%	0.236	0.225 \pm .000	0.418\pm.001	<u>0.403\pm.021</u>	0.406

Table 5: Mean annotator-level F1 on the multiclass SNT dataset under varying noise levels. Best is **bold**; second best is underlined. Unlike binary datasets, contrastive methods (AART, NRA-Embed) do not consistently outperform MichEmbed. MV = Majority Vote baseline.

A Supplementary Results on Multiclass Data

To assess generalization beyond binary classification, we evaluated on the Sentiment Analysis Dataset (SNT), introduced by [Díaz et al. \(2018\)](#). SNT is a sentiment classification resource aimed at addressing age-related biases in sentiment models, leveraging text from older adults’ blog posts containing age-related terms such as “old” and “young.” The dataset uses a 5-class Likert scale and contains 1,481 annotators averaging 41 examples each, with 4–12 annotators per sample (see Table 4).

A.1 Noise Injection for Multiclass Labels

For SNT, we applied symmetric noise injection following a similar procedure to the binary datasets. Each annotators labels were flipped independently according to the noise percentage, to one of the four other classes. As with our binary experiments, all noise was injected only in the training set; test labels remained clean.

A.2 Main Results on SNT

A.3 Effect of Parameter q on SNT

q	No Noise			20% Noise			40% Noise		
	0.5	0.75	1.0	0.5	0.75	1.0	0.5	0.75	1.0
SNT	0.494	0.495	0.502	0.471	0.463	0.464	0.400	0.396	0.400
	\pm .003	\pm .005	\pm .005	\pm .013	\pm .020	\pm .011	\pm .012	\pm .019	\pm .017

Table 6: Effect of robustness parameter q on SNT. Unlike binary datasets, the pattern is inconsistent: $q = 0.5$ performs best at 20% noise, while $q = 1.0$ is marginally better at 0% and 40% noise.

A.4 Discussion: Limitations of Contrastive Loss for Likert-Scale Data

Our main experiments focus on binary classification tasks, where the contrastive loss’s agree/disagree signal naturally aligns with the label

structure. The SNT results reveal important limitations when extending contrastive approaches to multiclass ordinal data.

Likert scales degrade contrastive loss performance.

On SNT, which uses Likert scale classification, we find that previously strong contrastive approaches like AART and our NRA-Embed method degrade in performance relative to MichEmbed. We hypothesize this is due to the fundamental mismatch between contrastive objectives and ordinal label structure. A contrastive loss treats labels “Strongly Agree” and “Moderately Agree” as a negative pair in exactly the same way it treats “Strongly Agree” and “Strongly Disagree”—both are simply instances of disagreement. This binary treatment discards the semantic similarity between adjacent ordinal categories, leading to suboptimal representation learning.

In contrast, MichEmbed, which relies on a combination of Annotator and Annotation Embeddings rather than contrastive loss, performs strongly on SNT. The annotation embeddings can implicitly capture relationships between label values, making this approach better suited for ordinal classification.

Multitask models struggle with sparse annotators.

The multitask model performed particularly poorly on SNT (Table 5). This is likely due to the large number of annotators (1,481) relative to their average annotation count (41 samples each), creating extremely sparse per-annotator training data (see Table 4). In contrast, multitask performed well on HSB, which has only 6 annotators who each labeled approximately 952 samples. This pattern aligns with findings from [Mokhberian et al. \(2023\)](#), who demonstrated that embedding-based approaches provide more equitable performance for sparse annotators.

Implications for future work.

These findings suggest that extending noise-robust contrastive learning to multiclass ordinal settings requires class-sensitive modifications. Future work should explore ordinal-aware contrastive objectives that weight negative pairs by label distance, for instance, penalizing “Strongly Agree” vs. “Strongly Disagree” more heavily than “Strongly Agree” vs. “Moderately Agree.”

978 B Per-Annotator Analysis

979 To examine how different models capture individual
980 annotator perspectives, we conduct a detailed
981 analysis on the HSB dataset, which contains anno-
982 tations from six crowdworkers (A1–A6) with di-
983 verse labeling patterns. Table 7 presents F1 scores
984 for each annotator across all methods and noise
985 conditions.

986 B.1 Annotator Heterogeneity

987 Even under clean conditions, substantial perfor-
988 mance variation exists across annotators. We ob-
989 serve two distinct groups: *high-agreement annota-*
990 *tors* (A1, A2, A3) who achieve F1 scores above
991 0.95 with embedding-based methods, and *low-*
992 *agreement annotators* (A4, A5, A6) who score be-
993 low 0.93. Notably, A5 consistently exhibits the
994 lowest F1 across all methods (0.79–0.89 at 0%
995 noise), suggesting this annotator’s labeling patterns
996 diverge most from learnable representations. This
997 heterogeneity reflects genuine differences in anno-
998 tator behavior rather than random variation, as the
999 ranking remains stable across methods.

1000 B.2 Cross-Annotator Equity

1001 The Gap metric (Max – Min F1) measures how
1002 equitably a model performs across annotators. At
1003 0% noise, NRA-Embed achieves the smallest Gap
1004 (0.07), followed by AART (0.08), indicating that
1005 contrastive embedding approaches produce more
1006 balanced representations than Multitask (0.09) or
1007 MichEmbed (0.10). This advantage persists at 20%
1008 noise, where NRA-Embed maintains a Gap of 0.08
1009 while Multitask’s Gap increases sharply to 0.15.
1010 The widening disparity for Multitask suggests that
1011 per-annotator classification heads are more suscep-
1012 tible to noise-induced overfitting for low-agreement
1013 annotators.

1014 B.3 Performance Under Noise

1015 At 40% noise, the pattern shifts. While NRA-
1016 Embed achieves the second-highest Mean F1 (0.78)
1017 among all methods, outperforming AART (0.76)
1018 and MichEmbed (0.75), its Gap increases to 0.16.
1019 Examining individual annotators reveals why: A5,
1020 the hardest annotator, degrades from 0.89 to 0.69
1021 for NRA-Embed ($\Delta = -0.20$), whereas Multitask
1022 shows smaller degradation for A5 ($\Delta = -0.15$).
1023 However, NRA-Embed shows the smallest degra-
1024 dation for A3 ($\Delta = -0.11$), the annotator where
1025 it achieves the best 40%-noise performance (0.85).

This suggests that robust contrastive learning helps
preserve representations for annotators whose pat-
terns align well with the learned embedding space,
but provides less protection for outlier annotators
under extreme noise.

1031 B.4 Summary

1032 This analysis reveals that: (1) annotator difficulty is
1033 consistent across methods, with A5 being hardest
1034 and A1–A3 being easiest; (2) NRA-Embed pro-
1035 vides the most equitable cross-annotator perfor-
1036 mance at low-to-moderate noise levels; and (3) at
1037 extreme noise, a trade-off emerges between over-
1038 all performance and cross-annotator equity. The
1039 performance of multitask at 40% noise may be
1040 attributed to the lack of sparse annotators that is
1041 unique to HSB (Akhtar et al., 2021).

1042 C Theoretical Analysis of Robustness in 1043 RINCE

1044 In this appendix, we provide the theoretical justifi-
1045 cation for utilizing the Robust InfoNCE (RINCE)
1046 loss (Chuang et al., 2022) in the context of sub-
1047 jective NLP tasks. We analyze and contrast the
1048 gradient dynamics of RINCE against the standard
1049 InfoNCE loss to demonstrate how the hyperparam-
1050 eter q controls the trade-off between learning from
1051 hard samples (exploration) and robustness to label
1052 noise (exploitation), and discuss why it’s a perfect
1053 fit for our problem requirements. For more details
1054 and the theoretical proofs, we refer to the RINCE’s
1055 original paper.

1056 C.1 Formulation

1057 Let s^+ denote the similarity score of a positive pair
1058 and $\{s_i^-\}_{i=1}^K$ denote the set of K negative similarity
1059 scores. We define the partition function $Z = e^{s^+} +$
1060 $\sum_{i=1}^K e^{s_i^-}$.

1061 The standard **InfoNCE** loss is defined as:

$$1062 \mathcal{L}_{\text{InfoNCE}}(s) = -\log\left(\frac{e^{s^+}}{Z}\right) = -s^+ + \log Z \quad (4)$$

1063 The **RINCE** loss utilizes a generalized Box-Cox
1064 transformation. For parameters $q \in (0, 1]$ and $\lambda >$
1065 0, the loss is defined as:

$$1066 \mathcal{L}_{\text{RINCE}}^{\lambda, q}(s) = -\frac{e^{qs^+}}{q} + \frac{\lambda}{q} Z^q \quad (5)$$

1067 Note that as $q \rightarrow 0$, Eq. (5) asymptotically recovers
1068 the behavior of the InfoNCE loss (up to a scaling
1069 factor).

Table 7: Per-annotator F1 scores on the HSB dataset (3 runs). **Bold** indicates best performance per annotator at each noise level. Mean averages across annotators; Gap = Max – Min F1 (lower indicates more equitable performance).

	0% Noise				20% Noise				40% Noise			
	Multi	MichE	AART	NRA	Multi	MichE	AART	NRA	Multi	MichE	AART	NRA
A1	.97	.98	.96	.96	.97	.95	.94	.93	.84	.81	.79	.80
A2	.97	.96	.95	.95	.96	.95	.94	.93	.82	.78	.77	.80
A3	.97	.97	.96	.96	.97	.96	.94	.94	.82	.81	.82	.85
A4	.89	.90	.89	.89	.86	.86	.87	.86	.76	.69	.75	.77
A5	.88	.88	.89	.89	.82	.85	.85	.85	.73	.67	.66	.69
A6	.92	.93	.92	.93	.89	.91	.90	.90	.76	.70	.75	.77
Mean	.93	.94	.93	.93	.91	.91	.91	.90	.79	.75	.76	.78
Gap	.09	.10	.08	.07	.15	.11	.10	.08	.11	.14	.16	.16

C.2 Gradient Analysis and Noise Robustness

To understand the robustness mechanism, we analyze the gradient of the loss with respect to the positive score s^+ . This gradient dictates how strongly the model pushes positive pairs together.

InfoNCE Dynamics ($q \rightarrow 0$). The gradient of the standard InfoNCE loss with respect to s^+ is:

$$\frac{\partial \mathcal{L}_{\text{InfoNCE}}}{\partial s^+} = -(1-p), \quad \text{where } p = \frac{e^{s^+}}{Z} \quad (6)$$

Here, the gradient magnitude is $|1-p|$. This creates a “hard-positive mining” effect:

- When the model is confident ($p \approx 1$), the gradient is near 0.
- When the model is incorrect ($p \approx 0$), the gradient magnitude is maximized (≈ 1).

While beneficial for clean data, this dynamic is detrimental for noisy or subjective NLP tasks. If a sample is mislabeled or highly subjective (a “noisy positive”), p will be small. InfoNCE will generate a large gradient, forcing the model to overfit this noise.

RINCE Dynamics ($q > 0$). Differentiating Eq. (5) with respect to s^+ yields:

$$\begin{aligned} \frac{\partial \mathcal{L}_{\text{RINCE}}}{\partial s^+} &= -e^{qs^+} + \lambda Z^{q-1} e^{qs^+} \\ &= -e^{qs^+} \left(1 - \lambda \left(\frac{e^{s^+}}{Z} \right)^{1-q} \right) \\ &= - \underbrace{e^{qs^+}}_{\text{Weighting Term}} \cdot \underbrace{\left(1 - \lambda p^{1-q} \right)}_{\text{Error Term}} \end{aligned} \quad (7)$$

The critical innovation in RINCE is the **weighting term** e^{qs^+} . Unlike InfoNCE, the gradient magnitude in RINCE is dependent on the absolute value of the similarity score s^+ :

- If s^+ is large (high confidence), e^{qs^+} is large, maintaining the learning signal.
- If s^+ is small (low confidence/noisy sample), $e^{qs^+} \rightarrow 0$.

Consequently, RINCE applies a “soft clipping” to the gradients. If a sample is too difficult (implying potential noise or extreme subjectivity), the e^{qs^+} term suppresses the gradient, effectively ignoring the outlier rather than corrupting the representation space.

C.3 The Exploration-Exploitation Trade-off

The parameter q governs a fundamental trade-off between exploration and exploitation, which is particularly relevant for subjective NLP:

1. **Exploration** ($q \rightarrow 0$, **InfoNCE**): The loss focuses on “hard” samples (tail of the distribution). This pushes the model to explore boundaries and learn fine-grained distinctions. In subjective tasks, hard samples often arise from annotator disagreement or label noise. Excessive exploration here can lead to overfitting to the noise.
2. **Exploitation** ($q \rightarrow 1$, **Robustness**): As q increases, the weighting term e^{qs^+} dominates. The model focuses on “easy” samples where it is already confident (head of the distribution). This represents *exploitation* of clear, agreed-upon signals while discarding ambiguous data points.

For subjective sentiment analysis, where ground truth is often a distribution rather than a single fact, the pure exploration of InfoNCE is dangerous. By selecting a moderate q (e.g., $q = 0.5$), RINCE enforces a necessary degree of exploitation, ensuring the model learns from high-consensus data while

1132 remaining robust to the inherent noise of subjective
1133 annotation.

1134 C.4 Suitability for Noisy, Subjective NLP

1135 Our problem setting involves inherently noisy,
1136 subjective NLP labels (e.g., sentiment, emotion,
1137 stance), where annotators frequently disagree and
1138 supervision is inconsistent. Contrastive pairs derived
1139 from such data are plagued by three distinct
1140 types of noise: (1) **False Positives** (semantically
1141 distinct examples labeled as similar), (2) **False Neg-**
1142 **atives** (semantically identical examples labeled as
1143 different), and (3) **Subjective Variance**. RINCE
1144 addresses these specific challenges through the fol-
1145 lowing mechanisms:

- 1146 • **Robustness to Noisy Positives via Exploita-**
1147 **tion.**

1148 Subjective datasets naturally produce weak
1149 or contradictory positives. As shown in the
1150 gradient analysis, the term e^{qs^+} acts as an au-
1151 tomated confidence gate. When the model
1152 encounters a “noisy positive” (a pair labeled
1153 similar but with low semantic alignment s^+),
1154 the gradient is down-weighted. This implicit
1155 *easy-positive mining* prevents the model from
1156 overfitting to specific annotator errors.

- 1157 • **Bounded Penalty for False Negatives.**

1158 In standard InfoNCE, a False Negative (a posi-
1159 tive pair mislabeled as negative) is treated as
1160 a “hard negative,” generating a gradient that
1161 grows without bound as the model tries to cor-
1162 rect the “error” by pushing the pair apart. This
1163 can destroy the semantic structure. RINCE,
1164 via the power function transformation (x^q),
1165 bounds the penalty for hard negatives. This
1166 ensures that while the model separates distinct
1167 classes, it remains stable in the presence of
1168 contradictory negative labels.

- 1169 • **Geometric Robustness via Wasserstein Dis-**
1170 **tance.**

1171 Theoretically, RINCE optimizes a lower
1172 bound on the Wasserstein Distance Maximiza-
1173 tion (WDM). This formulation ensures that
1174 the learning objective is focused on the global
1175 distributional alignment between representa-
1176 tions and labels, rather than sample-wise exact
1177 matching. In subjective tasks, where individ-
1178 ual labels are noisy samples from a latent sen-
1179 timent distribution, this geometric robustness

is crucial for capturing the underlying data
manifold.

- 1182 • **Implicit soft reweighting vs. explicit noise**
1183 **estimation.**

1184 Unlike methods that require estimating a noise
1185 transition matrix or setting a hard “forget
1186 rate” based on a known noise ratio (e.g., Co-
1187 teaching), RINCE handles noise implicitly via
1188 gradient suppression. While the hyperparam-
1189 eter q controls the degree of robustness, it
1190 modulates the curvature of the loss rather than
1191 setting a hard threshold, allowing the model
1192 to dynamically down-weight samples based
1193 on their individual learning difficulty (s^+ , s^-)
1194 rather than a global noise statistic.