

# Generative Calibration of Inconsistent Annotation for Label Distribution Learning

## supplementary material

### Overview

To demonstrate the adequacy of our work, we need to supplement the experimental results. In this supplementary material, we present the detailed experimental results of the experiments, including recovery experiment, prediction experiment, and ablation experiment.

### Experimental Results

#### Recovery Experiment

	P=4	GCIA	Kmeans	P=5	GCIA	Kmeans	P=6	GCIA	Kmeans	P=7	GCIA	Kmeans
ns												
Cheb (↓)	0.1102	0.2525	0.1299	0.1411	0.1681	0.1891	0.1778	0.1852	0.1808	0.1962	0.2003	<b>0.1945 •</b>
Clark (↓)	0.6897	2.4071	2.2345	0.8367	2.3613	2.3704	0.9834	2.2590	2.2360	1.0860	2.3200	2.0508
Canber (↓)	1.6856	6.5293	5.8222	2.0587	6.2172	6.2971	2.4418	5.9687	5.9669	2.7029	6.2083	5.4264
KL (↓)	0.1370	0.5738	0.1478	0.3224	<b>0.2853 •</b>	0.3238	0.5617	<b>0.3636</b>	<b>0.3237 •</b>	0.7258	<b>0.4235 •</b>	<b>0.4245</b>
Cosine (↑)	0.9467	0.8302	<b>0.9491 •</b>	0.9194	0.9243	<b>0.9016 •</b>	0.8879	<b>0.8940</b>	<b>0.8956 •</b>	0.8703	<b>0.8775 •</b>	<b>0.8730</b>
Intersec (↑)	0.8748	0.5950	0.8406	0.8388	<b>0.7800 •</b>	<b>0.7641</b>	0.7991	0.7781	0.7861	0.7753	0.7567	0.7713
Rho (↑)	0.7389	0.5520	<b>0.7474 •</b>	0.7002	0.7142	0.7076	0.6535	<b>0.6892</b>	<b>0.7109 •</b>	0.6166	<b>0.6643</b>	<b>0.6911 •</b>
spo5												
Cheb (↓)	0.3460	<b>0.2460 •</b>	<b>0.2475</b>	0.3654	<b>0.2590 •</b>	<b>0.2610</b>	0.3880	<b>0.2758 •</b>	<b>0.2767</b>	0.2233	<b>0.2016</b>	<b>0.2014 •</b>
Clark (↓)	0.8587	<b>0.5089 •</b>	<b>0.5166</b>	0.9124	<b>0.5334 •</b>	<b>0.5427</b>	0.9620	<b>0.5622 •</b>	<b>0.5695</b>	0.5398	<b>0.4594</b>	<b>0.4589 •</b>
Canber (↓)	1.2803	<b>0.7920 •</b>	<b>0.8007</b>	1.3694	<b>0.8355 •</b>	<b>0.8468</b>	1.4533	<b>0.8860 •</b>	<b>0.8932</b>	0.7896	<b>0.6840</b>	<b>0.6832 •</b>
KL (↓)	1.6647	<b>0.2163 •</b>	<b>0.2257</b>	1.9085	<b>0.2369 •</b>	<b>0.2480</b>	2.1544	<b>0.2639 •</b>	<b>0.2741</b>	0.5860	<b>0.2093</b>	<b>0.2069 •</b>
Cosine (↑)	0.7929	<b>0.8661 •</b>	<b>0.8647</b>	0.7790	<b>0.8558 •</b>	<b>0.8542</b>	0.7590	<b>0.8406 •</b>	<b>0.8399</b>	0.8865	<b>0.9017</b>	<b>0.9018 •</b>
Intersec (↑)	0.6540	<b>0.7540 •</b>	<b>0.7525</b>	0.6346	<b>0.7410 •</b>	<b>0.7390</b>	0.6120	<b>0.7242 •</b>	<b>0.7233</b>	0.7767	<b>0.7984</b>	<b>0.7986 •</b>
Rho (↑)	0.1257	<b>0.1545 •</b>	<b>0.1349</b>	0.1464	<b>0.1516 •</b>	<b>0.1512</b>	0.0909	<b>0.1070 •</b>	<b>0.0957</b>	0.2457	<b>0.2671</b>	<b>0.2717 •</b>

Table 1: Recovery performance on 2 datasets. Every three columns from left to right represent one group of experiments, with the same P-value used to generate noise. ‘P=4/5/6/7’ represents the difference between the unprocessed noise dataset and the ground-truth. ‘GCIA/KMeans’ represents the difference between the dataset calibrated by the calibrationn method and the ground-truth. Bold indicates that the calibrated dataset is superior to the noise dataset. A dot following the data indicates that the calibrationn method is better.

The experimental results of the recovery experiment were supplemented here. In addition to the 6 datasets mentioned in the main body, the experimental results of 2 additional datasets (No.7, and No.8) mentioned in the article are presented here. The experimental results met our expectations.

#### Prediction Experiment

In the main body, we presented the results of the ‘sj’ and ‘mov’ datasets. Here, we provide additional supplementary experimental results for other datasets.



Figure 1: Prediction performance. “Noise” refers to directly using a dataset that contains noise for training an LDL model. “GCIA/KMeans” represents using GCIA/KMeans as a preprocessing method for the dataset before utilizing it for model training. The processed dataset is then employed in the training process. “ $P=4/5/6/7$ ” indicates different values of professionalism  $P$  used to generate the noisy dataset.

## Ablation Experiment

In the main body, we presented experimental results for two metrics. Here, we provide additional supplementary results for the other two metrics in the ablation experiments.

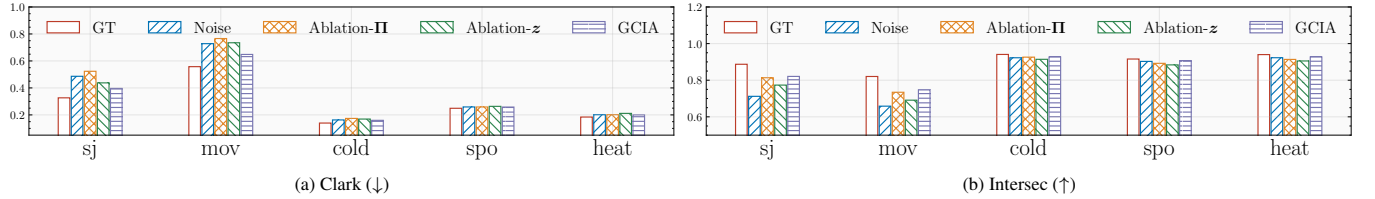


Figure 2: Ablation study. By setting  $K=1$ , we obtained the prediction experiment results without the categorical variable  $z$  (using SABFGS as an example). The figures depict the experimental results at  $P=4$ . ‘Ablation-’ indicates the removal of this component. ‘GCIA’ represents a complete model.