# FEATURE DISCRIMINATION ANALYSIS FOR BINARY AND TERNARY QUANTIZATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

In machine learning, quantization is widely used to simplify data representation and facilitate algorithm deployment on hardware. Considering the fundamental role of classification in machine learning, it is imperative to investigate the impact of quantization on classification. Current research primarily revolves around quantization errors, under the assumption that higher quantization errors generally lead to lower classification performance. However, this assumption lacks a solid theoretical foundation, and often contradicts empirical findings. For instance, some extremely low bit-width quantization methods, such as  $\{0, 1\}$ -binary quantization and  $\{0, \pm 1\}$ -ternary quantization, can achieve comparable or even superior classification accuracy compared to the original non-quantized data, despite exhibiting high quantization errors. To evaluate the classification performance more accurately, we propose to directly investigate the feature discrimination of quantized data, rather than analyze its quantization error. It is found that both binary and ternary quantization methods can surprisingly improve, rather than degrade, the feature discrimination of original data. This remarkable performance is validated through classification experiments on diverse data types, including images, speech and text.

025

004

010 011

012

013

014

015

016

017

018

019

021

#### 1 INTRODUCTION

031 Quantization has been widely applied in machine learning to simplify data storage and computation 032 complexities, while also catering to the requirements of algorithm deployment on digital hardware. 033 In general, this operation will lead to a decrease in classification accuracy (Baras & Dey, 1999; 034 Hoefler et al., 2021), due to reducing the precision of data or model parameters. To achieve a balance between complexity and accuracy, it is crucial to delve into the impact of quantization 035 on classification. Currently, the impact is mainly evaluated through quantization errors, with the assumption that larger quantization errors generally lead to decreased classification accuracy (Lin 037 et al., 2016a). However, this assumption lacks a solid theoretical basis (Lin et al., 2016a), as it merely adopts the quantization principle from signal processing (Gray & Neuhoff, 1998), which primarily focuses on data reconstruction fidelity rather than classification accuracy. In practice, it 040 seems challenging to accurately assess the classification performance solely based on quantization 041 errors. 042

For instance, it has been observed that some extremely low bit-width quantization methods, such 043 as the  $\{0,1\}$ -binary quantization and  $\{0,\pm1\}$ -ternary quantization, which have been successively 044 applied large-scale retrieval (Charikar, 2002) and deep network quantization (Qin et al., 2020; Gho-045 lami et al., 2022), can achieve comparable or even superior classification performance than their 046 full-precision counterparts (Courbariaux et al., 2015; Lin et al., 2016b; Lu et al., 2023), despite 047 suffering from high quantization errors. Apparently, the remarkable classification improvement re-048 sulting from quantization should not be attributed to the high quantization errors. This reveals the inadequacy of quantization errors in assessing the actual classification performance. Due to the absence of a theoretical explanation, the classification improvement induced by quantization has often 051 been regarded as incidental and received little attention. Instead of quantization errors, in the paper we demonstrate that this intriguing phenomenon can be reasonably explained by feature discrimi-052 nation. Following the Fisher's linear discriminant analysis (Fisher, 1936), we here refer to feature discrimination as the ratio between inter-class and intra-class scatters, and evaluate the classification performance based on the rule that the higher the feature discrimination, the easier the classification.
 To the best of our knowledge, this is the first study that exploits feature discrimination to analyze the
 impact of quantization on classification, although it is more direct and reasonable than quantization
 errors in evaluating classification performance. The scarcity of relevant research can be attributed to
 the nonlinearity of the quantization operation, which substantially increases the analytical complex ity of feature discrimination functions.

060 In the paper, it is demonstrated that the impact of the threshold-based binary and ternary quanti-061 zation on feature discrimination can be analyzed, when the data are appropriately modeled using 062 a Gaussian mixture model, with each Gaussian element representing one class of data. The Gaus-063 sian mixture model is chosen here based on two considerations. Firstly, the model has been well-064 established for approximating the distributions of real-world data (Torralba & Oliva, 2003; Weiss & Freeman, 2007) and their feature transformations (Wainwright & Simoncelli, 1999; Lam & Good-065 man, 2000). Secondly, the closure property of Gaussian distributions under linear operations can 066 simplify the analysis of the feature discrimination function. By analyzing the discrimination across 067 varying quantization thresholds, it is found that there exist certain quantization thresholds that can 068 improve the discrimination of original data, thereby yielding improved classification performance. 069 This finding is extensively validated through classification experiments both on synthetic and real data. 071

The related works are discussed as follows. As mentioned earlier, our work should be the first to 072 take advantage of feature discrimination to investigate the impact of quantization on classification. 073 In the filed of signal processing, there have been a few works proposed to reduce the negative impact 074 of quantization on signal detection or classification (Poor & Thomas, 1977; Oehler & Gray, 1995). 075 However, these studies did not employ feature discrimination analysis, distinguishing them from 076 our research in both methodology and outcomes. Specifically, in these studies the model design 077 accounts for both reconstruction loss and classification loss. The classification loss is primarily modeled in several ways, such as directly minimizing the classification error on quantized data 079 (Srinivasamurthy & Ortega, 2002), enlarging the inter-class distance between quantized data (Jana & Moulin, 2000; 2003), reducing the difference between the distributions of quantized data and 081 original data (Baras & Dey, 1999), as well as minimizing the discrepancy in classification before and after quantization (Dogahe & Murthi, 2011). Through analyses of these losses, the classification performance of quantized data can only approach, rather than surpass, the performance of original 083 data (Baras & Dey, 1999). 084

## 2 PROBLEM FORMULATION

In this section, we specify the feature discrimination functions for the original (non-quantized) and quantized data. Prior to this, we introduce the binary and ternary quantization functions, as well as the data modeling.

## 2.1 QUANTIZATION FUNCTIONS

The binary and ternary quantization functions are formulated as

$$f_b(x;\tau) = \begin{cases} 1, & \text{if } x > \tau \\ 0, & \text{otherwise} \end{cases}$$
(1)

and

085

087 088

089

090 091

092 093

098 099 100

$$f_t(x;\tau) = \begin{cases} 1, & \text{if } x > \tau \\ 0, & \text{if } -\tau \le x \le \tau \\ -1, & \text{if } x < -\tau \end{cases}$$
(2)

where the threshold parameter  $\tau \in (-\infty, +\infty)$  for  $f_b(x; \tau)$ , and  $\tau \in [0, +\infty)$  for  $f_t(x; \tau)$ . The two functions operate element-wise on a vector  $\mathbf{x} = [x_1, x_2, \cdots, x_n]^\top \in \mathbb{R}^n$ , namely  $f_b(\mathbf{x}; \tau) = (f_b(x_1; \tau), f_b(x_2; \tau), \cdots, f_b(x_n; \tau))^\top$  and the same applies to  $f_t(\mathbf{x}; \tau)$ .

#### 105 2.2 DATA DISTRIBUTIONS

Throughout the work, we denote each data sample using a vector. For the sake of generality, as discussed before, we assume that the data vector randomly drawn from a class is a random vector

108  $\mathbf{X} = \{X_1, X_2, \cdots, X_n\}^{\mathsf{T}}$ , with its each element  $X_i$  following a Gaussian distribution  $N(\mu_{1,i}, \sigma^2)$ ; 109 and similarly, for the random vector  $\mathbf{Y} = \{Y_1, Y_2, \cdots, Y_n\}^\top$  drawn from another class, we suppose 110 its each element  $Y_i \sim N(\mu_{2,i}, \sigma^2)$ , where  $\mu_{2,i} \neq \mu_{1,i}$ . Considering that the discrimination between 111 the two random vectors  $\hat{\mathbf{X}}$  and  $\hat{\mathbf{Y}}$  positively correlates with the discrimination between their each 112 pair of corresponding elements  $X_i$  and  $Y_i$ , we propose to analyze the discrimination at the element 113 level, specifically between  $X_i$  and  $Y_i$ , rather than between the entire vectors, X and Y. For notational convenience, without causing confusion, in the sequel we will omit the subscript "i" of  $X_i$ 114 and  $Y_i$ , and write their distributions as  $X \sim N(\mu_1, \sigma^2)$  and  $Y \sim N(\mu_2, \sigma^2)$ , where  $\mu_1 \neq \mu_2$ . Note 115 that we assume here that the two variables X and Y share the same variance  $\sigma^2$ . This assumption is 116 common in statistical research, as the data we intend to investigate are often drawn from the same 117 or similar scenarios and thus exhibit similar noise levels. 118

119 When standardization, a common practice in machine learning, is applied to the two variables, X120 and Y, their distributions will exhibit specific relationships. More precisely, in a binary classification 121 problem, the dataset we handle is a mixture, denoted as Z, comprising two classes of samples 122 drawn respectively from X and Y. Usually, the mixture Z is assumed to possess a balanced class 123 distribution, meaning that samples are drawn from X and Y with equal probabilities. Under this 124 assumption, when we perform standardization by subtracting the mean and dividing by the standard 125 deviation for each sample in Z, the distributions of X and Y (in Z) will become

127 128

129 130

131 132 133

$$\tilde{X} = \frac{X - E[Z]}{\sqrt{D[Z]}} \sim N\left(\frac{(\mu_1 - \mu_2)/2}{\sqrt{\sigma^2 + \frac{1}{4}(\mu_1 - \mu_2)^2}}, \frac{\sigma^2}{\sigma^2 + \frac{1}{4}(\mu_1 - \mu_2)^2}\right)$$
(3)

and

$$\tilde{Y} = \frac{Y - E[Z]}{\sqrt{D[Z]}} \sim N\left(\frac{-(\mu_1 - \mu_2)/2}{\sqrt{\sigma^2 + \frac{1}{4}(\mu_1 - \mu_2)^2}}, \frac{\sigma^2}{\sigma^2 + \frac{1}{4}(\mu_1 - \mu_2)^2}\right)$$
(4)

where E[Z] and D[Z] denote the expectation and variance of Z, which have expressions  $E[Z] = \frac{1}{2}(\mu_1 + \mu_2)$  and  $D[Z] = \sigma^2 + \frac{1}{4}(\mu_1 - \mu_2)^2$ .

From Equations (3) and (4), it can be seen that after standardization, the two classes of variables  $\tilde{X}$ and  $\tilde{Y}$  still exhibit Gaussian distributions, but showcase two interesting properties: 1) their means are symmetric about zero; and 2) they have the sum of the square of the mean and the variance equal to one. By the two properties, the distributions of two classes of standardized data are characterized in Property 1. In the paper, we will typically focus our study on the standardized data.

142 **Property 1** (The distributions of two classes of standardized data). The two classes of standardized 143 data we aim to study have their samples i.i.d drawn from  $X \sim N(\mu, \sigma^2)$  and  $Y \sim N(-\mu, \sigma^2)$ , 144 where  $\mu^2 + \sigma^2 = 1, \mu \in (0, 1)$ .

146 2.3 FEATURE DISCRIMINATION

Following the Fisher's linear discriminant rule, we define the discrimination between two classes of
data as the ratio of the expected inter-class distance to the expected intra-class distance, as specified
below.

**Definition 1** (Discrimination between two classes of data). For two classes of data with samples respectively drawn from the variables X and Y, the discrimination between them is defined as

$$D = \frac{E[(X_1 - Y_1)^2]}{E[(X_1 - X_2)^2] + E[(Y_1 - Y_2)^2]}$$
(5)

where  $X_1$  and  $X_2$  are i.i.d. samples of X, and  $Y_1$  and  $Y_2$  are i.i.d samples of Y.

In the sequel, we will utilize the above definition D to denote the discrimination between original (non-quantized) data; and for the binary and ternary quantized data, as detailed below, we adopt  $D_b$ and  $D_t$  to represent their discrimination.

**Definition 2** (Discrimination between two classes of quantized data). Following the discrimination specified in Definition 1, the discrimination between two binary quantized data  $X_b = f_b(X; \tau)$  and

153 154 155

157

145

 $Y_b = f_b(Y; \tau)$ , is formulated as

163 164

166

167

168

169 170 171

172 173

174 175

176

177

178 179

180 181

182

183

185

186

191

197

199

162

 $D_b = \frac{E[(X_{1,b} - Y_{1,b})^2]}{E[(X_{1,b} - X_{2,b})^2] + E[(Y_{1,b} - Y_{2,b})^2]}$ (6)

where  $X_{1,b}$  and  $X_{2,b}$  are i.i.d. samples of  $X_b$ , and  $Y_{1,b}$  and  $Y_{2,b}$  are i.i.d. samples of  $Y_b$ . Similarly, the discrimination between two ternary quantized data  $X_t = f_t(X;\tau)$  and  $Y_t = f_t(Y;\tau)$  is expressed as

$$D_t = \frac{E[(X_{1,t} - Y_{1,t})^2]}{E[(X_{1,t} - X_{2,t})^2] + E[(Y_{1,t} - Y_{2,t})^2]}$$
(7)

where  $X_{1,t}$  and  $X_{2,t}$  are i.i.d. samples of  $X_t$ , and  $Y_{1,t}$  and  $Y_{2,t}$  are i.i.d. samples of  $Y_t$ .

2.4 GOAL

The major goal of the paper is to investigate whether there exist threshold values  $\tau$  in the binary quantization  $f_b(x;\tau)$  and the ternary quantization  $f_t(x;\tau)$ , such that the quantization can improve the feature discrimination of original data, namely having  $D_b > D$  and  $D_t > D$ .

#### 3 **DISCRIMINATION ANALYSIS**

#### THEORETICAL RESULTS 3.1

**Theorem 1** (Binary Quantization). Consider the discrimination D between two classes of data  $X \sim$  $N(\mu, \sigma^2)$  and  $Y \sim N(-\mu, \sigma^2)$  as specified in Property 1, as well as the discrimination  $D_b$  between their binary quantization  $X_b = f_b(X;\tau)$  and  $Y_b = f_b(Y;\tau)$ . We have  $D_b > D$ , if there exists a quantization threshold  $\tau \in (-\infty, +\infty)$  such that

$$\beta - \alpha + \frac{\mu^2 (1 - 2\beta) - \mu \sqrt{\mu^2 + 4\beta (1 - \beta)}}{1 + \mu^2} > 0, \tag{8}$$

where  $\alpha = \Phi\left(\frac{\tau-\mu}{\sigma}\right)$  and  $\beta = \Phi\left(\frac{\tau+\mu}{\sigma}\right)$ , with  $\Phi(\cdot)$  denoting the cumulative distribution function of the standard normal distribution. 192

193 **Theorem 2** (Ternary Quantization). Consider the discrimination D between two classes of data  $X \sim N(\mu, \sigma^2)$  and  $Y \sim N(-\mu, \sigma^2)$  as specified in Property 1, as well as the discrimination  $D_t$ 194 between their ternary quantization  $X_t = f_t(X;\tau)$  and  $Y_t = f_t(Y;\tau)$ . We have  $D_t > D$ , if there 195 exists a quantization threshold  $\tau \in [0, +\infty)$  such that 196

$$\beta - \alpha + \frac{\mu^2 - \sqrt{\mu^4 + 8\mu^2 \beta}}{2} > 0, \tag{9}$$

where  $\alpha = \Phi\left(\frac{-\tau-\mu}{\sigma}\right)$  and  $\beta = \Phi\left(\frac{-\tau+\mu}{\sigma}\right)$ , with  $\Phi(\cdot)$  denoting the cumulative distribution function of the standard normal distribution. 200 201 202

**Remarks:** Regarding the two theorems, there are three issues worth discussing. 1) The two theo-203 rems suggest that both binary and ternary quantization methods can indeed improve the classification 204 performance of original data, if there exist quantization thresholds  $\tau$  that can satisfy the constraints 205 shown in Equations (8) and (9). The following numerical analysis demonstrates that the desired 206 threshold  $\tau$  does exist, when the two classes of data  $X \sim N(\mu, \sigma^2)$  and  $Y \sim N(-\mu, \sigma^2)$  are as-207 signed appropriate values for  $\mu$  and  $\sigma$ . This threshold  $\tau$  can be approximately estimated using the 208 bisection method. 2) Our theoretical analysis is based on the premise that the data vectors belonging to the same class have Gaussian distributions at the vectors' each coordinate. This condition should 210 hold true when two classes of data are readily separable, as in this case the data points within each 211 class should cluster tightly, allowing for Gaussian approximation. This explains the recent research 212 findings, that the binary or ternary quantization tends to achieve comparable or superior classifi-213 cation performance, when handling relatively simple datasets (Courbariaux et al., 2015; Lin et al., 2016b), or distinguishable features (Lu et al., 2023). 3) The conclusion we derive in Theorem 1 for 214  $\{0,1\}$ -binary quantization also applies to another popular  $\{-1,1\}$ -binary quantization (Qin et al., 215 2020), since the Euclidean distance of the former is equivalent to the cosine distance of the latter.



Figure 1: Consider two classes of data  $X \sim N(\mu, \sigma^2)$  and  $Y \sim N(-\mu, \sigma^2)$ , with  $\mu = 0.8$  and  $\sigma^2 = 0.36$ , as specified in Property 1. The values for the left and right sides of Equations (8) and (9) are provided in (a) and (c) for binary and ternary quantization, respectively; and the discrimination D,  $D_b$  and  $D_t$  statistically estimated with Equations (5), (6) and (7) are illustrated in (b) and (d) for binary and ternary quantization, respectively.

#### 3.2 NUMERICAL ANALYSIS

In this part, we conduct numerical analyses for two primary objectives. Firstly, we aim to prove the existence of the desired quantization threshold  $\tau$  that holds Equations (8) and (9), namely making the left sides of the two inequalities larger than their right sides (with values equal to zero). For this purpose, we compute the values of the left sides of Equations (8) and (9), through assigning specific values to  $\tau$ , as well as to the two variables X and Y's distribution parameters  $\mu$  and  $\sigma^2$ . Note that we here set  $\sigma^2 = 1 - \mu^2$ ,  $\mu \in (0, 1)$ , in accordance with Property 1. In Figure 1, we examine the case that fixes  $\mu = 0.8$  and  $\sigma^2 = 0.36$ , while varying the value of  $\tau$  with a step width 0.01. The results for binary quantization and ternary quantization are provided in Figures 1 (a) and (c), respectively. It can be seen that for the two quantization methods, their conditions shown in Equations (8) and (9) will hold when respectively having  $\tau \in [-0.2, 0.2]$  and  $\tau \in [0, 0.5]$ . This proves the existence of the desired quantization threshold  $\tau$  that can improve feature discrimination. For limited space, we here only discuss the case of  $\mu = 0.8$  (and  $\sigma^2 = 1 - \mu^2$ ) in Figure 1. By examining different  $\mu \in (0,1)$  in the same way, we can find that the quantization threshold  $\tau$  that holds Equations (8) and (9), is present when  $\mu \in (0.76, 1)$  and  $\mu \in (0.66, 1)$ , respectively; see Figures 7 and 8 for more evidences. This result implies two consequences. On one hand, ternary quantization has more chances to improve feature discrimination compared to binary quantization, as the former has a broader range of  $\mu$ . On the other hand, the improved discrimination tends to be achieved when  $\mu$  is sufficiently large, coupled with a correspondingly small  $\sigma$ , or when the discrimination between two classes of data is sufficiently high. Empirically, as depicted in Figure 17, the two specific ranges of  $\mu$  values are attainable for the commonly-used features of real data. 

The second goal is to verify that the quantization thresholds  $\tau$  we estimate with Equations (8) and (9) in Theorems 1 and 2, can indeed improve feature discrimination. To this end, it needs to prove that the ranges of  $\tau$  derived by Equations (8) and (9), such as the ones depicted in Figures 1 (a) and (c), are consistent with the ranges we can statistically estimate by the discrimination definitions D,  $D_b$  and  $D_t$ , as specified in Definitions 1 and 2. To estimate the discrimination D,  $D_b$  and  $D_t$ , we randomly generate 10,000 samples from  $X \sim N(0.8, 0.36)$  and  $Y \sim N(-0.8, 0.36)$ , respectively, and then statistically estimate them with Equations (5), (6) and (7), across varying  $\tau$  (with a step width 0.01). The results are provided in Figures 1 (b) and (d), respectively for binary quantization and ternary quantization. It can be seen that we have  $\tau \in [-0.2, 0.2]$  for  $D_b > D$ , and have  $\tau \in [0, 0.5] \text{ for } D_t > D. \text{ The results coincide with the theoretical results shown in Figures 1 (a) and}$ (c), validating the correctness of Theorems 1 and 2.

#### 4 EXPERIMENTS

274 275

273

Through previous theoretical and numerical analyses, we have demonstrated that binary and ternary quantization can improve feature discrimination between two classes of data, when the data vectors within each class exhibit Gaussian distributions across each coordinate point of their feature vectors. Given that improved feature discrimination should yield better classification performance, this section aims to validate this improvement by assessing classification results.

281 Our experiments will mainly investigate binary classification using two fundamental linear classi-282 fiers: the k-nearest neighbors (KNN) algorithm (with k = 5) (Peterson, 2009), employing both Euclidean and cosine distances as similarity metrics, and the support vector machine (SVM) (Cortes 283 & Vapnik, 1995), equipped with a linear kernel. We choose linear classifiers for binary classifica-284 tion based on two considerations. Firstly, linear binary classification can directly reflect the feature 285 discrimination between two classes, unlike more complex nonlinear classifiers that often involve 286 feature selection operations. Secondly, linear binary classification is a foundational concept in ma-287 chine learning. The insights gained from this analysis can be extended to multiclass and nonlinear 288 classifier-based classifications, as evidenced in the subsequent experiments. 289

To assess the robustness and generalizability of our theoretical findings, we will conduct classification evaluations on both synthetic and real data. Synthetic data can conform perfectly to the distribution conditions outlined in our theoretical analysis, whereas real data usually cannot.

293

295

296

4.1 SYNTHETIC DATA

4.1.1 Setting

297 In the simulation, we suppose that two classes of data have their samples i.i.d. drawn from two 298 different random vectors  $\mathbf{X} = \{X_1, X_2, \cdots, X_n\}^{\top}$  and  $\mathbf{Y} = \{Y_1, Y_2, \cdots, Y_n\}^{\top}$ , for which we set  $X_i \sim N(\mu_i, \sigma_i^2)$  and  $Y_i \sim N(-\mu_i, \sigma_i^2)$ , with  $\mu_i \in (-1, 0) \cup (0, 1)$  and  $\sigma_i^2 = 1 - \mu_i^2$ , in 299 300 accordance with the data distributions specified in Property 1. Considering the fact that the features 301 of real-world data usually exhibit sparse structures (Weiss & Freeman, 2007; Kotz et al., 2012), 302 we further suppose that the means  $\mu_i$  decay exponentially in magnitude, i.e.  $|\mu_{i+1}|/|\mu_i| = e^{-\lambda}$ ,  $\lambda \geq 0$ , and set  $\mu_1 = 0.8$  in the following simulation. It can be seen that with the increasing of 303  $\lambda$ , the mean's magnitude  $|\mu_i|$  (with i > 1) will become smaller, indicating a smaller data element 304  $X_i$  (in magnitude) and a sparser data structure. However, the data element  $X_i$  with smaller  $\mu_i$ , is 305 not favorable for quantization to improve feature discrimination, as indicated by previous numerical 306 analyses. The impact of data sparsity on quantization can be investigated by increasing the value of 307 the parameter  $\lambda$ . 308

With the data model described above, we randomly generate two classes of data, each class contain-309 ing 1000 samples. The dataset is split into two parts for training and testing, in a ratio of 4:1. Then 310 we evaluate the KNN and SVM classification on them. The classification accuracy is determined by 311 averaging the accuracy results obtained from repeating the data generation and classification process 312 100 times. The results for KNN with Euclidean distance are provided in Figures 2 and 3, and the 313 results for KNN with cosine distance and SVM are given in the appendix, Figures 9–12. It can be 314 seen that the three classifiers exhibit similar performance trends. For conciseness, we will focus 315 more on the results of KNN with Euclidean distance in the following discussion. 316

317

318

4.1.2 RESULTS

**Comparison between the data with different sparsity.** In Figure 2, we investigate the classification performance for the data generated with different parameters  $\lambda \in \{0, 0.01, 0.1, 1\}$ , namely with different sparsity levels. Recall that the larger the  $\lambda$ , the smaller the  $|\mu_i|$ , or say the smaller the data element  $X_i$  (in magnitude). By previous analyses, the data element  $X_i$  with smaller  $|\mu_i|$  is not conducive to enhancing feature discrimination through quantization. Nevertheless, empirically, the negative effect does not appear to be significant. From Figure 2, it can be seen that when increasing



Figure 2: KNN (Euclidean distance) classification accuracy for the 10,000-dimensional binary, ternary, and original data that are generated with the varying parameter  $\lambda \in \{0, 0.01, 0.1, 1\}$ , which controls the data sparsity.



Figure 3: KNN (Euclidean distance) classification accuracy for the binary, ternary, and original data generated with the parameter  $\lambda = 1$ , and with varying dimensions  $n \in \{1, 100, 10000\}$ .

 $\lambda$  from 0.1 to 1, there have been quantization thresholds  $\tau$  that can yield better classification performance than original data. In addition, it noteworthy that as  $\lambda$  increases, the overall classification accuracy of original data will decrease. This decreasing trend also impacts the absolute performance of the quantized data, even though it may outperform original data.

364 Comparison between the data with different dimensions. The impact of data dimensions  $n \in \{1, 100, 10000\}$  on classification is investigated in Figure 3, where the data are generated 366 with the exponentially decaying parameter  $\lambda = 1$ . It can be seen that with the increasing of data 367 dimension, the range of the quantization thresholds  $\tau$  that outperform original data tends to expand, 368 but the performance advantage declines. As previously discussed, the decline should be attributed 369 to the data element  $X_i$  with small means  $|\mu_i|$ , whose quantity will rise with the data dimension n, 370 particularly when the decay parameter  $\lambda$  of  $|\mu_i|$  is large. To alleviate this adverse effect, it is rec-371 ommended to choose a relatively smaller  $\lambda$  for high-dimensional data, indicating a structure that is not overly sparse. Conversely, when the high-dimensional data is highly sparse, we should reduce 372 its dimension to improve the classification performance under quantization. 373

374

356

357 358 359

360

361

362

363

336

337

338

375 **Comparison between binary quantization and ternary quantization.** From Figures 2 and 3, it can be seen that ternary quantization surpasses binary quantization by offering broader ranges 376 of quantization thresholds  $\tau$  that can yield higher classification accuracy than original data. This 377 observation is consistent with our previous theoretical and numerical analyses.

378 **Comparison between KNN and SVM.** Combining the results in Figures 2, 3, and 9–12, we can 379 say that both KNN and SVM enable quantization to improve the classification accuracy of original 380 data, within specific ranges of quantization thresholds  $\tau$ . If closely examining these ranges, it can be 381 observed that KNN often performs better when using Euclidean distance than using cosine distance. 382 This can be attributed to the advantage of Euclidean distance over cosine distance in measuring the distance between 0 and  $\pm 1$ . Also, KNN often outperforms SVM, such as the case of  $\lambda = 0.1$  shown in Figures 2 and 10. This is because the support vector of SVM relies on a few data points located 384 on the boundary between two classes, which may deteriorate during quantization. In contrast, KNN 385 depends on the high-quality data points within each class, making it resilient to quantization noise. 386

- 387
- 388 389

390

391

392

394

**Comparison between classification accuracy, feature discrimination and quantization error.** Figure 16 illustrates that the classification accuracy of quantized data across varying quantization threshold  $\tau$  can be reasonably reflected by feature discrimination, rather than quantization errors.

- 393 4.2 REAL DATA
- 395 4.2.1 SETTING

396 The classification is conducted on three different types of datasets, including the image datasets 397 YaleB (Lee et al., 2005), CIFAR10 (Krizhevsky & Hinton, 2009) and ImageNet1000 (Deng et al., 398 2009), the speech dataset TIMIT (Fisher et al., 1986), and the text dataset Newsgroup (Lang, 1995). 399 The datasets are briefly introduced as follows. YaleB contains face images of 38 persons, with about 400 64 faces per person. CIFAR10 consists of 60,000 color images from 10 different classes, with each 401 class having 6,000 images. ImageNet1000 consists of 1000 object classes, with 1,281,167 training 402 images, 50,000 validation images, and 100,000 test images. For the above three image datasets, 403 we separately extract their features using Discrete Wavelet Transform (DWT), ResNet18 (He et al., 404 2016) and VGG16 (Simonyan & Zisserman, 2014). For ease of simulation, the resulting feature vectors are dimensionally reduced by integer multiples, leading to the sizes of 1200, 5018, and 405 5018 respectively. From TIMIT, as in (Mohamed et al., 2011; Hutchinson et al., 2012), we extract 406 39 classes of 429-dimensional phoneme features for classification, totally with 1,134,138 training 407 samples and 58,399 test samples. Newsgroup comprises 20 categories of texts, with 11,269 samples 408 for training, and 7,505 samples for testing. The feature dimension is reduced to 5000 by selecting 409 the top 5000 most frequent words in the bag of words, as done in (Larochelle et al., 2012). 410

For each dataset, we iterate through all possible class pairs to perform binary classification. The 411 samples for training and testing are selected according to the default settings of the datasets. For 412 YaleB without prior settings, we randomly assign half of the samples for training and the remaining 413 half for testing. In the simulation, we need to test the classification performance of quantized data 414 across varying quantization threshold  $\tau$ . The value of  $\tau$  should correlate with the element scale of 415 the feature vectors, in the pursuit of improving classification over orginal data. To address the scale 416 varying of  $\tau$  across different data, we here suppose that  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average 417 magnitude of the feature elements (coordinates) in all the feature vectors used for classification, and 418  $\gamma$  is a scaling parameter. By adjusting  $\gamma$  within a narrow range, as illustrated later, we can derive the 419 desired  $\tau$  for various types of data.

420 To verify the generalizability of our feature discrimination analysis between two classes, we not 421 only evaluate binary classification using KNN and SVM, but also conduct multiclass classification, 422 as well as nonlinear classification using multilayer perceptron (MLP) (Rumelhart et al., 1986) and 423 decision trees (Quinlan, 1986). Due to space limitations, in the main body, we present the classifi-424 cation results of YaleB, Newsgroup, and TIMIT using KNN with Euclidean distance and SVM, as 425 illustrated in Figures 4 to 6. The results for other datasets, such as CIFAR10 and ImageNet1000, 426 and other classifiers, including KNN with cosine distance, MLP and decision trees, are provided in 427 the appendix, but briefly discussed within the main text.

428 429

- 4.2.2 RESULTS
- **Binary classification using KNN and SVM.** From Figures 4-6,13, 14 and 21, it can be seen that that within specific ranges of quantization thresholds  $\tau$ , both binary and ternary quantization



Figure 4: Classification accuracy for the binary, ternary, and original data by KNN (Euclidean distance) and SVM on YaleB. The parameter  $\gamma$  corresponds to a threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average magnitude of the feature elements in all feature vectors.



Figure 5: Classification accuracy for the binary, ternary, and original data by KNN (Euclidean distance) and SVM on TIMIT. The parameter  $\gamma$  corresponds to a threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average magnitude of the feature elements in all feature vectors.

can achieve superior or at least equivalent classification performance compared to the original data across five different datasets, although as shown in Figure 17, each data class does not adequately conform to the Gaussian distribution assumption underlying our theoretical analysis. Similarly as in the classification of synthetic data, we observe the following results. Firstly, when using Euclidean distance, KNN consistently identifies quantization thresholds  $\tau$  that improve the classification of original data across all datasets. Secondly, compared to cosine distance, Euclidean distance tends to allow KNN to encompass a wider range of  $\tau$  values that are beneficial for improving classi-fication. Thirdly, with SVM, quantization occasionally achieves comparable performance to the original data, rather than surpassing it, as exemplified in Figure 13. Fourthly, ternary quantization often outperforms binary quantization by providing a broader range of threshold  $\tau$  values that fa-cilitate classification improvement. The rationale behind these results has been elaborated in our



Figure 6: Classification accuracy for the binary, ternary, and original data by KNN (Euclidean distance) and SVM on Newsgroup. The parameter  $\gamma$  corresponds to a threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$ denotes the average magnitude of the feature elements in all feature vectors.

previous classification analysis of synthetic data. The consistent performance observed in both real and synthetic data underscores the broad applicability of our theoretical findings.

Multiclass and nonlinear classification. While our feature discrimination analysis is focused on linear, binary classification, experiments demonstrate that our results can also be extended to mul-ticlass and nonlinear classifications. For example, in multiclass classification on ImageNet1000, quantization thresholds  $\tau$  that improve the classification of original data have been identified, as shown in Figure 22. The extension from binary to multiclass classification may be explained by the fact that feature elements sharing a common coordinate (or feature attribute) across different classes tend to exhibit a binary state: strong or weak, as illustrated in Figure 18, which indicates the intensity of the feature attribute within a feature vector. This suggests that multiclass classifi-cation at each feature coordinate can be viewed as a binary classification problem. Figures 19 and 20 demonstrate that the desired thresholds  $\tau$  can also be obtained in nonlinear classifications using MLP and decision trees. The extension from linear to nonlinear classification may be attributed to the fundamental linear operations often involved in nonlinear classifiers, which assess the linear discrimination between features or model parameters. 

#### 5 CONCLUSION

In the paper, we have proposed utilizing feature discrimination to analyze the impact of quantiza-tion on classification. Unlike traditional analyses, which are primarily based on quantization errors, our feature discrimination-based approach offers a more direct and reasoned assessment of clas-sification performance. Through our analysis, we demonstrate that common binary and ternary quantization methods can improve the feature discrimination of original data, particularly when data vectors within the same class follow Gaussian distributions at each coordinate. This improved discrimination is validated through binary classification experiments on both synthetic and real data. While our feature discrimination analysis primarily focuses on linear, binary classification issues, our experiments indicate that the findings can be extended to multiclass and nonlinear classification scenarios. This underscores the broad applicability of our theoretical results. Importantly, our study challenges the traditional belief that larger quantization errors generally lead to lower classification performance, laying a theoretical foundation for developing better quantization methods.

# 540 REFERENCES

553

554

565

566

567

570

587

588

- John S Baras and Subhrakanti Dey. Combined compression and classification with learning vector
   quantization. IEEE Transactions on Information Theory, 45(6):1911–1920, 1999.
- Moses S Charikar. Similarity estimation techniques from rounding algorithms. In <u>Proceedings of</u>
   the thiry-fourth annual ACM symposium on Theory of computing, pp. 380–388, 2002.
- 547 Corinna Cortes and Vladimir Vapnik. Support-vector networks. <u>Machine learning</u>, 20(3):273–297, 1995.
- Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. <u>Advances in neural information processing systems</u>, 28, 2015.
  - J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In IEEE Conference on Computer Vision and Pattern Recognition, 2009.
- Behzad Mohammadi Dogahe and Manohar N Murthi. Quantization for classification accuracy in high-rate quantizers. In Digital Signal Processing and Signal Processing Education Meeting, pp. 277–282. IEEE, 2011.
- Ronald A Fisher. The use of multiple measurements in taxonomic problems. <u>Annals of eugenics</u>, 7 (2):179–188, 1936.
- William M. Fisher, George R. Doddington, and Kathleen M. Goudie-Marshall. The darpa speech recognition research database: Specifications and status. In <u>Proceedings of DARPA Workshop on</u> Speech Recognition, pp. 93–99, 1986.
  - Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In <u>Low-Power Computer</u> Vision, pp. 291–326. Chapman and Hall/CRC, 2022.
- Robert M. Gray and David L. Neuhoff. Quantization. <u>IEEE transactions on information theory</u>, 44 (6):2325–2383, 1998.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778, 2016.
- Torsten Hoefler, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity in deep
   learning: Pruning and growth for efficient inference and training in neural networks. Journal of
   Machine Learning Research, 22(241):1–124, 2021.
- Brian Hutchinson, Li Deng, and Dong Yu. A deep architecture with bilinear modeling of hidden representations: Applications to phonetic recognition. In <u>IEEE international conference on acoustics</u>, speech and signal processing, pp. 4805–4808. IEEE, 2012.
- Soumya Jana and Pierre Moulin. Optimal design of transform coders and quantizers for image classification. In <u>International Conference on Image Processing</u>, volume 3, pp. 841–844. IEEE, 2000.
- Soumya Jana and Pierre Moulin. Optimal transform coding of gaussian mixtures for joint classification/reconstruction. In <u>Data Compression Conference</u>, pp. 313–322. IEEE, 2003.
  - Samuel Kotz, Tomasz Kozubowski, and Krzystof Podgorski. <u>The Laplace distribution and</u> generalizations: a revisit with applications to communications, economics, engineering, and finance. Springer Science & Business Media, 2012.
- A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. <u>Master's</u>
   thesis, Department of Computer Science, University of Toronto, 2009.
- 593 Edmund Y Lam and Joseph W Goodman. A mathematical analysis of the dct coefficient distributions for images. IEEE transactions on image processing, 9(10):1661–1666, 2000.

594 595 596	Ken Lang. Newsweeder: Learning to filter netnews. In <u>Machine learning proceedings 1995</u> , pp. 331–339. Elsevier, 1995.
597 598 599	Hugo Larochelle, Michael Mandel, Razvan Pascanu, and Yoshua Bengio. Learning algorithms for the classification restricted boltzmann machine. <u>The Journal of Machine Learning Research</u> , 13 (1):643–669, 2012.
600 601	K. Lee, J. Ho, and D. Kriegman. Acquiring linear subspaces for face recognition under variable lighting. <u>IEEE Transactions on Pattern Analysis and Machine Intelligence</u> , 27(5):684–698, 2005.
603 604 605	Darryl Lin, Sachin Talathi, and Sreekanth Annapureddy. Fixed point quantization of deep con- volutional networks. In <u>International conference on machine learning</u> , pp. 2849–2858. PMLR, 2016a.
606 607	Zhouhan Lin, Matthieu Courbariaux, Roland Memisevic, and Yoshua Bengio. Neural networks with few multiplications. In International Conference on Learning Representations, 2016b.
608 609 610	Weizhi Lu, Mingrui Chen, Kai Guo, and Weiyu Li. Quantization: Is it possible to improve classification? In <u>Data Compression Conference</u> , pp. 318–327. IEEE, 2023.
611 612 613 614	<ul> <li>Abdel-rahman Mohamed, Tara N Sainath, George Dahl, Bhuvana Ramabhadran, Geoffrey E Hinton, and Michael A Picheny. Deep belief networks using discriminative features for phone recognition. In <u>IEEE international conference on acoustics</u>, speech and signal processing, pp. 5060–5063. IEEE, 2011.</li> </ul>
615 616 617 618	Karen L Oehler and Robert M Gray. Combining image compression and classification using vector quantization. <u>IEEE transactions on pattern analysis and machine intelligence</u> , 17(5):461–473, 1995.
619	Leif E Peterson. K-nearest neighbor. Scholarpedia, 4(2):1883, 2009.
620 621 622	H Poor and J Thomas. Applications of ali-silvey distance measures in the design generalized quan- tizers for binary decision systems. <u>IEEE Transactions on Communications</u> , 25(9):893–900, 1977.
623 624	Haotong Qin, Ruihao Gong, Xianglong Liu, Xiao Bai, Jingkuan Song, and Nicu Sebe. Binary neural networks: A survey. <u>Pattern Recognition</u> , 105:107281, 2020.
625 626	J. Ross Quinlan. Induction of decision trees. Machine learning, 1:81–106, 1986.
627 628	David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back- propagating errors. <u>nature</u> , 323(6088):533–536, 1986.
629 630 631	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <u>arXiv preprint arXiv:1409.1556</u> , 2014.
632 633	Naveen Srinivasamurthy and Antonio Ortega. Reduced complexity quantization under classification constraints. In <u>Data Compression Conference</u> , pp. 402–411. IEEE, 2002.
634 635 636	Antonio Torralba and Aude Oliva. Statistics of natural image categories. <u>Network: computation in</u> <u>neural systems</u> , 14(3):391–412, 2003.
637 638 639	Martin J. Wainwright and Eero P. Simoncelli. Scale mixtures of gaussians and the statistics of natural images. In <u>Proceedings of the 12th International Conference on Neural Information Processing</u> Systems, 1999.
640 641 642	Yair Weiss and William T Freeman. What makes a good model of natural images? In <u>IEEE</u> <u>Conference on Computer Vision and Pattern Recognition</u> , pp. 1–8. IEEE, 2007.
643	
644	
645 646	
647	

#### DETAILED PROOF А

#### A.1 PROOF OF THEOREM 1

Let  $X_1$  and  $X_2$  be i.i.d. samples of X, and  $Y_1$  and  $Y_2$  be i.i.d. samples of Y. Denote  $X_{i,b}$  and  $Y_{i,b}$  as the binary quantization of  $X_i$  and  $Y_i$ , i.e.  $X_{i,b} = f_b(X_i; \tau)$  and  $Y_{i,b} = f_b(Y_i; \tau)$ , where i = 1, 2. By the distributions of X and Y specified in Property 1 and the binary quantization function  $f_b(x;\tau)$ defined in Equation(1), the probability mass functions of  $X_{i,b}$  and  $Y_{i,b}$  can be derived as 

$$P(X_{i,b} = k) = \begin{cases} 1 - \alpha, & k = 1 \\ \alpha, & k = 0 \end{cases}$$
(10)

and

$$P(Y_{i,b} = k) = \begin{cases} 1 - \beta, & k = 1\\ \beta, & k = 0 \end{cases}$$
(11)

where  $\alpha = \Phi\left(\frac{\tau-\mu}{\sigma}\right)$  and  $\beta = \Phi\left(\frac{\tau+\mu}{\sigma}\right)$ . By the probability functions, it is easy to deduce that

$$E [(X_1 - X_2)^2] = 2\sigma^2, \qquad E [(X_{1,b} - X_{2,b})^2] = 2\alpha - 2\alpha^2, \\E [(Y_1 - Y_2)^2] = 2\sigma^2, \qquad E [(Y_{1,b} - Y_{2,b})^2] = 2\beta - 2\beta^2, \\E [(X_1 - Y_2)^2] = 2\sigma^2 + 4\mu^2, \qquad E [(X_{1,b} - Y_{1,b})^2] = \alpha + \beta - 2\alpha\beta.$$

With these equations, the discrimination D of original data, as specified in Definition 1, can be further expressed as

$$D = \frac{E[(X_1 - Y_1)^2]}{E[(X_1 - X_2)^2] + E[(Y_1 - Y_2)^2]} = \frac{\sigma^2 + 2\mu^2}{2\sigma^2},$$
(12)

and similarly, the discrimination  $D_b$  of binary quantized data, as specified in Definition 2, can be written as 

$$D_b = \frac{E[(X_{1,b} - Y_{1,b})^2]}{E[(X_{1,b} - X_{2,b})^2] + E[(Y_{1,b} - Y_{2,b})^2]} = \frac{\alpha - 2\alpha\beta + \beta}{(2\alpha - 2\alpha^2) + (2\beta - 2\beta^2)}.$$
 (13)

Next, we are ready to prove that  $D_b > D$  under the condition (8). By Equations (12) and (13), it is easy to see that  $D_b > D$  is equivalent to

$$(\sigma^2 + 2\mu^2)\alpha^2 - 2(\sigma^2\beta + \mu^2)\alpha + (\sigma^2 + 2\mu^2)\beta^2 - 2\mu^2\beta > 0.$$
(14)

This inequality can be viewed as a quadratic inequality in  $\alpha$ , which has the discriminant:

$$\Delta = 4\mu^4 + 16(1 - \beta)\mu^2\beta > 0$$

By the above inequality, the inequality (14) holds when  $\alpha \in (-\infty, \alpha_1) \cup (\alpha_2, +\infty)$ , where

$$\alpha_1 = \beta + \frac{\mu^2(1-2\beta) - \mu\sqrt{\mu^2 + 4\beta(1-\beta)}}{1+\mu^2}$$

and

$$\alpha_2 = \beta + \frac{\mu^2 (1 - 2\beta) + \mu \sqrt{\mu^2 + 4\beta(1 - \beta)}}{1 + \mu^2}.$$
(15)

Given (15), we can further derive  $\alpha_2 > \beta$ , since  $\mu^2(1-2\beta) + \mu\sqrt{\mu^2 + 4\beta(1-\beta)} > 0$ . However, this result contradicts the conclusion that  $\alpha < \beta$  we can derive with the probability mass functions shown in (10) and (11), mainly by the increasing property of  $\Phi(\cdot)$ . So the solution to the inequality (14) should be  $\alpha \in (-\infty, \alpha_1)$ , implying  $\beta - \alpha + \frac{\mu^2(1-2\beta)-\mu\sqrt{\mu^2+4\beta(1-\beta)}}{1+\mu^2} > 0$ . 

#### A.2 PROOF OF THEOREM 2

Let  $X_1$  and  $X_2$  be i.i.d. samples of X, and  $Y_1$  and  $Y_2$  be i.i.d. samples of Y. Denote  $X_{i,t} = f_t(X_i; \tau)$ and  $Y_{i,t} = f_t(Y_i; \tau)$ , where i = 1, 2. By the distributions of X and Y specified in Property 1 and

the ternary quantization  $f_t(x; \tau)$  defined in Equation (2), the probability mass functions of  $X_{i,t}$  and  $Y_{i,t}$  can be derived as

$$P(X_{i,t} = k) = \begin{cases} \beta, & k = 1\\ 1 - \alpha - \beta, & k = 0\\ \alpha, & k = -1 \end{cases}$$
(16)

$$P(Y_{i,t} = k) = \begin{cases} \alpha, & k = 1\\ 1 - \alpha - \beta, & k = 0\\ \beta, & k = -1 \end{cases}$$
(17)

712 where  $\alpha = \Phi(\frac{-\tau - \mu}{\sigma})$  and  $\beta = \Phi(\frac{-\tau + \mu}{\sigma})$ .

Then, by Definition 2, the discrimination  $D_t$  of ternary quantization can be derived as

$$D_t = \frac{E[(X_{1,t} - Y_{1,t})^2]}{E[(X_{1,t} - X_{2,t})^2] + E[(Y_{1,t} - Y_{2,t})^2]} = \frac{(\alpha + \alpha^2 - 2a\beta + \beta + \beta^2)}{2(\alpha - \alpha^2 + 2\alpha\beta + \beta - \beta^2)}.$$
 (18)

718 By Equations (12) and (18), it can be seen that  $D_t > D$  is equivalent to

$$\frac{(\alpha+\beta)+(\alpha-\beta)^2}{2(\alpha+\beta)-2(\alpha-\beta)^2} > \frac{\sigma^2+2\mu^2}{2\sigma^2},$$

722 which can simplify to

$$\alpha^{2} - (2\beta + \mu^{2})\alpha + \beta^{2} - \mu^{2}\beta > 0.$$
<sup>(19)</sup>

Clearly, (19) can be regarded as a quadratic inequality in  $\alpha$ , with its discriminant:

$$\Delta = \mu^4 + 8\mu^2\beta > 0.$$

This inequality implies that the inequality (19) holds when  $\alpha \in (-\infty, \alpha_1) \cup (\alpha_2, +\infty)$ , where

$$\alpha_1 = \beta + \frac{\mu^2 - \sqrt{\mu^4 + 8\mu^2\beta}}{2}$$

731 and

$$\alpha_2 = \beta + \frac{\mu^2 + \sqrt{\mu^4 + 8\mu^2\beta}}{2}.$$
(20)

<sup>734</sup> In (20), the term  $\mu^2 + \sqrt{\mu^4 + 8\mu^2\beta} > 0$ , implying  $\alpha_2 > \beta$ . In contrast, we will derive  $\alpha < \beta$  by <sup>735</sup> the probability functions shown in Equations (16) and (17), particularly by the increasing property <sup>736</sup> of  $\Phi(\cdot)$ . By this contradiction, we can say that  $D_t > D$  holds only under the case of  $\alpha \in (-\infty, \alpha_1)$ , <sup>737</sup> namely

$$\beta-\alpha+\frac{\mu^2-\sqrt{\mu^4+8\mu^2\beta}}{2}>0.$$

### **B** OTHER RESULTS

# 758 B.1 NUMERICAL ANALYSIS759



(a) Theoretical results for the data with distri- (b) Numerical results for the data with distribution parameters  $\mu = 0.99$  and  $\sigma^2 = 0.02$  bution parameters  $\mu = 0.99$  and  $\sigma^2 = 0.02$ 



(c) Theoretical results for the data with distri- (d) Numerical results for the data with distribution parameters  $\mu = 0.76$  and  $\sigma^2 = 0.42$  bution parameters  $\mu = 0.76$  and  $\sigma^2 = 0.42$ 

Figure 7: Consider the binary quantization on two classes of data  $X \sim N(\mu, \sigma^2)$  and  $Y \sim N(-\mu, \sigma^2)$  as specified in Property 1. For two kinds of data with distribution parameters ( $\mu = 0.99$ ,  $\sigma^2 = 0.02$ ) and ( $\mu = 0.76$ ,  $\sigma^2 = 0.42$ ), the values for the left and right side of Equations (8) are provided in (a) and (c) respectively; and their discrimination D and  $D_b$  statistically estimated with Equations (5) and (6) are illustrated in (b) and (d), respectively.



(a) Theoretical results for the data with distri- (b) Numerical results for the data with distribution parameters  $\mu = 0.99$  and  $\sigma^2 = 0.02$  bution parameters  $\mu = 0.99$  and  $\sigma^2 = 0.02$ 



(c) Theoretical results for the data with distri- (d) Numerical results for the data with distribution parameters  $\mu = 0.66$  and  $\sigma^2 = 0.56$  bution parameters  $\mu = 0.66$  and  $\sigma^2 = 0.56$ 

Figure 8: Consider the ternary quantization on two classes of data  $X \sim N(\mu, \sigma^2)$  and  $Y \sim N(-\mu, \sigma^2)$  as specified in Property 1. For two kinds of data with distribution parameters ( $\mu = 0.99$ ,  $\sigma^2 = 0.02$ ) and ( $\mu = 0.66$ ,  $\sigma^2 = 0.56$ ), the values for the left and right side of Equations (9) are provided in (a) and (c) respectively; and their discrimination D and  $D_t$  statistically estimated with Equations (5) and (7) are illustrated in (b) and (d), respectively.







Figure 9: KNN (Cosine) classification accuracy for the 10,000-dimensional binary, ternary, and original data that are generated with the varying parameter  $\lambda \in \{0, 0.01, 0.1, 1\}$ , which controls the data sparsity.



Figure 10: SVM classification accuracy for the 10,000-dimensional binary, ternary, and original data that are generated with the varying parameter  $\lambda \in \{0, 0.01, 0.1, 1\}$ , which controls the data sparsity.



Figure 11: KNN (Cosine) classification accuracy for the binary, ternary, and original data generated with the sparsity parameter  $\lambda = 1$ , and with varying dimensions  $n \in \{100, 10000\}$ .



Figure 12: SVM classification accuracy for the binary, ternary, and original data generated with the sparsity parameter  $\lambda = 1$ , and with varying dimensions  $n \in \{1, 100, 10000\}$ .



Figure 13: Classification accuracy for the binary, ternary, and original data by KNN (Euclidean distance) and SVM on CIFAR10. The parameter  $\gamma$  corresponds to a threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average magnitude of the feature elements in all feature vectors.





Figure 14: Classification accuracy for the binary, ternary, and original data by KNN (Cosine distance) on four different datasets. The parameter  $\gamma$  corresponds to a quantization threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average magnitude of the feature elements in all feature vectors.



Figure 15: The histogram (blue bar) of the element values on one coordinate of the feature vectors within a single class of samples selected from four different datasets, accompanied with a Gaussian fitting curve (red line).

# <sup>1134</sup> C RESPONSE TO REVIEWS



C.1 CLASSIFICATION ACCURACY VS. FEATURE DISCRIMINATION VS. QUANTIZATION ERROR

Figure 16: KNN (Euclidean distance) classification accuracy for the binary, ternary, and original synthetic data which are generated with the parameter  $\lambda = 1$ , and with data dimension equal to 1. For comparison, the feature discrimination values and quantization errors across different thresholds  $\tau$  are provided for both binary and ternary data. Comments: It can be observed that the changing trend of classification values across  $\tau$  can be reasonably reflected by feature discrimination, rather than by quantization errors.



The data distribution parameter  $\mu$  estimated with real data

(c) TIMIT (d) Newsgroup Figure 17: The histogram of the data distribution parameter  $\mu$  (defined in Property 1) for each ele-ment (coordinate) of the feature vectors used in binary classification. Comments: It can be seen that with certain probabilities, the  $\mu$  value of each feature element will fall within the regions of (0.76, 1)and (0.66, 1), which supports achieving improved classification by binary and ternary quantization. Despite the fact the probabilities are not large, namely the amount of feature elements falling within (0.76, 1) or (0.66, 1) is relatively few, as widely proved in our experiments, we can still ob-tain the desired thresholds  $\tau$  that support improving classification on these real data. This robustness should be attributed to the fact that classification performance is mainly determined by a few im-portant feature elements of large magnitude, such as the ones with absolute means  $\mu$  falling within 

 (0.76, 1) or (0.66, 1).

C.2



## 1242 C.3 THE BINARY ATTRIBUTE OF FEATURE ELEMENTS ACROSS MULTIPLE CLASSES

Figure 18: Two histograms are drawn, one for the feature elements with values less than zero (dark 1264 red) and the other for those greater than zero (dark blue). The feature elements are collected from a 1265 common coordinate of feature vectors across all classes. Both histograms are fitted with Gaussian 1266 curves. Comments: It can be seen that both histograms approximately exhibit Gaussian distribu-1267 tions, with their two means separable. This indicates the binary nature (strong and weak) of the fea-1268 ture elements at each coordinate of feature vectors, regardless of the number of classes from which 1269 the feature vectors are drawn. This property allows us to generalize our binary classification-based feature discrimination analysis to multiclass classification scenarios. The reason is as follows. Con-1270 sider a feature vector  $\mathbf{x} = [x_1, x_2, ..., x_n]^{\top}$  for a given sample, where each element  $x_i$  corresponds 1271 to a specific feature attribute, such as frequencies in DCT features, scale and spatial positions in 1272 DWT features, or filters in convolutional features. The value of  $x_i$  indicates the strength of the *i*-1273 th attribute present within the sample. The strength of  $x_i$  can characterized by two distinct states: 1274 strong and weak, which reflect the presence or absence of the i-th attribute in the sample. The two 1275 states are evidenced in our statistical analysis of the  $x_i$  values in real-data feature vectors x. The 1276 results, depicted in this figure, show that the large (>0) and small values (<0) both exhibit Gaussian 1277 distributions, with the means of theses distributions representing the strong and weak states, respec-1278 tively. Given this understanding, the classification of each attribute (coordinate)  $x_i$  in feature vectors 1279 x can be considered a binary classification problem, regardless of the number of classes from which 1280 the feature vectors  $\mathbf{x}$  are drawn. Consequently, we can conclude that the capability of quantization to improve binary classification can also be extended to multiclass classification, provided that the 1281 Gaussian distributions of the two attributes at each coordinate of feature vectors are sufficiently sep-1282 arable, as required in Theorems 1 and 2. 1283

1284

- 1285 1286
- 1287
- 1288
- 1289
- 1290
- 1291
- 1292
- 1293
- 1294



three different datasets. The parameter  $\gamma$  corresponds to a quantization threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average magnitude of the feature elements in all feature vectors. Comments: Despite the fact that our linear feature discrimination analysis on quantized data may not directly extend to nonlinear classification scenarios, experiments using classifiers MLP and decision trees demonstrate that binary and ternary quantization can achieve improved or at least comparable classification re-sults even with nonlinear classifiers. This should be attributed to the fact that nonlinear classifiers generally involve fundamental linear operations, that evaluate the linear discrimination among fea-tures or model parameters.

#### NONLINEAR CLASSIFIERS: MLP AND DECISION TREES C.4



Figure 20: Decision trees-based binary classification accuracy for the binary, ternary, and original data on three different datasets. The parameter  $\gamma$  corresponds to a quantization threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average magnitude of the feature elements in all feature vectors. Comments: Despite the fact that our linear feature discrimination analysis on quantized data may not directly extend to nonlinear classification scenarios, experiments using classifiers MLP and decision trees demonstrate that binary and ternary quantization can achieve improved or at least comparable clas-sification results even with nonlinear classifiers. This should be attributed to the fact that nonlinear classifiers generally involve fundamental linear operations, that evaluate the linear discrimination among features or model parameters. 



Figure 21: Binary classification accuracy for the binary, ternary, and original data in ImageNet1000, using the classifier KNN (Euclidean distance). The parameter  $\gamma$  corresponds to a quantization threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average magnitude of the feature elements in all fea-ture vectors. Comments: It is evident that there are quantization thresholds  $\tau$  that can improve the binary classification accuracy of ImageNet1000. Given the complexity of ImageNet1000, this vali-dates the generalizability of our findings. 



Figure 22: Multiclass classification accuracy for the binary, ternary, and original data in Ima-geNet1000, using the classifier KNN (Euclidean distance). The parameter  $\gamma$  corresponds to a quan-tization threshold  $\tau = \gamma \cdot \eta$ , where  $\eta$  denotes the average magnitude of the feature elements in all feature vectors. Comments: It can be seen that there are quantization thresholds  $\tau$  that can improve the multiclass classification accuracy of ImageNet1000. This validates that our feature discrimina-tion analysis, rooted in binary classification, can be extended to multiclass classification, owing to the binary state of the feature elements sharing a common coordinate across different classes. See Figure 18 for detailed explanations.