
Category coding with neural network application

Qizhi Zhang

qizhi.zqz@alibaba-inc.com

Kuang-chih Lee

kuang-chih.lee@alibaba-inc.com

Hongying Bao

hongying.bhy@alibaba-inc.com

Yuan You

youyuan.yy@alibaba-inc.com

Dongbai Guo

dongbai.gdb@alibaba-inc.com

Abstract

In many applications of neural network, it is common to introduce huge amounts of input categorical features, as well as output labels. However, since the required network size should have rapid growth with respect to the dimensions of input and output space, there exists huge cost in both computation and memory resources. In this paper, we present a novel method called category coding (CC), where the design philosophy follows the principle of minimal collision to reduce the input and output dimension effectively. In addition, we introduce three types of category coding based on different Euclidean domains. Experimental results show that all three proposed methods outperform the existing state-of-the-art coding methods, such as standard cut-off and error-correct output coding (ECOC) methods.

1 Introduction

In machine learning, many features are categorical, such as color, country, user id, item id, etc. In the multi-class classification problem, the labels are categorical too. The ordering relation doesn't exist among different values for these categories. Usually those categorical variables are represented by one-hot feature vectors. For example, red is encoded to 100, yellow to 010 and blue to 001. But if the number of categories are very huge, for example the user id and item id in e-commerce applications, the one-hot encoding scheme needs too many resources to compute classification results.

In the past years while SVM is widely used, ECOC (error-correct output coding) method is proposed for handling huge numbers of output class labels. The idea of ECOC is to reduce a multi-class classification problem of huge number of classes to some two-class classification problems using binary error-correct coding. But for the solution of handling huge number of input categorical features, the similar method doesn't exist, because the categories can not be separated by linear model, unless the one-hot encoding is used.

In recent year, the deep neural network has great improvement in terms of performance and speed. The coding method can be applied to deep neural network with some new beneficial reform.

In the classification problem, because the number of labels of a single neural network need not to be binary, if we use a deep learning network as a base learner, it is not necessary to limit the code to be binary. In fact, there is a trade-off between the class number of one base learner and the number of base learner used. According to information theory, if we use p classes classifiers as basic classifiers to solve a classification problem of N -class, we need at least $\lceil \log_p N \rceil$'s base learners. For example, if we need to solve a classifying problem of $1M$'s classes, and we use the binary classifier as base learners, we need at least 20 base learners. For some classical applications, for example, the CNN image classification, we need to build a CNN network for every binary classifier. It is huge cost for

computation and memory resources. But if we combine different base learners with 1000 classes, we need at least 2 base learners. We know that the number of parameters in a Deep neural network is usually big, hence using a small number of base learner benefits the reduction of the cost in computing and storage.

On the other hand, because the neural network has the ability of non-linear representation, we can use the encoding for categorical features too. Can we use classical error-correct coding for categorical features? We know that in machine learning, the sparsity is a basic rule to be satisfied, but the classical error-correct coding does not satisfy the sparsity. Hence we need to design a new sparse coding scheme for this application.

In this paper, we give some new encoding method, they can be applied to both label encoding and feature encoding and give better performance than classical method. In section 2, we give the definition of category coding (CC) and propose 3 classes of CC, namely Polynomial CC, Remainder CC and Gauss CC, which have good property. In section 3 we discuss the application of CC in label encoding. In section 4, we discuss the application of CC in feature encoding. Our main tool is finite field theory and number theory, which can refer to [1] and [5].

2 Category coding

For a N -class categorical feature or label, we define a category coding (CC) as a map

$$\begin{aligned} f : \mathbb{Z}/N\mathbb{Z} &\longrightarrow \prod_{i=1}^r \mathbb{Z}/N_i\mathbb{Z} \\ x &\mapsto (f_i(x))_i \end{aligned}$$

where each $f_i : \mathbb{Z}/N\mathbb{Z} \longrightarrow \mathbb{Z}/N_i\mathbb{Z}$ is called a ‘‘site-position function’’. category coding, for $i = 1, 2, \dots, r$.

Generally, N is a huge number, and N_i are some numbers of middle size.

We can reduce a N -classes classification problem to r ’s classification problems of middle size through a CC.

We can also use a r -hot ($\sum_{i=1}^r N_i$)-bit binary encoding instead of the one-hot encoding as the representation of the feature, i.e., use the composite of the CC map f and the nature embedding

$$\begin{aligned} \prod_{i=1}^r \mathbb{Z}/N_i\mathbb{Z} &\longrightarrow \prod_{i=1}^r \mathbb{F}_2^{N_i} = \mathbb{F}_2^{\sum_i N_i} \\ (x_i)_i &\mapsto (N_i \text{ bit one hot representation of } x_i)_i \end{aligned}$$

to get a r -hot encoding.

For a CC f , we call $\max_{x \neq y} \#\{i = 1, \dots, r \mid f_i(x) = f_i(y)\}$ the collision number of f , and denote $C(f)$. We have the following theorem.

Theorem 2.1. For a CC $f : \mathbb{Z}/N\mathbb{Z} \longrightarrow \prod_{i=1}^r \mathbb{Z}/N_i\mathbb{Z}$, where $N_1 \leq N_2 \leq \dots \leq N_r$, we have $C(f) \geq \min\{i = 1, \dots, r \mid N \leq \prod_{j=1}^i N_j\} - 1$.

Proof. Let $k := \min\{i = 1, \dots, r \mid N \leq \prod_{j=1}^i N_j\}$. Suppose $C(f) < k - 1$, i.e

$$\max_{x \neq y} \#\{i = 1, \dots, r \mid f_i(x) = f_i(y)\} < k - 1$$

Hence for any $x \neq y \in \mathbb{Z}/N\mathbb{Z}$, there are at most $k - 2$ same site-position value between $f(x)$ and $f(y)$. Hence $\mathbb{Z}/N\mathbb{Z} \longrightarrow \prod_{i=1}^{k-1} \mathbb{Z}/N_i\mathbb{Z}$ is an injection, and hence $N \leq \prod_{i=1}^{k-1} N_i$. It is a contradiction with the definition of k . \square

If a CC satisfying $C(f) = \min\{i = 1, \dots, r \mid N \leq \prod_{j=1}^i N_j\} - 1$, we call it has the **minimal collision** property. In both usage of label encoding and feature encoding, we wish the code has minimal collision property.

We give 3 classes of CC, i.e, Polynomial CC, Remainder CC and Gauss CC, which satisfies the minimal collision property.

2.1 Polynomial CC

For any prime number p , we can represent any non-negative integral number x less than p^k as the unique form $x = x_0 + x_1p + \dots + x_{k-1}p^{k-1}$ ($x_i \in \mathbb{Z}/p\mathbb{Z}$), which gives a bijection $\mathbb{Z}/p^k\mathbb{Z} \longrightarrow \mathbb{F}_p^k$, where \mathbb{F}_p is the Galois field (finite field) of p elements.

For the classification problem of N -classes and any small positive integral number k (for example, $k=2, 3$) and a small real number $\epsilon \in (0, 1)$, we take a prime number in $[N^{\frac{1}{k}}, N^{\frac{1}{k-\epsilon}}]$ (According to the Prime Number Theorem ([8], [6]), there are about $\frac{k(N^{\frac{1}{k-\epsilon}} - N^{\frac{1}{k}})}{\log N}$ such prime numbers.) , and get a injection $\mathbb{Z}/N\mathbb{Z} \rightarrow \mathbb{Z}/p^k\mathbb{Z} \rightarrow \mathbb{F}_p^k$ by p-adic representation.

Theorem 2.2. For r 's different elements x_1, x_2, \dots, x_r in \mathbb{F}_p , the code defined by the composite map f of the p-adic representation map and the map

$$\begin{aligned} \phi_1 : \mathbb{F}_p^k &\longrightarrow \mathbb{F}_p[x]_{deg < k} \\ (a_0, \dots, a_{k-1}) &\mapsto a_0 + a_1x + \dots + a_{k-1}x^{k-1} \end{aligned}$$

and the map

$$\begin{aligned} \phi_2 : \mathbb{F}_p[x]_{deg < k} &\longrightarrow \mathbb{F}_p^r \\ g(x) &\mapsto (g(x_1), \dots, g(x_r)) \end{aligned}$$

has the minimal collision property. ■

Proof. We need proof that $C(\phi) \leq \min\{i = 1, \dots, r | N \leq p^i\} - 1$. Because we know that $\min\{i = 1, \dots, r | N \leq p^i\} = k$, hence we need just prove $C(\phi) \leq k$, i.e for any $\alpha \neq \beta \in \mathbb{Z}/N\mathbb{Z}$, $\#\{i = 1, \dots, r | f_i(\alpha) = f_i(\beta)\} \leq k - 1$.

Because the p-adic representation map and is an injection, and the map ϕ_1 is a bijection, we need just to show that for any $g_1 \neq g_2 \in \mathbb{F}_p[x]_{deg < k}$, $\#\{i = 1, \dots, r | g_1(x_i) = g_2(x_i)\} \leq k - 1$. Suppose there are $g_1 \neq g_2 \in \mathbb{F}_p[x]_{deg < k}$ such that $\#\{i = 1, \dots, r | g_1(x_i) = g_2(x_i)\} > k - 1$, it means the polynomial $g_1 - g_2 \in \mathbb{F}_p[x]$ of degree at most $k - 1$ has at least k roots, it is a contradiction with the Algebraic Basic Theorem on fields. □

Remark. The composite map of ϕ_1 and ϕ_2 in above theorem is known as Reed-Solomon code also [7]. The Reed-Solomon code is a class of non-binary MDS (maximal distinct separate) code [11]. MDS property is a excellent property in error-corrected coding. But unfortunately, it has not find any nontrivial binary MDS code yet up to now. In fact, for some situation, the fact that there are not any nontrivial binary MDS code is proved. ([3] and Proposition 9.2 on p. 212 in [14]). This is an advantage of CC than ECOC in label encode also.

2.2 Remainder CC

For the original label's set $\mathbb{Z}/N\mathbb{Z}$, a small number k like 2, or 3, etc., and a small positive number $\epsilon \in (0, 1)$, select r 's pairwise co-prime numbers p_1, p_2, \dots, p_r in the domain $[N^{\frac{1}{k}}, N^{\frac{1}{k-\epsilon}}]$. (According to the Prime Number Theorem ([8], [6]), there are about $\frac{k(N^{\frac{1}{k-\epsilon}} - N^{\frac{1}{k}})}{\log N}$ prime and hence pairwise co-prime numbers in this domain.)

We define the remainder CC as

$$\begin{aligned} \mathbb{Z}/N\mathbb{Z} &\longrightarrow \prod_{i=1}^r \mathbb{Z}/p_i\mathbb{Z} \\ x &\mapsto f_i(x) \end{aligned}$$

where $f_i(x) = x \bmod p_i$, and $\{p_i\}$ is called its modules. Then we have the following proposition:

Theorem 2.3. The remainder CC has the minimal collision property.

Proof. We need only to show that, for any $x \neq y \in \mathbb{Z}/N\mathbb{Z}$, there are at most $k - 1$'s i such, that $f_i(x) = f_i(y)$.

Suppose there exist k 's different i such, that $f_i(x) = f_i(y)$, we can suppose that $f_i(x) = f_i(y)$ for $i = 1, 2, \dots, k$. Then we have $x \equiv y \bmod p_i$ for all $i = 1, 2, \dots, k$. Because $\{p_i\}$ are pairwise co-prime numbers, we have $x \equiv y \bmod \prod_{i=1}^k p_i$. But we know $x, y \in \{0, 1, \dots, N - 1\}$, which in $\{0, 1, \dots, \prod_{i=1}^k p_i - 1\}$, hence $x = y$. □

2.3 Gauss CC

We propose a CC based on the ring of Gauss integers [2] [5], and so called Gauss CC.

We write the ring of Gauss integers as $\mathbb{Z}[\sqrt{-1}] := \{a + b\sqrt{-1} \in \mathbb{C} | a, b \in \mathbb{Z}\}$. For a big integral number N , let t is the minimal positive real number such that the number of Gauss integers in the

closed disc $\overline{U_t(0)}$ is not less than N , i.e. $\#\overline{U_t(0)} \cap \mathbb{Z}[\sqrt{-1}] \geq N$ and $\#\overline{U_{t-\epsilon}(0)} \cap \mathbb{Z}[\sqrt{-1}] < N$ for any small $\epsilon > 0$. In general, we have $\#\overline{U_t(0)} \cap \mathbb{Z}[\sqrt{-1}]$ is about πt^2 , hence we can get such t about $\sqrt{N/\pi}$.

We can embed the original IDs to the Gauss integers in Gauss integers in the closed disc.

$$\mathbb{Z}/N\mathbb{Z} \hookrightarrow \overline{U_t(0)} \cap \mathbb{Z}[\sqrt{-1}]$$

Let k be a small positive integral number, like 2,3, and ϵ' be a small positive real number. Let p_1, p_2, \dots, p_r be r pairwise co-prime Gauss integral numbers satisfying $|p_i| \in [(2t)^{\frac{1}{k}}, (2t)^{\frac{1}{k-\epsilon'}})$ for $i = 1, 2, \dots, r$. We define the category mapping

$$\begin{array}{ccc} \overline{U_t(0)} \cap \mathbb{Z}[\sqrt{-1}] & \longrightarrow & \prod_{i=1}^r \mathbb{Z}[\sqrt{-1}]/(p_i) \\ z & \mapsto & (f_i(z))_i \end{array}$$

where (p_i) means the principle ideal of $\mathbb{Z}[\sqrt{-1}]$ generated by p_i , $f_i(z) = z \pmod{(p_i)}$. $\{p_i\}$ is called the modules of this Gauss CC, and we have the following theorem.

Theorem 2.4. *The Gauss CC has the minimal collision property.*

Proof. From the method to take $\{p_i\}$, we know $k = \min\{i = 1, \dots, r | N \leq \prod_{j=1}^i |\mathbb{Z}[\sqrt{-1}]/(p_j)|\}$. Hence we need only to show that, for any $x \neq y \in \overline{U_t(0)} \cap \mathbb{Z}[\sqrt{-1}]$, there are at most $k-1$'s i such, that $f_i(x) = f_i(y)$.

Suppose there exist k 's different i such, that $f_i(x) = f_i(y)$, we can suppose that

$$f_i(x) = f_i(y) \quad \text{for } i = 1, 2, \dots, k$$

Then we have $x - y \equiv 0 \pmod{(p_i)}$ for all $i = 1, 2, \dots, k$.

Because $\{p_i\}$ are pairwise co-prime Gauss integral numbers, hence $\{(p_i)\}$ are pairwise co-prime ideal of $\mathbb{Z}[\sqrt{-1}]$, and we have $x - y \in \prod_{i=1}^k (p_i)$. Hence $\mathbf{Nm}(x - y) \in \prod_{i=1}^k (\mathbf{Nm}(p_i))\mathbb{Z}$ i.e. $|x - y|^2 \in \prod_{i=1}^k |p_i|^2\mathbb{Z}$, and hence $|x - y| \equiv 0 \pmod{\prod_{i=1}^k |p_i|}$. But we know $x, y \in \overline{U_t(0)}$, hence $|x - y| \leq 2t$. On the other hand, we know $\prod_{i=1}^k |p_i| > 2t$, hence $|x - y| = 0$, and hence $x = y$. \square

3 Application for label encode

For a N -class classification problem, we use a CC

$$\begin{array}{ccc} f : \mathbb{Z}/N\mathbb{Z} & \longrightarrow & \prod_{i=1}^r \mathbb{Z}/N_i\mathbb{Z} \\ z & \mapsto & (f_i(z))_i \end{array}$$

to reduce a N -classes classification problem to r 's classification problems of middle size through a LM. Suppose the training dataset is $\{x_k, y_k\}$, where x_k is feature and y_k is label, then we train a base learner on the dataset $\{x_k, f_i(y_k)\}$ for every $i = 1, 2, \dots, r$. We call it the label encoding method.

A CC good for label encoding should satisfy the follow properties:

Classes high separable. For two different labels y, \tilde{y} , there should be as many as possible site-position functions f_i such that $f_i(y) \neq f_i(\tilde{y})$.

Base learners independence. When y are selected randomly uniformly from $\mathbb{Z}/N\mathbb{Z}$, the mutual information of $f_i(y)$ and $f_j(y)$ approximate to 0 for $i \neq j$.

The property ‘‘classes high separable’’ ensures that for any two different classes, there are as many as possible base learners are trained to separate them. The property ‘‘base learners independence’’ ensures that the common part of the information learned by any two different base learners is few.

Remark. These properties are the similar of the properties ‘‘Row separable’’ and ‘‘Column separable’’ of ECOC ([32]) in non-binary situation.

The minimal collision property ensure the CCs satisfy ‘‘Class high separable’’, we will show that they satisfy ‘‘Base learner independence’’ also.

3.1 Polynomial CC

We will prove that, the Polynomial CC satisfies the property ‘‘Base learners independence’’ also.

Theorem 3.1. *If u is a random variable with uniform distribution on $\mathbb{Z}/N\mathbb{Z}$, y_i and y_j are the i -site value and j -site value ($i \neq j$) of the codeword of u under the simplex LM described above, then the mutual information of y_i and y_j approach to 0 when N grows up.*

Proof.

For any u in $\mathbb{Z}/p^k\mathbb{Z}$, the i -th site value is $y_i = u_0 + u_1x_i + \cdots + u_{k-1}x_i^{k-1} \pmod p$, where u_0, u_1, \dots, u_{k-1} are the coefficients of the p -adic representation of u . We denote this map by $g_i : \mathbb{Z}/p^k\mathbb{Z} \rightarrow \mathbb{Z}/p\mathbb{Z}$.

Let $t = \lceil N/p \rceil$, consider the following commutative diagram:

$$\begin{array}{ccc} \mathbb{Z}/pt\mathbb{Z} & \longrightarrow & \mathbb{Z}/pt\mathbb{Z} \\ \downarrow g_i & & \downarrow g_i \\ \mathbb{Z}/p\mathbb{Z} & \longrightarrow & \mathbb{Z}/p\mathbb{Z} \end{array}$$

The horizontal arrow in up line is defined by $u_0 + u_1p + \cdots + u_{k-1}p^{k-1} \mapsto (u_0 + 1 \pmod p) + u_1p + \cdots + u_{k-1}p^{k-1}$, and the horizontal arrow in down line is defined by $y \mapsto (y + 1 \pmod p)$. The horizontal arrows are bijections, which shows that the numbers of the pre-images in $\mathbb{Z}/pt\mathbb{Z}$ of every element in $\mathbb{Z}/p\mathbb{Z}$ are same and hence equal to t .

On the other hand, we have the commutative diagram:

$$\begin{array}{ccccc} \mathbb{Z}/p(t-1)\mathbb{Z} & \longrightarrow & \mathbb{Z}/N\mathbb{Z} & \longrightarrow & \mathbb{Z}/pt\mathbb{Z} \\ & \searrow & \downarrow & \swarrow & \\ & & \mathbb{Z}/p\mathbb{Z} & & \end{array}$$

where the horizontal arrows are the natural embedding, and other arrows are the restriction of g_i .

But the number of pre-images in $\mathbb{Z}/pt\mathbb{Z}$ of every element in $\mathbb{Z}/p\mathbb{Z}$ is t , and the same logic shows that the number of pre-images in $\mathbb{Z}/p(t-1)\mathbb{Z}$ of every element in $\mathbb{Z}/p\mathbb{Z}$ is $t-1$. Therefore the number of pre-images in $\mathbb{Z}/N\mathbb{Z}$ of every element in $\mathbb{Z}/p\mathbb{Z}$ is t or $t-1$.

Hence if u is a random variable with uniformly distribution on $\mathbb{Z}/N\mathbb{Z}$, its probability at every point in $\mathbb{Z}/N\mathbb{Z}$ is $1/N$, then the probability of y_i at every point in $\mathbb{Z}/p\mathbb{Z}$ are $\frac{t}{N}$ or $\frac{t-1}{N}$. The same logic shows that the probability of y_j at every point in $\mathbb{Z}/p\mathbb{Z}$ are $\frac{t}{N}$ or $\frac{t-1}{N}$.

Let $s = \lceil N/p^2 \rceil$, we have the commutative diagram for any $(a, b) \in \mathbb{F}_p^2$:

$$\begin{array}{ccc} \mathbb{Z}/p^2s\mathbb{Z} & \longrightarrow & \mathbb{Z}/p^2s\mathbb{Z} \\ \downarrow (g_i, g_j) & & \downarrow (g_i, g_j) \\ \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z} & \longrightarrow & \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z} \end{array}$$

where the up horizontal arrow is defined by $u_0 + u_1p + \cdots + u_{k-1}p^{k-1} \mapsto (u_0 + a \pmod p) + (u_1 + b \pmod p)p + \cdots + u_{k-1}p^{k-1}$, and the down horizontal arrow is defined by $(y_i, y_j) \mapsto (y_i + a + bx_i \pmod p, y_j + a + bx_j \pmod p)$. Both the horizontal arrows are bijections.

Because $x_i \neq x_j$ we know that when (a, b) runs over all the pairs in $\mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}$ the down horizontal map maps $(0, 0)$ to all the pairs in $\mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}$. Therefore all the number of pre-images in $\mathbb{Z}/p^2s\mathbb{Z}$ of any element in $\mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}$ are same, and hence equal to s .

A similar method shows that if u is a random variable with uniformly distribution on $\mathbb{Z}/N\mathbb{Z}$, the joint probability of (y_i, y_j) at every point in $\mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}$ are $\frac{s}{N}$ or $\frac{s-1}{N}$.

We know that the mutual information of y_i and y_j is $I(Y_i; Y_j) = \sum_{(y_i, y_j) \in \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}} P_{i,j}(y_i, y_j) \log \frac{P_{i,j}(y_i, y_j)}{P_i(y_i)P_j(y_j)}$.

a.) When $k = 2$, i.e. $p < N \leq p^2$, we know $s = 1$ and $p_{i,j}(y_i, y_j) = \frac{1}{N}$ on N 's point in $\mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}$ and 0 on other points. Hence we have

$$\begin{aligned} I(Y_i; Y_j) &\leq \sum_{(y_i, y_j) \in \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}} p_{i,j}(y_i, y_j) \log \frac{p_{i,j}(y_i, y_j)}{\left(\frac{t-1}{N}\right)^2} = N \times \frac{1}{N} \log \frac{1/N}{\left(\frac{t-1}{N}\right)^2} = 2 \log \frac{N}{t-1} - \log N \\ &\leq 2 \log \frac{N}{N/p-1} - \log N = 2 \log p - 2 \log(1 - \frac{p}{N}) - \log N = 2 \log p + 2O(\frac{p}{N}) - \log N \end{aligned}$$

However, $p \in [N^{\frac{1}{2}}, N^{\frac{1}{2}-\epsilon}]$ implies that $p = N^{\frac{1}{2}}(1 + o(1))$, hence we have

$$I(Y_i; Y_j) = \log N + 2 \log(1 + o(1)) + 2O(N^{-\frac{1}{2}}) - \log N = o(1) \rightarrow 0 \text{ as } N \rightarrow \infty$$

b.) When $k > 2$, i.e. $N > p^2$, we have

$$\begin{aligned} I(Y_i; Y_j) &= \sum_{(y_i, y_j) \in \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}} p_{i,j}(y_i, y_j) \log p_{i,j}(y_i, y_j) - \sum_{(y_i, y_j) \in \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}} p_{i,j}(y_i, y_j) (\log p_i(y_i) + \log p_j(y_j)) \\ &= \sum_{(y_i, y_j) \in \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}} p_{i,j}(y_i, y_j) \log p_{i,j}(y_i, y_j) - \sum_{y_i \in \mathbb{Z}/p\mathbb{Z}} p_i(y_i) \log p_i(y_i) - \sum_{y_j \in \mathbb{Z}/p\mathbb{Z}} p_j(y_j) \log p_j(y_j) \\ &\leq p^2 \frac{s}{N} \log(\frac{s}{N}) - 2p \frac{t-1}{N} \log \frac{t-1}{N} \end{aligned}$$

Because $(s-1)p^2 < N \leq sp^2$ and $(t-1)p < N \leq tp$, we have

$$\begin{aligned} I(Y_i; Y_j) &< (1 + \frac{p^2}{N}) \log(\frac{1}{p^2} + \frac{1}{N}) - 2(1 - \frac{p}{N}) \log(\frac{1}{p} - \frac{1}{N}) = \log \frac{1 + \frac{p^2}{N}}{(\frac{1}{p} - \frac{1}{N})^2} + \frac{p^2}{N} \log(\frac{1}{p^2} + \frac{1}{N}) + 2\frac{p}{N} \log(\frac{1}{p} - \frac{1}{N}) \\ &= \log \frac{1 + \frac{p^2}{N}}{(1 - \frac{p}{N})^2} + \frac{p^2}{N} (\log(1 + \frac{p^2}{N}) - 2 \log p) + 2\frac{p}{N} (\log(1 - \frac{p}{N}) - \log p) < \log \frac{1 + \frac{p^2}{N}}{(1 - \frac{p}{N})^2} + \frac{p^2}{N} \log(1 + \frac{p^2}{N}) \\ &= O(\frac{p^2}{N}) + O(\frac{p}{N}) + \frac{p^2}{N} O(\frac{p^2}{N}) = O(\frac{p^2}{N}) \end{aligned}$$

However, $p \in [N^{\frac{1}{k}}, N^{\frac{1}{k}-\epsilon}]$ implies that $p = N^{\frac{1}{k}}(1 + o(1))$, hence we have

$$I(Y_i; Y_j) = O(N^{\frac{2}{k}-1}) \rightarrow 0 \text{ as } N \rightarrow \infty$$

□

3.2 Remainder CC and Gauss CC

The theorem 2.3, 2.4 tells us that the Remainder CC and Gauss CC satisfies the ‘‘Classes high separable’’ property. In fact, they satisfy the property ‘‘Base learners independence’’ also.

Theorem 3.2. *Let $f : \mathbb{Z}/N\mathbb{Z} \rightarrow \prod_{i=1}^r \mathbb{Z}/N_i\mathbb{Z}$ be a Remainder CC, and x be uniformly randomly selected from $\mathbb{Z}/N\mathbb{Z}$, we have that for any $i \neq j$, the mutual Information of $f_i(x)$ and $f_j(x)$ approximate 0.*

Proof.

Let $t_i := \lceil \frac{N}{p_i} \rceil$ and $s_{ij} = \lceil \frac{N}{p_i p_j} \rceil$ for every i, j . We have that the probabilities of $f_i(x)$ at every point in $\mathbb{Z}/p_i\mathbb{Z}$ are $\frac{t_i}{N}$ or $\frac{t_i-1}{N}$ and the probabilities of $(f_i(x), f_j(x))$ at every point in $\mathbb{Z}/p_i\mathbb{Z} \times \mathbb{Z}/p_j\mathbb{Z}$ are $\frac{s_{ij}}{N}$ or $\frac{s_{ij}-1}{N}$ by using the similar method in the proof of Theorem 3.1.

We know that the mutual information of $y_i = f_i(x)$ and $y_j = f_j(x)$ is

$$I(Y_i; Y_j) = \sum_{(y_i, y_j) \in \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}} p_{i,j}(y_i, y_j) \log \frac{p_{i,j}(y_i, y_j)}{p_i(y_i)p_j(y_j)}$$

a.) When $k = 2$, we have $N < p_i p_j$ and hence $s = 1$ and $p_{i,j}(y_i, y_j) = \frac{1}{N}$ on N 's point in $\mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}$ and 0 on other points. Hence we have

$$\begin{aligned} I(Y_i; Y_j) &\leq \sum_{(y_i, y_j) \in \mathbb{Z}/p\mathbb{Z} \times \mathbb{Z}/p\mathbb{Z}} p_{i,j}(y_i, y_j) \log \frac{p_{i,j}(y_i, y_j)}{\frac{1}{N^2}} \\ &= N \times \frac{1}{N} \log \frac{1/N}{\frac{1}{N^2}} \\ &= \log N - \log(t_i - 1) - \log(t_j - 1) \\ &< \log N - \log(\frac{N}{p_i} - 1) - \log(\frac{N}{p_j} - 1) \\ &\leq \log N - 2 \log(\frac{N^{\frac{1}{2}}}{N^{2-\epsilon}} - 1) \\ &= -2 \log(\frac{1}{N^{2-\epsilon}} - N^{-\frac{1}{2}}) \\ &\rightarrow 0 \text{ as } N \rightarrow \infty \end{aligned}$$

b.) When $k \geq 3$, we have $p_i p_j < N^{\frac{2}{k-\epsilon}} < N$, and

$$\begin{aligned}
& I(Y_i; Y_j) \\
&= \sum_{(y_i, y_j) \in \mathbb{Z}/p_i\mathbb{Z} \times \mathbb{Z}/p_j\mathbb{Z}} p_{i,j}(y_i, y_j) \log p_{i,j}(y_i, y_j) \\
&\quad - \sum_{(y_i, y_j) \in \mathbb{Z}/p_i\mathbb{Z} \times \mathbb{Z}/p_j\mathbb{Z}} p_{i,j}(y_i, y_j) (\log p_i(y_i) + \log p_j(y_j)) \\
&= \sum_{(y_i, y_j) \in \mathbb{Z}/p_i\mathbb{Z} \times \mathbb{Z}/p_j\mathbb{Z}} p_{i,j}(y_i, y_j) \log p_{i,j}(y_i, y_j) \\
&\quad - \sum_{y_i \in \mathbb{Z}/p_i\mathbb{Z}} p_i(y_i) \log p_i(y_i) - \sum_{y_j \in \mathbb{Z}/p_j\mathbb{Z}} p_j(y_j) \log p_j(y_j) \\
&\leq p_i p_j \frac{s_{ij}}{N} \log\left(\frac{s_{ij}}{N}\right) - p_i \frac{t_i-1}{N} \log\left(\frac{t_i-1}{N}\right) - p_j \frac{t_j-1}{N} \log\left(\frac{t_j-1}{N}\right)
\end{aligned}$$

Because

$$\begin{aligned}
(s_{ij} - 1)p_i p_j &< N \leq s_{ij} p_i p_j \\
(t_i - 1)p_i &< N \leq t_i p_i \\
(t_j - 1)p_j &< N \leq t_j p_j
\end{aligned}$$

We have

$$\begin{aligned}
& I(Y_i; Y_j) \\
&< \left(1 + \frac{p_i p_j}{N}\right) \log\left(\frac{1}{p_i p_j} + \frac{1}{N}\right) \\
&\quad - \left(1 - \frac{p_i}{N}\right) \log\left(\frac{1}{p_i} - \frac{1}{N}\right) - \left(1 - \frac{p_j}{N}\right) \log\left(\frac{1}{p_j} - \frac{1}{N}\right) \\
&= \log\left(\frac{\frac{1}{p_i} + \frac{1}{N}}{\left(\frac{1}{p_i} - \frac{1}{N}\right)\left(\frac{1}{p_j} - \frac{1}{N}\right)} + \frac{p_i p_j}{N}\right) + \frac{p_i p_j}{N} \log\left(\frac{1}{p_i p_j} + \frac{1}{N}\right) \\
&\quad + \frac{p_i}{N} \log\left(\frac{1}{p_i} - \frac{1}{N}\right) + \frac{p_j}{N} \log\left(\frac{1}{p_j} - \frac{1}{N}\right) \\
&\leq \log\left(1 + \frac{p_i p_j}{N}\right) - \log\left(1 - \left(\frac{1}{p_i} + \frac{1}{p_j}\right) \frac{p_i p_j}{N} + \frac{1}{N^2}\right) \\
&\leq \log\left(1 + \frac{p_i p_j}{N}\right) - \log\left(1 - \left(\frac{1}{p_i} + \frac{1}{p_j}\right) \frac{p_i p_j}{N}\right) \\
&= O\left(\frac{p_i p_j}{N}\right) + O\left(\left(\frac{1}{p_i} + \frac{1}{p_j}\right) \frac{p_i p_j}{N}\right) \\
&= O\left(\frac{p_i p_j}{N}\right) \\
&= O\left(N^{\frac{2}{k-\epsilon}} - 1\right) \\
&= O\left(N^{\frac{2+\epsilon-k}{k-\epsilon}}\right) \rightarrow 0 \quad \text{as } N \rightarrow \infty
\end{aligned}$$

■

This theorem tells us that, the Remainder CC satisfies the property ‘‘Base learners independence’’.

Similarly, we have

Theorem 3.3. *Let $f : \mathbb{Z}/N\mathbb{Z} \rightarrow \prod_{i=1}^r \mathbb{Z}/N_i\mathbb{Z}$ be a Gauss CC, and x be uniformly randomly selected from $\mathbb{Z}/N\mathbb{Z}$, we have that for any $i \neq j$, the mutual Information of $f_i(x)$ and $f_j(x)$ approximate 0.*

□

This theorem tells us that, the Gauss CC satisfies the property ‘‘Base learners independence’’ also.

3.3 Decode Algorithm

Suppose we used the LM $f_i : \mathbb{Z}/N\mathbb{Z} \rightarrow \mathbb{Z}/N_i\mathbb{Z}$ ($i = 1, 2, \dots, n$) to reduce a classification problem of class number N to the classification problems of class number N_i 's, and trained n base learner for every f_i , the output of every base learner i is a distribution P_i on $\mathbb{Z}/N_i\mathbb{Z}$. Now, for a input feature data, how we collect the output $\{P_i : i = 1, 2, \dots, n\}$ of every base learner to get the predict label?

In this paper, we search the $x \in \mathbb{Z}/N\mathbb{Z}$ such that $\sum_i \log P_i(f_i(x))$ is maximal, and let such x be the decoded label. (In fact, $\sum_i \log P_i(f_i(a)) = -\sum_i KL(f_{i*} \delta(x-a) || P_i)$, where $\delta(x-a)$ is the Delta distribution at $a \in \mathbb{Z}/N\mathbb{Z}$, and $f_{i*} \delta(x-a)$ is the marginal distribution of $\delta(x-a)$ induced by f_i .)

3.4 Numeric Experiments

We use the Inception V3 network and LM on the dataset ‘‘CJK characters’’. CJK is a collective term for the Chinese, Japanese, and Korean languages, all of which use Chinese characters and derivatives (collectively, CJK characters) in their writing systems. The data set ‘‘CJK characters’’ is the grey-level image of size 139x139 of 20901 CJK characters (0x4e00 ~ 0x9fa5) in 8 fonts.

We use 7 fonts as the train set, and other one font as the test set. We use inception v3 network as base learner, and train the networks using batch size=128 and 100 batch per an epoch.

We use three CCs as follows, and get the performance like in Table 1.

a. The polynomial CCs with $k=2$ and $p=181$. These Polynomial CCs are defined by $f : \mathbb{Z}/N\mathbb{Z} \rightarrow \mathbb{F}_p^r$, where $N = 21901$, and $f_i(x) = ((x \bmod p) + \text{floor}(x/p)i) \bmod p$, and $r=2$ or $r=6$.

b. The Remainder CCs with $k=2$ and $p_i \in \{173, 191, 157, 181, 193, 199\}$. These Remainder CCs are defined by $f : \mathbb{Z}/N\mathbb{Z} \rightarrow \prod_{i=1}^r \mathbb{Z}/p_i\mathbb{Z}$, where $N = 21901$, $f_i(x) = x \bmod p_i$, and $r = 2$ or 6 .

c. the Gauss CCs with $k=2$ and $p_i \in \{10 \pm 9\sqrt{-1}, 13 \pm 2\sqrt{-1}, 12 \pm 7\sqrt{-1}\}$. These Gauss CCs are defined by $f : \overline{U_{82}(0)} \cap \mathbb{Z}[\sqrt{-1}] \rightarrow \prod_{i=1}^r \mathbb{Z}[\sqrt{-1}]/(p_i)$, where $N = 21901$, and $f_i(x) = x \bmod (p_i)$, and $r=2$ or $r=6$.

d. ECOC of 15 bit.

ep.	ECOC of 15 bit	Poly. CC of 2 sites	Rem. CC of 2 sites	Gauss CC of 2 sites	Poly. CC of 6 sites	Rem. CC of 6 sites	Gauss CC of 6 sites
20	0.0069	0.0118	0.0081	0.0017	0.0640	0.0459	0.0230
40	0.0795	0.6657	0.6130	0.4308	0.9878	0.9667	0.9760
60	0.3660	0.8172	0.7629	0.8436	0.9968	0.9962	0.9966
80	0.5740	0.8684	0.8757	0.9195	0.9988	0.9983	0.9985
param. num (10^7)	2.18×15	2.21×2	2.21×2	2.21×2	2.21×6	2.21×6	2.21×6

Table 1: Comparing of ECOC and CCs

We can see, even when the base learner number 2 of CCs is much less than the base learner number 15 of ECOC, the performance of CCs are better than the ECOC which trainable number of parameters of networks bigger than CCs.

4 Application for feature encode

For a categorical feature take value in $\mathbb{Z}/N\mathbb{Z}$, where N is a huge integral number, we can use the composite mapping of a CC $\mathbb{Z}/N\mathbb{Z} \rightarrow \prod_{i=1}^r \mathbb{Z}/N_i\mathbb{Z}$ and the nature embedding

$$\begin{aligned} \prod_{i=1}^r \mathbb{Z}/N_i\mathbb{Z} &\longrightarrow \prod_{i=1}^r \mathbb{F}_2^{N_i} = \mathbb{F}_2^{\sum_i N_i} \\ (x_i)_i &\mapsto (N_i \text{ bit one hot representation of } x_i)_i \end{aligned}$$

to get a r -hot encoding. We use this r -hot encoding as feature encoding.

Apart from the CC feature encoding, the more natural ideas for feature encoding are

COO. Cut off of one-hot encoding. We call a n -bit binary code the 'Cut off of one-hot', if the $n - 1$ most frequently used ID's are one-hot encoded in the front $n - 1$ bits, and all the other ID's are encoded to the code '0...01'.

RMP. Using a code frequently used in error-correct encoding. For example, a Reed-Muller code [13] with punch by a random subset of bits. For a binary code $\{f_i\}_{i \in \mathbb{Z}/n\mathbb{Z}} : C \hookrightarrow \mathbb{F}_2^m$ and a subset $Q \subset \mathbb{Z}/n\mathbb{Z}$ of m elements, the punch of f by Q means the code $\{f_i\}_{i \in \mathbb{Z}/n\mathbb{Z} \setminus Q} : C \hookrightarrow \mathbb{F}_2^{n-m}$.

We will show that, the performance of our Polynomial CC, Remainder CC and Gauss CC are better than both the code COO and RMP.

4.1 Numeric Experiments

We use the dataset "Movie Lends" ([4]), which has the columns UserID, MovieID, Rating and Timestamp. The UserIDs range between 1 and 6040, and MovieIDs range between 1 and 3952, ratings are made on a 5-star scale, timestamp is represented in seconds. Each user has at least 20 ratings. We use only the column UserID, MovieID and Rating. and use a DNN with an embedding layer and two full-connected layers. In the embedding layer, the User code and Movie code are embedded to real vectors of dimension 32 respectively, the dimension of the output the two full-connected layers are 64 and 1 respectively. After the first full-connected layer we use 'RELU', after the second full-connected layer we use $x \mapsto 4 * \text{sigmoid}(x) + 1$. We use this network as a regression model, and train it by minimize MSE. The ratio between train data and validation data is 8:2. We compare the validation loss of the following methods:

1. 582 bit cut off of the one-hot code for UserID, and 474 bit cut off of the one-hot code for MovieID.

2. 582 bit random punch of RM(12,1) for UserID, and 474 bit random punch of RM(11,1) for MovieID.
3. 582 bit 6-hot Polynomial code based on finite field \mathbb{F}_{97} for UserID, and 474 bit 6-hot Polynomial code based on finite field \mathbb{F}_{73} for MovieID.
4. 582 bit Remainder code with modules $\{83, 89, 97, 101, 103, 109\}$ for UserID, and 474 bit Remainder code with modules $\{67, 71, 73, 79, 83, 101\}$ for MovieID.
5. 582 bit Gauss code with modules $\{8 \pm 5\sqrt{-1}, 9 \pm 4\sqrt{-1}, 10 + \sqrt{-1}, 10 + 3\sqrt{-1}\}$ for UserID, and 474 bit Remainder code with modules $\{67, 71, 73, 79, 83, 101\}$ for MovieID.

The validation losses are like in Table 2. We see that the performance of Polynomial CC, Remainder CC and Gauss CC are better than the one-hot cut and RM code with punch of same length significantly. Moreover, the performance of Gauss CC is best, and then the Remainder CC.

4.2 Theoretical analysis for feature coding

We see the performance of Polynomial CC, Remainder CC and Gauss CC are good for feature coding, but we don't know how to choose the non-zero bit number r in the coding. More generally, how to study the performance of codes without experiments? In the theory of error-correcting code, we know the Hamming distance is an important metric for codes. In general, if the original IDs and length of coding is fixed, the error-correcting codes with big Hamming distance have good performance. But for feature coding, Hamming distance is not a good metric. For example, we compare the performance of Method 2 introduced in the previous subsection and the anti-Method 2.

The codings used in anti-Method 4 and Method 4 have the relationship: $x \mapsto 1 - x$. The corresponding pair of codes in the two method has same Hamming distance, but the performance is difference (in Table 2). Hence the Hamming distance is not a good choose for metric of feature encoding.

For a binary r -hot codeword c of length n , we can view $\frac{1}{r}c$ as a distribution on $\mathbb{Z}/n\mathbb{Z}$, and call it **the reduced distribution of \mathbf{x}** , write it as $\text{dist}(c)$. The **average minimal KL-divergence (AMKL)** of a code $I \rightarrow C$ is defined as $\sum_i \min_j \text{KL}(\text{dist}(c_i) || \text{dist}(c_j)) p_i$. We propose that use **AMKL** as the metric of code, and give the conjecture:

Conjecture 4.1. *The feature code with bigger AMKL has better performance.*

To examine the conjecture 4.1, we give a lemma to compute the **AMKL** firstly:

Lemma 4.2. *For a n bit r -hot code $I \rightarrow \mathbb{F}_2^n$, if for any codeword c_i the maximal common non-zero bit number between c_i and any other codeword in C is r , the **AMKL** equal to $(1 - \frac{r}{n})\infty$.*

Proof. For any $i \in I$, let x_i denote the codeword of i . For any $i \neq j$ in I , the reduced distribution of x_i, x_j are $\text{dist}(x_i) = \frac{1}{r}x_i, \text{dist}(x_j) = \frac{1}{r}x_j$ respectively. Hence the KL-divergency of $\text{dist}(x_i), \text{dist}(x_j)$ is $\text{KL}(\text{dist}(x_i) || \text{dist}(x_j)) = \frac{1}{r} \log \frac{1/r}{0} \times (r - r) + \frac{1}{r} \log \frac{1/r}{1/r} \times r = (1 - \frac{r}{r}) \log \infty$. Hence $\sum_i \min_j \text{KL}(\text{dist}(c_i) || \text{dist}(c_j)) p_i = \sum_i (1 - \frac{r}{r}) p_i \log \infty = (1 - \frac{r}{r}) \log \infty$. \square

We use the some numeric experiments to examine the conjecture 4.1. We use the following encoding Method on dataset "Movie Lends", and their **AMKL** and performance is like in table 3. We see that the **AMKL** has positive effect to performance. Moreover, the performance of Gauss CC > Remainder CC > Polynomial CC with same length and **AMKL**.

Method 1. 582 bit Remainder code with modules $\{289, 293\}$ for UserIDs, and 474 bit Remainder code with modules $\{235, 239\}$ for MovieIDs.

Method 2. 582 bit Remainder code with modules $\{193, 194, 195\}$ for UserIDs, and 474 bit Remainder code with modules $\{157, 158, 159\}$ for MovieIDs.

Method 3. 582 bit 6-hot Polynomial code based on finite field \mathbb{F}_{97} for UserIDs, and 474 bit 6-hot Polynomial code based on finite field \mathbb{F}_{73} for MovieIDs.

Method 4. 582 bit Remainder code with modules $\{83, 89, 97, 101, 103, 109\}$ for UserIDs, and 474 bit Remainder code with modules $\{67, 71, 73, 79, 83, 101\}$ for MovieIDs.

Method 5. 582 bit Gauss code with modules $\{8 \pm 5\sqrt{-1}, 9 \pm 4\sqrt{-1}, 10 + \sqrt{-1}, 10 + 3\sqrt{-1}\}$ for UserIDs, and 474 bit Remainder code with modules $\{67, 71, 73, 79, 83, 101\}$ for MovieIDs.

Method 6. 582 bit Remainder code with modules {19, 23, 25, 27, 29, 31, 32, 37, 41, 43, 47, 49, 53, 59, 67 } for UserIDs, and 473 bit Remainder code with modules {17, 19, 23, 25, 27, 29, 31, 32, 37, 41, 43, 47, 49, 53} for MovieIDs.

ep.	Method 1	Method 2	Method 3	Method 4	Method 5	anti-Meth. 4
1	1.044	0.960	0.938	0.936	0.932	1.235
5	1.036	0.862	0.857	0.857	0.854	1.301
9	1.034	0.857	0.853	0.851	0.849	1.177
13	1.029	0.855	0.853	0.850	0.846	1.141

Table 2: Comparing of Coding Methods

method	u. avg min KL	i. avg min KL	MSE ep.1	MSE ep. 2	MSE ep.3
1 (two hot)	0.5	0.5	1.0335	0.951	0.913
2 (three hot)	0.667	0.667	0.992	0.912	0.887
3 (six hot)	0.833	0.833	0.938	0.886	0.873
4 (six hot)	0.833	0.833	0.936	0.883	0.867
5 (six hot)	0.833	0.833	0.932	0.883	0.866
6 (15, 14 hot)	0.867	0.857	0.908	0.875	0.864

Table 3: avg min KL

5 Conclusion

We propose three classes of category coding (CC) with minimal collision property. They are Polynomial CC, Remainder CC and Gauss CC.

In the application for label coding in the classification problem with huge labels number using CNN, we prove that they have good theoretical properties and show that they have good performance in numerical experiments.

In the application for feature coding in collaborative filtering using DNN, we show that their performance is better than cut-off method and classical binary coding method. Moreover, we give a metric “AMKL” of feature coding, and show it has positive effect to the performance. In additional, we show that the performance of Gauss CC > Remainder CC > Polynomial CC with same length and AMKL.

References

- [1] Nathan Jacobson. Lectures in Abstract Algebra III: Theory of Fields and Galois Theory. Springer-Verlag New York, 1964.
- [2] https://en.wikipedia.org/wiki/Gaussian_integer
- [3] E. Guerrini and M. Sala. A classification of MDS binary systematic codes. BCR preprint 2006. www.bcricri.ucc.ie/FILES/PUBS/BCRI_57.pdf,
- [4] <http://grouplens.org/datasets/movielens/1m/>
- [5] Lang, Serge. Algebraic Number Theory. Springer-Verlag New York, 1994.
- [6] Apostol T. M. Introduction to Analytic Number Theory. Springer-Verlag, 1976, New York.
- [7] Irving S. Reed and Gustav Solomon. Polynomial codes over certain finite fields. JSIAM volume 8(2), jun of 1960, p300–304. [url=http://links.jstor.org/sici?sici=0368-4245%28196006%298%3A2%3C300%3APCOCFF%3E2.0.CO%3B2-2](http://links.jstor.org/sici?sici=0368-4245%28196006%298%3A2%3C300%3APCOCFF%3E2.0.CO%3B2-2)
- [8] Bernhard Riemann. Ueber die Anzahl der Primzahlen unter einer gegebenen Grosse. Monatsberichte der Berliner Akademie, November 1859.
- [9] https://en.wikipedia.org/wiki/Prime_number_theorem.
- [10] Irving S. Reed and Gustav Solomon. Polynomial codes over certain finite fields. J. SIAM, 8:300-304, 1960.
- [11] Richard C. Singleton. Maximum distance q-nary codes. IEEE Transactions on Information Theory, 10(2):116–118, April 1964.
- [12] Sejnowski T.J., Rosenberg C.R.(1987).Parallel networks that learn to pronounce english text. Journal of Complex Systems,1(1), 145-168.
- [13] Shu Lin; Daniel Costello (2005). Error Control Coding (2 ed.). Pearson. ISBN 0-13-017973-6. Chapter 4.

- [14] L. R. Vermani. Elements of Algebraic Coding Theory. CRC Press, 1996.
- [15] E. Guerrini and M. Sala. A classification of MDS binary systematic codes. BCRI preprint, www.bcric.ucc.ie 56, UCC, Cork, Ireland, 2006.
- [16] T. G. Dietterich and G. Bakiri. Solving multiclass learning problems via error-correcting output codes. *Journal of artificial intelligence research*, pp. 263–286, 1995.
- [17] E. L. Allwein, R. E. Schapire, and Y. Singer. Reducing multiclass to binary: A unifying approach for margin classifiers. *The Journal of Machine Learning Research*, 1:113–141, 2001.
- [18] A. Passerini, M. Pontil, and P. Frasconi. New results on error correcting output codes of kernel machines. *Neural Networks, IEEE Transactions on*, 15(1):45–54, 2004.
- [19] Langford, J., and Beygelzimer, A. 2005. Sensitive error correcting output codes. In COLT.
- [20] G. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006
- [21] S. Escalera, O. Pujol, and P. Radeva. Ecoc-one: A novel coding and decoding strategy. In ICPR, volume 3, pp. 578–581, 2006.
- [22] O. Pujol, P. Radeva, and J. Vitria. Discriminant ECOC: a heuristic method for application dependent design of error correcting output codes. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 28(6):1007–1012, 2006.
- [23] O. Pujol, S. Escalera, and P. Radeva. An incremental node embedding technique for error correcting output codes. *Pattern Recognition*, 41(2):713–725, 2008.
- [24] S. Escalera, O. Pujol, and P. Radeva. Separability of ternary codes for sparse designs of error-correcting output codes. *Pattern Recognition Letters*, 30(3):285–297, 2009.
- [25] G. Zhong, K. Huang, and C.-L. Liu. Joint learning of error-correcting output codes and dichotomizers from data. *Neural Computing and Applications*, 21(4):715–724, 2012.
- [26] G. Zhong and M. Cheriet. Adaptive error-correcting output codes. In IJCAI, 2013.
- [27] G. Zhong and C.-L. Liu. Error-correcting output codes based ensemble feature extraction. *Pattern Recognition*, 46(4):1091–1100, 2013.
- [28] Yang, Luo, Loy, Shum, Tang. Deep Representation Learning with Target Coding. Twenty-Ninth AAAI Conference on Artificial Intelligence, 2015 :3848-3854.
- [29] Berger, A.: Error-Correcting Output Coding for text classification. In: IJCAI(1999)
- [30] Ghani, R.: Using error-correcting codes for text classification. *Proceedings of ICML-00, 17th International Conference on Machine Learning* (pp. 303–310). Stanford, US: Morgan Kaufmann Publishers, San Francisco, US.
- [31] Ghani, R. Using Error-Correcting Codes for Efficient Text Classification with a Large Number of Categories. KDD Lab Project Proposal.
- [32] T. G. Dietterich and G. Bakiri. Solving multiclass learning problems via error-correcting output codes. *Journal of artificial intelligence research*, 1995, p263-286.