# Data preprocessing in MySQL

Group 6

**CHING HUSAN LIU***
Data Science Degree Program
National Taiwan University
Taipei, Taiwan
r10946013@ntu.edu.tw

**YU SYUAN CHUANG**
Department of Electrical
Engineering
National Taiwan University
Taipei, Taiwan
r10921095@ntu.edu.tw

**YAN TING YE**
Communication Engineering
National Taiwan University
Taipei, Taiwan
r10942153@ntu.edu.tw

**CHIAO YIN TENG**
Bioenvironmental Systems Engineering
National Taiwan University
Taipei, Taiwan
b07602001@ntu.edu.tw

## ABSTRACT

Data mining and data analysis usually cost lots of effort and are time-consuming. When we train models in machine learning, data pre-processing has a great influence on the performance of the model. It is time-consuming to write different data pre-processing programs for different datasets.

In this article, we present practical issues and common solutions when preparing and transforming datasets with SQL language. We implement common data pre-processing, such as filling blank values, removing duplicate data, statistical applications (quartiles, percentiles), one-hot encoding, image augmentation and feature correlation. Users only need to add specified parameters in the query field to get the pre-processed data directly.

## KEYWORDS

data pre-processing, image augmentation, feature correlation, encoding, MySQL.

## 1 Data pre-processing

Figure 1 is the dataset we implement null value filling. You can see that the bottom column has the missing value "NULL".



Figure 1: dataset of the missing value

### 1.1 Missing value

First, we calculate the prices of the other three types of book, and then update the missing value to the average. The result shows in Figure 2.



Figure 2: implement of the average

For maximum value, we take the maximum price of the other three types of book, and then update

the value to the maximum. The result shows in Figure 3.



Figure 3: implement of the maximum

For calculation of mode, we set a count to calculate the price of occurrences of the other three types of book, and then update the value to the mode. The result shows in Figure 4.



Figure 4: implement of the mode

The random value in the range is to find the maximum and minimum values of the price, and then use the RAND() function to find the random value and then update the value to the missing value. The result shows in Figure 5.



Figure 5: implement of the random value

## 1.2    **Statistical applications**

We designed a new dataset to find the quartile of the data and made a statistical table. Figure 6 is the dataset we implement quartile.



Figure 6: dataset of the quartile

However, when implementing the GROUP BY, an error message may occur: "Error Code: 1055." or "Error Code: 1140." The solution is to delete the function that is enabled by default in MySQL and reset the "SQL_mode" to solve the problem. The results of the final implementation of quartiles are shown in Figure 7, and the quartiles of each category are calculated according to the category of books.

| Category | Book_ID | AssessedValue | Quartile |
|----------|---------|---------------|----------|
| 童書 | 12 | 220 | 1 |
| 漫畫 | 3 | 230 | 1 |
| 童書 | 9 | 289 | 1 |
| 童書 | 7 | 295 | 2 |
| 漫畫 | 5 | 300 | 2 |
| 童書 | 11 | 313 | 2 |
| 童書 | 10 | 331 | 3 |
| 童書 | 8 | 379 | 3 |
| 漫畫 | 1 | 441 | 3 |
| 漫畫 | 2 | 447 | 4 |
| 漫畫 | 4 | 496 | 4 |
| 漫畫 | 6 | 525 | 4 |

12 rows in set (0.00 sec)

| Category | Minimum | 1Quartile | Median | 3Quartile | Maximum | Count |
|----------|---------|-----------|--------|-----------|---------|-------|
| 漫畫 | 230 | 230 | 300 | 441 | 525 | 6 |
| 童書 | 220 | 289 | 313 | 379 | 379 | 6 |

2 rows in set (0.00 sec)

Figure 7: implement of the quartile

For percentile, users can set the percentage by themselves, and finally output the percentile of the category of books. The result shows in Figure 8.

Figure 8: implement of the percentile

## 1.3 **One-hot Encoding**

We implement the one-hot encoding with coalesce to display 1 for the data that exists in this column, and display 0 for the rest. The result shows in Figure 9.



Figure 9: implement of one-hot encoding

## 2 **Image Augmentation**

Image augmentation is a common technique of data preprocessing for image tasks. In this section, we propose several image augmentation functions in MySQL:

1. Horizontal / Vertical Shift

2. Horizontal / Vertical Flip

3. Crop

These functions are implemented purely using MySQL language. We implement them as stored functions. Once functions are stored in a MySQL server, users can call these functions to get the desired augmented images.

Here we only deal with grayscale images for simplicity.

### 2.1 **Image in MySQL**

Before writing these functions, we first need to understand how images are stored in MySQL. MySQL supports directly reading a PNG / JPG file and storing it as a BLOB object. However, in such a way we can not access pixel values of images directly. A PNG / JPG decoder and encoder would be necessary for decoding images and modifying them. But a PNG / JPG decoder and encoder in MySQL is not what we want to discuss, so we choose to store images in another way.

A grayscale image can be treated as a 2-D numeric array, with all entries are uint8. We can serialize such an array and then store it as a BLOB object. A BLOB object can be viewed as a binary string in MySQL. We can access and modify any element in a binary string.



Fig. 7: Images are treated as arrays and then serialized in order to be stored in MySQL as BLOB datatype.

### 2.2 **Function Implementation**

As we mention above, any element in a BLOB object can be arbitrarily modified, which means any pixel value can be changed to the value we want. Theoretically, any kind of image augmentation can be realized in MySQL. However, we encounter some difficulties due to the lack of vector and array algorithms in

MySQL. Thus we implement shift, flip, and crop first.

The shift value, corp range are function arguments that the user can specify.

Utilizing the SUBSTRING query in MySQL to reuse the consecutive elements in original images is the main idea. Here we use right shift as an example:

RIGHT_SHIFT(image, w, h, s)

```
1  a_image = ''
2  pad = '\x00' repeat s times
3  FOR i=1..h
4     r = SUBSTRING(image, 0, w-s)
5     row = CONCAT(pad, r)
6     a_image = CONCAT(a_image, row)
7  RETURN a_image
```

The padding can be reused so we construct it first and then reuse it in every iteration.

Other functions are implemented in a similar way. We will discuss their complexity later.

## 2.3 Time Complexity Analysis

For time complexity analysis, we use the number of query calls as the measure, i.e. one query call takes one unit of time.

Here we assume that the height and the width of images are both N for simplicity. (It is reasonable because typically we have width = $\Theta$(height) for most images.)

Horizontal shift with shift value s:

It takes O(N) to construct the padding row. We need N iterations, 2 CONCAT and 1 SUBSTRING in each iteration. The overall time complexity is $\Theta$(N).

Vertical shift with shift value s:

It takes $\Theta$(N) to construct the padding row. We need s CONCAT and 1 SUBSTRING. Overall time complexity is $\Theta$(s) (or O(N)).

Horizontal flip:

We need N iterations, N CONCAT and N SUBSTRING in each iteration. Overall time complexity is $\Theta(N^2)$.

Vertical flip:

We need N iterations, 1 CONCAT and 1 SUBSTRING in each iteration. Overall time complexity is $\Theta(N)$.

Corp:

It takes $\Theta$(N) to construct the right, left, top, and bottom padding row. We need N iterations, 1~3 CONCAT and 0~1 SUBSTRING in each iteration. Overall time complexity is $\Theta(N)$.

## 2.4 Result

We can see in Fig. 8 all functions work correctly.



(a) original      (b) horizontal / vertical shift



(c) horizontal / vertical flip      (d) corp
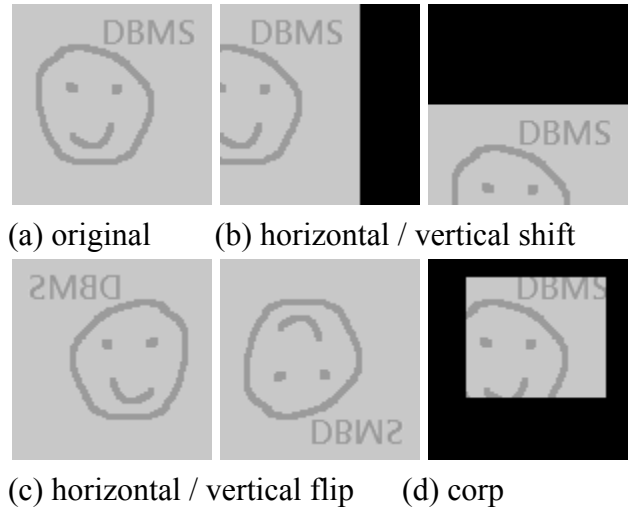
Fig. 8: Augmented images

To compare the execution time of our image augmentation and existing one, we execute each query 1000 times and measure the execution time. (we compare it to Pytorch's implementation.)

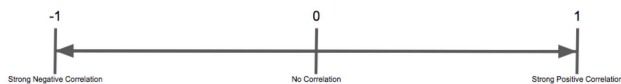|      | h_s  | v_s  | h_f  | v_f  | c    |
|------|------|------|------|------|------|
| our  | 8.02 | 2.06 | 700  | 7.69 | 4.89 |
| PT   | 0.14 | 0.14 | 0.18 | 0.14 | 0.19 |

Table 1: The 1000 times execution time of each query. The h_s, v_s stands for horizontal and vertical shift; h_f, v_f stands for horizontal and vertical flip; c stands for corp. PT means using Pytorch's image augmentation implementation.

## 3 Feature Correlation in MySQL

The correlation coefficient is a widely used method to determine the strength of the relationship between two variables or two sets of numbers. This coefficient is a number between -1 and 1, which 1 means the strongest possible positive correlation , -1 means the strongest possible negative correlation and 0 means there is no correlation.

A positive correlation means that as one number increases the second number will also increase. A negative correlation means that as one number increases the second number decreases.

The example of possible positive correlation is : sales number and marketing spend. The example of possible negative correlation is : product price and its sales.



### 3.1 Correlation v.s. Causation

The possible positive correlation might not mean there is causation between two variables. For example, the strong correlation between two numbers might be caused by randomness. Also, while there might be a causation between two variables and there is a strong correlation between them, the correlation coefficient may not tell us anything about the causation between the two variables.

### 3.2 the Pearson Correlation Coefficient

The Pearson Correlation Coefficient is a widely used method to calculate the correlation coefficient. The original equation is:

$$r = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n\Sigma x^2 - (\Sigma x)^2][n\Sigma y^2 - (\Sigma y)^2]}}$$
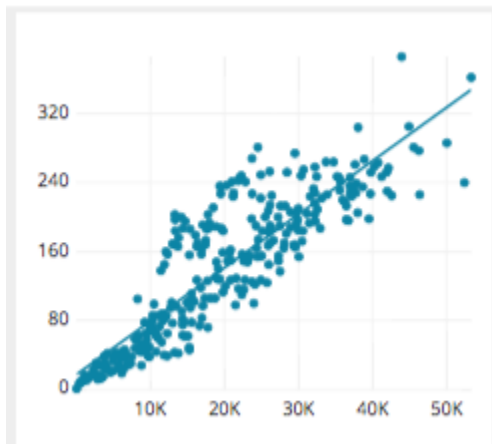
when use in MySQL:

```sql
SELECT
    (
    (tot_sum - (amt_sum * act_sum / _count)) /
    sqrt((amt_sum_sq - pow(amt_sum, 2.0) / _count) *
    (act_sum_sq - pow(act_sum, 2.0) / _count))
    )
    AS "Corr Coef Using Pearson"
```

### 3.3 Implement

```sql
SELECT
    ((tot_sum - (amt_sum * act_sum / _count)) /
    sqrt((amt_sum_sq - pow(amt_sum, 2.0) / _count) * (act_sum_sq - pow(act_sum, 2.0) / _count)))
    AS "Corr Coef Using Pearson"


FROM(
SELECT
    sum("Amount") AS amt_sum,
    sum("Activities") AS act_sum,
    sum("Amount" * "Amount") AS amt_sum_sq,
    sum("Activities" * "Activities") AS act_sum_sq,
    sum("Amount" * "Activities") AS tot_sum,
    count(*) as _count

FROM(
SELECT
    DATE_TRUNC('day', p.payment_date)::DATE AS "Day",
    SUM(p.amount) AS "Amount",
    COUNT(DISTINCT a.activity_id) AS "Activities"
FROM
    public.payments p
    INNER JOIN public.subscriptions s ON p.subscription_id = s.subscription_id
    INNER JOIN public.users u ON s.user_id = u.user_id
    INNER JOIN public.activity a ON a.user_id = u.user_id

GROUP BY 1) as a

) as b

GROUP BY tot_sum, amt_sum, act_sum, _count, amt_sum_sq, act_sum_sq
```

When compared to corr() , a function in PostgreSQL which can calculate the correlation coefficient directly , they have the same correlation coefficient 0.9 and its figure below.

## ACKNOWLEDGMENTS

Insert paragraph text here. Insert paragraph text here. Insert paragraph text here. Insert paragraph text here. Insert paragraph text here. Insert paragraph text here. Insert paragraph text here. Insert paragraph text here. Insert paragraph text here. Insert paragraph text here. Insert paragraph text here.

## REFERENCES

[1] Patricia S. Abril and Robert Plant, 2007. The patent holder's dilemma: Buy, sell, or troll? *Commun. ACM* 50, 1 (Jan, 2007), 36-44. DOI: https://doi.org/10.1145/1188913.118891 5.

[2] Sten Andler. 1979. Predicate path expressions. In *Proceedings of the 6th. ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages (POPL '79)*. ACM Press, New York, NY, 226-236. DOI:https://doi.org/10.1145/567752.567774

[3] Ian Editor (Ed.). 2007. *The title of book one* (1st. ed.). The name of the series one, Vol. 9. University of Chicago Press, Chicago.DOI:https://doi.org/10.1007/3-540-09 237-4.

[4] David Kosiur. 2001. *Understanding Policy-Based Networking* (2nd. ed.). Wiley, New York, NY.

[5] Ordonez, C. (2011). Data set preprocessing and transformation in a database system. *Intelligent Data Analysis*, *15*(4), 613-631.

[6] Ni, J., Xie, S., Li, Z., & Jia, C. (2018, August). Optimization design method of cache extension in MySQL database. In *2018 IEEE International Conference on Mechatronics and Automation (ICMA)* (pp. 2295-2299). IEEE.

[7] JOSHI, A. P., & PATEL, B. V. (2020). Data Preprocessing: the Techniques for Preparing Clean and Quality Data for Data Analytics Process.

## TEAM MEMBER WORKS

劉馨瑄: Data preprocessing、code cleaning
鄧喬尹: Encoding、video clip、written report
葉彥廷: image augmentation、written report
莊于萱: feature correlation、written report