

Metode Normalisasi

Ada berbagai macam metode normalisasi, seperti MinMax, Zscore, Decimal Scaling, Sigmoid, dan Softmax. Pemakaiannya tergantung pada kebutuhan dari dataset dan jenis analisa yang dilakukan.

MinMax

Metode Min-Max merupakan metode yang cukup bisa dibayangkan karena termasuk metode normalisasi yang bersifat linier dengan data aslinya. Namun, metode ini bisa menyebabkan out of bound pada beberapa kasus.

Min - Max

Formula :

$$\text{new_data} = \frac{(\text{current_data} - \text{min}) * (\text{new_max} - \text{new_min})}{(\text{max} - \text{min})} + \text{new_min}$$

Dimana :

New_data = data baru hasil normalisasi

Current_data = data terkini yang akan dinormalisasi

Min = nilai terkecil dari satu kolom baris data

Max = nilai terbesar dari satu kolom baris data

New_min = batas nilai terkecil dari normalisasi

New_max = batas nilai terbesar dari normalisasi

Kenapa bisa terjadi out of bound? Out of Bound terjadi apabila ada data baru masuk, dimana data tersebut melebihi nilai maksimal atau nilai minimal dari data yang sudah ada. Secara otomatis, perhitungan yang berlaku pada data yang sudah diperoleh tadi harus diulangi lagi semuanya dengan data baru yang masuk atau data baru yang mempunyai nilai maksimal/minimum yang melebihi tadi tidak bisa diproses. Karena kekurangan inilah MinMax tidak cocok untuk analisa real time / evolving system. Dimungkinkan dalam kasus-kasus terjadi kasus out of bound pada MinMax.

MinMax sangat dianjurkan untuk kasus-kasus berbasis time frame analisis dan forecasting. Perhitungan dari metode ini cukup mengurangi data yang asli dengan nilai minimal dari fitur tersebut, kemudian hasil tersebut dikalikan dari hasil pengurangan nilai maksimal yang baru dengan nilai minimal yang baru dan kemudian dibagi dengan nilai max dan min data di setiap fitur terakhir ditambah dengan nilai min yang baru.

Z-Score

Zscore adalah metode yang sering digunakan dalam berbagai penelitian berbasis data mining atau data science. Z-score merupakan metode normalisasi yang berdasarkan mean (nilai rata-rata) dan standard deviation (deviasi standar) dari data. Kenapa Z-Score sangat populer? Selain tidak banyak variabel yang diset dalam perhitungannya. Z-Score sangat dinamis dalam melakukan perhitungan normalisasi. Kelemahan dari Z-Score adalah prosesnya akan terulang lagi jika ada data baru yang masuk. Selain itu elemen yang dibutuhkan untuk perhitungan Z-Score juga membutuhkan proses yang cukup lama baik standar deviation ataupun rata-rata dari setiap kolom.

Z-Score

Formula :

$$\text{new_data} = \frac{(\text{current_data} - \text{mean of columns})}{\text{standar_deviation of columns}}$$

Dimana :

New_data = data baru hasil normalisasi

Current_data = data terkini yang akan dinormalisasi

Mean of columns = rata-rata dari setiap kolom

Standar_deviation of columns = standar deviasi dari setiap kolom

Decimal Scaling

Decimal Scaling

Formula :

$$\text{new_data} = \frac{\text{current_data}}{10^n}$$

Dimana :

New_data = data baru hasil normalisasi

Current_data = data terkini yang akan dinormalisasi

n = pangkat untuk pembagi

Softmax

Softmax merupakan metode normalisasi pengembangan transformasi secara linier. Output range-nya adalah 0-1. Metode ini sangat berguna pada saat data yang ada melibatkan data outlier.

Softmax

Formula

$$\text{new_data} = \frac{1}{(1+e^{(-\text{transfdata})})}$$

Dimana :

New_data = data baru hasil normalisasi

transfdata = $(\text{current_data} - \text{mean}) / (x * (\text{standar_deviasi} / (2 * 3.14)))$

Current_data = data terkini yang akan dinormalisasi

X = respon linier pada standar deviasi (range 0 – 1)

e = 2.718281828

Sigmoid

Sigmoidal merupakan metode normalization melakukan normalisasi data secara nonlinier ke dalam range -1 s/d 1 dengan menggunakan fungsi sigmoid. Metode ini sangat berguna pada saat data yang ada melibatkan data outlier. Data outlier adalah data yang keluar jauh dari jangkauan data lainnya

Sigmoid

Formula :

$$\text{new_data} = \frac{(1-e^{(-x)})}{(1+e^{(-x)})}$$

Dimana :

New_data = data baru hasil normalisasi

e = 2.718281828 (eksponensial)

$x = \frac{(\text{current_data} - \text{mean of columns})}{\text{standar_deviation of columns}}$