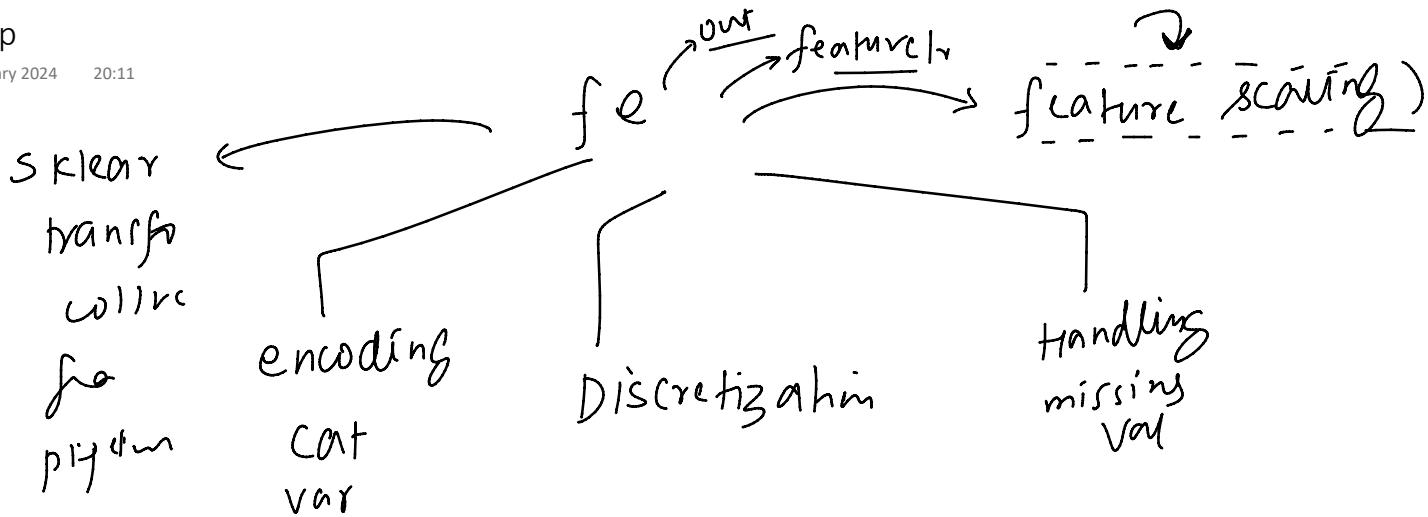


Recap

29 February 2024 20:11



What is Feature Scaling

29 February 2024 14:38

feature

Feature scaling is a method used in data pre-processing for machine learning and data science that involves scaling the range of variables or features of data. The goal of feature scaling is to normalize the range of independent variables or features of data within a specific range, such as between -1 and 1 or 0 and 1, or to standardize the features to have a mean of 0 and a standard deviation of 1. This process is important because features on different scales can distort the distance between data points in some algorithms and can also affect the performance of algorithms negatively.

Why do we need feature scaling?

1. Improves the performance of Distance based algorithms
2. Improves in optimization techniques
3. ~~deep learning~~

Which algorithms are affected if the features are not scaled?

1. Gradient Descent based algorithms
 - a. Linear Regression
 - b. Logistic Regression
 - c. Neural Network
2. SVM
3. Distance based algorithms
 - a. KNN
 - b. K-Means
4. Dimensionality Reduction Techniques
 - a. PCA
 - b. LDA
5. Regularization Techniques
 - a. Ridge Regression
 - b. Lasso Regression

Which algorithms are not affected?

1. Tree Based Algorithms
 - a. CART ✓
 - b. Random Forest ✓
 - c. AdaBoost ✓
 - d. Gradient Boosting ✓
 - e. XGBoost ✓
2. Naïve Bayes

Types of Feature Scaling Techniques

- 1) Standardization
- 2) Normalization
 - + min max scalar (0-100)
 - + Robust scalar
 - + Max abs scalar

distance
y

Kmeans

0-100

51, 20000

concept of distance

manhattan

Knn

compar

optimization

J

scaled data / unscaled

salary

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

$$\hat{w}_2 = w_2 - \eta \left(\frac{\partial L}{\partial w_2} \right)$$

sign

1. Standardization \rightarrow scaling

29 February 2024 15:06

$$x'_i = \frac{x_i - \mu}{\sigma}$$

$$\text{age} \quad \mu = 30 \quad \sigma = 10$$

$$\sigma^2 = 1$$

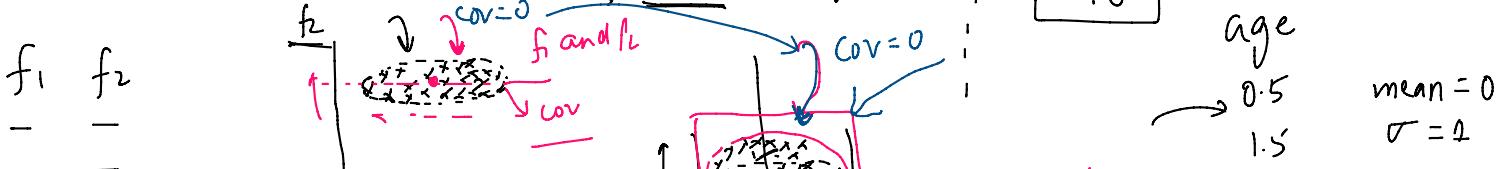
Description: Scales features to have zero mean and unit variance, $z = \frac{(x - \mu)}{\sigma}$ where μ is the mean and σ is the standard deviation.

$$\text{age} \quad \mu = 30 \quad \sigma = 10$$

$$35 \quad 45 \quad \text{Standard} \quad \frac{35 - 30}{10} = \frac{5}{10} = 0.5$$

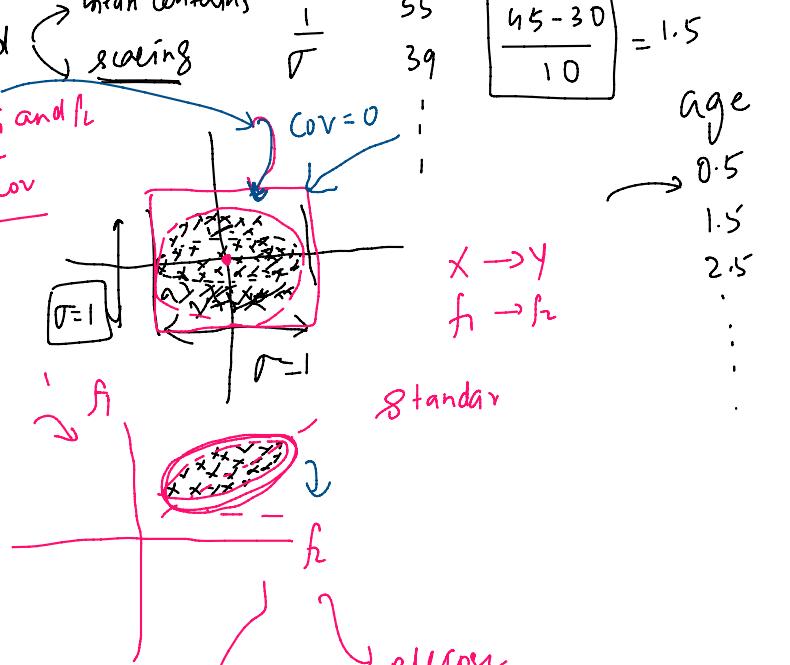
$$55 \quad 39 \quad \frac{45 - 30}{10} = 1.5$$

Geometric Intuition:



$$\text{mean} = 0$$

$$\sigma = 10$$



Advantages:

1. Simple and efficient to implement

Disadvantages:

1. Does not work with algs that assume non negative values.

algo \rightarrow negative
 ↳ transform
 ↳ Normalization

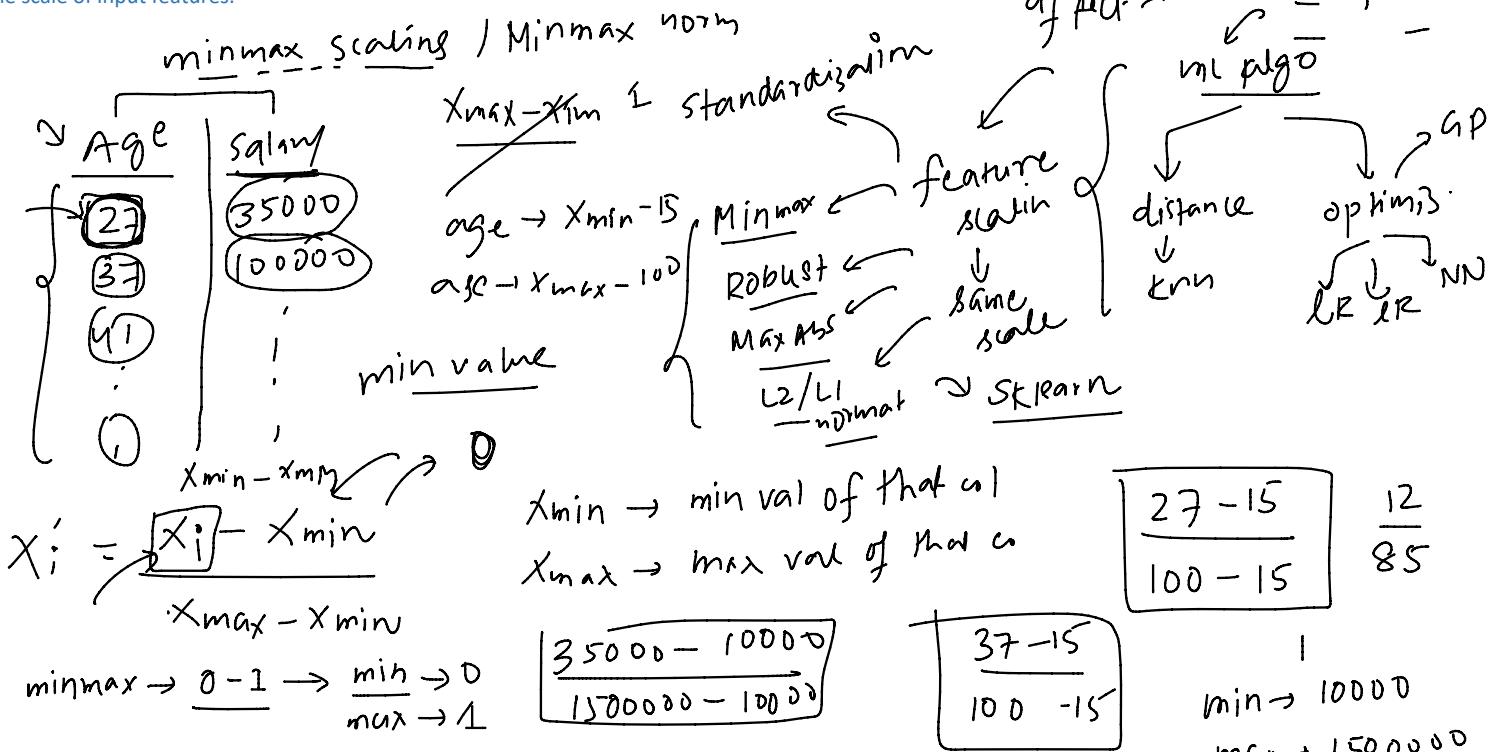
multinomial naive bayes
 - - - - -
 deep \rightarrow relu

circular \rightarrow ellipses

2. Minmax Scaling

29 February 2024 17:06

Normalization is a data pre-processing technique used in machine learning to adjust the values of numeric columns in the dataset to a common scale, without distorting differences in the ranges of values or losing information. Unlike standardization, which rescales data to have a mean of 0 and a standard deviation of 1, normalization typically rescales the data into the range between 0 and 1. This process is important for modeling algorithms that are sensitive to the scale of input features.



1. Fixed Range: After applying Min-Max scaling, all feature values are transformed to lie within a predefined range, typically [0, 1]. If X is a feature, the transformed value 'X' will be in the range where the minimum value of the feature becomes 0, and the maximum value becomes 1.

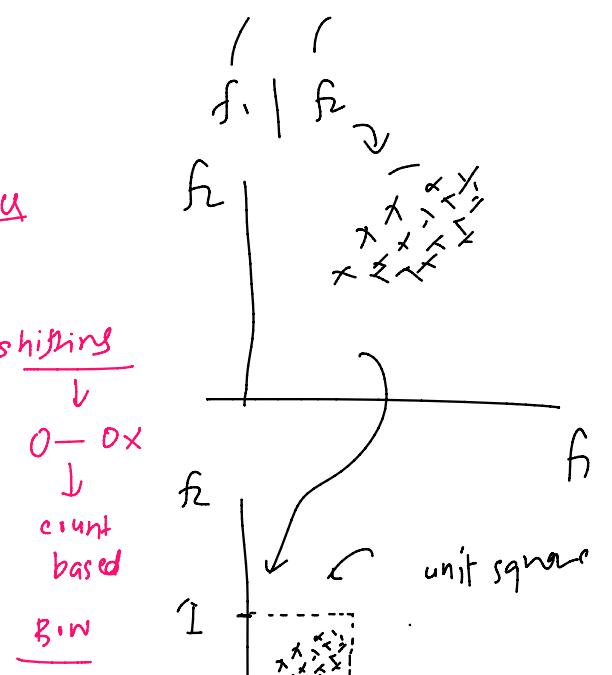
2. Relative Distances Preserved: The relative distances between values in each feature are preserved. While the absolute values change to fit within the new range, the proportions relative to the minimum and maximum of the original data are maintained.

3. No Impact on Shape of the Data Distribution: Min-Max scaling does not alter the shape of the feature's distribution. If the original data is normally distributed, skewed, or has any other distribution, the scaled data will maintain that distribution, albeit within the new range.

4. Dependency on Min and Max Values: The scaled values are highly sensitive to the minimum and maximum values in the data. If there are outliers, they will compress the rest of the data into a very small portion of the range, potentially reducing the effectiveness of the scaling.

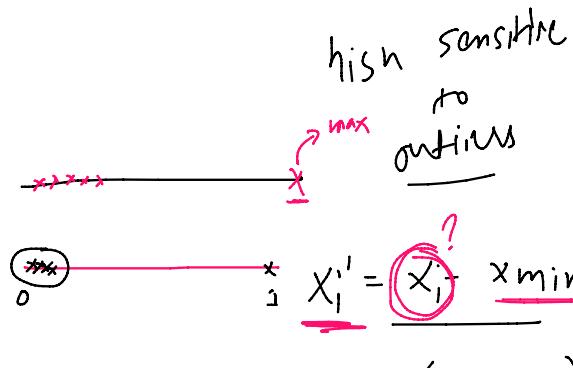
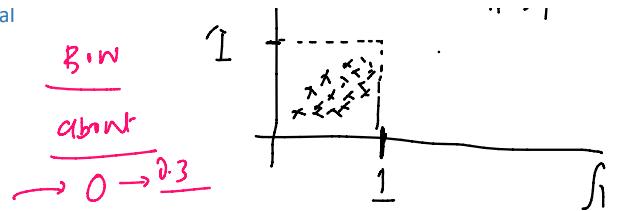
5. Zero Values Transformation: If the original data contains zero values, they may no longer be zero after scaling unless the minimum value in the original data is 0. The transformation shifts all values proportionally between the new minimum and maximum.

6. Impact on Algorithms: Min-Max scaling can significantly impact the performance of machine learning algorithms that are sensitive to the magnitude of the input values, such as gradient descent-based algorithms, nearest neighbors, and neural networks. The scaling ensures that features contribute equally to the model's learning process, preventing features with larger scales from dominating the learning.

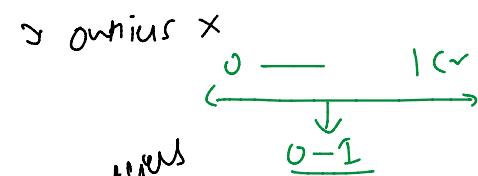


values, such as gradient descent-based algorithms, nearest neighbors, and neural networks. The scaling ensures that features contribute equally to the model's learning process, preventing features with larger scales from dominating the learning.

7. Reversibility: The Min-Max scaling process is reversible. Given the scaled value, along with the original minimum and maximum, you can calculate the original value. This is important for interpreting model outputs in the original scale.



outliers
more less



$Z - \text{standard}$

$$x_i = \frac{x_i - \mu}{\sigma}$$

$x_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$

activation \rightarrow $0-1$

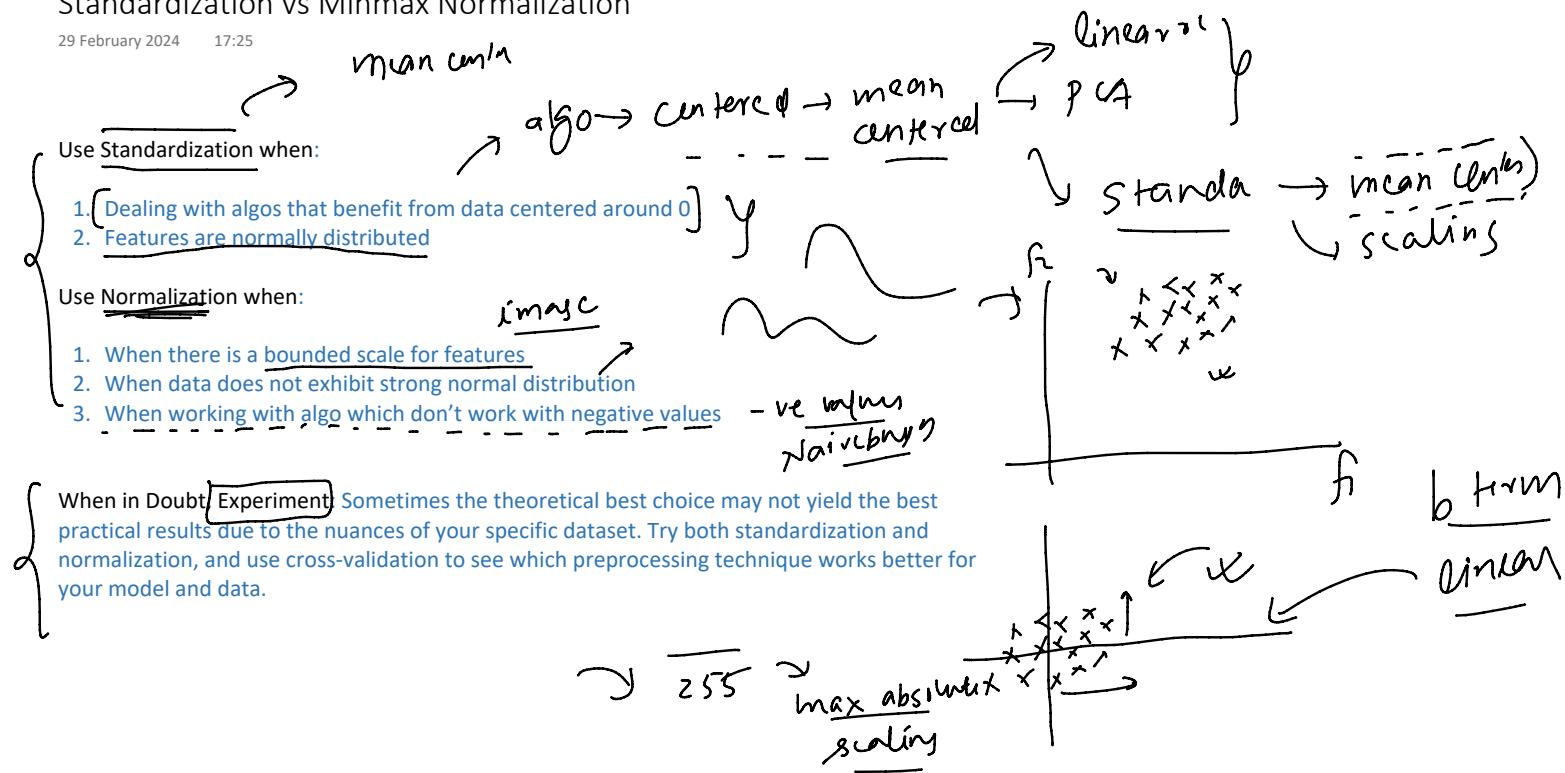
\rightarrow bounded image

scaling $f_1 f_2 \dots$
 $f_1 \rightarrow$ did normal
normal-like

$$\frac{0-255}{0-1}$$

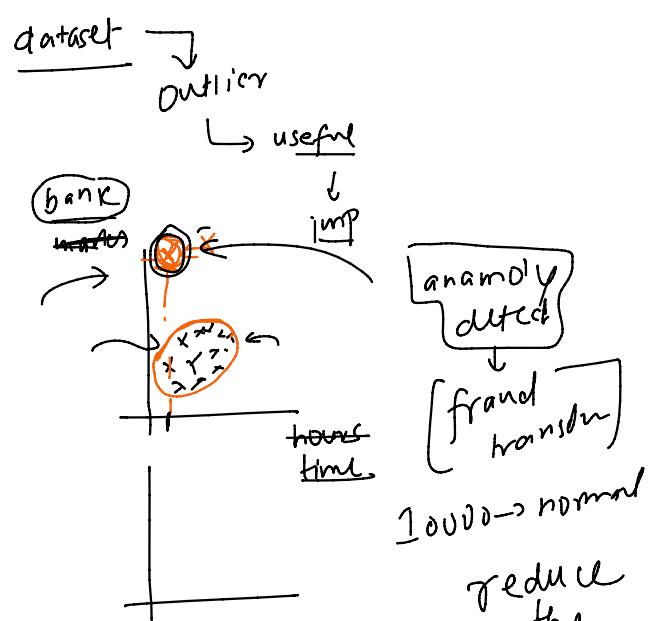
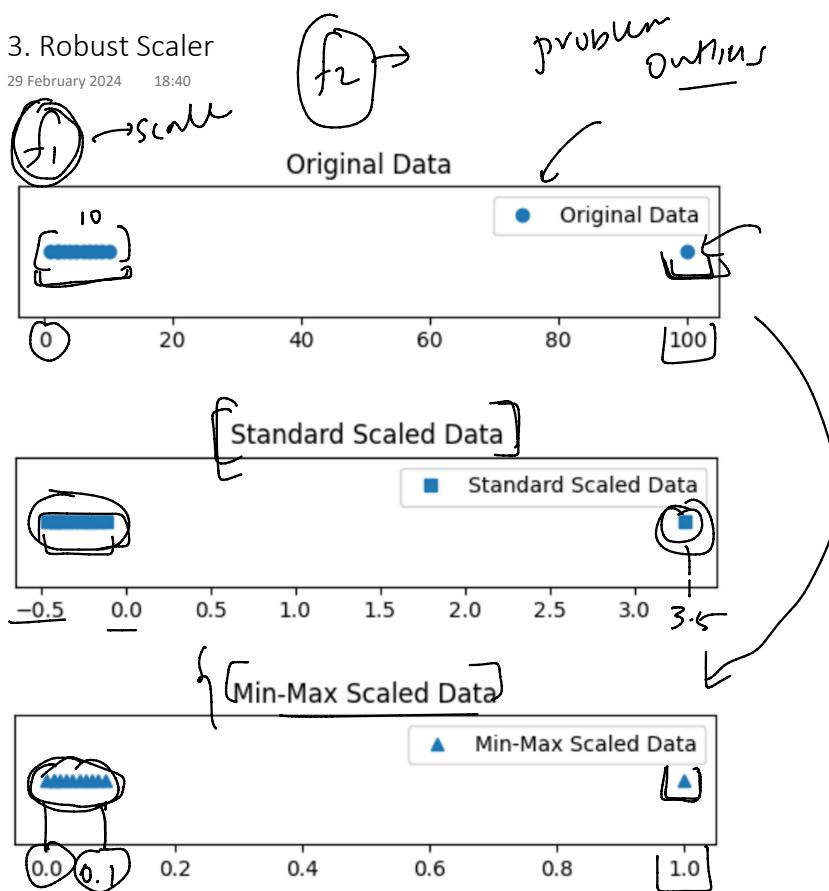
Standardization vs Minmax Normalization

29 February 2024 17:25



3. Robust Scaler

29 February 2024 18:40



Robust scale

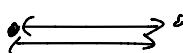
→ Robust to outliers

Robust → Stand

$x'_i = \frac{x_i - \text{mean}}{\text{std}}$

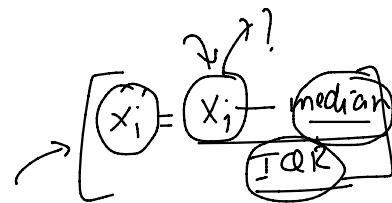
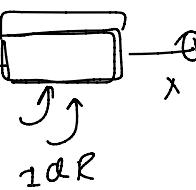
$$x'_i = \frac{x_i - \text{mean}}{\text{std}}$$

middle



Robust

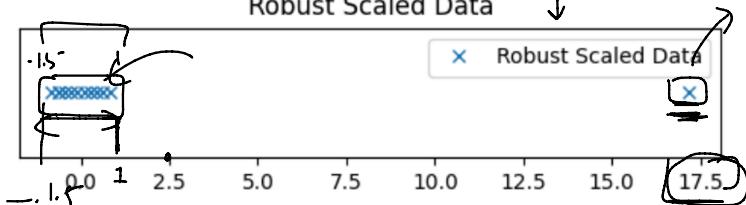
→ reduce outliers



25pm

75pe - 25m

Robust Scaled Data



- Median Centering:** The median of each feature will be shifted to zero. Since the robust scaler subtracts the median from each data point, the center of the scaled data for any feature will be its median, aligning it at 0 in the scaled dataset.
- IQR Scaling:** The IQR (the range between the 25th and 75th percentiles) becomes the scaling reference, set to 1. This means that the difference between the first quartile (Q1) and the third quartile (Q3) in the scaled data will be 1. Consequently, data points that were at Q1 and Q3 in the original data will be at -1 and 1, respectively, in the scaled data.
- Outlier Influence Minimization:** The scaling is less sensitive to outliers compared to methods like standardization or min-max scaling. Since the median and IQR are more robust against outliers, extreme values have less influence on the scaling process, ensuring that the core data is not skewed by anomalous points.
- Data Distribution Shape:** The overall shape of the data distribution is maintained. While the scale and location of the data change, the relative positioning of data points within the distribution remains constant.
- No Fixed Range for Scaled Data:** Unlike min-max scaling, the robust scaler does not confine the data to a fixed range. Scaled values can be greater than 1 or less than -1, especially for data points that lie outside the interquartile range.
- Impact on Learning Algorithms:** Scaling data with a robust scaler can improve the performance of machine learning algorithms sensitive to outliers. By ensuring that the feature's core data is effectively normalized, models can focus on the more representative part of the data.
- Reversibility:** The transformation is reversible, provided that you retain the original median and IQR values. Knowing these parameters allows you to reconstruct the original data values from the scaled data.

→ median
center



4. Max Absolute Scaler

29 February 2024 18:40

max_0_wit

$$\left\{ \begin{array}{c} \text{age} \\ 27 \\ 31 \\ 42 \end{array} \right. \rightarrow \frac{27}{42} =$$

The formula for max absolute scaling a value X_i in a feature is:

$$X'_i = \frac{X_i - \min(X)}{\max(X)} \quad \text{scaling}$$

$y \in [-1, +1]$

Range: The scaled data will be within the range $[-1, 1]$. If the original data contains both positive and negative values, the scaled data will span this entire range. If the original data is all non-negative or non-positive, the scaled data will lie within $[0, 1]$ or $[-1, 0]$, respectively.

Preservation of Zero: If a value is 0 in the original data, it remains 0 in the scaled data, as scaling does not introduce a shift.

Outlier Sensitivity: Like min-max scaling, max absolute scaling is sensitive to outliers since extreme values will dictate the scaling factor for the entire feature.

Why \rightarrow

$$\left\{ \begin{array}{c} \text{BOW} \rightarrow \text{Sentiment} \\ \text{---} \\ \text{---} \end{array} \right. = \frac{\text{BOW}}{\text{BOW}} =$$

1. Sign Preservation: The scaling maintains the sign of each value. Positive values remain positive, and negative values stay negative, ensuring that the relative relationships in terms of direction (positive or negative) are preserved.
2. Zero Values: Zero values in the data remain unchanged. This characteristic is crucial for datasets where the zero point carries specific meaning or indicates the absence of a feature.
3. Relative Distances: The relative distances between values in each feature are preserved. While the scale changes, the proportional relationships between data points remain consistent.
4. Outlier Sensitivity: The scaling process is sensitive to outliers. Extreme values can dominate the scaling factor, affecting how the rest of the data is scaled.
5. No Centering: The data is not centered around the mean or median; only the scale changes. The center of the distribution remains the same as in the original data.
6. Dimensional Consistency: Each feature is scaled independently, ensuring that the scaling does not distort the relationships between different features in multivariate data.
7. Sparse Data Suitability: The method is particularly suitable for sparse data, as it

$\xrightarrow{\text{pos/neg}}$ $-1, 1$
 $\xrightarrow{\text{pos}}$ $0, 1$
 $\xrightarrow{\text{neg}}$ $-1, 0$

$\boxed{\text{sparse data}}$

BOW \downarrow scaling

How | about | movie | . . . | sentiment
 1 | | 3 | . . . | 0
 2 | | 2 | . . . | !

Standar \downarrow min \downarrow robust \downarrow 0 sparse \downarrow predictor
 \downarrow \downarrow \downarrow \downarrow
 ∞ ∞ ∞ -1

sparse $\rightarrow 0, 0, 0 \rightarrow 0 \text{ has some meaning}$

max_abs_Scaler

$$\left\{ \begin{array}{c} 0 \rightarrow \text{blank} \\ 255 \rightarrow 255 \\ 213 \rightarrow 213 \end{array} \right.$$

$0 \rightarrow 255$
 $b \rightarrow b$
 $1 \rightarrow 0$
 $0-4 \rightarrow 1$

Scaling does not distort the relationships between different features in multivariate data.

7. Sparse Data Suitability: The method is particularly suitable for sparse data, as it does not alter zero entries, maintaining the sparsity pattern of the dataset.

8. Scale Interpretation: Post-scaling, the interpretation of a unit change in the data is consistent across features, facilitating comparisons and analyses that rely on the magnitude of values.

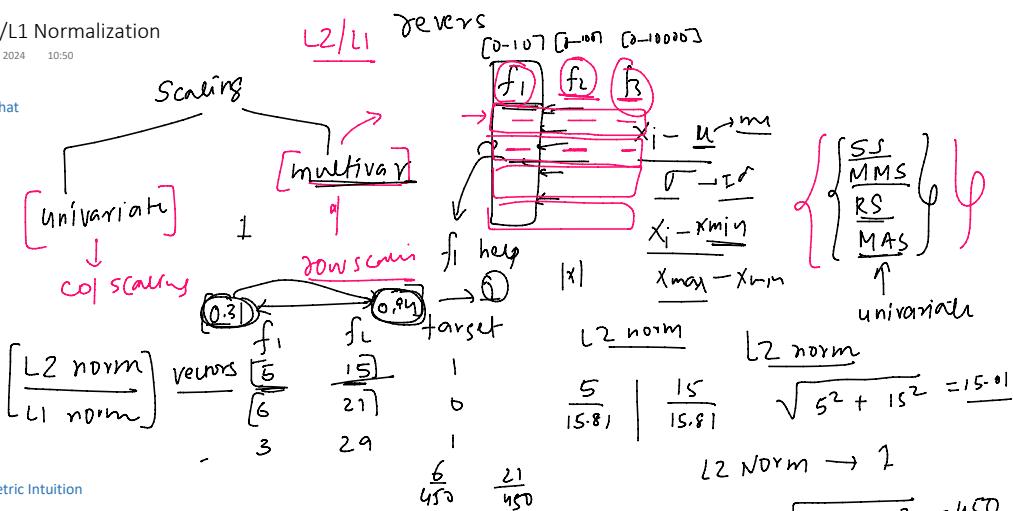
$X \rightarrow \text{Scaling}$

9. Reversibility: The scaling transformation is reversible. Knowing the original maximum absolute value allows you to rescale the data back to its original range, preserving the integrity of the data for inverse transformations if needed.

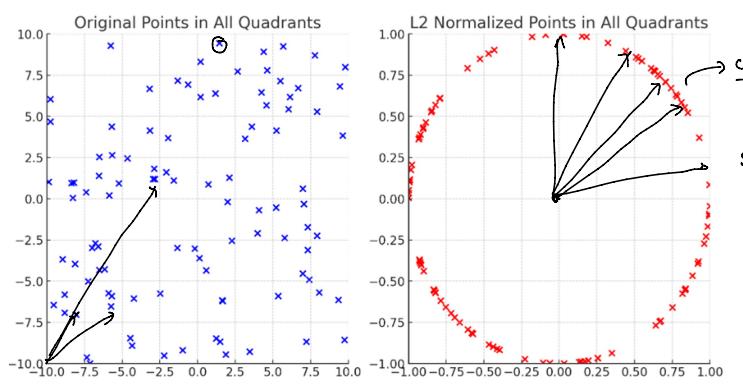
5. L2/L1 Normalization

05 March 2024 10:50

The What



Geometric Intuition



$$L2 \text{ norm} = 1$$

$$\text{origin} = 2$$

L1 Normal

$$\begin{matrix} f_1 \\ 5 \\ 6 \\ 3 \end{matrix}$$

$$\begin{matrix} f_2 \\ 15 \\ 21 \\ 29 \end{matrix}$$

L2 Norm → L1 Norm

$$|5| + |15| = 20$$

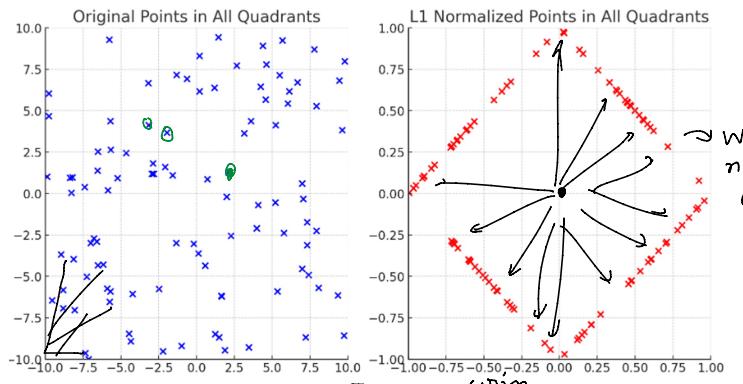
$$\boxed{0.25} \quad \boxed{\frac{5}{20}} \quad \boxed{\frac{15}{20}} \quad \boxed{0.75}$$

$$L1 \text{ norm} = 1$$

$$\begin{matrix} 6 \\ 27 \\ 0.23 \end{matrix}$$

$$\begin{matrix} 21 \\ 27 \\ 0.77 \end{matrix}$$

$$|6| + |21| = 27$$



Kim

h

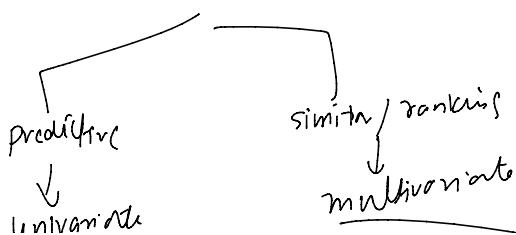
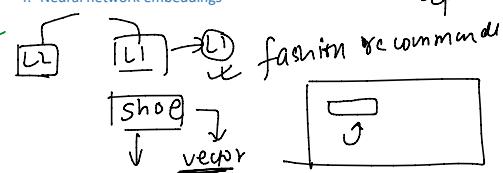
256 norm

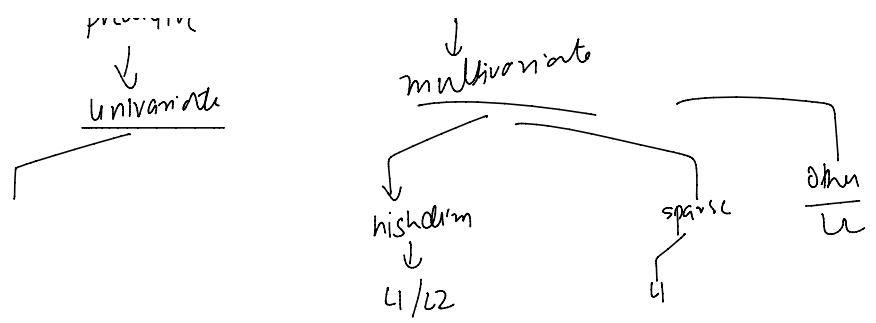


1000 D
→ vectors

Use-case

- 1. Text Preprocessing - sparse features → L1
- 2. Similarity based algos [cosine similarity]
- 3. Image processing based applications
- 4. Neural network embeddings





Comparison

02 March 2024 17:35

Universität

Parameter/Criteria	Standardization	Min-Max Scaling	Robust Scaling	Max Absolute Scaling
<u>Formula</u>	$\frac{X-\mu}{\sigma}$	$\frac{X-\text{Min}}{\text{Max}-\text{Min}}$	$\frac{X-\text{Median}}{\text{IQR}}$	$(\frac{X}{\text{Max}} - \frac{\text{Min}}{\text{Max}})$
<u>Centering</u>	Centers data around <u>mean (0)</u>	No centering; shifts data to start at 0	Centers data around the <u>median</u>	No centering; retains original data center
<u>Scaling</u>	Scales data by standard deviation	Scales data to a specified range, typically [0, 1]	Scales data by the interquartile range	Scales data relative to the maximum absolute value
<u>Sensitivity to Outliers</u>	Yes, outliers can significantly affect mean and standard deviation	Yes, outliers affect the minimum and maximum values ↓	No, uses median and IQR which are less influenced by outliers	Yes, the maximum absolute value determines the scaling factor

<u>Range of Scaled Data</u>	No fixed range; typically within [-3, 3] for data close to normal distribution	Fixed and known range, typically [0, 1] or [-1, 1]	No fixed range; values can vary widely based on the IQR	Typically [-1, 1], but not strictly bounded if there are outliers
<u>Impact on Distribution Shape</u>	Maintains the shape of the distribution	Maintains the shape of the distribution	Maintains the shape of the distribution	Maintains the shape of the distribution
<u>Suitability for Sparse Data</u>	Not particularly suitable as it shifts zero values	Not suitable as it changes zero values unless the range starts at 0	Not particularly suitable as it can shift zero values	Suitable, especially when zero values are meaningful
<u>Preservation of Zero</u>	Zero values are shifted unless the mean is 0	Zero values are shifted unless the minimum value is 0	Zero values are shifted unless the median is 0	Zero values are preserved

<u>Common Use Cases</u>	Data with a Gaussian distribution, linear models, clustering, PCA	Neural networks, image processing, when data needs to be in a [0, 1] range	Data with outliers, robust machine learning models, financial data	Data where maximum value is crucial, text processing, sparse data
<u>Interpretability</u>	Interpretation in terms of how many standard deviations from the mean	Direct interpretability due to fixed range	Interpretation in terms of median and IQR	Interpretation in terms of a proportion of the maximum value
<u>Impact on Feature Importance</u>	Equalizes features' importance based on variance	Equalizes features' importance based on their range	Equalizes features' importance based on their spread (IQR)	Normalizes features' influence based on their maximum absolute value

<u>Algorithmic Suitability</u>	Algorithms assuming features with Gaussian distribution, linear models	Algorithms sensitive to feature scale, neural networks	Algorithms robust to outliers, median-based clustering	Algorithms where magnitude and sign are crucial, preserving data structure
<u>Invariance to Transformations</u> X	Not invariant to multiplicative transformations of input	Invariant to linear transformations of input	Not invariant to multiplicative transformations of input	Invariant to scaling transformations of input
<u>Mathematical Properties</u> X	Linear transformation	Linear transformation	Non-linear transformation (if outliers are present)	Linear transformation
<u>Parameter Dependency</u>	Depends on <u>mean</u> and <u>standard deviation</u>	Depends on <u>minimum</u> and <u>maximum</u> values	Depends on <u>median</u> and <u>IQR</u>	Depends on the <u>maximum absolute value</u>

<u>Effect on Data Ordering</u>	Maintains the order of data	Maintains the order of data	Maintains the order of data	Maintains the order of data
<u>Typical Value Range Post-Scaling</u>	Mostly within [-3, 3] but can have outliers	[0, 1] or a user-specified range like [-1, 1]	Varies widely; central data points fall between [-1, 1]	[-1, 1], but extremes can exceed this range if outliers are present
<u>Ease of Reversibility</u>	Easy to reverse with knowledge of original mean and std	Easy to reverse with original min and max values	Reversible with original median and IQR	Easy to reverse with the original max absolute value

multivariate

Parameter/Criteria	L2 Normalization	L1 Normalization
Formula	$\frac{X}{\ X\ _2}$ where $\ X\ _2 = \sqrt{\sum x_i^2}$	$\frac{X}{\ X\ _1}$ where $\ X\ _1 = \sum x_i $
Centering	No centering; does not alter the mean of the data	No centering; does not alter the mean of the data
Scaling	Scales the feature vector so its Euclidean norm (L2 norm) is 1	Scales the feature vector so its Manhattan norm (L1 norm) is 1
Sensitivity to Outliers	Moderately sensitive; outliers affect the norm, thus the scaling	Less sensitive; since it sums absolute values, the impact of outliers is diluted
Range of Scaled Data	No fixed range; the vector's length is 1, but individual elements can vary	No fixed range; elements are scaled relative to the sum of absolute values
Impact on Distribution Shape	Changes the magnitude but not the direction of data points in feature space	Changes the magnitude but not the direction of data points in feature space
Suitability for Sparse Data	Suitable; does not alter zero entries	Suitable; does not alter zero entries

Preservation of Zero	Zero values remain zero	Zero values remain zero
Common Use Cases	Text classification, clustering, image processing where direction matters	Text processing, creating sparse representations, feature selection
Interpretability	Normalizes data points to unit length; useful in cosine similarity	Scales data points by their aggregate absolute magnitude; useful in sparsity
Impact on Feature Importance	Normalizes overall influence of a feature vector, not individual features	Balances the overall influence by the total magnitude, promoting sparsity
Algorithmic Suitability	Algorithms where the magnitude of vectors isn't as important as their direction	Algorithms benefiting from sparsity and robustness to outliers or scale
Invariance to Transformations	Not invariant to additive transformations, invariant to scalar multiplications	Not invariant to additive or scalar multiplicative transformations
Mathematical Properties	Non-linear transformation (based on Euclidean norm)	Non-linear transformation (based on Manhattan norm)
Parameter Dependency	Dependent on the entire feature vector's norm	Dependent on the entire feature vector's norm
Effect on Data Ordering	Does not alter the order of features but normalizes their magnitude	Does not alter the order of features but normalizes their magnitude
Typical Value Range Post-Scaling	Each data point is scaled such that the sum of squares is 1	Each data point is scaled such that the sum of absolute values is 1
Ease of Reversibility	Not directly reversible without the original vector's norm	Not directly reversible without the original vector's norm