

## Feature Scaling - Normalization

→ Goal of Normalization is to change values in numeric columns in the data set to use common scale, w/o distorting differences in ranges of values or losing information.

You can columns with different units of measurement Kg, cm, inch, meter, grams, °C, °F etc all together.

We ignore these measurement and try to bring <sup>all</sup> values under a common scale.

### Techniques for Normalization

- ① Min Max Scaling (used 90% of time)
- ② Mean Normalization
- ③ Max Abs Scaling
- ④ Robust Scaling,

### Min Max Scaling

for transforming value we use value,  $x_i$  value need to be transformed

$$x_i' = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}} \rightarrow \text{min value is variable/column}$$

$$\downarrow \quad [x_{\max}] - x_{\min}$$

Transformed  $\rightarrow$  max value of Data Value

Now after all values are transformed,  
the Max value of these transformed  
values now is  $["1"]$  and min value  
is  $["0"]$ .

Geometric Intuition

So min max scaling squeeze the values  
under 1 unit of square/rectangle in  
2D data.

1 unit of cube for 3D data. 1 unit of  
hyper cube for n-dimensional data.

Remember the distribution remains same  
but after transformation the values gets  
squeezed under 1 unit of measurement.

### Mean Normalization

Its kind of like standardization where  
data is centered towards mean

$$x_i' = \frac{x_i - \bar{x}_{\text{mean}}}{x_{\text{max}} - x_{\text{min}}}$$

Gives value in the range of  $-1 \leq 1$   
Extremely zero to use.

People prefer standardization over this  
cause scikit learn has no support of this

## 3 Max Absolute Scaling

used for sparse data  $\Rightarrow$  where we have lots of zeros.

$$x_i' = \frac{x_i^0}{|x_{\max}|} \rightarrow \text{Abs value}$$

## 4 Robust Scaling

$$x_i' = \frac{x_i^0 - \text{median}}{\text{IQR} \quad (75^{\text{th}} \text{percentile} - 25^{\text{th}} \text{percentile})}$$

### Robust Scalar Class

This thing is robust to Outliers because median " doesn't gets effected by outliers like in case of mean.

The Robust scalar can be utilized when we have lots of outliers in our data.

## 5 Normalization VS Standardization

90% we use this  $\leftarrow y$

i) Is Feature Scaling Required?

$\Rightarrow$  we need to understand what algo require feature scaling.

a) Most of time we get better results

with standardization.

2) Most Normalization (min max scale)  
is used where we before hand know  
the min & max value of data.  
Mostly used in CNN. (0-255) values.

Q Impact of Outlier on Min Max Scaling?

As value gets squeezed under 0 & 1,  
the outlier value also now gets  
reduced under 1 & between 0.  
Now this can effect the algorithm.

2)  
But to avoid this we can use robust  
Scaling if we have lots of outliers in  
our data.

