

Definition

Dimensionality reduction, or dimension reduction, is the transformation of data from a high-dimensional space into a low-dimensional space so that the low-dimensional representation retains some meaningful properties of the original data, ideally close to its intrinsic dimension.

- 2D & 3D → Scatter Plot and other related plots
- 4D , 5D , & 6D → Pair Plots
- nD → Dimensionality Reduction
  - PCA (Principal Component Analysis)
  - t-SNE (t-distributed Stochastic Neighborhood Embedding)

```
In [ ]:
```

Row and Column Vector

- A row vector is a row of entires. It has 1 row and n columns.

a = [a\_1, a\_2, a\_3, ..., a\_n]

- A column vector is a column of entries. It has 1 column and n rows.

```
x = [[x1],
     [x2],
     [x3],
     ...,
     [xn]]
```

Note

- By default, when someone says a vector, it means that it is a column vector.
- The transpose of column vector is called a row vector.
- Please refer to [wiki](#) article.

```
In [ ]:
```

Dataset representation

A dataset is represented as  $D = \{x_i, y_i\}$  where  $X = x_i$  (independant variables or features) and  $Y = y_i$  (target variable or dependant variable).

```
For example
-----
```

Iris data ⇒

```
[ [PL],
  [PW],
  [SL],
  [SW] ]
```

which are features and

```
[species]
```

represents target variable.

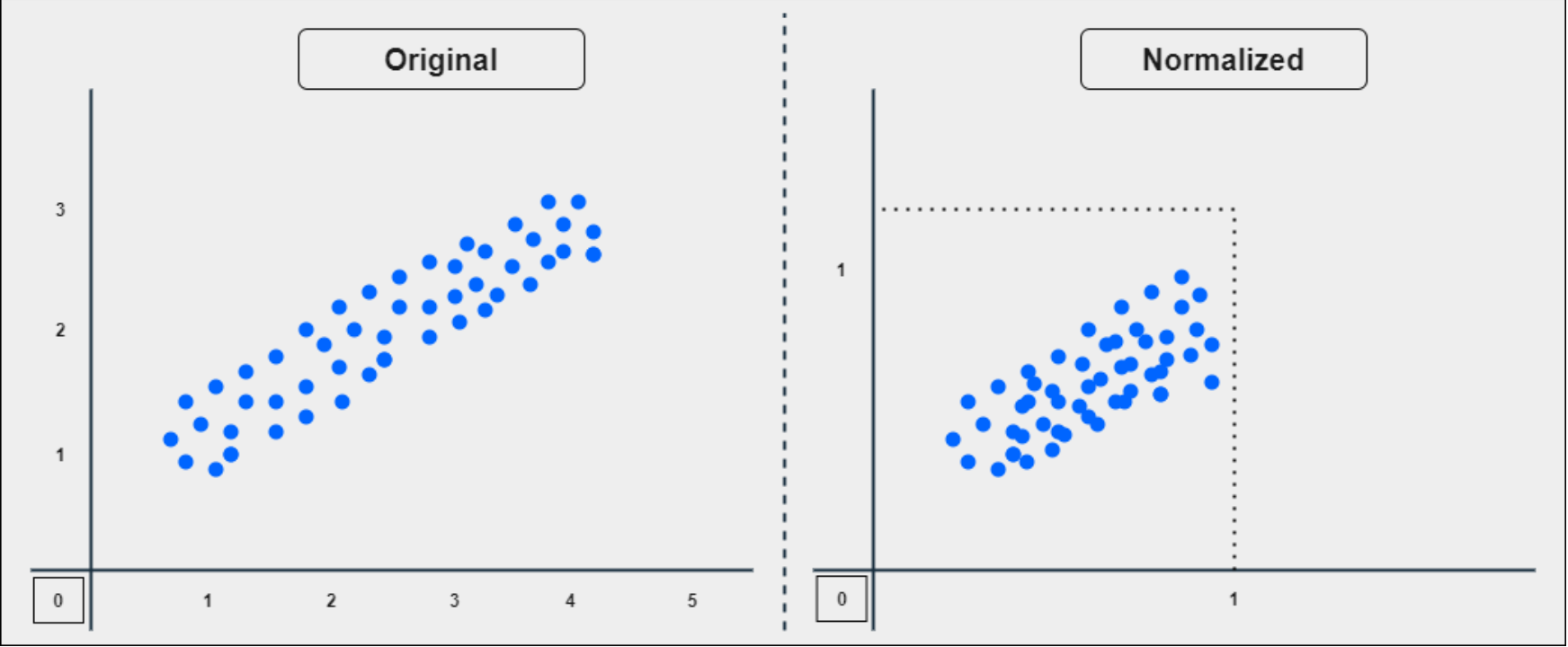
```
In [ ]:
```

Data Preprocessing

- Column Normalization
  - Consider each col from the dataset and find out col\_min and col\_max.
  - Compute

col\_i^1 = (col\_i - col\_min) / (col\_max - col\_min)

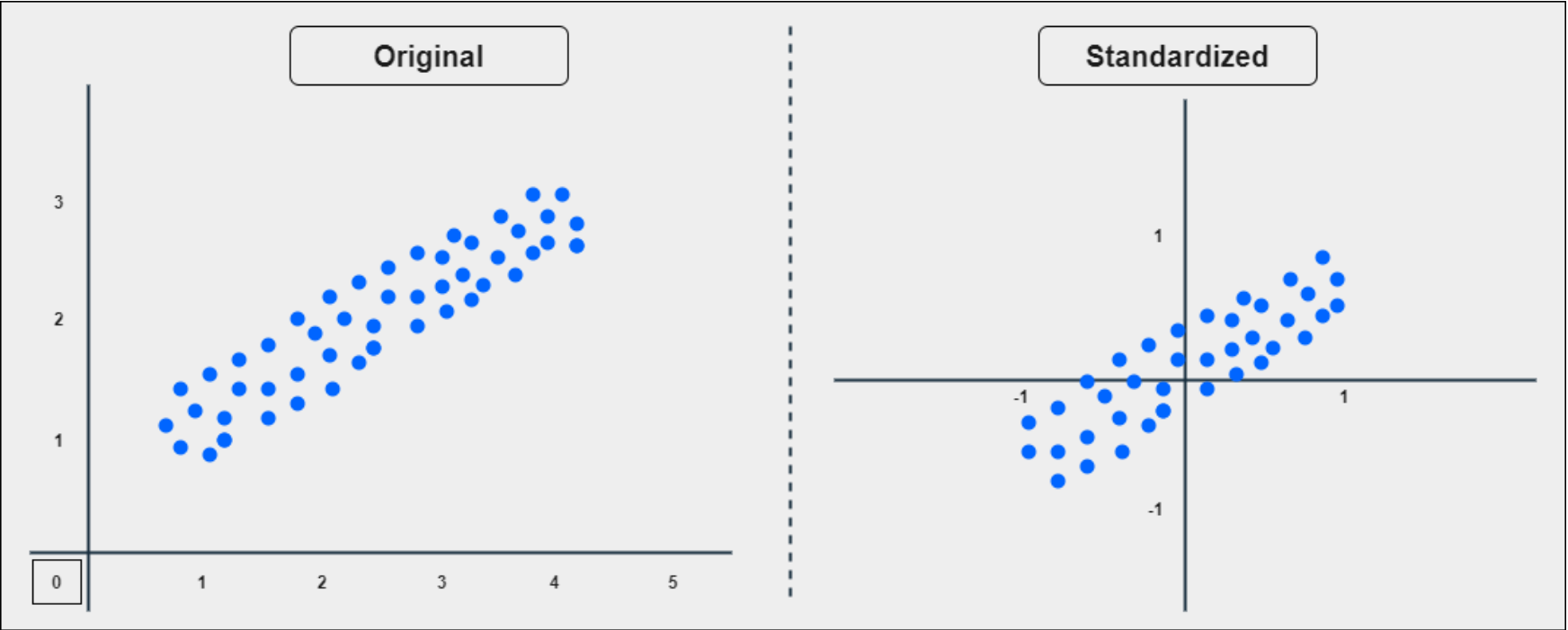
- All the value of col will be in the range of 0 and 1, col\_i^1 ∈ [0, 1]
- This helps to get rid off scale measurement.
- Squishes the data into one unit measurement.



- Column Standardization
  - Consider each col from the dataset and find out μ\_col and σ\_col.
  - Compute standard normal variate for the col such as

col\_z = (col\_i - μ\_col) / σ\_col

- The mean of col\_z is 0 and standard deviation is 1.
- It is also called as mean centering, i.e., mean is at origin and scaling is done by standard deviation (1).



```
In [ ]:
```

Co-Variance Data Matrix (Symmetric Matrix)

Let A be a matrix where A\_ij = A\_ji then A is known as symmetric matrix.

```
A = [[2, 1, 2],
     [1, 1, 5],
     [2, 5, 3]]
```

Co-Variance Data Matrix is a symmetric matrix.

- Cov(X, Y) = 1/n ∑(x\_i - μ\_x)(y\_i - μ\_y)
- Cov(X, X) = Var(X)
- Cov(X, Y) = Cov(Y, X)

Let f\_1 and f\_2 are two features which are column standardized . The Cov(f\_1, f\_2) is written as -

⇒ 1/(n-1) ∑ f\_1 f\_2

⇒ f\_1^T f\_2 / (n-1)

**Note** - We consider (n - 1) so as to make sure we get an unbiased estimator.

If X is a dataset irrespective of target variable. Assuming X is column standardized , we get covariance matrix as -

S\_ij = X^T X / (n-1) = f\_i^T f\_j / (n-1)

where

- f\_i and f\_j are features.

```
In [ ]:
```