

col. std. $\rightarrow \begin{cases} \rightarrow \text{moving the mean-vector to origin.} \\ \rightarrow \text{squishing/expanding s.t std-dev for any feature is 1.} \end{cases}$

col. stand. :- mean-centering \rightarrow origin

+ scaling $\rightarrow \text{std-dev}_{\text{for all feature}} = 1$

Co-Variance Matrix

$\xrightarrow{\text{col-vector } j^{\text{th}} \text{ feature}}$

$$X = \begin{bmatrix} x_1 & & & & & & \\ x_2 & & & & & & \\ x_3 & & & & & & \\ \vdots & & & & & & \\ x_i & \xleftarrow{x_i^T} & \begin{matrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_i \\ \vdots \\ a_n \end{matrix} & \xrightarrow{x_i^T} & & & \\ \vdots & & & & & & \\ x_n & & & & & & \end{bmatrix}_{n \times d} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

def: ; $S = \text{cov-matrix of } X$

$S_{ij} = \text{ele in } S$

$S_{ij} = \text{col } i^{\text{th}} \text{ row } j^{\text{th}}$

$S_{ij} = \text{feature for } i^{\text{th}} \text{ data-point}$

$S_{ij} = \text{cov}(f_i, f_j)$

$i: 1 \rightarrow d$

$j: 1 \rightarrow d$

\uparrow Square matrix.

$$S_{ij} = \text{cov}(f_i, f_j)$$

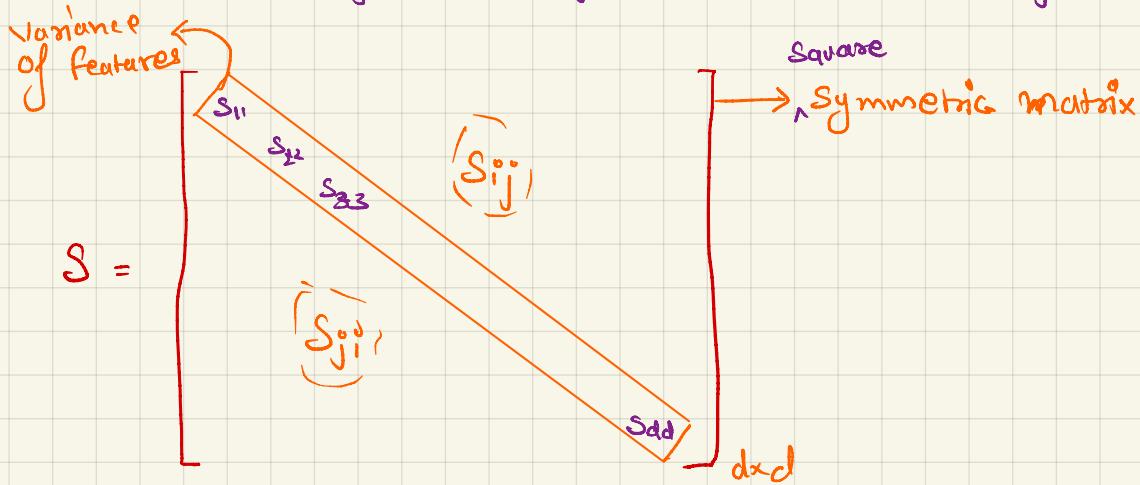
$$S_{ij} = \text{col } i^{\text{th}} \text{ row } j^{\text{th}}$$

$$S_{ij} = \text{cov}(f_i, f_j)$$

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{cov}(x, x) = \text{var}(x) \quad -\textcircled{1}$$

$$\text{cov}(f_i, f_j) = \text{cov}(f_j, f_i) \quad -\textcircled{2} \rightarrow S_{ij} = S_{ji}$$



$$X = \begin{bmatrix} 1 & f_1 & f_2 & f_3 & \dots & f_d \\ 2 & x_{11} & & & & \\ 3 & x_{21} & & & & \\ \vdots & & & & & \\ n & x_{n1} & & & & \end{bmatrix}_{n \times d}$$

Let, X : col. standardization $\Rightarrow \text{mean}\{f_i\} = 0$
 $\Rightarrow \text{std-dev}\{f_i\} = 1$

$\rightarrow \text{mean}(f_i)$

$$\text{cov}(f_1, f_2) = \frac{1}{n} \sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) \quad \rightarrow \text{mean}(f_i)$$

$$= \frac{1}{n} \sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)$$

$$\text{cov}(f_1, f_2) = \underbrace{(f_1^T f_2)}_{f_1 \cdot f_2} * \frac{1}{n}$$

$$X = \begin{bmatrix} 1 & f_1 & f_2 & f_3 & \dots & f_d \\ 2 & \checkmark & \checkmark & & & \\ 3 & \checkmark & \checkmark & & & \\ \vdots & & & & & \\ n & x_{11} & x_{12} & & & \end{bmatrix}_{n \times d}$$

If f_1, f_2 have been standardized,

$$\text{cov}(f_1, f_2) = \frac{f_1^T f_2}{n}$$

conclusion

$$S_{d \times d} = \frac{1}{n} (X^T X)_{d \times n} = ()_{d \times d} \rightarrow \text{here I am assuming that 'X' has been col. standardized.} \rightarrow \text{data-matrix}$$

$$S_{ij} = \text{cov}(f_i, f_j) = \frac{\mathbf{f}_i^T \mathbf{f}_j}{n}$$

$$(i, j) = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_n \\ f_0 \\ f_1 \\ \vdots \\ f_d \end{bmatrix} \left[\begin{array}{|c|c|c|c|} \hline \bullet & \checkmark & * & \dots \\ \hline \end{array} \right] \star \left[\begin{array}{|c|c|c|c|} \hline \bullet & \checkmark & * & \dots \\ \hline \end{array} \right]$$

$$S_{ij} = \frac{f_i^T f_j}{n}$$

$$S_{d \times d} = X_{d \times n}^T X_{n \times d} \quad \text{if } x \text{ has been col. std.}$$

MNIST dataset (784 features)

Let, \rightarrow 60k training data points

10k Test datapoints.

$$\text{data points} \leftarrow D = \{x_i, y_i\}_{i=1}^{60K}$$

$$y_i \in \{0, \cancel{1}, 2, 3, 4, \dots, 9\}$$

$$X^T \quad X$$

Biased Estimator

A biased estimator systematically overestimates or underestimates the parameter it is trying to estimate. In the context of covariance:

- Covariance Matrix (Biased Estimator)

$$\text{Cov}(X) = \frac{1}{n} X^T X$$

where n is the number of samples

- The biased estimator uses n in the denominator. This makes the calculation slightly more efficient but introduces bias, meaning the estimate might not be centered around the true population parameter, especially for small sample sizes.

Unbiased Estimator

An unbiased estimator, on average, hits the true parameter it estimates, meaning the expected value of the estimator is equal to the true value of the parameter.

- Covariance Matrix (Unbiased Estimator):

$$\text{Cov}(X) = \frac{1}{n-1} X^T X \longrightarrow \text{num} \beta$$

where n is the number of samples.
 - The unbiased estimator uses $n - 1$ in the denominator. This adjustment corrects the bias that occurs when estimating the population parameter from a sample. It ensures that the average of the

by default uses this.

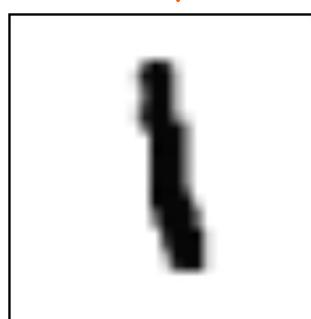
Obj :- Given a new image, determine whether a user has written {0,1,2...9}.
OR
Classify the written character into one of the '10' numeric characters.

$$x_i : [\quad]$$

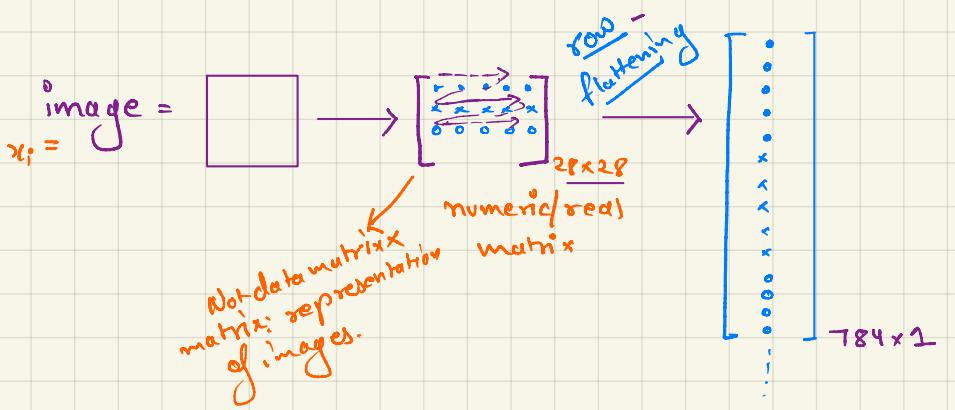
→ How do I convert a [°] image into it?

Every MNIST data point, every image, can be thought of as an array of numbers describing how dark each pixel is. For example, we might think of  as something like:

example, we might think



Since each image has 28 by 28 pixels, we get a 28x28 array. We can flatten each array into a $28 * 28 = 784$ dimensional vector. Each component of the vector is a value between zero and one describing the intensity of the pixel. Thus, we generally think of MNIST as being a collection of 784-dimensional vectors.



Overall dataset



Code to LOAD Mwist

CSV

comma separate value

f_1 f_2 f_3 f_4 → header &
y: 2, 2.1, 2.3, 2.4

2

63:

printing
first five

rows.

o index

```

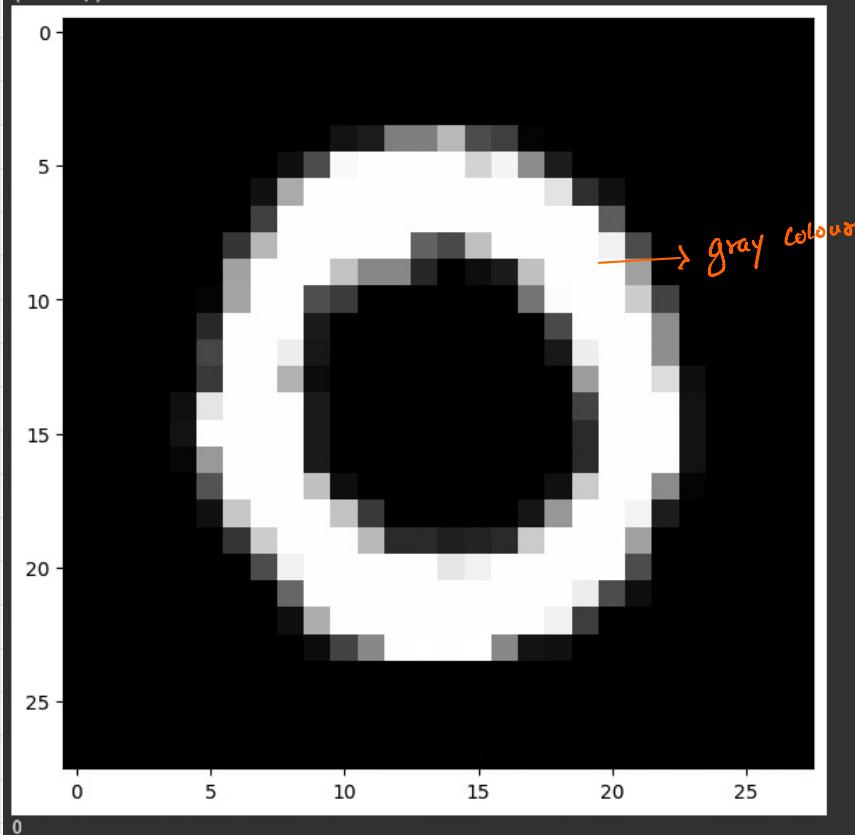
print(d.shape) → (42000, 784)
print(l.shape) → (42000,)

# display or plot a number.
plt.figure(figsize=(7,7))
idx = 1

grid_data = d.iloc[idx].to_numpy().reshape(28,28) # reshape from 1d to 2d pixel array
plt.imshow(grid_data, interpolation = "none", cmap = "gray")
plt.show() ↳ image show
print(l[idx])

(42000, 784)
(42000,)

```



get the data of '1' row &
convert it into '28x28' matrix.

→ Series.to_numpy()

↓

it is used to convert
'Series' to numpy array.
or
index

head() in pandas

DataFrame.head(n)

`DataFrame.head(n=5)` [source]

Return the first *n* rows.

This function returns the first *n* rows for the object based on position. It is useful for quickly testing if your object has the right type of data in it.

For negative values of *n*, this function returns all rows except the last $|n|$ rows, equivalent to `df[:-n]`.

If *n* is larger than the number of rows, this function returns all rows.

Parameters: `n : int, default 5`
Number of rows to select.

Returns: `same type as caller`
The first *n* rows of the caller object.

See also
`DataFrame.tail`
Returns the last *n* rows.

Examples

```
>>> df = pd.DataFrame({'animal': ['alligator', 'bee', 'falcon', 'lion',
...                                'monkey', 'parrot', 'shark', 'whale', 'zebra']})
... df
   animal
0  alligator
1      bee
2    falcon
3     lion
4    monkey
5    parrot
6     shark
7     whale
8     zebra
```

Viewing the first 5 lines

```
>>> df.head()
   animal
0  alligator
1      bee
2    falcon
3     lion
4    monkey
```

Viewing the first *n* lines (three in this case)

```
>>> df.head(3)
   animal
0  alligator
1      bee
2    falcon
```

For negative values of *n*

```
>>> df.head(-3)
   animal
0  alligator
1      bee
2    falcon
3     lion
4    monkey
5    parrot
```

DataFrame.drop()

pandas.DataFrame.drop

```
DataFrame.drop(labels=None, *, axis=0, index=None, columns=None,
level=None, inplace=False, errors='raise')  
[source]  
Drop specified labels from rows or columns.  
  
Remove rows or columns by specifying label names and corresponding axis, or by specifying  
directly index or column names. When using a multi-index, labels on different levels can be  
removed by specifying the level. See the user guide for more information about the now  
unused levels.  
  
Parameters: labels : single label or list-like  
Index or column labels to drop. A tuple will be used as a single label and  
not treated as a list-like.  
  
axis : {0 or 'index'; 1 or 'columns'}, default 0  
Whether to drop labels from the index (0 or 'index') or columns (1 or  
'columns').  
  
index : single label or list-like  
Alternative to specifying axis ( labels, axis=0 is equivalent to  
index=labels ).  
  
columns : single label or list-like  
Alternative to specifying axis ( labels, axis=1 is equivalent to  
columns=labels ).  
  
level : int or level name, optional  
For MultiIndex, level from which the labels will be removed.  
  
inplace : bool, default False  
If False, return a copy. Otherwise, do operation inplace and return None.  
  
errors : {'ignore', 'raise'}, default 'raise'  
If 'ignore', suppress error and only existing labels are dropped.  
  
Returns: DataFrame or None  
DataFrame without the removed index or column labels or None if  
inplace=True .
```

Examples

```
>>> df = pd.DataFrame(np.arange(12).reshape(3, 4),  
...                   columns=['A', 'B', 'C', 'D'])  
>>> df  
   A  B  C  D  
0  0  1  2  3  
1  4  5  6  7  
2  8  9  10 11  
  
Drop columns  
  
>>> df.drop(['B', 'C'], axis=1)  
   A  D  
0  0  3  
1  4  7  
2  8  11  
  
>>> df.drop(columns=['B', 'C'])  
   A  D  
0  0  3  
1  4  7  
2  8  11  
  
Drop a row by index  
  
>>> df.drop([0, 1])  
   A  B  C  D  
2  8  9  10 11
```

Drop columns and/or rows of MultiIndex DataFrame

```
>>> midx = pd.MultiIndex(levels=[[['lama', 'cow', 'falcon'],  
...                                ['speed', 'weight', 'length']],  
...                            [[0, 0, 0, 1, 1, 2, 2, 2],  
...                             [0, 1, 2, 0, 1, 2, 0, 1, 2, 1]]],  
...                            codes=[[45, 30], [200, 100], [1.5, 1], [320, 250],  
...                                    [250, 150], [1.5, 0.8], [320, 250],  
...                                    [1, 0.8], [8.3, 0.2]])  
>>> df  
   big  small  
lama  speed  45.0  30.0  
      weight 200.0 100.0  
      length  1.5  1.0  
cow   speed  30.0  20.0  
      weight 250.0 150.0  
      length  1.5  0.8  
falcon speed 320.0 250.0  
      weight  1.0  0.8  
      length  0.3  0.2
```

Drop a specific index combination from the MultiIndex DataFrame, i.e., drop the combination
'falcon' and 'weight', which deletes only the corresponding row

```
>>> df.drop(index='falcon', 'weight'))  
   big  small  
lama  speed  45.0  30.0  
      weight 200.0 100.0  
      length  1.5  1.0  
cow   speed  30.0  20.0  
      weight 250.0 150.0  
      length  1.5  0.8  
falcon speed 320.0 250.0  
      length  0.3  0.2
```

```
>>> df.drop(index='cow', columns='small')  
   big  
lama  speed  45.0  
      weight 200.0  
      length  1.5  
falcon speed 320.0  
      weight  1.0  
      length  0.3
```