## Exercise Set 12

# 1 Exam-style Python questions

a) What output do you expect from the following code?

```
import numpy as np
myArray = np.array([3,2,7,5])
sortedArray = np.sort(myArray)
print(sortedArray[1])
```

b) In the next Jupyter cell, the code contines as follows. What output do you expect?

```
newArray = myArray+10
print(newArray)
```

c) In the next Jupyter cell, the code contines as follows. What output do you expect?

```
print(np.shape(myArray))
reshapedArray = np.reshape(myArray,(2,2))
print(np.shape(reshapedArray))
```

d) In the next Jupyter cell, the code contines as follows. What output do you expect and why?

```
print(myArray+np.array([1,2,3,4,5]))
```

- e) You would like to calculate the unbiased estimator for the population standard deviation from the sample/data contained in myArray using the function np.std. The documentation (as found by running the code np.std? or similarly the numpy.org webpage) is given in the last page. How would you execute the function (i.e. with which arguments?)?
- f) How would you find the minimum value of myArray (there are many possible solutions, any of them that works is fine)? (Write down the code but do not worry about precise syntax).
- g) Name (at least) one difference between a standard Python list and a numpy array (for example how they respond differently to a certain operation).

# 2 Proof of the variance of a sum [normal]

Prove that

$$var(X + Y) = var(X) + var(Y) + 2cov(X, Y)$$

You can use the fact that  $\operatorname{var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$  where  $\mathbb{E}$  denotes the expectation value and  $\mu_X \equiv \mathbb{E}(X)$ . You can also use the definition of the covariance  $\operatorname{cov}(X,Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}(XY) - \mu_X \mu_Y$ .

# 3 Skiing and Fondue [normal]

A large Swiss cheese company thinks about marketing its product on ski hills. The idea is that people who enjoy outdoor activities on the mountains in also enjoy cheese more, and thus one could sell Fondue to them better than to the general population. To test this, a marketing institute sampled 36 people (9 in each group) who ate/didn't eat fondue and went/didn't go skiing. Each person should rate how much they enjoyed their day on a scale between 1 and 10.

- a) First, we want to investigate the effects of the individual factors. Create a table of the means for all four possible conditions, including marginal/partial means, and the total mean.
- b) Create a two-factor ANOVA table, assuming the factors are independent. Can you say ( using a level of significance of  $\alpha=0.05$ ), that each factor (skiing, and eating fondue) has a statistically significant effect? How different are the two effects?
- c) Test if people who went skiing and ate fondue on the same day gained significantly more enjoyment from this combined action than just the sum of the individual actions. Use a level of significance of  $\alpha = 0.05$ .

No skiing, no fondue	no skiing, fondue	skiing, no fondue	skiing, fondue
8	4	3	8
6	8	6	9
5	6	4	6
4	6	9	8
1	4	8	9
3	9	8	6
4	8	5	5
3	5	10	9
4	6	4	9

# Documentation of the np.std function

As printed when running np.std? after the code above.

```
Signature:
np.std(
    a,
    axis=None,
    dtype=None,
    out=None,
    ddof=0,
    keepdims=<no value>,
    where=<no value>,
                 np.std(*args, **kwargs)
Call signature:
Type:
                 _ArrayFunctionDispatcher
                 <function std at 0x734d0c46b6a0>
String form:
                 /opt/jlab-env/lib/python3.12/site-packages/numpy/core/fromnumeric.py
File:
Compute the standard deviation along the specified axis.
```

Returns the standard deviation, a measure of the spread of a distribution, of the array elements. The standard deviation is computed for the flattened array by default, otherwise over the specified axis.

#### Parameters

-----

#### a : array\_like

Calculate the standard deviation of these values.

axis : None or int or tuple of ints, optional

Axis or axes along which the standard deviation is computed. The default is to compute the standard deviation of the flattened array.

.. versionadded:: 1.7.0

If this is a tuple of ints, a standard deviation is performed over multiple axes, instead of a single axis or all the axes as before.

dtype : dtype, optional

Type to use in computing the standard deviation. For arrays of integer type the default is float64, for arrays of float types it is the same as the array type.

out : ndarray, optional

Alternative output array in which to place the result. It must have the same shape as the expected output but the type (of the calculated values) will be cast if necessary.

ddof : int, optional

Means Delta Degrees of Freedom. The divisor used in calculations is ''N - ddof'', where ''N'' represents the number of elements. By default 'ddof' is zero.

keepdims : bool, optional

If this is set to True, the axes which are reduced are left in the result as dimensions with size one. With this option, the result will broadcast correctly against the input array.

If the default value is passed, then 'keepdims' will not be passed through to the 'std' method of sub-classes of 'ndarray', however any non-default value will be. If the sub-class' method does not implement 'keepdims' any exceptions will be raised.

where : array\_like of bool, optional

Elements to include in the standard deviation.

See '~numpy.ufunc.reduce' for details.

.. versionadded:: 1.20.0

### Returns

-----

standard\_deviation : ndarray, see dtype parameter above.

If 'out' is None, return a new array containing the standard deviation, otherwise return a reference to the output array.

See Also

\_\_\_\_\_

```
var, mean, nanmean, nanstd, nanvar
:ref:'ufuncs-output-type'
```

### Notes

----

The standard deviation is the square root of the average of the squared deviations from the mean, i.e., "std = sqrt(mean(x))", where "x = abs(a - a.mean())\*\*2".

The average squared deviation is typically calculated as "x.sum() / N", where "N = len(x)". If, however, 'ddof' is specified, the divisor "N - ddof" is used instead. In standard statistical practice, "ddof=1" provides an unbiased estimator of the variance of the infinite population. "ddof=0" provides a maximum likelihood estimate of the variance for normally distributed variables. The standard deviation computed in this function is the square root of the estimated variance, so even with "ddof=1", it will not be an unbiased estimate of the standard deviation per se.

Note that, for complex numbers, 'std' takes the absolute value before squaring, so that the result is always real and nonnegative.

For floating-point input, the \*std\* is computed using the same precision the input has. Depending on the input data, this can cause the results to be inaccurate, especially for float32 (see example below). Specifying a higher-accuracy accumulator using the 'dtype' keyword can alleviate this issue.

### Examples

```
-------
>>> a = np.array([[1, 2], [3, 4]])
>>> np.std(a)
1.1180339887498949 # may vary
>>> np.std(a, axis=0)
array([1., 1.])
>>> np.std(a, axis=1)
array([0.5, 0.5])
```

In single precision, std() can be inaccurate:

```
>>> a = np.zeros((2, 512*512), dtype=np.float32)
>>> a[0, :] = 1.0
>>> a[1, :] = 0.1
>>> np.std(a)
0.45000005
```

Computing the standard deviation in float64 is more accurate:

```
>>> np.std(a, dtype=np.float64)
0.44999999925494177 # may vary
```

Specifying a where argument:

```
>>> a = np.array([[14, 8, 11, 10], [7, 9, 10, 11], [10, 15, 5, 10]])
>>> np.std(a)
2.614064523559687 # may vary
>>> np.std(a, where=[[True], [True], [False]])
2.0
```