

## Module 3

### Variance (Part 1)

(Describes the spread of data around the mean value)

Let's look at 2 data sets  $D_1$  and  $D_2$

$D_1$  is represented by (blue dots): 1, 3, 4, 5

$D_2$  is represented by (red squares): -1, 3, 7

(see pgs.

~~1 and 2~~

$D_1$  and  $D_2$  has the same mean, which is 3. Set the data points in  $D_2$  are less concentrated around the mean than the data points in  $D_1$ .

Remember, the mean value is the data that you expect on average, but to describe the concentration of data around the mean, <sup>value</sup> we can use the Concept of Variance

The variance is used to characterize the variability or spread of data points in dataset.

In 1D we can look at the average square distance of datapoint from mean value of the dataset.

Let's do it for  $D_1$  and  $D_2$

$$D_1 = \{1, 2, 4, 5\} \quad E[D_1] = 3$$

and expected mean value was 3

$$D_2 = \{-1, 3, 7\} \quad E[D_2] = 3$$

Now we want to compute the average square distance (from ~~the~~) of  $D_1$  from mean and from  $D_2$  from the same mean.

Let's do it for  $D_1$  first:

$$D_1: \frac{(1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2}{4}$$

$$= \frac{4+1+1+4}{4} = \frac{10}{4} \quad (B)$$

$$D_2: \frac{(-3)^2 + (3-3)^2 + (7-3)^2}{3} = \frac{16+0+16}{3} = \frac{32}{3} \text{ (A)}$$

- (A) is now bigger than (B), which means that the average square distance of  $D_2$  from mean value, is bigger than average square distance of  $D_1$  from mean value.

Which indicates that the spread of data is higher in  $D_2$  than in  $D_1$ .

So how can we formalize what we have done?  
we can define the average square distance as follows:

$\{x_1, \dots, x_N\} \Rightarrow$  define this as a dataset  $X$ .

$$\therefore \text{variance} : \text{Var}[x] = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2$$

where  $\mu$  is mean value of data set  $X$ .

$$\therefore \mu = E[X]$$



What we have done here is exactly the same as we did before with  $D_1$  and  $D_2$ , we computed an average square distance of data points in dataset from mean value ( $\mu$ ) of dataset.

Now, we can also make some statements about it.

First, the variance as defined here, can never be negative, as we just sum up square values that also means, we can take square root of variance, and this is called the standard deviation.

The standard deviation is expressed in the same units as the mean value, whereas the variance is unfortunately expressed in squared unit.

So comparing them is quite difficult.

Therefore when we talk about spread of data, we ⑤.  
usually look at the standard deviation

So, in this session, we looked at variances of  
1D data set, in next session, we will  
(generalize) generalize it to higher Dims.