

Table of Contents

Introduction

Aims and Overview	3
Defining Efficiency	3
Efficiency Notation	4
Formal Definitions	5

Cases

Mean	9
Variance	13
Standard Deviation	19

Conclusion

Evaluation of Success	21
Guidelines for Efficiency	21
Competing goals of	22
Equations	

Bibliography	23
--------------	----

Introduction

Aims and Overview

This Internal Assessment aims to find ways to make formulas provided by the International Baccalaureate more efficient in computation. I will be using big O notation to find the efficiency of statistical calculations (mean, variance, and standard deviation), then making modifications to the formulas to make them more efficient. First, the IA will begin by defining efficiency and efficiency notation. It will then progress through examples from the formula booklet. Finally, I will draw conclusions and general guiding principles to increase the efficiency of programs.

I chose this IA topic because it is related to computer science, a passion of mine. I hope to study computer science in university, and to this end I have been studying some coding. Big-O notation has been a subject that has interested me in these studies. Furthermore, this topic allows me to examine the thinking behind how math is presented to students, and how that may differ from the way math should be presented to computers. Finally, it also allows for an exercise in algebra when transforming equations into more efficient versions.

Defining Efficiency

Efficiency can be defined in different ways. It is possible to define efficiency in terms of the number of recalls to recorded data or the number of calculation steps. In this IA, I will be focusing on the number of calculation steps to quantify the complexity of an algorithm. To clarify, in this IA, “checking” if values are equal or finding values will not be a computing step. Computing steps are limited to an algebraic operation such as adding 1 to a value or squaring a value. This can be found in terms of the input size n to find the efficiency of the program.

Efficiency can be defined in terms of best, worst, and average cases, with best case meaning the shortest possible computation time for a given input, best case being the longest, and average case being the average computation time. Big O notation most commonly refers to the worst case because this is the limiting factor on the efficiency of a program. In programming it is often most useful to plan for the worst case (the upper

limit for computation time), to ensure that no matter the input one can guarantee that the program will finish in a certain time. In this IA, I may also refer to the best case and average case if those become deciding factors between different equations for which the worst case is equal. To be more efficient is to require fewer steps for each input size n . In other words, $f(n)$ is more efficient than $g(n)$ if:

$$O(f(n)) < O(g(n))$$

For this IA, each step that must be repeated n times in the worst-case scenario will be bolded. Each step that must be performed once in the worst-case scenario will be underlined. Lines of pseudocode that do not represent real actions will be left as plain text. For example, as defined above, actions like “check if...” and “move on” will not be considered computing actions.

Efficiency Notation

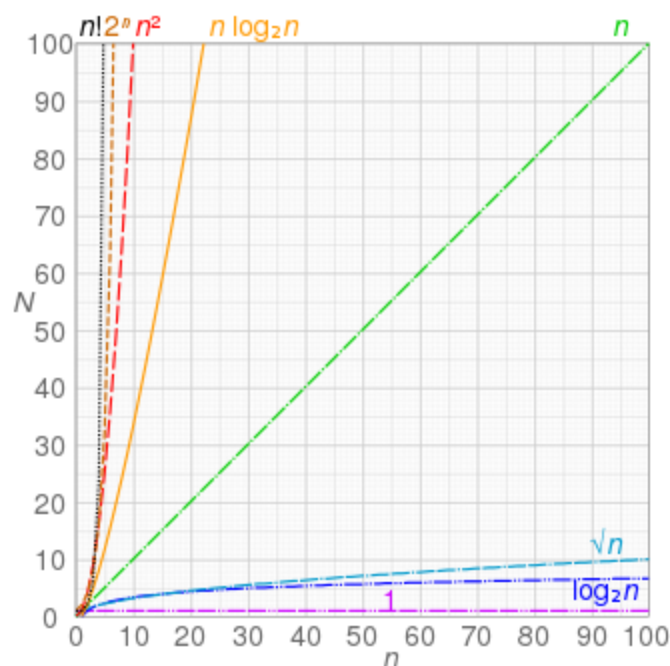
Efficiency notation can be used to understand and estimate the computing power or time required to find the solution to a formula in relation to the size of the input. This can be found in reference to the number of computing steps required in a process, number of calls to memory, or other statistics. More efficient programs grow in necessary computing power more slowly. The efficiency of these programs can be measured in big O notation such as $O(n)$. For example, if a program's calculation time doubles with a doubling of the input, it would be efficient according to $O(n)$. On the other hand, a program that has a processing speed which quadruples with a doubling of the input has efficiency $O(n^2)$. There are a few important rules to note when using big O notation:

1. Only the highest order term (the term that grows most quickly) is considered, because as the function grows in time, it will dominate the processing power required. As n approaches infinity, lower order terms come to matter less, because the highest order term grows the fastest. $O(n^2+7n+12)$ can be represented as $O(n^2)$.
2. Constants are dropped from big O notation because their effect is insignificant compared to different orders of operations. Furthermore, constants do not change the ratio of operations between different input sizes, and different

systems' computing times often vary linearly, so coefficients change between systems. $O(17n^2)$ can be represented as $O(n^2)$.

It is important to note that although big O notation does not formally take into account things like coefficients and lower order terms, in the context of this IA, I may take these factors into account to differentiate between different versions of algorithms if the simplified big O notation is the same. These differences may be small, but they are notable.

Formal Definitions



Example of different big O notations.

The big O notation of a function can be understood as the asymptotic upper bound of a function (Black, 2019). Formally, $f(n)$ can be considered the big O notation for $g(n)$ if there exists a positive real coefficient (m) for which:

$$mf(n) \geq g(n) \text{ for all } n \geq k$$

where k is a specific positive real value of n . It can be trivially demonstrated that, for example:

$$n = O(17n + 2)$$

This is because there exists a positive real coefficient (such as 18) for which:

$$mn \geq 17n + 2 \text{ for all } n \geq k$$

Solving this inequality algebraically and substituting 18 for m gives:

$$n \geq 2$$

As a result, it can be said that:

$$18f(n) \geq g(n) \text{ for all } n \geq 2$$

On the other hand, it can also be demonstrated that higher order functions cannot have a lower order big O notation. For example:

$$n \neq O(n^2)$$

This is because there is no positive real number for which $f(x)$ will always be greater than $g(x)$. In other words:

$$mn \geq n^2$$

cannot be true for sufficiently large values because solving this equality yields:

$$m \geq n$$

It can be concluded that:

$$mf(n) \geq g(n) \text{ for all } n \leq m$$

but this restricted n to *less than or equal to* m rather than *more than or equal to*, meaning that the condition for big O notation is not met, because when n is greater than m , $f(n)$ is no longer greater than $g(n)$.

Another case that can be explored is logarithmic and exponential functions. Logarithmic functions can be expressed as the inverse of exponential functions. For two functions where:

$$f(n) \geq g(n)$$

these functions produce points:

$$(n_i, y_1); (n_i, y_2), y_1 \geq y_2$$

with $f(n)$ producing the point on the left and $g(n)$ producing the point on the right. As they are both positively sloping functions (as logarithmic and exponential functions are), it also stands to reason that for one y value one point can be produced by each equation where:

$$(n_1, y_i); (n_2, y_i) , x_1 \leq x_2$$

because a larger function produces the same y value with smaller input. Using the nature of inversions, where the n and y coordinates of a function are switched, the the inverse points are:

$$(n_i, y_1); (n_i, y_2) , y_1 \leq y_2$$

By noting that this is the same as the explanation of the inequality expressed above, and that $f(n)$ remains on the left while $g(n)$ remains of the right, it can be deduced that:

$$f^{-1}(n) \leq g^{-1}(n)$$

Next, it can be shown that a logarithmic relationship is just the inverse of an exponential relationship because:

$$\ln(n) = y$$

By switching the n and y inputs:

$$\ln(y) = n$$

$$e^n = y$$

By applying these above principles, if:

$$m \log(n) \geq \ln(n) \text{ for all } n \geq k$$

then:

$$me^n \geq 10^n \text{ for all } n \geq k$$

It is not true that:

$$me^n \geq 10^n \text{ for all } n \geq k$$

because:

$$10^n = \left(\frac{10^n}{e^n}\right)e^n$$

by substitution:

$$me^n \geq \left(\frac{10^n}{e^n}\right)e^n$$

$$me^n \geq \left(\frac{10}{e}\right)^n e^n$$

$$m \geq \left(\frac{10}{e}\right)^n$$

Because m is not greater than a constant, but rather an expression with a variable, it does not satisfy the condition of being a single positive real number. This proves that

$$e^n \neq O(10^n)$$

or, more generally that two exponential functions with different bases cannot be the big O functions of each other. As described above, if logarithmic functions with different bases could be the big O notation of each other, then exponential functions of different bases could too. Since we know that exponential functions with different bases cannot be big O notations of each other, neither are logarithmic functions. For example:

$$\log(n) \neq O(\ln(n))$$

Cases

Mean

Mean can be expressed in the following equation according to the IB formula booklet:

$$\mu = \frac{\sum_{i=1}^k f_i x_i}{n}$$

With f_i being the number recurrences of each different value of x (x_i) in a data set and n being the number of total values. I will call this equation *mean a*. In other words, this formula finds the sum of all values and divides them by the total number of values (n).

Mean a can be represented in the following steps:

1. For every value of x
 - a. **Add 1 to the total number of recurrences of that value (f_i)**
 - b. **Add 1 to the total number of values in the data set (n)**
2. For each value of x check if it matches with any previous values:
 - a. If yes, move to the next value
 - b. If no
 - i. **Multiply it by the total number of that value (f_i)**
 - ii. **Add the product to the running sum for the data set**
3. Divide the final sum by the number of values (n)

The efficiency notation of this would be $O(4n+1)$ because in the worst case scenario where each value is different, steps 1a and 1b would take place once for every value or n number of times, in addition to steps 2bi and 2bii, because there would be no repeated values. Finally, step 3 is necessary, but is only performed once regardless of the input size n . In the worst case scenario, 4 steps would be repeated n times and one step would be repeated once giving $O(4n+1)$. This simplified according to the rules of big O notation is $O(n)$.

In the best case scenario where all values are equal, this program would require $2n+3$ steps. This is because steps 1a and 1b would take place n number of times, but step 2bi and 2bii would only need to take place once. All other values could be skipped over because they match the first value. Finally, step 3 must still be performed. As 2 steps must be repeated n times and 3 steps must be repeated once, this gives $O(2n+3)$, or simplified to $O(n)$.

To simplify *mean a*, it would be useful to eliminate step 1a because it does not deliver any benefit in the worst-case scenario where values must be added up individually. To remove this step, one would need to remove f_i from the equation, as this step is used to define f_i . By the basic properties of multiplication:

$$x_1 + x_2 + x_2 + x_2 + \dots + x_2 = x_1 + f_i(x_2)$$

where the number of x_2 's added together is equal to f_i . Because f_i is definitionally equal to the number of recurrences of any given number (x_i), adding together each of the values is equal to multiplying of recurrences of each different value (f_i) by that value (x_i) and then adding these together. This can be represented as:

$$\sum_{i=1}^k f_i x_i = \sum_{i=1}^k x$$

where x_i is each different value, but x is every value even if there are repeats. Dividing both sides by the number of values n gives:

$$\mu = \frac{\sum_{i=1}^k f_i x_i}{n} = \frac{\sum_{i=1}^k x}{n}$$

In other words the following equation for mean is possible:

$$\mu = \frac{\sum_{i=1}^k x}{n}$$

I will call this equation *mean b*. To test if this equation is really more efficient than *mean a*, it is necessary to represent it in terms of computing steps as below.

1. For every value of x (x)
 - a. **Add 1 to the total number of values in the data set (n)**
 - b. **Add the value (x) to the running sum for the data set**
2. Divide the final sum by the number of values (n)

The efficiency notation of this would be $O(2n+1)$ because no matter the input, steps 1a and 1b are necessary for every value of x . The best and worst case scenarios are the same because there are no *if* statements. It is then necessary to complete step 2 only once no matter the input size (n). As there are 2 steps that must be completed n times and 1 step that must be completed once, the efficiency of this is $O(2n+1)$, which is, simplified, $O(n)$.

Mean b is clearly more efficient than *mean a* because while the simplified big O notation of both equations is $O(n)$, the unsimplified notation of *mean b* is $O(2n+1)$ which is better than the best case for *mean a* ($O(2n+3)$) and much better than the worst case for the *mean a* ($O(4n+1)$).

I tested these conjectures by creating a program for *mean a* and *mean b*, then using python's "timeit" function to time the function for differently sized inputs of numbers between 1 and 100. For each input, the function was run 10,000 times and the average time was taken. *Mean a* consistently outperformed *mean b*, especially as the input size increased. It is important to note that as these are randomly generated values, they more closely resemble an average case rather than a worst case. Here are the images of my code, where x is a list of numbers as the input, and the graph:

```

1 import timeit
2 print(timeit.timeit("""
3 l=[x]
4 diction={}
5 total=0
6 n=0
7 z=0
8 for a in l:
9     if a in diction:
10         diction[a]+=1
11     else:
12         diction[a]=1
13     n+=1
14 for a in l:
15     z+=1
16     if a in diction:
17         number=a*diction[a]
18         total+=number
19     del diction[a]
20 mean=total/n
21 print(mean)
22 """, number=10000)/10000))

```

mean a

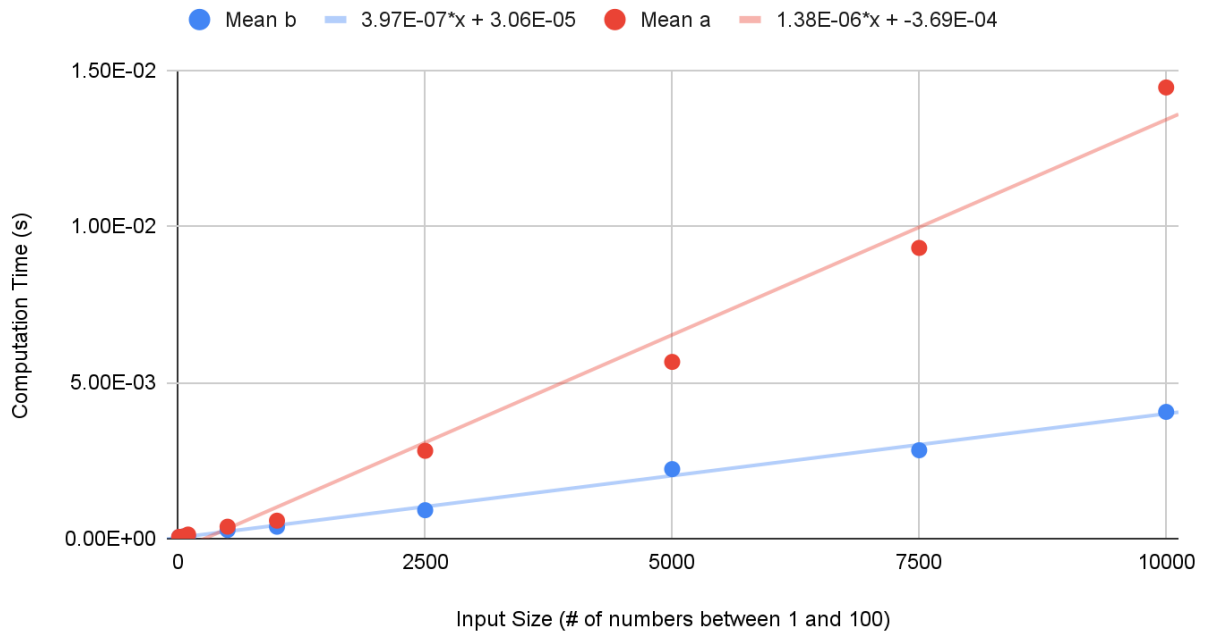
```

1 import timeit
2 print(timeit.timeit("""
3 l=[x]
4 n=0
5 total=0
6 for q in l:
7     total+=q
8     n+=1
9 mean=total/n
10 print(mean)
11 """, number=10000)/10000))
--

```

mean b

Calculation Time for Mean a and Mean b Depending on Input Size



Graph of the efficiency of *mean a* and *mean b*.

The discrepancy between the formula presented by IB (*mean a*) and the most computationally efficient program (*mean b*) can be explained by the differences between the way humans and computers compute numbers. For a computer, just counting up the number of times each value recurs (f_i) is as demanding as just adding together all of the numbers because for simple forms of arithmetic, time is constant regardless of the size of numbers being added together (Paul, 2016). On the other hand, for humans, counting up the number of recurrences of each value is far easier than adding all of the values together. Because *mean a* represents the system easier for humans, it is recommended by the IB. In the rest of the IA, when I can replace *mean a* with *mean b* as part of other equations, I will, because *mean b* is more efficient for computers.

Variance

According to the IB, variance can be expressed in the following equation:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i(x_i - \mu)^2}{n}$$

With f_i being the number recurrences of each different value of x (x_i) in a data set, n being the number of total values and μ being the mean. Importantly, this variance formula requires knowledge of the mean, but a computer given a full set of data would not already know the mean. For this reason, a substitution of μ for the mean equation is necessary. In this case I will use *mean b* as it has been demonstrated to be more efficient. This results in the following equation:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i(x_i - \frac{1}{n} \sum_{i=1}^k x)^2}{n}$$

I will call this equation *variance a*. This can be represented in the following steps. In this equation I will first find the mean and then use it for all following steps rather than finding the mean anew for each value.

1. For every value of x (x)
 - a. **Add 1 to the total number of recurrences of that value (f_i)**
 - b. **Add 1 to the total number of values in the data set (n)**
 - c. **Add the value (x) to the running sum for the data set**
2. Divide the final sum by the number of values (n) to find the mean (μ)
3. For each value of x check if it matches with any previous values:
 - a. If yes, move to the next value
 - b. If no
 - i. **Subtract the mean (μ)**
 - ii. **Square the result**
 - iii. **Multiply it by the total number of that value (f_i)**
 - iv. **Add the product to the running sum for the data set**
4. Divide the sum by the number of values (n)

In the worst case where there are no repeated values, this equation is $O(7n+2)$. Steps 1a, 1b, and 1c must be repeated for every value, and steps 3bi, 3bii, 3biii, and 3biv must each be repeated, as no values are repeated and can be skipped. Finally, steps 2 and 4 must be completed once each no matter the size of the input. As 7 steps must be repeated for every value and 2 must be completed once, the efficiency notation of this equation is $O(7n+2)$, or in simplified form $O(n)$.

In the best case scenario, where each value is the same, the efficiency of this program is $O(3n+6)$. This is because steps 1a, 1b, and 1c will still be run once for every value, but now steps 3bi, 3bii, 3biii, and 3biv only need be run once, as all following values can be skipped past as they are equal. Steps 2 and 4 must also be completed once. 3 steps must be completed once for every value and 6 steps must be completed once, giving $O(3n+6)$, or in simplified form $O(n)$.

The largest source of inefficiency in this equation is the necessity of subtracting the mean from every value before squaring the result. It is possible to eliminate this inefficiency by removing the subtraction from the summation part of the equation, thus making it possible to just do this step once. (In this explanation I will use μ for the mean instead of *mean b* for the sake of readability)

$$\sigma^2 = \frac{\sum_{i=1}^k f_i(x_i - \mu)^2}{n}$$

First, expand the square:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i(x_i^2 - 2x_i\mu + \mu^2)}{n}$$

Next by multiplying each value in the brackets by f_i :

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2 - 2f_i x_i \mu + f_i \mu^2}{n}$$

By separating different terms:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \frac{\sum_{i=1}^k 2f_i x_i \mu - f_i \mu^2}{n}$$

By removing common factors:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \mu \frac{\sum_{i=1}^k 2f_i x_i - f_i \mu}{n}$$

By distributing the summation through the equation:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \mu \frac{\sum_{i=1}^k 2f_i x_i - \sum_{i=1}^k f_i \mu}{n}$$

And removing constants from the summations:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \mu \frac{2 \sum_{i=1}^k f_i x_i - \mu \sum_{i=1}^k f_i}{n}$$

Because f_i is equal to the recurrences of each value, the summation of all f_i 's for each value must equal the total number of values n . Applying this in the equation gives:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \mu \frac{2 \sum_{i=1}^k f_i x_i - \mu n}{n}$$

Then substituting *mean a* for μ and simplifying gives:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \frac{\sum_{i=1}^k f_i x_i}{n} * \frac{2 \sum_{i=1}^k f_i x_i - \frac{\sum_{i=1}^k f_i x_i}{n} n}{n}$$

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \frac{\sum_{i=1}^k f_i x_i}{n} * \frac{2 \sum_{i=1}^k f_i x_i - \sum_{i=1}^k f_i x_i}{n}$$

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \frac{\sum_{i=1}^k f_i x_i}{n} * \frac{\sum_{i=1}^k f_i x_i}{n}$$

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \left(\frac{\sum_{i=1}^k f_i x_i}{n} \right)^2$$

Finally, substituting *mean a* with *mean b* to improve efficiency:

$$\sigma^2 = \frac{\sum_{i=1}^k f_i x_i^2}{n} - \left(\frac{\sum_{i=1}^k x}{n} \right)^2$$

I will call this second formula presented by IB *variance b*. It successfully removes the mean subtraction from the summation, simplifying the equation. The steps for computing through this formula are as follows:

1. For every value of x (x)
 - a. **Add 1 to the total number of recurrences of that value (f_i)**
 - b. **Add 1 to the total number of values in the data set (n)**
 - c. **Add the value (x) to the running sum for the data set**
2. Divide the final sum by the number of values (n) to find the mean (μ)
3. For each value of x check if it matches with any previous values:
 - a. If yes, move to the next value
 - b. If no
 - i. **Square the result**
 - ii. **Multiply it by the total number of that value (f_i)**
 - iii. **Add the product to the running sum for the data set**
4. Divide the sum by the number of values (n)
5. Square the mean

6. Subtract mean from the quotient from step 4

In the worst-case scenario where every value is different, this equation is computed according to the efficiency $O(6n+4)$. This is because steps 1a, 1b, and 1c must still be completed once for every data point. Similarly, 3bi, 3bii, and 3biii must be computed once for every value because every value is different. Finally, steps 2, 4, 5, and 6 must be computed once no matter the input size. As 6 steps must be completed once per input and 4 must be completed once, the efficiency of the worst-case scenario would be $O(6n+4)$, or in simplified form $O(n)$.

For the best-case scenario where all values are equal, *variance b* is actually slower than *variance a* and requires steps $O(3n+7)$. This is because steps 1a, 1b, and 1c must still be performed once per input, but steps 3bi, 3bii, and 3biii only need be performed once. Steps 2, 4, 5, and 6 will also be performed once each. Because 3 steps need to be performed n times and 7 steps need to be repeated once, the efficiency is $O(3n+7)$, or in simplified form $O(n)$.

The gain in efficiency between *variance a* and *variance b* for the worst-case scenario but loss of efficiency in the best-case scenario can be explained by the replacement of one step that may need to be repeated n times (step 3bi in *variance a*) by two steps that need to be repeated once no matter the input (steps 5 and 6 in *variance b*). In the worst case, where step 3bi in *variance a* needs to be repeated n times, it is clearly more efficient in almost all cases to just replace it with 2 steps that only need to be repeated once. On the other hand, in the best case, when step 3bi in *variance a* only needs to be repeated once because all the values are the same, it is less efficient to replace it with two steps. In general, *variance b* is more efficient, because the efficiency gain in the worst case is n but the efficiency loss in the best case is only 1. Furthermore, in deciding between two equations' efficiency, one must generally consider the worst case (Wynne-McHardy, 2018).

Finally, it may be more efficient to follow a similar process as completed with the mean for the leftmost term of *variance b*. Knowing that *mean a* and *mean b* are equal:

$$\frac{\sum_{i=1}^k f_i x_i}{n} = \frac{\sum_{i=1}^k x}{n}$$

Substituting x^2 for x in the mean equation gives the following result:

$$\frac{\sum_{i=1}^k f_i x_i^2}{n} = \frac{\sum_{i=1}^k x^2}{n}$$

Substituting this second equation into *variance b* gives:

$$\sigma^2 = \frac{\sum_{i=1}^k x^2}{n} - \left(\frac{\sum_{i=1}^k x}{n} \right)^2$$

Combining the fractions gives:

$$\sigma^2 = \frac{n \sum_{i=1}^k x^2 - \left(\sum_{i=1}^k x \right)^2}{n^2}$$

I will call this formula *variance c*. Importantly, in this equation x is each value, not every different value. The steps taken to compute this can be represented as follows:

1. For every value of x (x)
 - a. **Add 1 to the total number of values in the data set (n)**
 - b. **Add the value (x) to the running sum for the data set**
2. Square the sum found in step 1
3. For each value of x :
 - a. **Square the result**
 - b. **Add the product to the running sum for the data set**
4. Multiply the sum found in step 2 by n
5. Subtract result of step 2 from the sum found in step 4
6. Square the number of values (n)
7. Divide difference found in step 5 by the result of step 6

This equation has an efficiency in the worst- and best-case scenario of $O(4n+5)$. There are no *if* statements, so the efficiency is constant no matter the input. There are 4 steps that must be taken n times (1a, 1b, 3a, and 3b), and 5 that must be taken 1 time (2, 3, 4, 5, and 6), giving it an efficiency $O(4n+5)$, or in the simplified form $O(n)$.

The worst-case efficiency for *variance c* is much better than the worst-case for *variance a* and *variance b*. This is because those formulas require a counting of the number of recurrences of each value, which is useless in the worst-case situation where all values are different.

The best-case efficiency is still best with *variance a* because in *variance c* it is necessary to square every value individually. In the best case where all values are equal, it is more efficient to square the repeated value once and then multiply it by the number of recurrences (f_i).

As noted above, the worst case is considered for determining the efficiency of programs, so *variance c* is the most efficient of the 3 variance equations. This equation is likely not taught because it bears the least resemblance to the intuitive variance equation (*variance a*) and because it is relatively easy to find the mean using a calculator. Using a hand calculator, *variance b* is most efficient because solving for the mean is relatively easy, and finding the number of recurrences of values is easier than adding all the values together. Despite this being a case, for a computer *variance c* remains the best equation.

Standard Deviation

According to the IB, standard deviation can be expressed in the following equation:

$$\sigma = \sqrt{\frac{\sum_{i=1}^k f_i(x_i - \mu)^2}{n}}$$

With f_i being the number recurrences of each different value of x (x_i) in a data set, n being the number of total values and μ being the mean. Standard deviation can also just be expressed as the square root of variance.

$$\sigma = \sqrt{\sigma^2}$$

For this reason, the steps required to compute it can be expressed as:

1. Find variance
2. Take the square root of variance

Although “find variance” is not one step, this formulation is useful, because it shows the efficiency of any notation to find standard deviation is equal to the efficiency of variance plus one. In other words, $O(\text{standard deviation}) = O(\text{variance} + 1)$.

According to this, the efficiency of *standard deviation a* in the worst case would be $O(7n+3)$ because *standard deviation a* uses *variance a* which has $O(7b+2)$ for the worst case. Similarly, the best case would have $O(3n+7)$ because the best-case scenario for *variance a* has efficiency $O(3n+6)$.

To find the most efficient standard deviation equation, one need merely substitute the most efficient variance equation (*variance c*) for variance. This gives the following equation:

$$\sigma = \sqrt{\frac{n \sum_{i=1}^k x^2 - \left(\sum_{i=1}^k x\right)^2}{n^2}}$$

Conclusions

Evaluation of Success

In this IA I was able to marginally improve the mean, variance, and standard deviation equations. Although I improved the fully described efficiency of each over the formula provided in the formula booklet, I was unable to improve the simplified efficiencies of any formula. The simplified big O notation of each formula remained at $O(n)$. Given that the formulas presented by IB all had efficiency $O(n)$, it was impossible to improve the efficiency, because the only improvement would be $O(1)$ where the input size does not matter. Achieving $O(1)$ with a summation is impossible as at the least every value must be added together or recorded in some way. I believe for each of these equations I have found the most efficient variation.

Guidelines for Efficiency

There were a few major guidelines to make equations more efficient to be drawn from this investigation. These are: minimize the number of variables, move constants outside summations, and use previously efficient formulas.

Minimizing the number of variables minimizes the analysis necessary of the data. For example, gathering information about the number of recurrences of each value is costly in terms of time, and may not deliver any benefit in the worst-case scenario. By removing extraneous variables, even if the data is arranged so as to be very unfavourable for efficiency, there is no loss of speed.

Moving constant values outside summations means that there is one fewer step that must be done n times. This may result in less efficiency in the best-case scenario, but in the worst-case scenario, it prevents wastefully subtracting a value over and over. Instead, by taking the mean out of the summation, for example, it is possible to execute the subtraction only once. This makes the formula far more efficient, especially given larger inputs.

Finally, using previously efficient formulas in modifying new ones can help to eliminate unnecessary steps. For example, by using the most efficient mean formula to find mean in the variance and standard deviation equations, a few steps are saved. This

also makes work more efficient by preventing redundant work removing the same steps from multiple equations.

Competing Goals of Equations

The reason many formulas presented by IB aren't already optimized is because they are designed to be easy for humans to calculate and understand. For example, it is easier for people to count up the number of recurrences of a value than to add those values together. For a computer, this is not a consideration. Similarly, to promote an understanding of variance, the least efficient formula is presented because it clearly represents the squaring of the difference between mean and each value. It is useful to have both efficient and intuitive formulas presented to students so that they can understand the operation and apply it most effectively. It is for this reason that the formula booklet includes both *variance a* and *variance b*.

References

- Black, P. E. (2019, September 6). *Big-O Notation*. National Institute of Standards and Technology. <https://xlinux.nist.gov/dads/HTML/bigOnotation.html>
- Paul. (2016, March 8). *does a computer take more time to multiply, divide, subtract, add two big numbers than smaller number*. Stack Overflow.
<https://stackoverflow.com/questions/35868686/does-a-computer-take-more-time-to-multiply-divide-subtract-add-two-big-number>
- Wynne-McHardy, R. (2018, November 26). *Big O for Beginners*.
<https://hackernoon.com/big-o-for-beginners-622a64760e2>