

## Data Bias in HR analytics

**Author: Tatenda Emma Matika . November 2020**

The rise in the use of data analytics in the workplace has resulted in companies adopting new methods of working with their data. This is a progressive step for companies that have managed to do so and now, they no longer give reports based on intuition, but based on what the data is saying.

This is all good, however, putting too much trust in data will lead to problems. In all predictive models that are created using data, there is always a bias of some sort. But we all know that all models are wrong, but there are some that are useful. The useful ones are those that use data that represents reality well. But representing reality well may not be enough when it comes to HR issues such as diversity and gender balance.

We use past data when creating predictive models. A company's past data may be showing a larger proportion of males in leadership roles compared to females. This then means that using this data to predict future leaders or to select a leader will result in a model that is biased towards males. This does not mean that the data is wrong, it is correct, but skewed towards a certain group. Now, as HR practitioners, we do not want such a bias as it will likely result in grievances from employees.

An important step that is then used to overcome such bias is to make sure that the data is balanced. In real life situations, it is difficult to have balanced data, but there are work-around methods that can be used. The simplest methods are

1. Up-sampling the Minority Class and
2. Down-sampling the Majority Class.

The first method involves creating more data values that represent the minority class such that we have a high number of records in the class. This is a relatively difficult task since it may require adding fictional data. There are methods that are used to do this but when it comes to explaining this to other departments, it might become challenging.

The second method involves removing some of the records from the majority class. If there are 800 males and 200 females in the data, we may want to reduce the number of males to a value closer to 200. But this also becomes problematic as it will result in loss of important

information. Another question would be how to determine which ones to remove and which ones to keep.

This bias does not only occur with gender, it can also occur with values such as race, schools people obtained their qualifications, age and so on. A key to overcoming such problems is to focus on people's personal traits and skills/talent. AT IPC, we consider psychometric tests. These show a person's

personality and other things such as attention to detail. These say more about a person than the gender, age or schools which were traditionally considered to be important.

This step of determining which data to use is should be a critical step in data analytics processes. The models that we create cannot make these decisions for themselves. They only take what we give them and give output depending on that. A common phrase for this is garbage-in garbage-out. We need to make sure that we do not feed garbage to our models. Maybe in future we might have algorithms that can check whether they have been fed garbage or not, but for now, we need our human intelligence.

This bias that we have discussed is about unbalanced data, and this is what is usually found in HR data. However, we also have some other types of data bias. These are:

### **1. Selection bias**

This occurs the data is selected wrongly, for example selecting a specific subset of employees, rendering the sample unrepresentative of the whole company.

### **1. Recall bias**

This occurs when the data depends on people's memory. We might ask for information that people do not remember clearly and they give biased responses.

### **1. Omitted variable bias**

This occurs when an important variable is omitted in the data. For example, age might be an important variable for the model, and if it is omitted, it affects the model.

### **1. Cause-effect bias**

This occurs when we assume that a correlation between variables means a causation. We might introduce a policy and immediately, employees start to perform poorly. There is need to investigate why employees are performing poorly, before rushing to conclude that the policy is the cause. There are many examples of correlations that are not related in any way on [the spurious correlations website](#). This will give insight into the cause-effect bias.

There are other types of bias, but these are some that are important in HR analytics. We need to make use of our human intelligence to identify bias before feeding data into a model for better results.

**Tatenda Emma Matika is a Business Analytics Trainee at Industrial Psychology Consultants (Pvt) Ltd, a management and human resources consulting firm.**

## **References**

<https://www.hrtechnologist.com/articles/digital-transformation/experts-on-ai-bias-in-hr/>

<https://www.analyticsinhr.com/blog/ai-bias-diversity-by-design/>

<https://data36.com/statistical-bias-types-explained/>

<https://www.tylervigen.com/spurious-correlations>

<https://thehumancapitalhub.com/articles/Data-Bias-In-HR-Analytics>