# Doppelgä ngers effects confound ML model

**Introduction**

Doppelgä ngerss effects happen when biomedical data sets have duplicate or nearly identical records. This can throw off machine learning models by introducing bias or skewing the results. Doppelgä ngerss effects can be caused by mistakes in data entry or data management, or they can be made on purpose for a certain reason (such as to test a machine learning model).

In general, Doppelgä ngerss effects can be severe for machine learning because they can change the way data is distributed and trick the model, which can lead to wrong or unreliable results. For example, if a machine learning model is trained on a data set with duplicated records, it may learn to over-represent or under-represent certain patterns or relationships that don't exist in the real underlying data. This can make it hard to generalise and make it hard to do well on data we haven't seen.

Before using data sets for machine learning, it is important to carefully clean and pre-process them to reduce the effects of Doppelgä ngerss effects. This could mean finding duplicate records and removing them, as well as fixing any other mistakes or inconsistencies in the data. Also, machine learning models should be carefully validated and tested to make sure they are strong and can work well with data they haven't seen before.

**Doppelgä ngers effects in different types of datasets**

Doppelgä ngers effects are not unique to biomedical data, and they can occur in any type of data set where there is the potential for duplicate or nearly identical records. For example, there are many potential examples of doppelgä ngers effects in different types of data. Here are a few examples:

1. In patient medical records, doppelgä ngers effects can occur when a patient is accidentally entered into the system multiple times, either with slightly different personal information or with completely identical information. This can lead to multiple records for the same patient, which can confuse machine learning models trying to predict outcomes or analyze trends.
2. In financial data, doppelgä ngers effects can occur when transactions are accidentally recorded multiple times, either due to errors in data entry or due to fraud. This can lead to duplicate records that mislead machine learning models trying to predict fraud or identify patterns in spending.
3. In social media data, doppelgä ngers effects can occur when users create multiple accounts, either to evade bans or to manipulate the system. This can lead to duplicate records that mislead machine learning models trying to predict user behavior or identify trends in social media activity.
4. In meteorological data, doppelgä ngers effects can occur when weather stations report the same data multiple times, either due to errors in data transmission or due to intentional manipulation of the data. This can lead to duplicate records that mislead machine learning models trying to predict

weather patterns or analyze climate data.

Overall, Doppelgä ngerss effects can occur regardless of the type of data and can have a wide range of effects on machine learning models, depending on the specific situation and data set.

**Doppelgä ngers effects in biomedical types of datasets**
Because biomedical data is complicated and sensitive, Doppelgä ngerss effects can be especially hard to deal with. For example, mistakes or duplications in patient medical records can have serious effects on care and treatment, and they can also confuse machine learning models that are used to predict patient outcomes or look for trends in healthcare data. Because of this, it is especially important to carefully clean and pre-process biomedical data sets before using them for machine learning, to make sure that Doppelgä ngerss effects will not introduce bias or skew results.

1. **In imaging data,** doppelgä ngers effects can occur when an image is accidentally acquired or processed multiple times, either due to errors in data management or due to intentional duplication for testing purposes. This can lead to duplicate records in the data set, which can mislead machine learning models trying to predict diseases or analyze trends in imaging data.
2. **In gene sequencing data,** doppelgä ngers effects can occur when the same genetic sample is sequenced multiple times, either due to errors in data management or due to intentional duplication for testing purposes. This can lead to duplicate records in the data set, which can mislead machine learning models trying to predict diseases or analyze trends in genetic data.
3. **In metabonomics data,** doppelgä ngers effects can occur when the same metabolic sample is analyzed multiple times, either due to errors in data management or due to intentional duplication for testing purposes. This can lead to duplicate records in the data set, which can mislead machine learning models trying to predict diseases or analyze trends in metabolic data.

In each of these cases, doppelgä ngers effects can distort the distribution of data and mislead machine learning models, leading to inaccurate or unreliable results. As such, it is important to carefully clean and pre-process these types of data sets before using them for machine learning, to ensure that doppelgä ngers effects do not introduce bias or skew results.

**Identification of data doppelgä ngerss**
Data doppelgä ngerss refer to two or more sets of data that are identical or nearly identical. These doppelgä ngerss can arise due to a variety of reasons, such as data entry errors, duplication of data, or intentional manipulation of data.

**MD5 or SHA-1**
Use checksum algorithms: Use checksum algorithms such as MD5 or SHA-1 to calculate the checksum for each data set. If the checksums are identical, it is likely that the data sets are identical. The algorithms are as followed:

*MD5:*

1. Initialize a message digest buffer to the initial value specified in the MD5 specification.
2. Pad the message so that its length is a multiple of 512 bits.
3. Divide the padded message into 512-bit blocks.
4. For each block, apply the MD5 compression function to the message digest buffer and the block.
5. After all blocks have been processed, the final message digest is the hash value.

*SHA-1*

1. Choose a checksum algorithm. There are several algorithms available, including SHA-1, SHA-256, and MD5.
2. Split the data set into fixed-size blocks. The size of the blocks will depend on the chosen algorithm. For example, SHA-1 processes data in 512-bit blocks.
3. Calculate the checksum for each block by applying the chosen algorithm to the data.
4. Concatenate the checksums for all blocks to create the final checksum value.

As for using checksum algorithms to identify the "doppelgä ngers effect," it is not clear how checksum algorithms would be relevant to this phenomenon. The doppelgä ngers effect refers to the phenomenon of people having perceived "lookalikes" or "twins," which is typically studied using statistical measures such as Pearson's correlation coefficient or Spearman's rank correlation coefficient. These measures are calculated based on the values of the variables being compared, whereas checksum algorithms create a hash based on the data itself, regardless of the values of the variables. Therefore, it is not appropriate to use checksum algorithms to identify the doppelgä ngers effect.

*Pearson's correlation coefficient*

PPCC: To use Pearson's correlation coefficient to identify the "doppelgä ngers effect," we would need to have two sets of data that are related to the phenomenon we are studying. For example, if we are studying the doppelgä ngers effect in terms of people's perceptions of their own physical attractiveness, we might have one set of data that represents people's self-reported attractiveness ratings and another set of data that represents the attractiveness ratings given to them by others. We would then use the following algorithm and function to calculate the Pearson's correlation coefficient between these two sets of data:

1. Standardize both sets of data by subtracting the mean and dividing by the standard deviation. This will ensure that the variables are on the same scale and will allow we to compare them directly.
2. Calculate the covariance of the two sets of data by multiplying the standardized values for each individual and summing the products.
3. Divide the covariance by the product of the standard deviations of the two sets of data to obtain the Pearson's correlation coefficient.

*Spearman's rank correlation coefficient*

To use Spearman's rank correlation coefficient to identify the doppelgä ngers effect, we would need to have data on the perceived attractiveness ratings of individuals as well as their own self-perceived attractiveness ratings. We would then calculate the correlation coefficient between these two variables to determine the strength and direction of the relationship. A positive correlation would indicate that individuals who perceive themselves as more attractive tend to be rated as more attractive by others, while a negative correlation would indicate the opposite relationship. The steps for calculating Spearman's rank correlation coefficient is as follows:

1. Sort the values of the two variables in ascending order.
2. Assign ranks to the values, with the lowest value receiving a rank of 1, the second lowest value receiving a rank of 2, and so on. If there are ties, assign the same rank to all tied values and adjust the ranks of the remaining values accordingly. For example, if there are three tied values with ranks 2, 3, and 4, the next value would receive a rank of 5 instead of 4.
3. Calculate the difference between the ranks for each pair of values.
4. Calculate the sum of the squared differences.
5. Calculate the Spearman's rank correlation coefficient using the following formula:
6. $r = \frac{1-(6 \times sum\ of\ \text{squared differences})}{n(n^2-1)}$ where n is the number of data points.

The resulting coefficient will range from -1 to 1, with a value of 1 indicating a perfect positive correlation, a value of -1 indicating a perfect negative correlation, and a value of 0 indicating no correlation.

In summary, data doppelgä ngerss can be identified by comparing data sets, using checksum algorithms, using data deduplication tools, and using data visualization tools.

**Improve the doppelgä ngerss effect in ML model**
1. Data preprocessing: Data preprocessing involves cleaning and formatting the data before it is used to train a machine learning model. This includes identifying and removing duplicate records, which can help in improving the doppelgä ngerss effect.
2. Data quality checks: Data quality checks can be implemented to ensure the accuracy and completeness of the data. This can include checks such as verifying the data against external sources or using checksum algorithms to identify duplicate records.
3. Data governance: Data governance policies and procedures can be put in place to ensure the integrity and reliability of the data. This can include establishing strict protocols for data entry and handling, as well as implementing access controls to prevent unauthorized manipulation of the data.
4. Data annotation: Data annotation involves adding additional information or context to the data. This can help in identifying and resolving data doppelgä ngerss as the annotations can differentiate between similar data sets.

In summary, the doppelgä ngerss effect in machine learning models can be improved by implementing data pre-processing, data quality checks, data governance, and data annotation measures.

**Reference:**
[1] Wang, L. R., Wong, L., & Goh, W. W. B. (2022). How doppelgä ngers effects in biomedical data confound machine learning. Drug discovery today, 27(3), 678–685. https://doi.org/10.1016/j.drudis.2021.10.017
[2] Gama, J., Zliobaite, I., Bifet, A., Pechenizkiy, M., & Bouchachia, A. (2014). A survey on concept drift adaptation. ACM Computing Surveys (CSUR), 46(4), 44.
[3] Hassanien, A. E., & Zomaya, A. Y. (2018). Handling concept drift and class imbalance in data streams. In Big Data Analytics (pp. 69-92). Springer, Cham.
[4] Bifet, A., & Gama, J. (2010). Adaptive learning from evolving data streams. In Machine Learning and Knowledge Discovery in Databases (pp. 19-46). Springer, Berlin, Heidelberg.