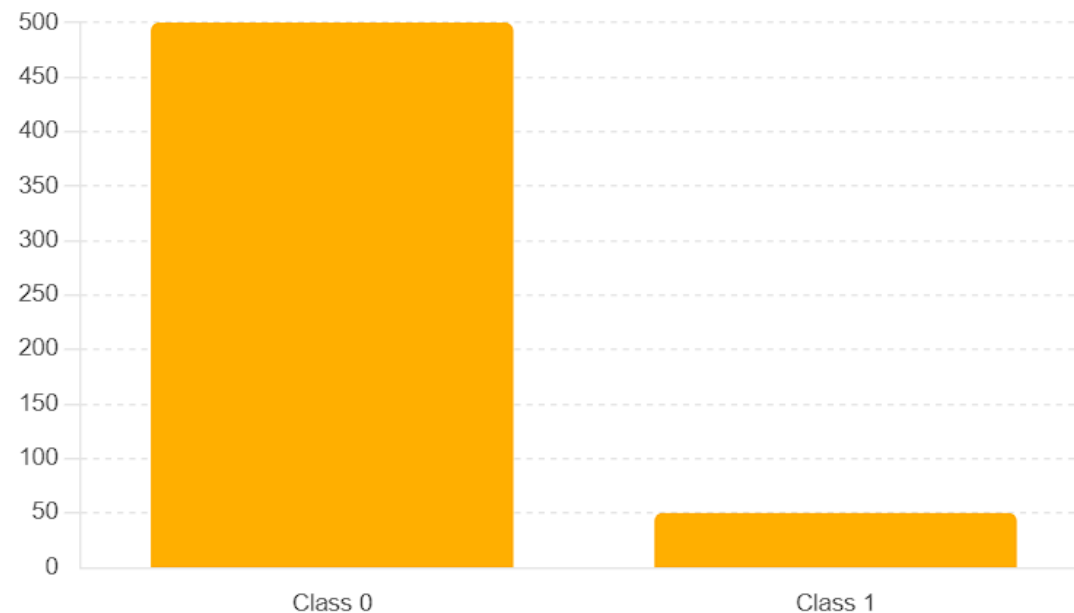
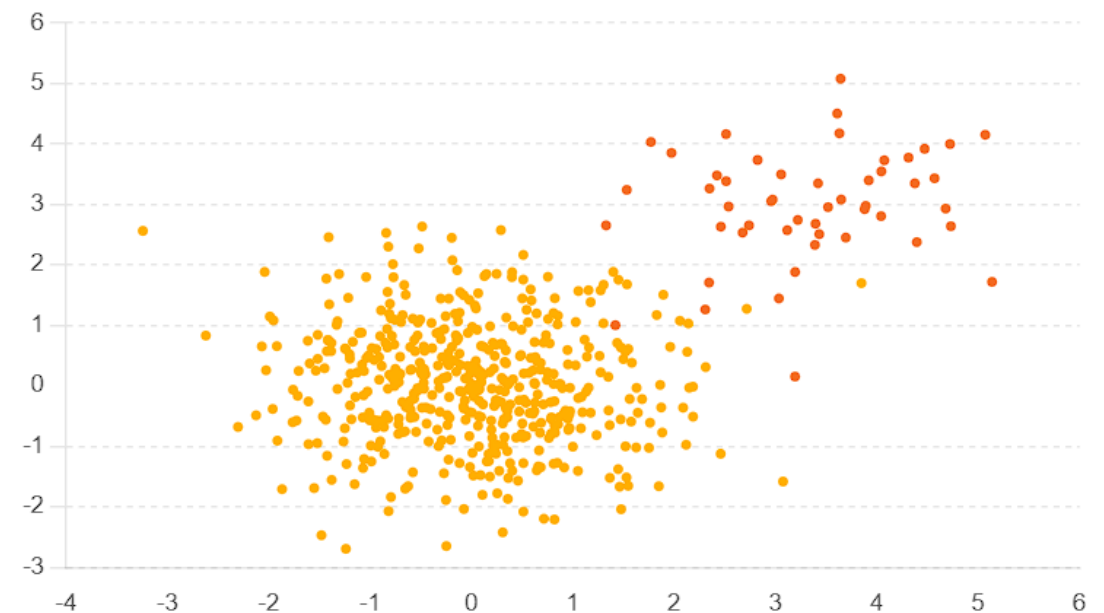


An imbalanced dataset is a dataset where the classes are not represented equally. In other words, one class has significantly more instances than the other(s). This imbalance can pose significant challenges for machine learning models, as they may become biased towards the majority class and perform poorly on the minority class. This is especially problematic in classification tasks where the minority class is often the class of most interest, such as fraud detection, medical diagnosis, or anomaly detection.

Y Count by X Class



Y Feature 2 by X Feature 1 for Class 0 and Class 1



- ❖ **Disproportionate Class Distribution:** One class has a much higher number of instances compared to the other(s).
- ❖ **Model Bias:** Machine learning models tend to be biased towards the majority class, leading to poor predictive performance on the minority class.
- ❖ **Evaluation Challenges:** Standard evaluation metrics such as accuracy can be misleading. Metrics that focus on the performance of the minority class, such as precision, recall, F1-score, and the area under the ROC curve (AUC-ROC), become more relevant.

- ❖ **Biased Predictions:** Models trained on imbalanced data tend to predict the majority class more often, leading to high false-negative rates for the minority class.
- ❖ **Misleading Metrics:** Overall accuracy can be high even if the model performs poorly on the minority class, masking poor performance.
- ❖ **Training Difficulties:** Standard algorithms may struggle to learn the decision boundary correctly due to the skewed class distribution.

No. There isn't a strict rule or universally agreed-upon threshold for what constitutes an imbalanced dataset. It's often context-dependent and can vary based on the specific problem, the nature of the data, and the goals of the analysis or model.

However, a commonly used heuristic is the **80-20 rule**, where a dataset is considered imbalanced if the class distribution is roughly 80% to 20% or worse. In this case, **the majority class would have around 80% of the samples and the minority class around 20%.**

What's considered imbalanced can vary widely. In some cases, a class distribution of **60-40 might be considered imbalanced**, especially if the minority class is critical or costly to misclassify. In other cases, a distribution of **90-10 might be considered balanced**, especially if the classes are naturally imbalanced in the real-world scenario to which the model will be applied.

Ultimately, it's important to consider the **specific domain**, the implications of misclassifications, and the goals of the analysis when determining whether a dataset is imbalanced. Additionally, the choice of threshold might be influenced by practical considerations and domain expertise.

❖ Resampling Techniques:

- ❖ **Oversampling the Minority Class:** Duplicate or generate synthetic instances of the minority class using methods like SMOTE (Synthetic Minority Over-sampling Technique).
- ❖ **Undersampling the Majority Class:** Reduce the number of instances in the majority class to balance the class distribution.

❖ Algorithmic Approaches:

- ❖ **Cost-sensitive Learning:** Assign higher misclassification costs to the minority class to make the algorithm more sensitive to these errors.
- ❖ **Ensemble Methods:** Use ensemble techniques like Balanced Random Forest or EasyEnsemble, which are designed to handle class imbalance.

Let's Implement with Python