| | from imblearn. over_sampling import | Description / Steps |
|---|---|---|
| **Random Oversampling** | RandomOverSampler | - extract at random observation of the majority class until a certain balancing ratio is reached<br>- it is a naive technique (requires no assumption) |
| **SMOTE** | SMOTE | 1) Isolate the minority class<br>2) Determine how many new samples need to be generated and select from which original samples the new one will be generated<br>3) For each chosen minority observation, find k-nearest neighbors<br>4) Determine the L2-distance between the observation and k-nearest neighbors<br>5) The distance has to be multiplied by a factor (a random number [0;1])<br><br>formula:<br>new sample = or_smpl - f * (or_smpl - minor_nbr) |
| **SMOTE-NC (Nominal Continuous)** | SMOTENC(<br>    categorical_features=[n, ...]<br>) | - extands the functunality of SMOTE to categorical variables (ADASYN cannot do this)<br><br>having a categorical feature (column):<br>1) from all numerical features, calculate median(sum(std(all numerical features)))<br>2) proceed with SMOTE putting this median value as an L2-distance between feature values<br>3) when putting the new-generated sample to the place, assign the most frequent categorical feature values to it |
| **Borderline SMOTE** | BorderlineSMOTE | - create new sample only from the original observations that are the closets to the borderline with the majority class<br><br>1) fit KNN with *all* dataset<br>2) find and ignore observations from the minority class which K-ns belongs to the majority class (noise and irrelevant)<br>3) find and ignore observations from the minority class if most of their neighbors are from the minority class (safe and easy to classify)<br>4) Select the observations of the minority class if most of their neighbors are from the majority class<br>5) fit KNN to the minority class observations<br><br>Next - division into 2 variants |
| **Borderline SMOTE (variant 1)** | BorderlineSMOTE(<br>    kind='borderline-1',<br>) | 6) as a regular SMOTE: a new sample to be between original observations of the *minority* class |

| | | |
|---|---|---|
| **Borderline SMOTE (variant 2)** | BorderlineSMOTE(<br>    kind='borderline-2',<br>) | 6) a new sample to be between original observations of the *majority* class |
| **K-Means SMOTE** | KMeansSMOTE | - for clusters<br><br>1) determine clusters -- K-means algorithms to the whole dataset<br>2) select clusters where the % of the minority classes is above a threshold (typically 0.5)<br>3) weight the cluster -- how many new samples to create in each cluster<br>    L2mean = Mean L2 between minority observations<br>    density = (number of minority observations / L2mean) * number of feature<br>    sparsity = 1 / density<br>    cluster sparsity = sparsity / sum()<br>4) calculate the number of the synthetic samples to be generated for each cluster:<br>    $g(i) = cs(i)*G$, where<br>                cs(i) = cluster sparsity,<br>                G = total num of samples to generate,<br>                g(i) = num of samples to generate from cluster i |
| **ADASYN** | ADASYN | - synthetic data is more generated from *all* observations that are harder to classify<br>(this is the main difference beetwen this one and SMOTE)<br><br>1) determine the balancing ratio: X(minority) / X(majority)<br>2) determine the number og samples to generate:<br>    G = ( X(majority) - X(minority) ) * factor<br>3) train KNN using the *entire* dataset to find closest K-ns for each observation of minority class<br>4) determine the weighting r: D/K, where D = neighbors from the majority class, K = neighbors<br>5) normalize r: r/sum(r-s)<br>6) calculate the number of the synthetic samples to be generated for each observation of minority class:<br>    $g(i) = r(i) * G$, where<br>                r(i) = weight for observation i,<br>                G = total num of samples to generate,<br>                g(i) = num of samples to generate from observation i<br>7) for each minority class x(i) generate g(i) synthetic samples<br><br>formula:<br>new sample = or_smpl - f * (or_smpl - *any*_nbr) |