**Imbalanced Data**

fraud
detection

output → category → rare disease → 1000

dumb model
- - - - -
NO

| | |
|---|---|
| 995 NO | 5 Yes |

99.5 %

Problem → [bias for majority] → minority class

Problem → metrics: accuracy

Recall
Precision

Solution:
- Undersampling ① (1:1)
  - Monday
- Oversampling ②
  - Random oversampling
  - SMOTE
  - variants
- Hybrid
- Ensemble methods
- Cost sensitive learning (L)
- Class weighting
- threshold tuning

995 → 50

$\begin{bmatrix} 0.5 \end{bmatrix} \nearrow 1 \searrow 0 \rightarrow 0.6 \nearrow 1 \searrow 0$

Thu
  └→ [Interview questions on Projects] → [Rohan]   AVP - BoB

# Random Oversampling

SMOTE (why?)

Imbalanced data

↓

Biased for majority class

→ Random Oversampling

Overfitting → bad results

imb

↓ balanced
duplicate

green
increase

train test split

X train        X_test ↗ imbal

↓

resampling
over

↳ X_resampled    algorithm

b → (9) (10)

**Synthetic minority oversampling Technique**

oversampling → [minority → up sample]

[K=25]

interpolation

[K=5]

a) Class Imbalance

b) SMOTE

New synthetic data given by:
s1 = x1 + (rand(0,1) * x2-x1)

KNN → [k=5]

maj = min

- Train a KNN on minority class observations - find each observation's 5 closest neighbours

$0-1$

0.5

**To create the new synthetic data**:

- Select examples from the minority class at random.
- Select a neighbour of each example at random (for the interpolation)
- Extract a random number between 0 and 1
- Calculate the new examples as = (original sample (factor) * (original sample - neighbour)
- The final dataset consists of the original dataset + the newly created examples

800

| maj | min |
|-----|-----|
| 900 | 100 |

variants

$X (2,2)$
$(1,1)$ $(1.5, 1.5)$

④

sample − 0.5 [sample − neighbor]

$(1,1) - 0.5 [(1,1) - (2,2)]$

$(1,1) - 0.5 [(-1,-1)]$

$(1,1) - (-0.5, -0.5)$

$(1.5, 1.5)$

Disadvantages

1. Does Not Handle Categorical Data Well

   numerical

   SMOTE assumes that the data is numeric and continuous. Therefore, it struggles with categorical data because it relies on linear interpolations between samples. Applying SMOTE directly to categorical features can lead to the generation of nonsensical instances, as interpolation might produce values that do not correspond to any valid category.

2. Computational Complexity

   high − distances

   The method is computationally expensive, especially with large datasets or datasets with a high number of dimensions (features). This is because SMOTE calculates distances between all pairs in the minority class to find the k-nearest neighbors, which can be computationally intensive.

3. Dependency on the Choice of Neighbors
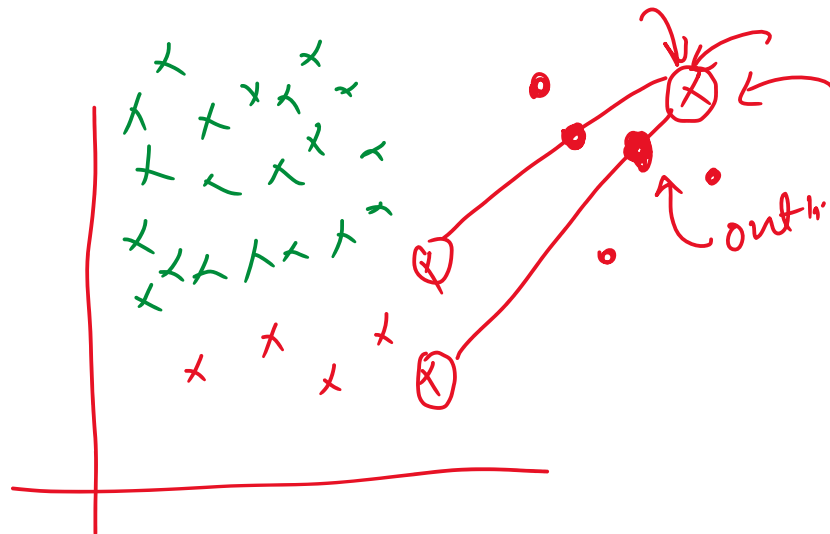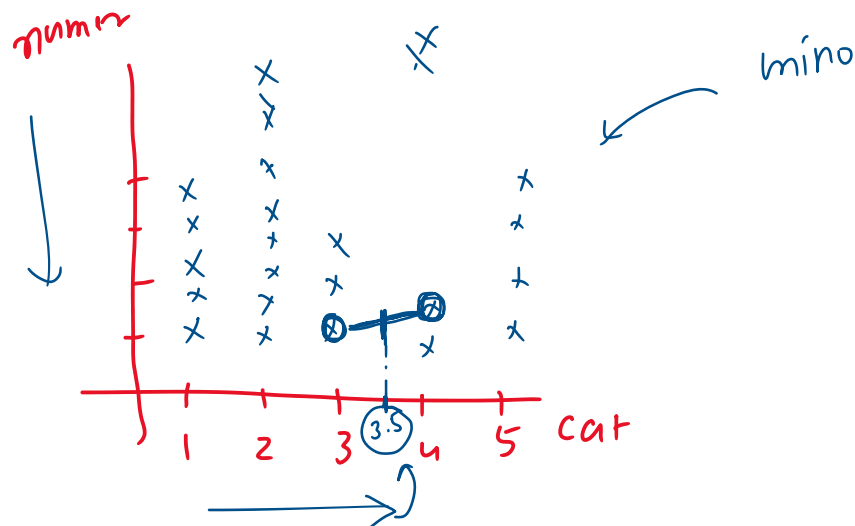
   [K=5] [K=3] [K=25]

   The performance of SMOTE heavily depends on the number of nearest neighbors (k-neighbors) chosen for generating the synthetic samples. Choosing too few neighbors can limit the diversity of the synthetic samples, while choosing too many might include distant, less relevant samples for interpolation, leading to less accurate synthetic data.

4. Sensitive to Outliers

SMOTE can be sensitive to outliers because it uses these as bases to create more samples. If outliers are present in the minority class, SMOTE will generate more outliers, which can skew the model and lead to poor generalization on unseen data.
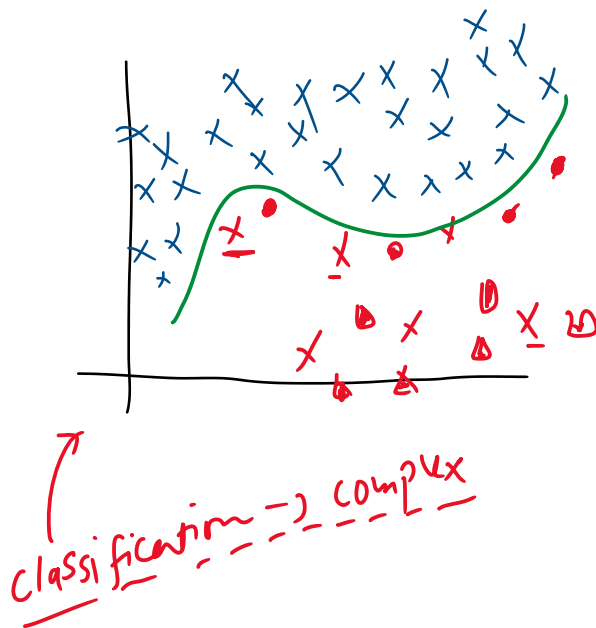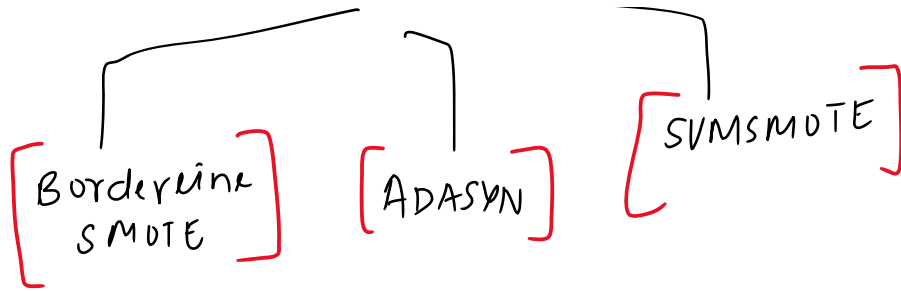
5. Balance Achieved May Not Reflect True Nature

Balancing the classes via synthetic data does not mean that the resulting dataset reflects the true nature of the underlying problem. This artificial balance can sometimes lead to models that perform well on balanced training data but less effectively on real-world, imbalanced datasets.
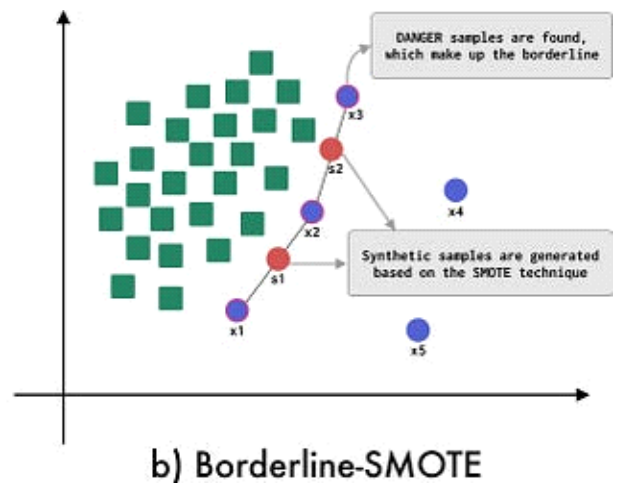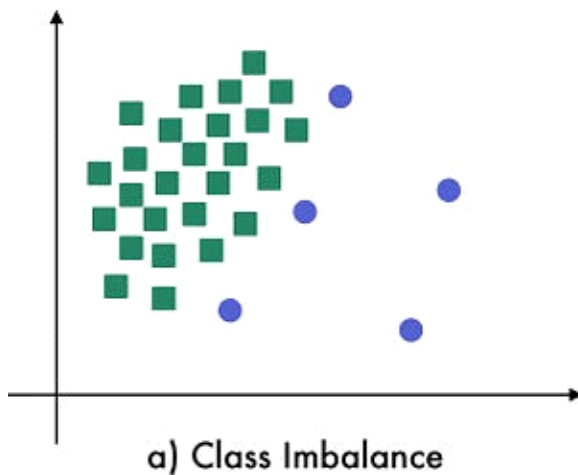




Variants

[ Borderline SMOTE ] [ ADASYN ] [ SVMSMOTE ]



classification → complex

A → more pts near the

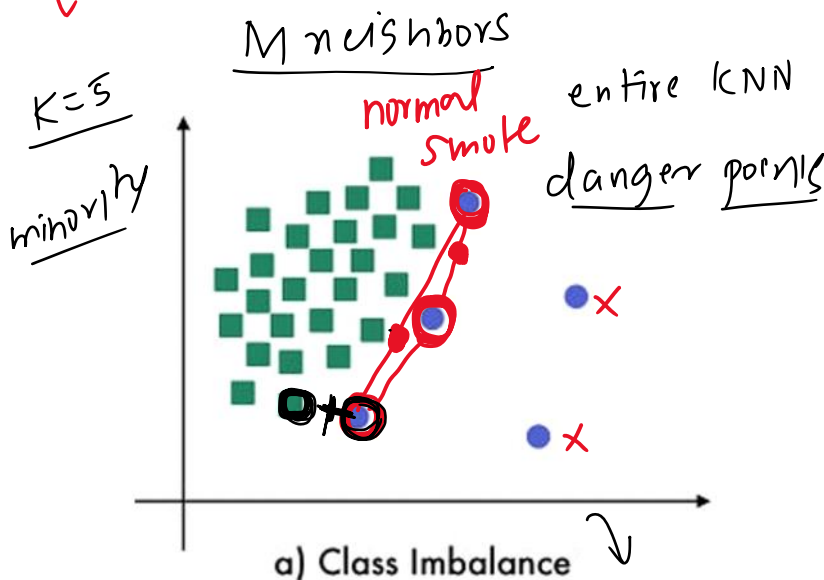B → far from decision boundary

a) Class Imbalance

b) Borderline-SMOTE

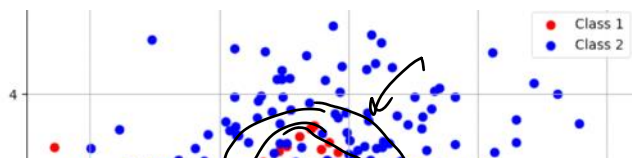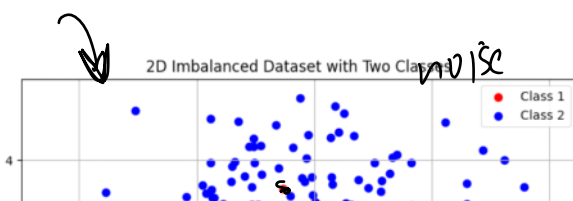Core Principle - To improve the decision boundary

- Train a KNN on entire dataset
- Find the M closest neighbours to each observation from the minority group
- If most, but not all, neighbours belong to a different class, add the observation to a DANGER group

- Train another KNN only on minority group, find each DANGER group observation's closest K neighbours
- Interpolate as in SMOTE from templates in DANGER group to minority neighbours

$K = 5$

$\dfrac{majority}{minority}$

minority



M neighbors

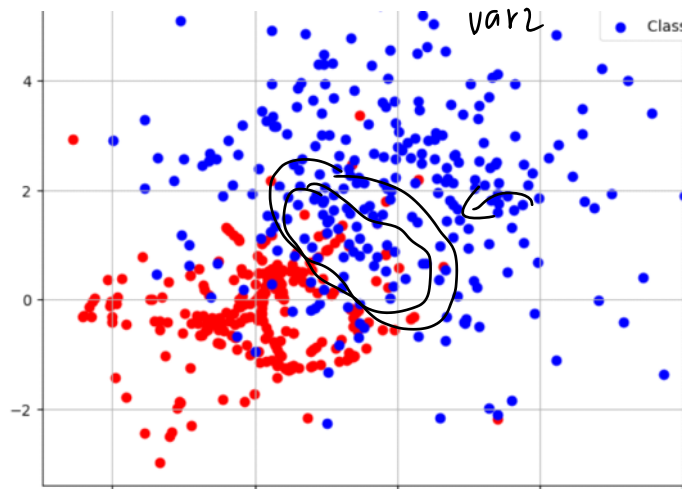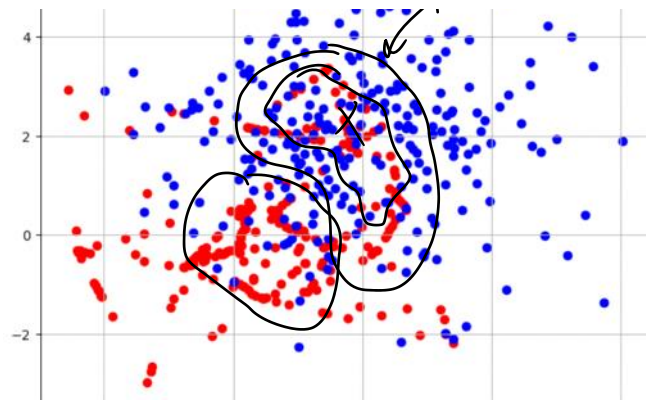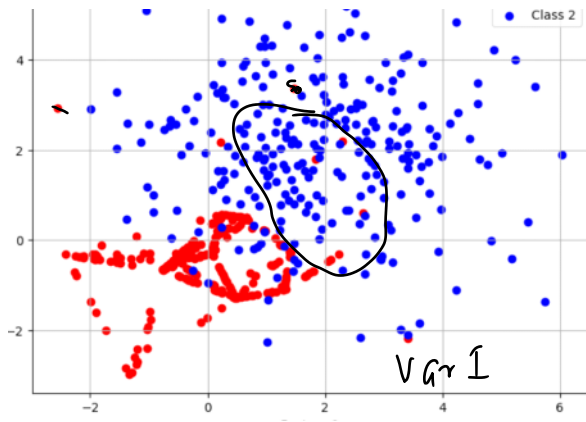normal smote

entire KNN

danger points

a) Class Imbalance

$K = 5$

minority

danger data points

3 type
↓
minority

① → all the neighbors are majority value

② → all datapoint neighbor minority



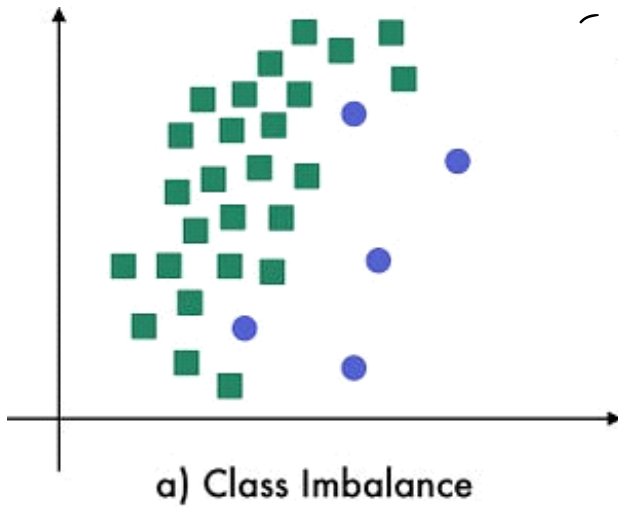2D Imbalanced Dataset with Two Classes

noise

Class 1
Class 2

Class 1
Class 2

Var I



var2

$factor [0-0.5]$

$[variant 2] \longrightarrow \overline{unclear}$

Adaptive Synthetic → Classification

01 May 2024    08:59

easy points

Those samples located in a low density zone are found.

Hard data points

Synthetic samples are generated based on the SMOTE technique

a) Class Imbalance

b) ADASYN

[Smote]    ●    [borderly smote]

1. Determine the balancing ratio = X(minority)/X(majority)
2. Determine the number of new examples to create: G = (Xmaj - Xmin) * factor (the factor is 1 to attain a balancing ratio of 1)
3. Train a KNN on entire dataset
4. Find the K closest neighbours to each example from the minority class
5. Determine a weighting factor for each observation of Xmin: ri = D / K, where K is the number of neighbours and D is the number of neighbours that do not belong to Xmin
6. Normalise ri: rnorm = ri / sum(r)
7. Determine how many observations should be created from each observation of Xmin: Gi = ri * G (G was determined in step 2, ri on step 6)
8. Select a neighbour of each example at random (for the interpolation, can be Xmin or Xmaj)
9. Extract a random number between 0 and 1
10. Calculate the new examples as = original sample - factor * (original sample - neighbour)

11.

minority

0.3

5 mino

→ 0.3

weighting

$\frac{0.2D}{K} = \boxed{\frac{4}{5}}$

0.1

0.1

①

$G = 00$

$(800)$  factor = 1

maj 900   min 100

balance = $\frac{100}{900} = \frac{1}{9}$

$(900 - 100)$

$(800) \times$ factor = 1

$\frac{P}{K} = \boxed{\frac{3}{5}}$

weight = $\frac{P}{K}$

weight ↑ ①

$\frac{}{0 \to}$

● The final dataset consists of the original dataset + all the newly created examples

$800 \times 0.3 = 240$

$0.1 \times 800 = 80$

$(160)$

# SVM SMOTE

01 May 2024    15:02

- Train a SVM on entire dataset
- Find the support vectors of the minority class, these will be the templates
- Train a KNN on entire dataset, find the 10 closest neighbours of the support vectors.
- Decide between inter and extrapolation: if most of the neighbours are from majority class, interpolate, otherwise, extrapolate
- Train another KNN, this time only on minority group, find the 5 closest neighbours to the support vectors
- Create the synthetic examples by inter or extrapolation between templates and their neighbours
- Note that the neighbours are not chosen at random, but instead from the closer to the furthest to create the synthetic data
- The final dataset consists of the original dataset + all the newly created examples

$$new\_sample = support\_vector + \lambda \times (neighbor - support\_vector)$$

# Which one to use when?

01 May 2024      15:07

1. Class Imbalance Severity:

   - If the imbalance is severe, more aggressive methods like SMOTE or its variants might be necessary. For milder imbalances, simpler techniques like random oversampling might suffice without introducing much noise.

2. Data Complexity and Overlap:

   - For datasets with complex decision boundaries or significant overlap between classes, techniques like Borderline SMOTE or ADASYN (Adaptive Synthetic Sampling) can be more effective because they focus on the regions where the classifier is most likely to make errors.

3. Model Sensitivity:

   - Consider how sensitive your model is to noise. Techniques like SMOTE might introduce synthetic outliers which can be problematic for models that are sensitive to outliers, such as k-nearest neighbors or linear models. In such cases, less aggressive techniques or a combination of oversampling with undersampling might be preferred.

4. Feature Type:

   - Check whether features are categorical or continuous. Some oversampling techniques work best with continuous features since they involve interpolation. Techniques like SMOTE and its variants typically assume continuous features. For categorical data, you might need to adapt these methods or choose techniques designed for categorical features.

5. Computational Resources:

   - Some oversampling methods are computationally intensive. If you're dealing with very large datasets or limited computational resources, you might prefer simpler or more efficient methods like random oversampling.

6. Domain-Specific Considerations:

   - Some domains might have specific requirements or constraints. For example, in medical diagnostics, false negatives might be more critical than false positives, affecting the choice of oversampling techniques.

7. Experimentation:

- Ultimately, the best way to determine the most effective oversampling technique for your dataset is through experimentation. Trying out different methods and tuning their parameters while monitoring performance on a validation set is often the most practical approach.