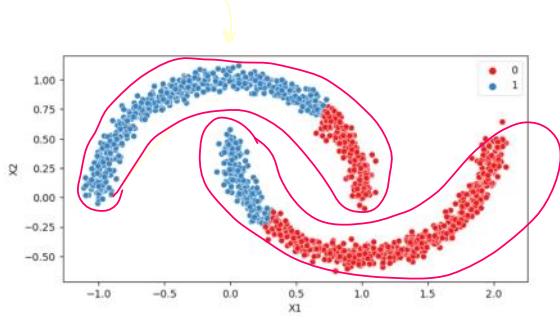
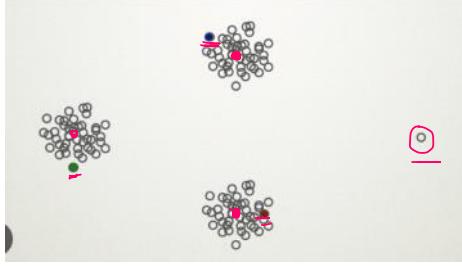
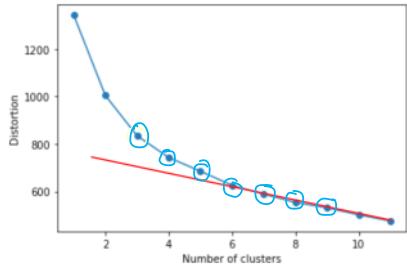


Lecture: 131 DBSCAN अंतर्गत Clustering Algorithm | पुस्तकाने प्राप्ति K-mean Clustering अंतर्गत |
-कैसे DBSCAN प्राप्ति? अंतर्गत K-mean के किसी limitation अवश्यक नहीं।

Why DBSCAN?

16 December 2023 15:10



Disadvantage of K-Mean Clustering :-

- कंप्यूटर elusion रहे लेटो बने दितु इम्‌। Elbow-method दितु यथा
- रूपूटर elusion रहे अवैतो वैष्ण नाले याणे एका उपयोग मतो वापरले
- कंप्यूटर elusion रहे लेटो असाम नासृ रक्षेश्वरी।

Outliers තුළු k-Mean clustering හෝ perform කුණු නා

 K-Mean diye karanje cluster का गुणकीय ढान shape दे रहे। यदि, clusters को complex हैं, स्पर्शरेत्रे बिंदु
मतज, तो K-mean करने पर perform करने मिलते हैं।

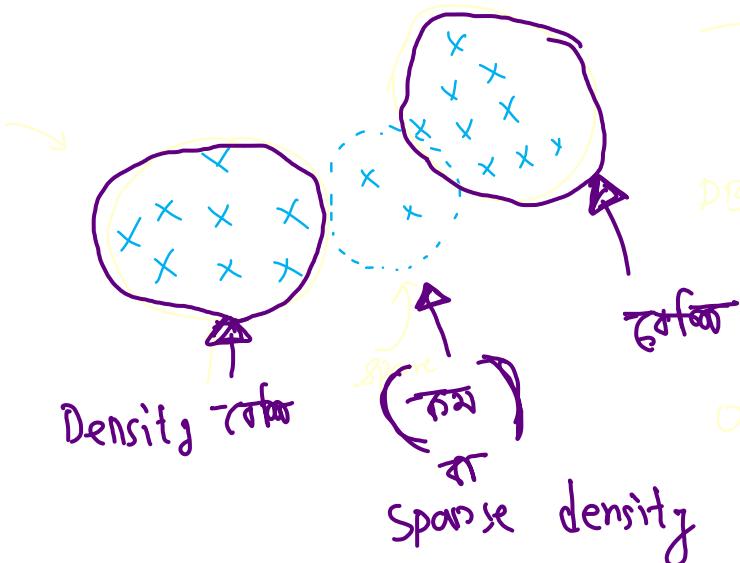
Clustering algorithms

- Centroid based clustering
- Density Based Clustering Algorithm



What is Density Based Clustering

16 December 2023 15:12



Kmean
↓
→ centroid
clustering

DBSCAN → Density Based Spatial
Clustering of Application
with Noise

OPTICS

Sparse density গুরুত্বের সমান ক্ষেত্রে, DBSCAN
অনিয়ন্ত্রিত ক্ষেত্রে, OPTICS, 2 টি মুক্ত এলসেন হতে পারে।

তাই,
DBSCAN → Density Based Spatial
Clustering of Application
with Noise

* Density Based আলগোরিদম আছে যার
নাম OPTICS

Hyperparameter, (Minpts & Epsilon) in DBSCAN

MinPts & Epsilon

16 December 2023 15:13

Density

density
hyperparameter

(E)epsilon ← {radius, #points (threshold)}

Minpts ← {#points}

x x x
x x x
dataset

1st

density
dense → ⑤
③

epsilon
↳ radius of
the neighbor

minPts
↳ threshold

draw
circle
with
 $\epsilon = 1$

$\epsilon = 1$ unit
sparse

2nd

epsilon → 1 unit

$\epsilon = 1$ unit
minPts = 3

1st ৰ, গোমাদ্বৰু dataset দেওয়া আছে । এটা $\epsilon = 1$ and minPts=3
দেওয়া আছে । প্রথম dataset ($\epsilon = 1$) র সম্মত radius
নিয়ে circle করা হি, তাৰপৰ যদি circle রে minPts
গুৰু নমান না হৈলে point থকে তাৰপৰ অন্তি Dense বলা
হৈবু না থকলে Sparse বলা ।

3 types of point in DBSCAN

Core Points, Border Points & Noise Points

16 December 2023

15:13

Border: ϵ द्वारा नेहबर हो पाये, total data point
MinPts का अधिक कम। इसका बीच $< \epsilon$ है। याहूं, ϵ के नेहबर होने की ओर बोर्डर होता है। याहूं
 ϵ के नेहबर का शाफ़ जास्ति होता है।

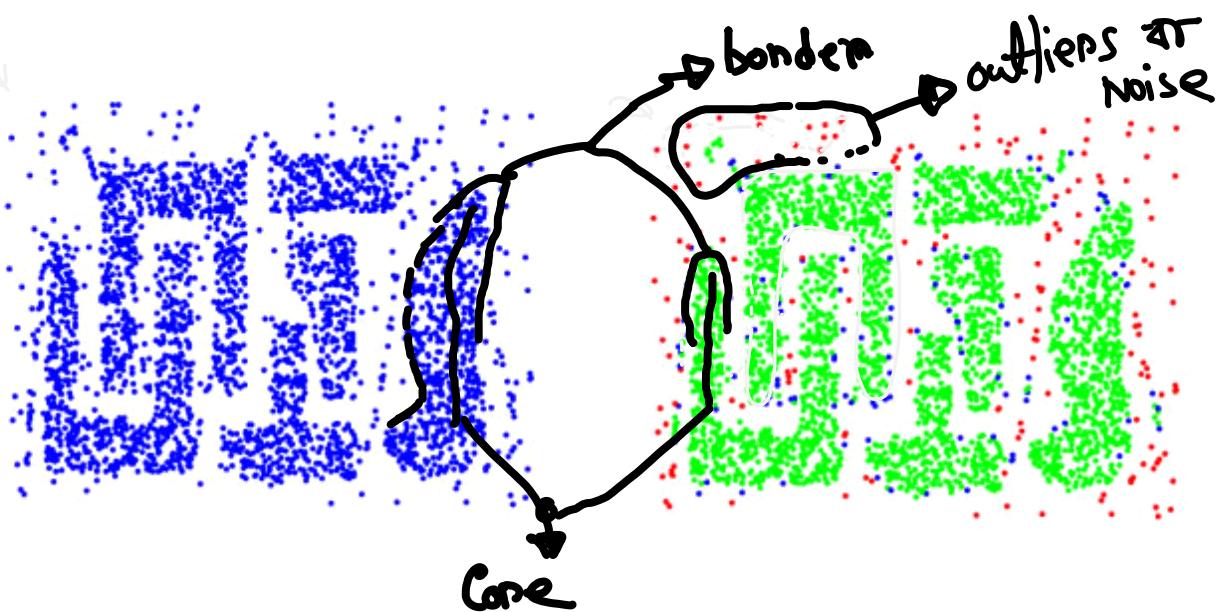
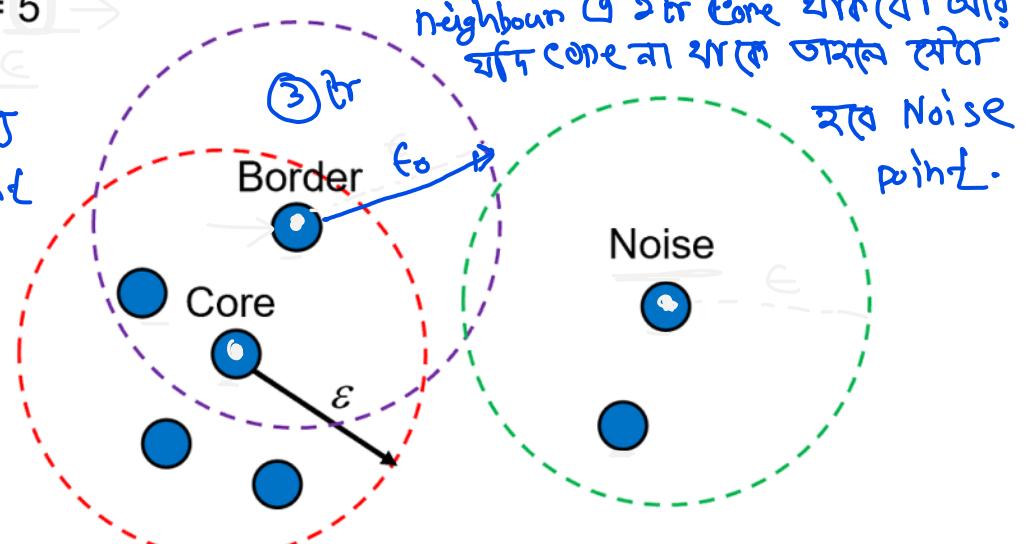
→ MinPts = 5 →

Core:

एक data point एवं उसके ϵ द्वारा नेहबर data point

जो मात्र minPts का मात्रा
हो तो उसे एक Core
point कहता है।

→ मात्र एक circular बनाता है।



Eps = 10, MinPts = 4

Density Connected Points

16 December 2023 15:13

Density Connected Points

ଏଥିର, A B ରୁଟ୍ଟେ point Density connected ହେଉଛି କାହିଁବେ ସେବଳେ ?

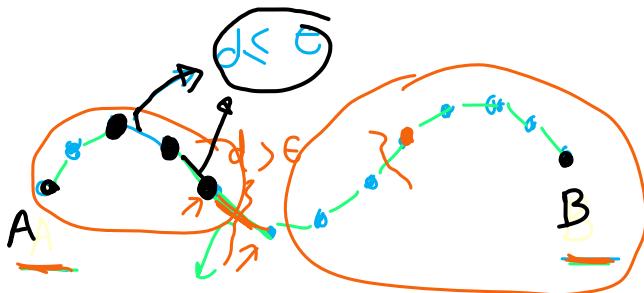


simply mean ଯାଏ, 2ଟେ point

A,B Density Connected ହେଲେ
ଏହିକୁ ବନ୍ଦର୍ମୁଖ କହ୍ନେ ପାଇବା
ଏହି cluster ହେଲାମିତି

ଫେଲାନ୍ତି ନାହିଁ ।

?



ଆଜି, A ରୁକ୍ଷେ B ରୁକ୍ଷେ d<epsilon ଏବଂ point ଏହି
ଆଧୁନିକୁ ହେଲାମିତି
(ମରିଗଲାଣ୍ଡ)

[Now let's learn the DBSCAN Algorithm]

DBSCAN Algorithm

16 December 2023 15:14

step 0 - We decide the value of E & minPTS:

Step 1 - Identify all points as either core point, border point or noise point

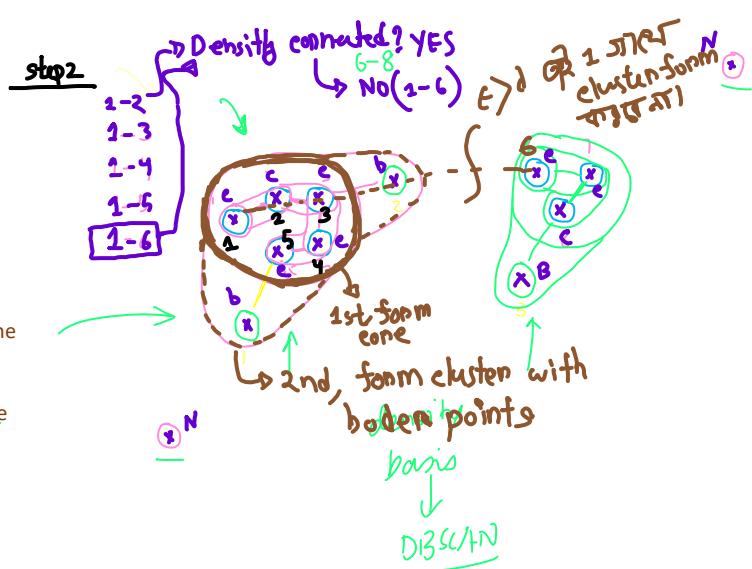
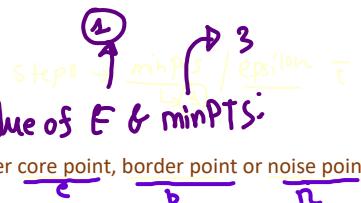
Step 2 - For all of the unclustered core points

Step 2a - Create a new cluster

Step 2b - add all the points that are unclustered and density connected to the current point into this cluster

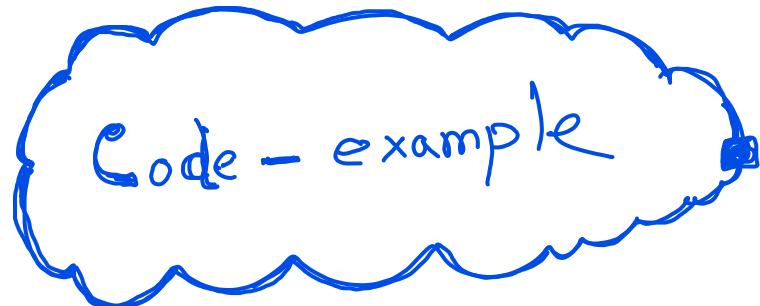
→ Step 3 - For each unclustered border point assign it to the cluster of nearest core point

Step 4 - Leave all the noise points as it is.



Code

16 December 2023 15:14



(Advantage & Disadvantage of DBSCAN)

Advantage

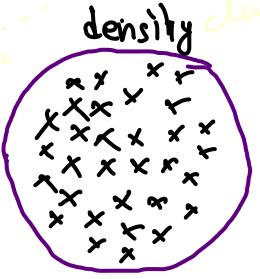
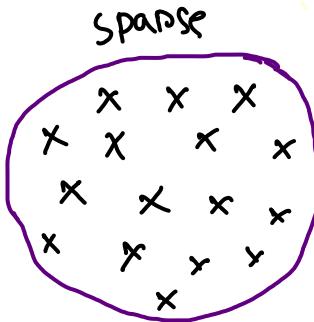
1. Robust to outliers
2. No need to specify clusters
3. Can find arbitrary shaped clusters
4. Only 2 hyperparameters to tune

জো অনলম্ব ডেক্টেশন (ব্যবহার করতে পারি)।

Disadvantage

1. Sensitivity to hyperparameters
2. Difficulty with varying density clusters
3. Does not predict

prediction prediction



এখান, ২টা cluster আছে, ফিল্টারো দেখি
রুটি, অবস্থা কম, কিন্তু, আমাদের E কৃত মান তা নেইস্থ ইন্দু
cluster কুচে হোগে। এই ইন্দুর data তে DBSCAN কাজে
perform করে না।

অবশ্যে মজার বিষয় ইটু, DBSCAN এ predict করে - কোথাও
নেই। কাবু, DBSCAN নমুনা Data দ্বাৰা স্পষ্ট prediction
কৰতে পারে না।

Visualization

16 December 2023 15:14

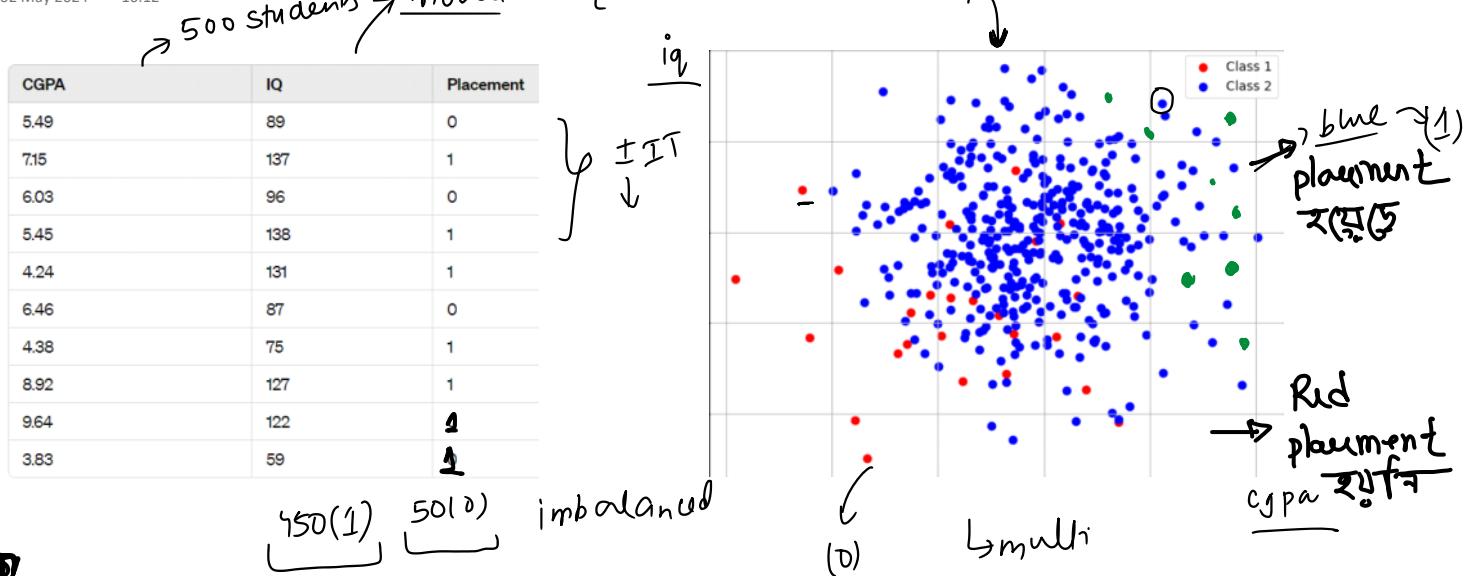


Imbalanced Data in mL
Lecture: 132

What is Imbalanced Data?

02 May 2024

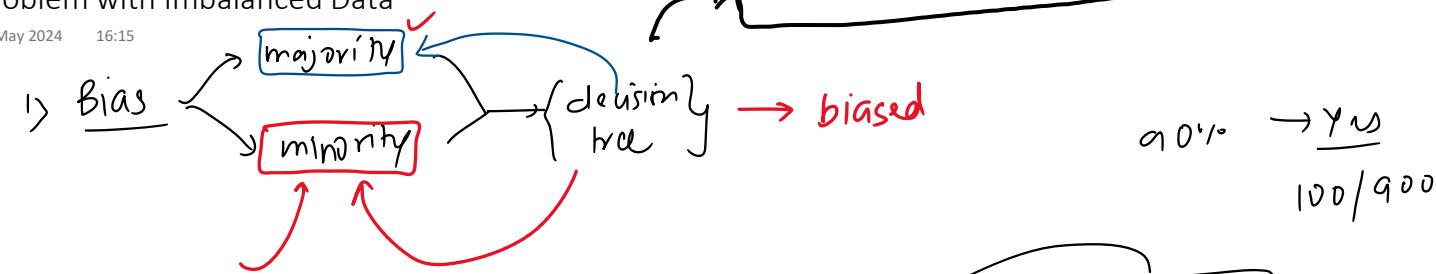
16:12



→ imbalance data বলতু এমন data কে হিসাবে for example
কৃষ্ণজি গোত্র placement হয়েছে ① আবু না হয়ে রয়ে আছে ②। এখন,
1 দুই ঘোড়া ০ দুই মুঠো অনেক কম, তাই হী dataset দিয়ে
model train করান, ① দিয়ে biased model পাওয়ে । এই, প্রশ্নের
dataset -কে imbalance dataset বলে ।

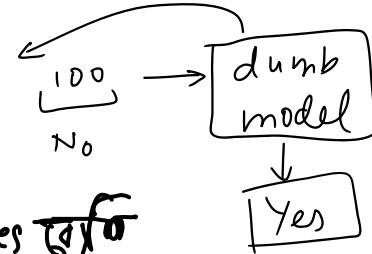
Problem with Imbalanced Data

02 May 2024 16:15



- 1> Bias → majority → decision tree → Biased
2> Metrics are not reliable → Accuracy →

$\frac{90}{10}$
Yes



কোনো যেস দ্রুতি

No কম, test করুন,

১০০ থেকে ৭০০ বলুন মাত্র আস দিবে

আরু ১০০ বলুন মাত্র, তা হলু, accuracy
৭০%। f1, F1, Recall, Precision

দ্বিতীয় স্বাক্ষর কে এখন কোনো কী ?

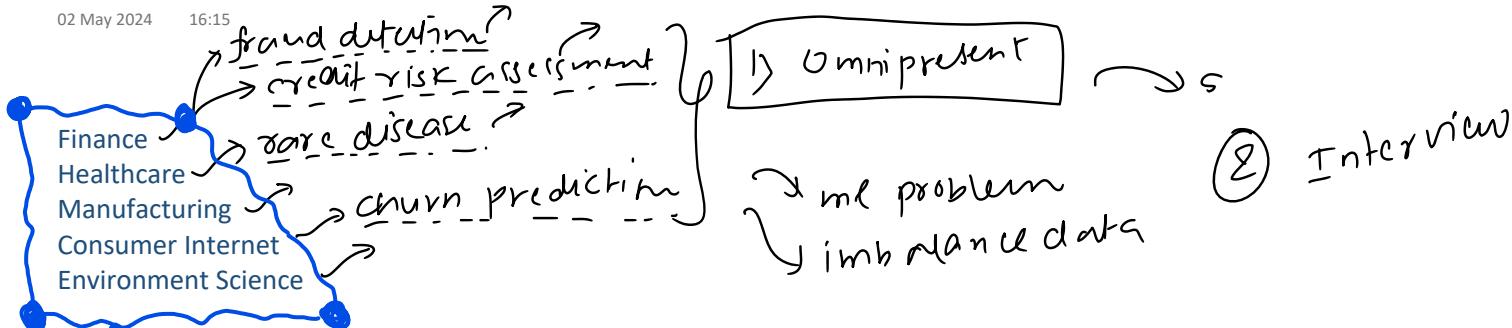
Code Example

(কোন imbalanced Data প্রক্রিয়া মাত্র) ?

Why studying Imbalanced Data is important?

02 May 2024

16:15



এই field গুলোতে এখনো ML ক্ষেত্রে ব্যবহৃত হয়। আবু এখানে
imbalanced data আকর্তু।

- Fraud Detection
- Credit Risk Assessment
- Rare Disease
- Churn Prediction
- Earthquake

→ এখান, শাজাহান সর্বোচ্চ ক্রিয়া পদক্ষেপ করে।
জনৈ পথ load দিয়ে ফেরত আবে না আম নোলাম করা।

Rare disease তাত্ক্ষণ্য কর মানুষের ইব।

ইস কর সংগ্রহ করে কোর্স platform ইস্টে।

যাবু, কোর্স এ দেশ region এ সেখ কীভ হো রয়।

অবু, এখানের data ইস্টে imbalanced।

→ Interview কর করে আগুন।

Some Technique to deal with imbalanced data.

i) Undersampling

ii) Oversampling

iii) SMOTE

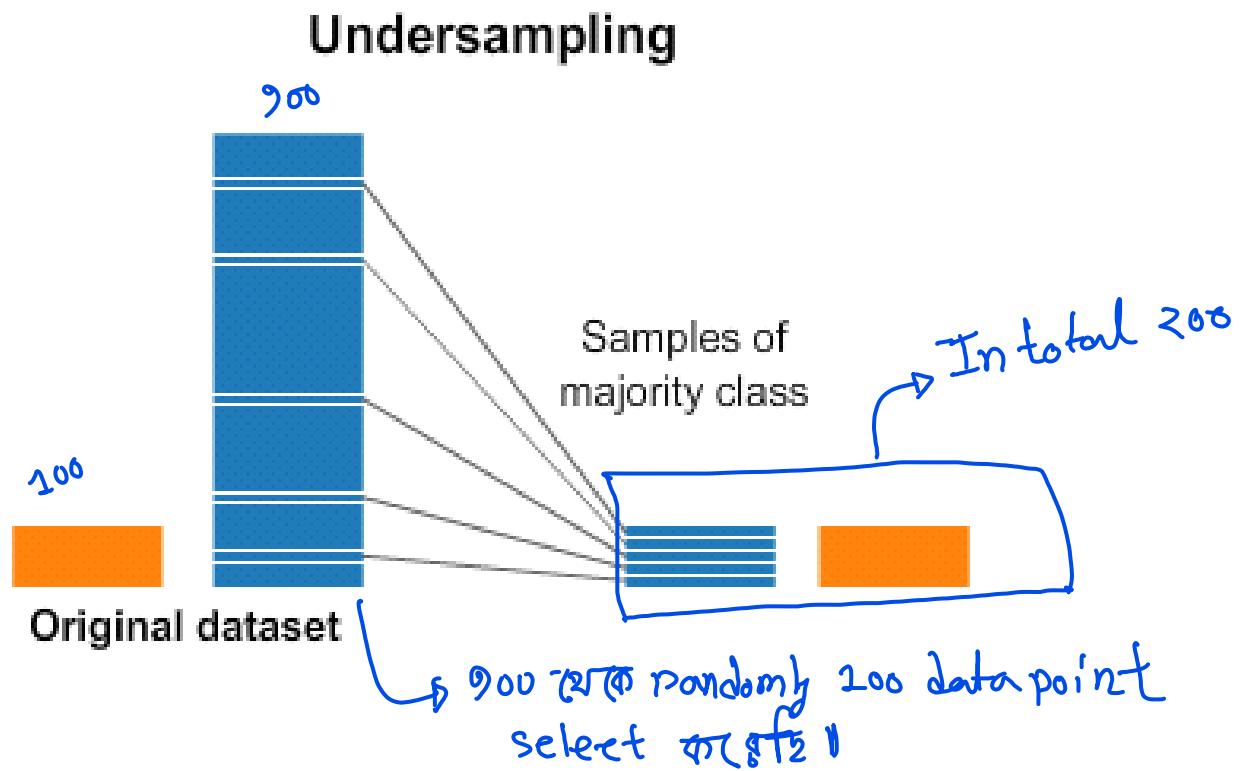
iv) Ensemble Learning

v) Cost Sensitive Learning

vi) Others

1. Undersampling

02 May 2024 16:17



✳️ Code example ✳️

Advantages :

1. Reduction in bias
2. Faster training

Disadvantages :

1. Information loss leading to underfitting
2. Sampling Bias

↳ Randomly data sample pick করার পরে

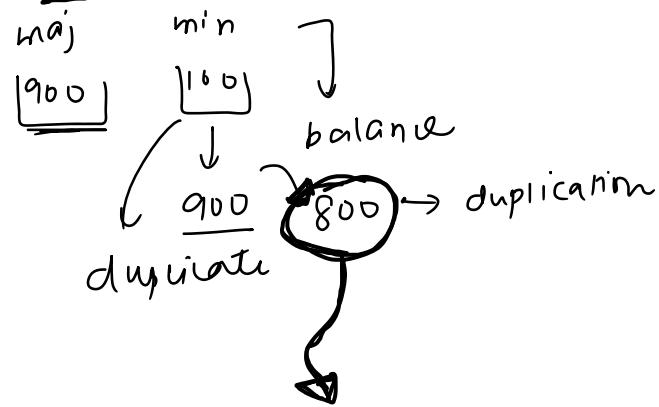
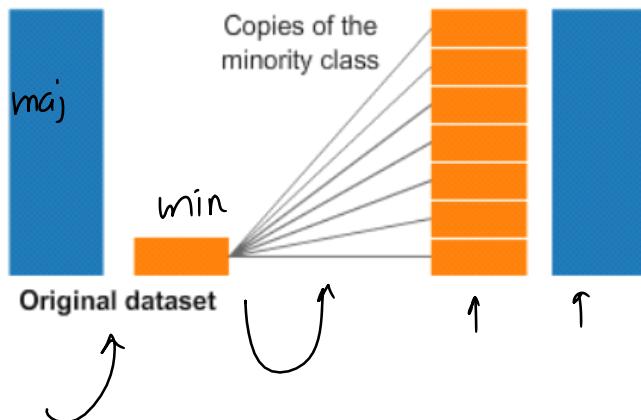
Bias হ্রাস করা।

2. Oversampling

03 May 2024 01:13

ପ୍ରଥାତି, datapoint ଅଦ୍ଧକଣ କରିବି

Random
Oversampling



ଏହାଟି 800 data duplicate
କରିବାକୁ ମାଧ୍ୟମେ achieve କରିବାକୁ

Code Example

Advantage

1. Reduced Bias ✓

Disadvantage

1. Increased size
2. Duplication of data may cause overfitting

ml ଓଲ୍ଡୁ
imp → overfitting
↓
dealing w/ ROS → overfitting

like, DT ଏବଂ ଗ୍ରେଟ୍ କିମ୍ବା
-ଫୁଲ, overfitting ହେଉ ଦେଖିବାକୁ

Synthetic Minority Over-sampling Technique:

3. SMOTE

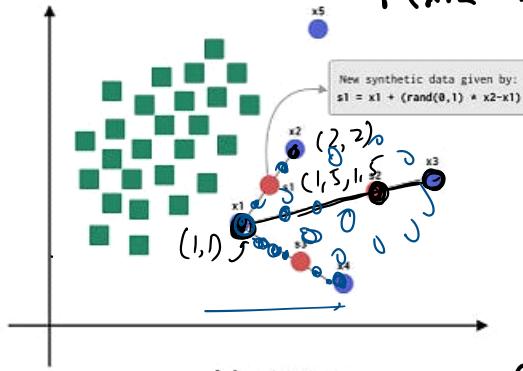
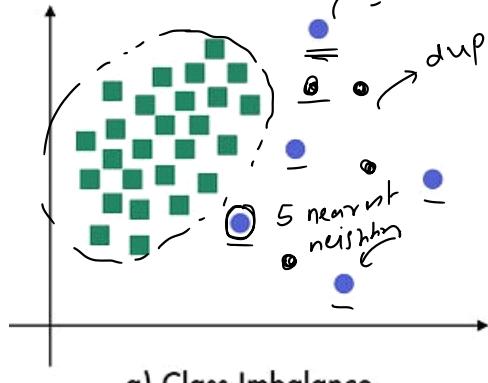
popular

Oversampling

→ min

oversampling

প্ৰযোগ, min value দ্বাৰা oversampling
কৰি, interpolation কৰিবলৈ
oversampling এ দ্বাৰা, duplicate কৰা
হৈলৈ।



Algorithm:

5 KNN

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

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

Code example

Disadvantages

$K=5$

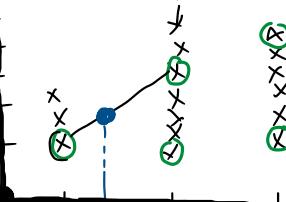
- Does Not Handle Categorical Data Well
- Computational Complexity (KNN-এর কালুজ)
- Dependency on the Choice of Neighbors
- Sensitive to Outliers
- Balance Achieved May Not Reflect True Nature

KNN এ, $K=1$ নিলৈ data point কৰা
generate কৰে, K সামৰ কৰা data point
কৰি generate কৰে। K গো optimum value
হৈবলৈ কৱিতু কৰিব, তাৰ মাধ্যমাবলী মত
হৈচাঁক হৈবলৈ কৰিব।

$0 \rightarrow K \rightarrow K=1$

$K=5$

Data generate কৰতু থাকি ঘোষণ না পৰ্যন্ত both
Sample কৰাৰ রেছে।



Categorical data

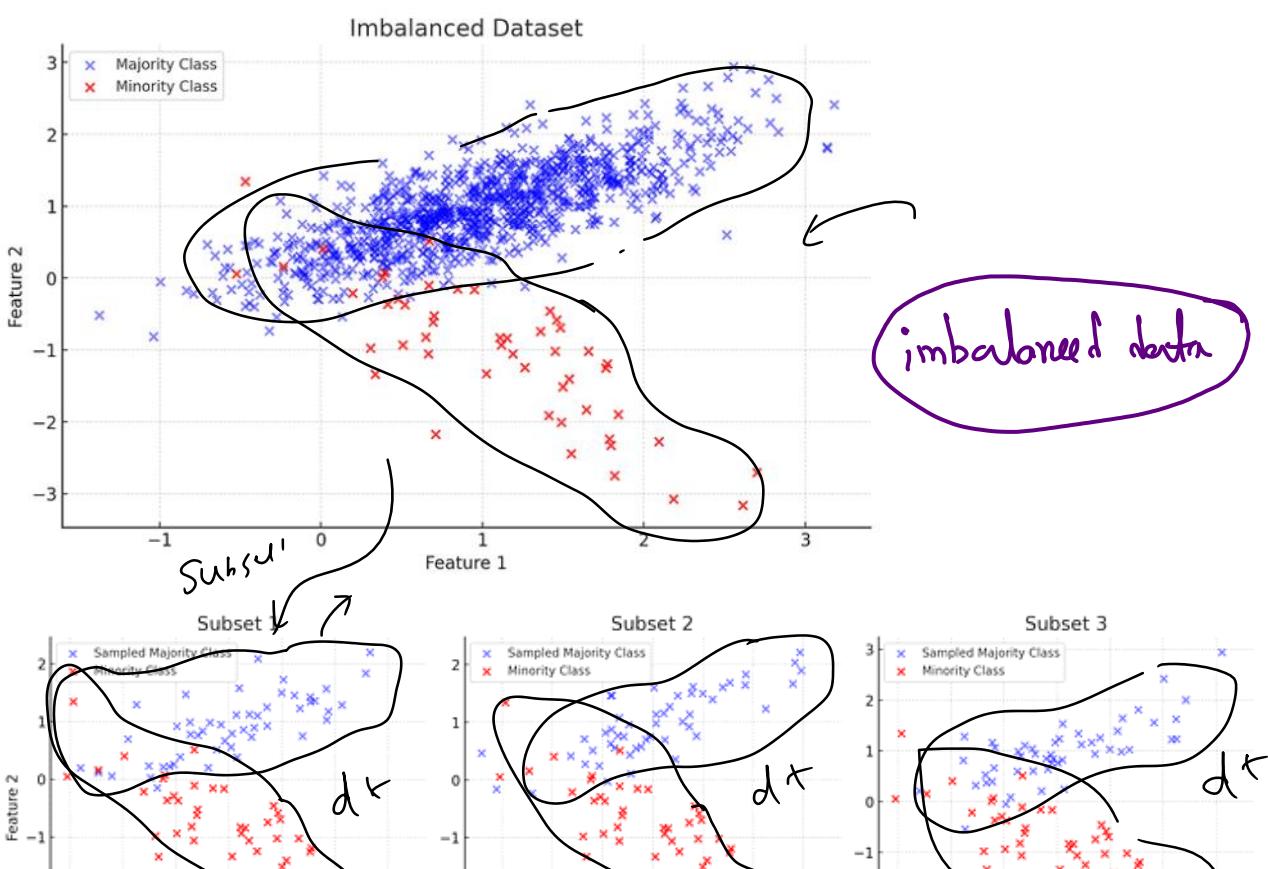
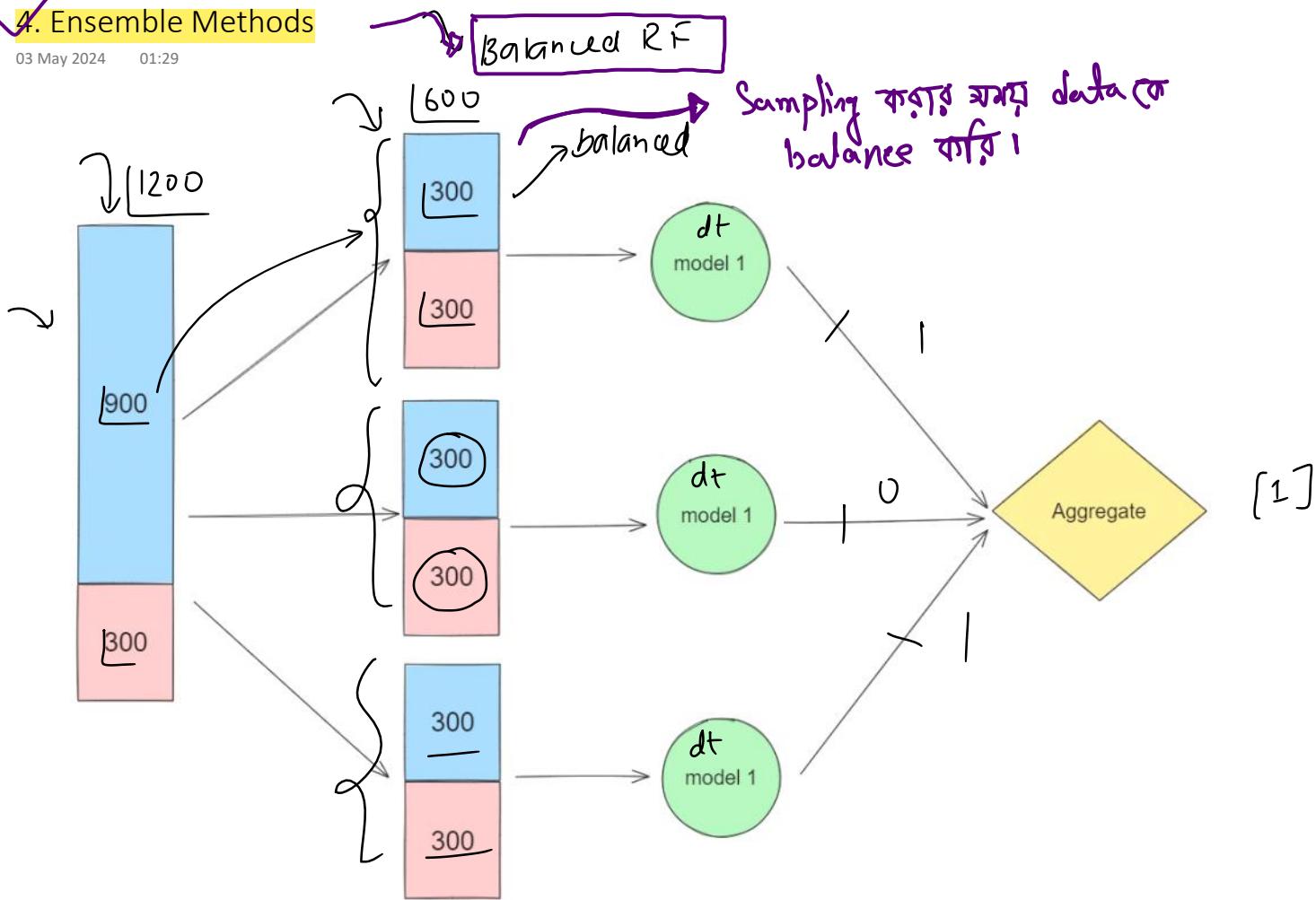
Interpolation (প্ৰযোগ কৰি) কৰা মাধ্যমে গ্ৰেডে data point
Generate কৰে এৰ দ্বাৰা গ্ৰেডে general value পৰিবৰ্তন না।

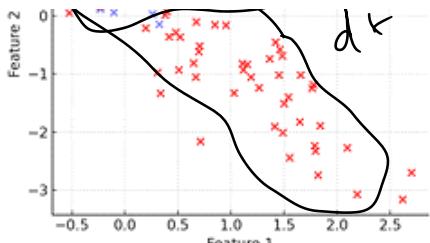
outliers
Generated Synthetic
point
data
Synthetic point দৃঢ়ভাৱে কৰতো perform
দ্বাৰা।
Outliers কৰি কৰে new generated
Synthetic point দৃঢ়ভাৱে কৰতো perform
দ্বাৰা।

→ Balance achieve কৰতো, এটা বাধুণ্যমূলক নাও হতো।

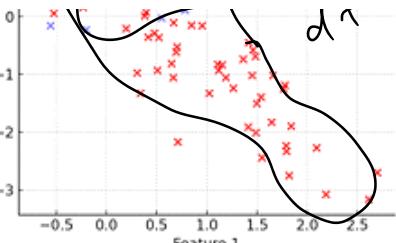
4. Ensemble Methods

03 May 2024 01:29

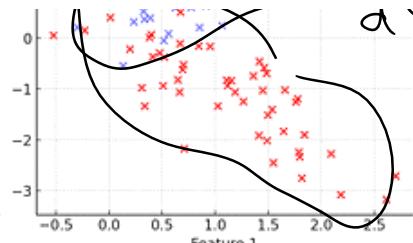




sample-1



Sample-2



Sample-3

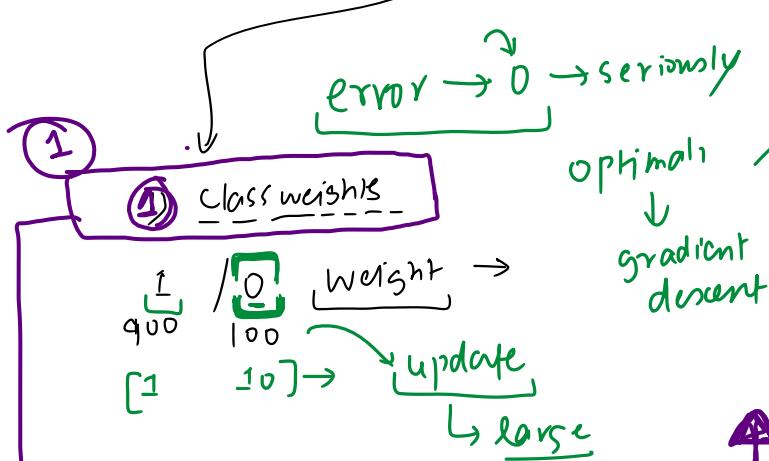
☞Boosting, ensemble learning এর এই algorithm
স্টোক ডেটা কে sample করি, তাহলে আমরা imbalanced
ডেটা কে handle করতে পারব।

5. Cost Sensitive Learning

03 May 2024 01:30

model का training process उत्तरांक रखने का chance नहीं
याकू फालि imbalance data handle करना !

learning process → imbalance data



याकू, class weight प्रियोप्ति।
Majority class के कम आये minority class के बढ़वा weightage प्रियोप्ति।
जरु, training process में यहां minority class के लिए error बाढ़ावा, इसे जाकू बढ़ावा देवा।

sk-learn में अनुकूल algorithm में,
class_weight वाले parameter की याद।

- i) LOR
- ii) SVM
- iii) DT

But, KNN, Naive Bayes एवं गुणात्र
प्राप्त नहीं।

2. custom loss fn design
- याकू, यहां: GB, xgbm एवं
custom loss fn व्यवहार करा
जाए।

List of all the techniques

03 May 2024 15:06

Imbalance learn library ର ଡୁକ୍ମେନ୍ସେୟୁୱନ୍ ଫୋର୍ସ୍ କାମ୍ପ୍ଯୁଟିଙ୍ଗ
follow କରୁଥିଲୁ | ଏହାର ବାବୀ ଅନ୍ତର୍ଭାବରେ
ଆହେ imbalance data handle କରୁଥିଲୁ ।