

## Machine Learning worksheet - 8

- ① (B)
- ② (A)
- ③ (C)
- ④ (C)
- ⑤ (D)
- ⑥ (B)
- ⑦ (C)
- ⑧ (B) & (C)
- ⑨ (C) & (D)
- ⑩ (B)

⑪ . The disadvantage of One Hot Encoding is that for high cardinality, the feature space can really blow up quickly and you started fighting with curse of dimensionality.

• In these cases One hot encoding is employed followed by PCA for dimensionality reduction.

(17) Techniques for handling imbalanced datasets:

→ Under Sampling.

- Random Under Sampling aims to balance class distribution by randomly eliminating majority class examples. This is done until majority & minority classes are balanced out.

→ Over Sampling.

- It increases the number of instances in minority classes by randomly replicating them in order to present a higher representation of the minority class in sample.

→ Cluster based Over Sampling

↳ Cluster based over sampling, K means algorithm is independently applied to minority and majority class instances.

~~Subsequently~~ Subsequently, each cluster is oversampled such that all clusters of the same class have an equal number of instances & all classes have same size.



- (13) • The key difference between ADASYN and SMOTE is that the former uses a density distribution as a criteria to automatically decide 'the number of synthetic samples that must be generated for each minority sample by adaptively changing the weights of the minority samples to compensate for the skewed distributions
- The latter generates the same number of synthetic samples for each original minority sample.