

1 Proposed Approach

Consider two datasets: $\mathcal{D}^a = (\mathbf{X}^a, \mathbf{y}^a)$ and $\mathcal{D}^b = \mathbf{X}^b$, where only \mathcal{D}^a was labeled. Both \mathbf{X}^a and \mathbf{X}^b represent feature vectors in the same space (\mathbb{R}^m) and may have different sizes in terms of the number of instances. From \mathcal{D}^a , we trained a model, an ensemble of *naïve* Bayes, to evaluate the importance of attributes in \mathbf{X}^a . The same model predicts labels for \mathbf{X}^b , and the predicted labels ($\hat{\mathbf{y}}^b$) are then utilized to evaluate the importance of features in \mathcal{D}^b .

For testing purposes, we utilized \mathcal{D}^b with known labels. However, it's important to note that the models did not have access to \mathbf{y}^b during the training phase. This information was solely used in the production of the method's quality measure (Δ_{acc}).

The model could provide two feature importance measures. They were named Difference Between Conditional Probabilities (DBCP) and Minimal Sufficient Set (MSS) [?], and were employed to compute four distinct strategies:

1. **DBCP**(\mathcal{D}^b): The DBCP method calculates the feature importances for \mathcal{D}^b .
2. **DBCP**(\mathcal{D}^b) – **DBCP**(\mathcal{D}^a): This strategy calculates the difference in DBCP feature importances between \mathcal{D}^b and \mathcal{D}^a .
3. **MSS**(\mathcal{D}^b): The MSS method calculates the feature importances for \mathcal{D}^b .
4. **MSS**(\mathcal{D}^b) – **MSS**(\mathcal{D}^a): This strategy calculates the difference in MSS feature importances between \mathcal{D}^b and \mathcal{D}^a .

The result of each strategy was min-max normalized to ensure values are between 0 and 1. The normalized values are then transformed into probability distributions by dividing each element by the sum of the absolute values of the same vector, following the L1 normalization. Each resulting probability distribution (Λ_k), where k represents a specific strategy, satisfies $\sum_{j=1}^m \Lambda_{kj} = 1, \forall k \in \{1, 2, 3, 4\}$. These distributions are then utilized as bias vectors in a Biased Random Subspace model.

Algorithm 1

```

1: function GETDBCPb( $\mathbf{X}^a, \mathbf{y}^a, \mathbf{X}^b$ )
2:   MODEL1TRAIN( $\mathbf{X}^a, \mathbf{y}^a$ )
3:    $\hat{\mathbf{y}}^b \leftarrow$  MODEL1PREDICT( $\mathbf{X}^b$ )
4:   MODEL2TRAIN( $\mathbf{X}^b, \hat{\mathbf{y}}^b$ )
5:   return MODEL2FEATUREIMPORTANCES
6: end function

```

Algorithm 2

```
1: function GETDBCPb – DBCPa( $\mathbf{X}^a, \mathbf{y}^a, \mathbf{X}^b$ )
2:   MODEL1TRAIN( $\mathbf{X}^a, \mathbf{y}^a$ )
3:   DBCPa  $\leftarrow$  MODEL2FEATUREIMPORTANCES
4:   DBCPb  $\leftarrow$  GETDBCPb( $\mathbf{X}^a, \mathbf{y}^a, \mathbf{X}^b$ )
5:   return PROBABILITYDISTRIBUTION(DBCPb – DBCPa)
6: end function
```

Algorithm 3

```
1: function GETMSSb( $\mathbf{X}^a, \mathbf{y}^a, \mathbf{X}^b$ )
2:   MODELTRAIN( $\mathbf{X}^a, \mathbf{y}^a$ )
3:   return MODELMINIMALSUFFICIENTSET( $\mathbf{X}^b$ )
4: end function
```

Algorithm 4

```
1: function GETMSSb – MSSa( $\mathbf{X}^a, \mathbf{y}^a, \mathbf{X}^b$ )
2:   MODELTRAIN( $\mathbf{X}^a, \mathbf{y}^a$ )
3:   MSSa  $\leftarrow$  MODELMINIMALSUFFICIENTSET( $\mathbf{X}^a$ )
4:   MSSb  $\leftarrow$  MODELMINIMALSUFFICIENTSET( $\mathbf{X}^b$ )
5:   return PROBABILITYDISTRIBUTION(MSSb – MSSa)
6: end function
```

Algorithm 5

```
1: function PROBABILITYDISTRIBUTION( $\mathbf{v}$ )
2:    $\mathbf{v} \leftarrow$  MINMAXNORMALIZATION( $\mathbf{v}$ )
3:   return  $\frac{\mathbf{v}}{|\mathbf{v}|_1}$  // return  $\mathbf{v}$  normalized by L1 norm
4: end function
```

Algorithm 6

```
1: function DELTAACC( $\mathbf{X}^a, \mathbf{y}^a, \mathbf{X}^b, \mathbf{y}^b, \Lambda_k$ )
2:    $\mathbf{X}_{train}, \mathbf{y}_{train}, \mathbf{X}_{test}, \mathbf{y}_{test} \leftarrow$  TRAINTESTSPLIT( $\mathbf{X}^a, \mathbf{y}^a, 10\%$ )
3:   MODEL1TRAIN( $\mathbf{X}_{train}, \mathbf{y}_{train}$ )
4:    $\hat{\mathbf{y}}_{ID} \leftarrow$  MODEL1PREDICT( $\mathbf{X}_{test}$ ) // ID means in distribution
5:   MODEL2(biasedSubspaces=True,  $\Lambda_k$ )
6:   MODEL2TRAIN( $\mathbf{X}^a, \mathbf{y}^a$ )
7:    $\hat{\mathbf{y}}_{OOD} \leftarrow$  MODEL2PREDICT( $\mathbf{X}^b$ ) // OOD means out of distribution
8:   return ACC( $\hat{\mathbf{y}}_{OOD}, \mathbf{y}^b$ ) – ACC( $\hat{\mathbf{y}}_{ID}, \mathbf{y}_{train}$ )
9: end function
```
