

Projeto - Estratégias de Sampling (Métricas)

1. (10 points) Implemente em R:

1. Uma função *houdout* que receba como parâmetro um dataset e o limiar para divisão. Esta função deverá selecionar aleatoriamente as amostras para compor o dataset. Como saída a função deverá retornar o dataset para teste.
2. Uma função *houdoutStratified* que receba como parâmetro um dataset e o limiar para divisão. Esta função deverá selecionar aleatoriamente as amostras para compor o dataset resultante de modo que seja garantidas amostras para treinamento e teste considerando a divisão. Como saída a função deverá retornar o dataset para teste.
3. Uma função chamada *accuracyFromMatrix* que receba uma Matriz de Confusão e retorne a acurácia da modelagem.
4. Uma função chamada *precisionFromMatrix* que receba uma Matriz de Confusão e retorne a precisão da modelagem.
5. Uma função chamada *recallFromMatrix* que receba uma Matriz de Confusão e retorne a revocação da modelagem.
6. Uma função chamada *fScoreFromMatrix* que receba uma Matriz de Confusão e retorne o valor de F-score da modelagem.

2. (10 points) Usando a base “iris” como base para discussão dos resultados e validação dos experimentos, faça uma comparação entre os algoritmos RandomForest (pacote RandomForest) e SVM (pacote e1071) para a classificação com diferentes estratégias de amostragem para avaliação dos resultados, usando:

- Houdout;
- Houdout com Estratificação;
- Crossvalidation 2 folds;
- Crossvalidation 7 folds;
- Crossvalidation 12 folds;
- Leave one out;

3. Realize 30 iterações dos métodos comentados no exercício anterior e por meio de um Boxplot, discuta os resultados de desempenho e estabilidade. Calcule a métrica de MSE (Mean Square Error) para auxiliar a discussão.

Leitura recomendada:

Japkowicz, Nathalie, and Mohak Shah. Evaluating learning algorithms: a classification perspective. Cambridge University Press, 2011.