

O processo completo de descoberta de um novo fármaco envolve um alto custo financeiro, além de um longo período para seu desenvolvimento. Visando a otimização desse processo e a redução dos custos, inclui-se a essa abordagem uma etapa *in silico* para a predição e avaliação de interações entre pequenas moléculas (ligantes) candidatas e proteínas-alvo (receptores). Essa predição pode ser realizada utilizando algoritmos de docagem molecular que buscam prever a melhor conformação/orientação de um ligante em um sítio de ligação de um receptor e estimar a afinidade desse complexo. Para realizar essas estimativas e prover um resultado confiável, muitas funções de escore tem sido desenvolvidas, obtendo computacionalmente a afinidade de ligação entre o complexo proteína-ligante. As funções de escore são divididas em categorias de acordo com a metodologia utilizada: baseadas em Física, empíricas, *Knowledge-based* e baseadas em aprendizado de máquina (AM). As funções de escore baseadas em AM dependem fortemente do conjunto de dados que são utilizados no seu treinamento. Estes conjuntos de dados normalmente são compostos por um grande número de características e instâncias, ocasionando em um grande fluxo de dados e um alto custo computacional para realização dos treinamentos. Além disto, em conjuntos de dados com um grande número de atributos, o desempenho do estimador tende a mostrar resultados menos precisos a partir de um certo número de atributos, mesmo que estes sejam úteis. Este fenômeno é conhecido como “A Maldição da Dimensionalidade”, e deve-se ao fato de que a densidade de um conjunto de dados exponencialmente integrado aumenta à medida que o número de recursos usados aumenta. Para reduzir a dimensionalidade podem ser utilizados métodos de seleção de atributos buscando um conjunto de atributos mais relevantes para o modelo com o objetivo de aumentar a qualidade dos resultados e reduzir o tempo de treinamento. O presente trabalho demonstra o impacto da redução de dimensionalidade em um conjunto de dados de dimensão 4,152 (instâncias) \times 502 (atributos), utilizando o *Random Forest* (RF) como algoritmo de seleção dos atributos. Nesse conjunto obtido do PDBBind2018, as instâncias correspondem a complexos proteína-ligante e os atributos são descritores do receptor, ligante ou do complexo. Além disto, a função de escore obtida neste trabalho é gerada a partir do treinamento de um modelo com redes neurais. Os resultados mostram a qualidade dos modelos gerados com diferentes dimensões, fornecendo assim as informações referentes a onde ocorre a curva de perca de qualidade dos resultados para este conjunto de dados de treinamento.