

Aprendizagem Automática 2

Trabalho Prático

-

Mestrado em Engenharia Informática
Universidade do Minho

Grupo nº 11

PG41080	João Ribeiro Imperadeiro
PG41081	José Alberto Martins Boticas
PG41091	Nelson José Dias Teixeira
PG41851	Rui Miguel da Costa Meira

6 de Maio de 2020

Conteúdo

1	Introdução	3
2	Implementação	4
2.1	Estrutura	4
2.1.1	Problema	4
2.1.2	Algoritmos	4
2.1.2.1	Classificação	5
2.1.2.2	Regressão	5
2.1.3	Métricas	5
2.1.3.1	Classificação	6
2.1.3.2	Regressão	6
2.1.4	Otimização	6
3	Retrospectiva e Prospetiva	7
4	Conclusão	8
A	Observações	9

Lista de Figuras

Capítulo 1

Introdução

No âmbito da unidade curricular (UC) *Aprendizagem Automática II* (AA2), foi requerida a realização de um trabalho prático para avaliação. Tal como foi proposto a 23 de abril, o grupo escolheu a opção relativa ao "desenvolvimento de algoritmos/*software* no âmbito da Aprendizagem Máquina". Mais especificamente, optou-se pelo desenvolvimento de uma *framework* de *AutoML*, com o objetivo de obter o melhor modelo para problemas de *supervised learning* e *unsupervised learning*, de forma automática e com a menor intervenção possível por parte do programador. À *framework* idealizada foi atribuído o nome *UnicornML*.

Por sugestão do docente da UC, foram postos de parte os problemas de *unsupervised learning*, pela sua complexidade e menor atenção dada durante as aulas. Assim, sobram apenas os problemas de *supervised learning* que podem ser divididos em duas categorias: classificação e regressão. Mais à frente serão abordadas as duas categorias em pormenor.

Atualmente, a *UnicornML* é capaz de proceder à identificação e distinção entre problemas de regressão e de classificação, sendo que existe também a possibilidade de o utilizador fornecer essa informação, através de uma *flag*, poupando esse trabalho extra. Para além disto, o utilizador pode ainda indicar, da mesma forma, quais os algoritmos e métricas que pretende que sejam testados e, ainda, qual das otimizações prefere (*randomized search*, otimização bayesiana ou nenhuma). Todos estes detalhes são devidamente explicados nos próximos capítulos.

Capítulo 2

Implementação

2.1 Estrutura

A *UnicornML* tem uma estrutura simples. Existe uma classe principal, com o mesmo nome da *framework*, que permite treinar dados supervisionados para problemas de classificação e regressão. Neste nível, são processadas todas as informações iniciais, que podem incluir:

- **Problema** - se se trata de um problema de regressão ou classificação; caso não seja fornecida essa informação, é identificado o tipo do problema em questão.
- **Algoritmos** - quais os algoritmos, dentro dos disponibilizados, que o utilizador pretende que sejam testados (passado em forma de lista);
- **Métricas** - métricas a avaliar (passado em forma de lista).

Todas estas informações são passadas pelo terminal e são opcionais. Caso não sejam fornecidas, serão testadas todas as hipóteses oferecidas pela *framework*. O utilizador pode indicar o problema e os algoritmos, só os algoritmos, só o problema ou nenhum dos dois, mas deverá sempre ter em conta as limitações e consequências de cada uma das escolhas tomadas.

2.1.1 Problema

Os problemas de aprendizagem supervisionada podem ser divididos em dois conjuntos: problemas de regressão e problemas de classificação. Como tal, é importante perceber-se qual dos dois problemas enfrentamos, de forma a que se possa poupar tempo e recursos de computação na procura do melhor modelo. Para isso, foi pensada uma forma de identificar o tipo do problema. No entanto, tal como referido na proposta já entregue, esta não é uma prioridade, pelo que o método para já utilizado é simples e identifica apenas a presença de inteiros ou *floats* para fazer esta distinção.

No entanto, este processamento é evitado se o utilizador indicar qual dos problemas os seus dados representam. Esta indicação é dada através de uma opção, sendo passada uma de duas *strings*: *Regression* ou *Classification*.

2.1.2 Algoritmos

A *UnicornML* oferece diversos algoritmos para cada um dos tipos de problemas. O utilizador pode escolher, dentro dos algoritmos disponíveis, quais os que quer que sejam testados. No entanto, os algoritmos só serão testados se estiverem disponíveis para o tipo de problema identificado pela *framework* ou indicado pelo mesmo.

Caso o utilizador não indique quais os algoritmos que prefere que sejam testados, a *framework* testará todos os algoritmos disponíveis para o tipo de problema identificado pela mesma ou indicado pelo utilizador.

2.1.2.1 Classificação

Os algoritmos disponíveis para problemas de classificação são os seguintes:

- Regressão logística;
- *K-Nearest Neighbors* (KNN);
- *Support Vector Classifier* (SVC) - uma *Support Vector Machine* (SVM) para classificação;
- *kernel SVM* - uma SVM com uma função *kernel*, que permite a classificação em espaços de dimensão superiores;
- Classificadores Bayesianos - família de classificadores baseados na teoria de Bayes. Foram implementados quatro algoritmos diferentes: *Gaussian*, *Multinomial*, *Bernoulli* e *Complement*;
- Árvore de decisão;
- *Random Forest* - operam construindo uma multitude de árvores de decisão.

Estes algoritmos encontram-se na classe `Regression`, em `unicornML/regression/_-__init__.py`.

2.1.2.2 Regressão

Os algoritmos disponíveis para problemas de regressão são os seguintes:

- Regressão linear;
- Regressão polinomial - os polinómios testados variam entre grau 2 e grau igual ao número de colunas da base de dados;
- *Support Vector Regressor* (SVR) - uma SVM para regressão;
- Árvore de decisão;
- *Random Forest* - operam construindo uma multitude de árvores de decisão.

Estes algoritmos encontram-se na classe `Classification`, em `unicornML/classification/_-__init__.py`.

2.1.3 Métricas

As métricas permitem avaliar o desempenho de um certo modelo. O utilizador também pode escolher as métricas que irão ser tomadas em consideração e, posteriormente, apresentadas. Mais uma vez, isso está limitado às métricas disponíveis para cada tipo de problema. De realçar que nem todas as métricas podem estar disponíveis num determinado momento.

2.1.3.1 Classificação

As métricas disponíveis para problemas de classificação são as seguintes:

- *Accuracy*;
- *F1*;
- *Precision*;
- *Recall*.

2.1.3.2 Regressão

As métricas disponíveis para problemas de regressão são as seguintes:

- *R-squared* (R^2);
- *Mean Square Error* (MSE).

2.1.4 Otimização

A classe `Model`, em `unicornML/model/__init__.py`, é o coração de toda a *framework*. A mesma foi pensada de forma a simplificar o restante código e reduzir duplicações do mesmo. É, ainda, onde é feita a procura do melhor modelo, segundo as opções escolhidas pelo utilizador. Esta classe é utilizada pelas classes `Regression` e `Classification`.

Encontramos aqui mais uma opção que será futuramente disponibilizada ao utilizador: a escolha do método de otimização - *grid search* (*randomizedSearch*) ou otimização bayesiana (Bayes). Por defeito, o método de otimização utilizado será `randomizedSearch`. Futuramente será também possível optar pela otimização bayesiana, embora esta ainda não esteja implementada.

Capítulo 3

Retrospectiva e Prospetiva

O trabalho realizado até agora incidiu muito no planeamento da UnicornML, de forma a que a mesma fosse o mais simples possível de implementar, perceber e modificar. Numa primeira fase, muitos dos esforços se dedicaram à estrutura que a *framework* hoje apresenta e que serviu de base para todo o trabalho até agora desenvolvido e todo o que virá ainda a realizar-se. Estes esforços de planeamento e pensamento traduziram-se num código estruturado e fácil de compreender e desenvolver, pelo que se perspectiva que não existirão grandes dificuldades no trabalho que falta desenvolver.

Importa realçar que tudo o que foi descrito neste relatório está já implementado, embora nem tudo esteja totalmente operacional, o que encaramos como absolutamente normal.

Para o futuro próximo e até ao prazo limite, prevê-se o seguinte:

- Implementação de redes neuronais, tanto para problemas de classificação como de regressão;
- Verificar a existência ou não de *overfitting* e sua resolução;
- Disponibilização de todas as métricas apresentadas;
- Disponibilização do método de otimização bayesiana;
- Correção de eventuais *bugs*;
- Otimizações gerais.

Para o relatório final, será dado um maior detalhe no que toca a cada algoritmo desenvolvido e a cada escolha tomada.

Capítulo 4

Conclusão

Será um pouco prematuro tentar tirar ilações de um trabalho incompleto, pelo que esta conclusão permite apenas refletir sobre as perspectivas existentes, tendo em conta o trabalho desenvolvido até agora. No entanto, é possível perceber que a *UnicornML* está no bom caminho, tendo em conta os resultados obtidos para os métodos já implementados e a ideia do que será possível ainda desenvolver.

Apêndice A

Observações

- Documentação *Python 3*:
<https://docs.python.org/3/>
- Documentação *scikit-learn* - API:
<https://scikit-learn.org/stable/modules/classes.html>
- Documentação *scikit-learn* - *supervised learning*:
https://scikit-learn.org/stable/supervised_learning.html#supervised-learning