

## Redes neuronais com protótipos para quantificar os determinantes do desempenho acadêmico: evidências de um país europeu

A. Beatriz-Afonso (NOVA INFORMATION MANAGEMENT SCHOOL)<sup>1</sup>

F. Cruz-Jesus (NOVA INFORMATION MANAGEMENT SCHOOL)<sup>2</sup>

C. Nunes (NOVA INFORMATION MANAGEMENT SCHOOL)<sup>3</sup>

M. Castelli (NOVA INFORMATION MANAGEMENT SCHOOL)<sup>4</sup>

T. Oliveira (NOVA INFORMATION MANAGEMENT SCHOOL)<sup>5</sup>

### Resumo

Desde os anos 50 do século passado que o desempenho acadêmico tem sido foco de interesse por parte de investigadores e decisores políticos. No entanto, apenas recentemente os métodos de ciência de dados começaram a ser aplicados de forma mais sistemática a este tema. Este trabalho utiliza os dados dos exames nacionais de matemática e português da população portuguesa no ano letivo 2018/2019 para, através de redes neuronais, avaliar e comparar quais os fatores que afetam os resultados desses exames, e de que forma. Além disso, uma nova abordagem é apresentada para lidar com o dilema da "caixa negra" dos métodos de ciências de dados mais avançados. Esta abordagem passa pela criação de um conjunto de protótipos através de Redes Neuronais, fornecendo uma estimativa de quanto cada potencial impacta o desempenho acadêmico.

*Palavras-chave: Educação; Sucesso escolar; Ciência de dados*

### Abstract

Since the 1950s, academic performance has been the focus of interest by researchers and policymakers. However, only recently have data science methods begun to be applied more systematically to this topic. This work uses data from national mathematics and Portuguese exams of the Portuguese population in the 2018/2019 school year to, through neural networks, evaluate and compare which factors affect the results of these exams and in what way. Furthermore, a new approach is presented to deal with the "black box" dilemma of more advanced data science methods. This approach involves creating a set of prototypes through Neural Networks and estimating how much each potential impacts academic performance.

*Keywords: Education, Academic Achievement, Data Science*

---

1 Contato: [aafonso@novaims.unl.pt](mailto:aafonso@novaims.unl.pt)

2 Contato: [fjesus@novaims.unl.pt](mailto:fjesus@novaims.unl.pt)

3 Contato: [cnunes@novaims.unl.pt](mailto:cnunes@novaims.unl.pt)

4 Contato: [mcastelli@novaims.unl.pt](mailto:mcastelli@novaims.unl.pt)

5 Contato: [toliveira@novaims.unl.pt](mailto:toliveira@novaims.unl.pt)

## Introdução

A educação é o primeiro pilar dos direitos sociais na União Europeia (EUROPEAN COMMISSION, 2017), visto que, está intrinsecamente relacionada a praticamente todos os aspectos da vida e do bem-estar. Por meio de uma melhor educação, os indivíduos têm acesso a melhores empregos e são mais propensos a obter maiores rendimentos, o que, por sua vez, está associado a melhores condições materiais e de saúde. Além dos aspectos financeiros, a educação também proporciona melhor acesso à cultura e ao conhecimento, tornando os indivíduos mais aptos a lidar com as complexidades e decisões de suas próprias vidas. Adicionalmente, os governos também podem beneficiar, arrecadando mais impostos e reduzindo subsídios relacionados a questões socioeconômicas, como pobreza e desemprego, diminuindo a dependência do apoio financeiro e impulsionando a economia (MÜNICH; PSACHAROPOULOS, 2018). Além disso, cidadãos educados trabalharão para um mundo mais democrático e sustentável, tomando decisões informadas e responsáveis, por exemplo, exercendo seu direito de voto (COUNCIL OF EUROPE, 2018; MÜNICH; PSACHAROPOULOS, 2018). Essencialmente, a educação é a base para um mundo em paz, justo e inclusivo. A Organização das Nações Unidas (ONU) adotou a diretriz "ninguém será deixado para trás" (UNITED NATIONS, 2015, p. 5) como meta educacional, onde o objetivo é que todos os jovens, independentemente de seu gênero, deficiência ou situação, conclua o ensino primário e secundário com resultados práticos e habilidades relevantes até 2030 (UNITED NATIONS, 2015).

O desempenho acadêmico (DA) é o resultado da educação. Ao compreender o amplo espectro de fatores subjacentes que levam a disparidades em DA, os decisores políticos podem implementar estratégias mais eficazes e precisas (CRUZ-JESUS *et al.*, 2020; MIGUÉIS *et al.*, 2018), seja em nível escolar ou governamental, resultando na melhoria do desempenho dos alunos e na redução das lacunas educacionais (POKROPEK; BORGONOV; JAKUBOWSKI, 2015; VANDELANNOTE; DEMANET, 2020).

Embora estes fatores sejam um foco de interesse para muitos investigadores, a grande maioria ainda usa métodos estatísticos clássicos em amostras de alunos (CRUZ-JESUS *et al.*, 2020). O uso de métodos de Inteligência Artificial (AI) para entender o DA só começou a crescer recentemente (MUSSO; HERNÁNDEZ; CASCALLAR, 2020). Apesar da reconhecida superioridade em termos de desempenho, esses métodos possuem sobre os clássicos (por exemplo, regressão), o seu uso no contexto do DA ainda é limitado (COSTA-MENDES *et al.*, 2020), especialmente para fins de interpretação e não apenas de previsão. Esse facto pode ser explicado pelo problema da “caixa negra” das técnicas de IA e ML de

última geração (ver, por exemplo, Lipton (2018)). Contribuímos para a literatura fornecendo o que é, tanto quanto sabemos, uma abordagem inovadora neste contexto. Treinamos modelos de última geração com dados de praticamente todos os alunos matriculados numa escola pública portuguesa do ensino secundário, e depois aplicamos esses modelos a um conjunto de alunos “artificiais” (protótipos). Posteriormente avaliámos como as previsões dos modelos mudam ao longo dos diferentes cenários representados por cada protótipo. Assim, com este trabalho, aproveitámos ao máximo os métodos de IA e ML de última geração para esclarecer os impulsionadores do DA, quantificando o peso que cada um dos impulsionadores mais conhecidos da literatura (por exemplo, alunos, pais, escolas, e características dos professores) tem nos exames nacionais de matemática e língua materna dos alunos do ensino secundário. Em particular, respondemos às seguintes questões de pesquisa:

1. Quais são os principais determinantes do DA no ensino secundário?
2. Qual a influência de cada uma das principais variáveis (determinantes) nas notas dos exames nacionais de Matemática e Português?

Para responder a estas questões, o resto do documento está estruturado da seguinte forma: na segunda seção, realizamos uma revisão da literatura sobre DA; de seguida, há a apresentação dos métodos utilizados; na quarta seção, são apresentados os resultados; a discussão dos resultados está contida na quinta seção; e, por fim, são delineadas as conclusões e limitações do trabalho.

## 2. Revisão de literatura

O desempenho académico (DA) é definido como os resultados e habilidades dos alunos entre diferentes disciplinas de estudo que permitem que os alunos sejam bem-sucedidos profissional e academicamente (GENESE *et al.*, 2006). O DA é um fator essencial para moldar o curso da vida, determinando muitas oportunidades, como acesso ao ensino superior e oportunidades de emprego (SALMELA-ARO; TYNKKYNEN, 2012). Como o DA é um conceito amplo, existem várias formas para o medir. Entre estes estão a transição (ou não) do ano letivo (CRUZ-JESUS *et al.*, 2020), as notas finais numa disciplina específica (BERTHELON *et al.*, 2019), média de notas (GPA) (COHN *et al.*, 2004; JAYANTHI *et al.*, 2014; MARKER; GNAMBS; APPEL, 2018), entre outras.

### 2.1. Revisão de Literatura em Sucesso Académico Determinantes

Nos últimos 50 anos, os investigadores têm procurado identificar os determinantes do DA e sua magnitude. Coleman (1969) procurou entender quais os fatores que causam variações no sucesso dos alunos, com o intuito de reduzir as disparidades. Este autor

forneceu um ponto de partida para muitos outros estudos que tentaram identificar os fatores subjacentes ao sucesso educacional.

Da literatura encontrada, é notório o poder das habilidades cognitivas e do comportamento acadêmico passado em comparação com os demais fatores. Notas de anos académicos anteriores são constantemente apontadas como os mais fortes preditores do DA, provando assim que os alunos tendem a ser homogêneos com suas notas durante a vida académica (ASIF *et al.*, 2017). Portanto, quem já falhou um ano está mais propenso a falhar novamente, criando um “efeito bola de neve” (COSTA-MENDES *et al.*, 2020; CRUZ-JESUS *et al.*, 2020).

Apesar de Miguéis *et al.* (2018) afirmar que as variáveis socioeconómicas são fracos estimadores em comparação com as notas anteriores dos alunos, os fatores socioeconómicos têm desempenhado um papel preponderante na previsão do DA. O estatuto socioeconómico (SES) influencia fortemente o desempenho educacional dos alunos, indicando que a pobreza, de facto, afeta negativamente as capacidades para aprender (ARCHIBALD, 2006; DELEN, 2010; OPDENAKKER; VAN DAMME, 2001). As crianças cujas famílias precisam de apoio financeiro são mais propensas a ter notas mais baixas (ARCHIBALD, 2006; DELEN, 2010). Está provado também que o DA é impactado negativamente em estudantes sobrecarregados com dificuldades financeiras (ARCHIBALD, 2006; OPDENAKKER; VAN DAMME, 2001; VANDELANNOTE; DEMANET, 2020). No entanto, e sem surpresa, o efeito do SES é menor nos países mais ricos. Um sistema de alta qualidade pode oferecer aos alunos oportunidades para derrubar barreiras impostas por diferentes níveis de rendimento (POKROPEK; BORGONOV; JAKUBOWSKI, 2015). O SES de um aluno ou família pode ser medido de diferentes formas. O rendimento é uma variável evidente para avaliar os níveis socioeconómicos (FISCHBEIN, 1990; MENSAH; KIERNAN, 2010).

O nível de educação dos pais é um proxy comum e relevante para o SES dos alunos. A escolaridade dos pais afeta positivamente a relação dos filhos com a escola (STEINMAYR *et al.*, 2010). Pais com formação educacional mínima lutam para valorizar a educação, dar o suporte necessário e fornecer os recursos necessários para melhorar o DA dos alunos (COLEMAN, 1968; TESFAGIORGIS *et al.*, 2020). A escolaridade dos pais também pode ser considerada um proxy para a motivação (STEINMAYR; DINGER; SPINATH, 2010), cuja importância será mencionada mais adiante. Além disso, o envolvimento da família na educação, particularmente a figura dos pais, pode medir o desempenho de um aluno (FAN; CHEN, 2001; HILL; TAYLOR, 2004). Fan e Chen (2001) destacam o impacto das expectativas dos pais na educação dos filhos. A investigação de

Hill e Taylor (2004) sustenta que o envolvimento dos pais tem suas particularidades de acordo com a demografia e formação cultural de cada aluno, onde as famílias de classe baixa voltam a estar em desvantagem, apresentando inúmeras lutas, seja por falta de conhecimento, motivação, tempo ou recursos quando se trata de ser ativo na educação de seus filhos.

As características intrínsecas às escolas que afetam o DA, e sua magnitude, são também um objeto de estudo há muitos anos. Tal como acontece com outros aspetos, o nível geral de pobreza de uma escola está relacionado com DA. Além disso, o tamanho da escola, que tem sido uma variável de grande estudo entre os académicos, e tem apresentado resultados contraditórios. Enquanto Archibald (2006) afirma que o impacto é negativo, o estudo de Costa-Mendes *et al.* (2020) apoiou efeitos positivos. No entanto, os autores acreditavam que isso poderia refletir as características sociodemográficas do local onde a escola estava inserida. Outros fatores essenciais que devem ser destacados ao nível escolar são as repercussões negativas das drogas e da violência (ABAD; LÓPEZ, 2017) e o efeito positivo do capital social no desempenho dos alunos (VANDELANNOTE; DEMANET, 2020). Além disso, embora alguns países europeus tenham discutido os custos e benefícios da redução do tamanho das turmas, estudos anteriores evidenciam que essa redução não beneficia igualmente os alunos. Geralmente, o impacto é genuinamente pequeno (BOSWORTH, 2014; RIVKIN; HANUSHEK; KAIN, 2005). Há algumas evidências de que o tamanho da turma só prejudica o aproveitamento dos alunos quando esta tem mais de 30 alunos (CRUZ-JESUS *et al.*, 2020). Opdenakker e Van Damme (2001) provaram que todos os alunos se beneficiam de aulas e escolas de alta habilidade em relação aos sistemas de rastreamento educacional. Logo, apenas alunos de alto nível ficam beneficiados, baixo SES e alunos de baixa habilidade são negligenciados, e as diferenças acabam a ser maximizadas (OPDENAKKER; VAN DAMME, 2001). O estudo de Coleman (1968) sustentou que as salas de aula devem ser diversificadas em termos de habilidades e desempenho para que alunos de baixa habilidade ou baixo SES se sintam desafiados.

As características dos professores também podem influenciar o DA dos alunos (AARONSON; BARROW; SANDER, 2007; ARCHIBALD, 2006; COLEMAN, 1968). Alguns exemplos destas características são o seu percurso académico, a experiência, o número de certificações (AARONSON; BARROW; SANDER, 2007), ou escalão profissional do mesmo (ARCHIBALD, 2006). Da mesma forma, a relação entre professores e alunos revela resultados positivos significativos (OPDENAKKER; VAN DAMME, 2001; VANDELANNOTE; DEMANET, 2020), o que pode ser particularmente benéfico se os professores e alunos compartilharem a mesma formação, raça ou género (AARONSON; BARROW; SANDER,

2007; HILL; TAYLOR, 2004). No espectro oposto, as reuniões entre professores e pais estão associadas a mau comportamento e notas mais baixas (VANDELANNOTE; DEMANET, 2020).

### **Aplicação de Machine Learning em DA**

Nos últimos anos, a investigação ampliou as suas técnicas e métodos para uma melhor compreensão dos determinantes do DA, estendendo os métodos estatísticos clássicos para a aplicação da IA. Em termos de desempenho, é notório que, embora o peso dos métodos estatísticos clássicos seja ainda bastante considerável, as técnicas de ML apresentam melhores resultados.

A partir do estudo de Musso, Hernández e Cascallar (2020), onde um conjunto de testes e questionários foi aplicado e informações académicas foram extraídas de 655 estudantes de universidades privadas de Buenos Aires para prever GPA, retenção académica e conclusão de graduação. Redes neurais artificiais com algoritmos de retropropagação foram utilizadas e apresentaram-se como um método extremamente eficiente. No entanto, deve ser dada ênfase aos dados relatados, que levam a resultados imprecisos (MARKER; GNAMBS; APPEL, 2018).

Şen, Uçar e Delen (2012) usaram dados de 5.000 alunos aleatórios da 8ª série na Turquia para prever a pontuação dos testes de aferição, aplicando e comparando o desempenho de árvores de decisão, redes neurais artificiais, máquinas de vetores de suporte e regressão logística. As técnicas de ML provaram a sua maior força preditiva, sendo as árvores de decisão a técnica mais precisa. Abad e López (2017) entrevistaram 18.935 estudantes do ensino básico de 99 instituições de ensino no México para identificar quais fatores estão ligados ao DA. Foram utilizadas árvores de decisão, dada a transparência que este método proporciona, na compreensão das condições subjacentes que conduzem ao resultado e apoiam a criação de análises mais precisas. Da mesma forma, Asif *et al.* (2017) estudaram um conjunto de notas durante a licenciatura de 210 alunos de uma universidade de engenharia no Paquistão para compreender a variação nas notas finais. Os autores optaram por usar árvores de decisão em favor de outros métodos de ML pelas suas vantagens. Além disso, através de uma técnica de agrupamento provaram a homogeneidade das notas durante os quatro anos da licenciatura. Adicionalmente, Kotsiantis (2012) comparou o desempenho das árvores modelo com outros métodos de ML e técnicas estatísticas clássicas em 354 alunos de uma universidade aberta na Grécia para prever o teste do exame final e concluiu que as árvores modelo têm os resultados mais precisos. Mesmo quando Delen (2010) usou redes neurais artificiais, árvores de decisão, máquinas de vetores de suporte, regressão logística, bagging, boosting, e fusão de

informações, para comparar quais modelos preveem a evasão com mais precisão, em 16.066 alunos de quatro anos diferentes numa universidade pública nos Estados Unidos. Árvores de decisão que forneceram apenas os segundos melhores resultados, dadas as vantagens mencionadas anteriormente, foram vistas como a técnica mais adequada. No entanto, Delen (2010) defendeu que quando os métodos tradicionais de ML são comparados com os de ensemble, as técnicas de ensemble são mais potentes e trazem a vantagem de fornecer sistemas de previsão robustos.

Da mesma forma, Miguéis *et al.* (2018) usaram florestas aleatórias, árvores de decisão, máquinas de vetores de suporte, Naïve Bayes, bagged trees, e adaptive boosting trees em 2.459 estudantes de uma universidade pública europeia para entender o DA. Uma técnica de validação cruzada de dez vezes permitiu a avaliação dos diferentes métodos de ML. Embora todas tenham apresentado resultados promissores, as florestas aleatórias superaram os demais, com precisão acima de 96%, comprovando a força dos modelos ensemble.

É essencial mencionar que todos os estudos referidos acima foram baseados em amostras de alunos. Apenas dois estudos de base populacional foram encontrados, ambos usaram a população de alunos portugueses em escolas públicas de ensino secundário e demonstraram o poder das técnicas de ensemble Machine Learning. Quando Cruz-Jesus *et al.* (2020) compararam o desempenho de redes neuronais artificiais, árvores de decisão, árvores extremamente aleatórias, florestas aleatórias, máquinas de vetor de suporte, K-Nearest Neighbours e regressão logística para entender os fatores que levam ao fracasso ou não de um ano letivo, florestas aleatórias provaram ser o procedimento mais preciso. Em contraste, Costa-Mendes *et al.* (2020) aplicaram redes neuronais artificiais, florestas aleatórias, máquinas de vetor de suporte, um modelo de regressão multilinear e uma máquina de impulso de gradiente extremo, e confirmaram que esta última técnica detinha os melhores resultados.

### 3. Metodologia

Esta seção enquadra a metodologia utilizada neste trabalho. Está separada em duas subseções, a primeira uma visão geral de Redes Neuronais, enquanto a segunda se concentra num problema típico de ciência de dados, a seleção de variáveis.

#### 3.1. Redes Neuronais

As redes neuronais artificiais, ou apenas redes neuronais (RN), imitam, em certa medida, o cérebro humano para resolver problemas complexos compostos por um conjunto

de neurónios e conexões entre eles. As RN podem ser usadas em Machine Learning supervisionado, não supervisionado ou híbrido, que possui diferentes arquiteturas e algoritmos associados (JAIN; MAO, 1996, p. 14).

As redes neuronais com multicamadas são compostas pela camada de entrada, um número variável de camadas ocultas e a camada de saída. Todas essas camadas são formadas por neurónios e conectadas por meio de canais. Primeiro, um peso aleatório é atribuído a cada canal e um enviesamento para os neurónios das camadas ocultas. Adicionalmente, é escolhida uma função de ativação, que recebe o valor do neurónio anterior e converte-o numa possível entrada para o próximo, o que acrescenta a habilidade de não linearidade na rede. Na camada de saída, as diferenças entre o resultado previsto e o desejado são calculadas para encontrar o erro. A segunda etapa é a retro propagação, que se inicia na camada de saída e passa por todas as camadas em direção à camada de entrada. Durante este processo, atualiza os pesos com uma regra delta para minimizar a função de custo. Por fim, para completar uma iteração, o primeiro passo é repetido para os novos parâmetros. Este processo pode ser repetido até que o erro esteja abaixo de um limite ou o número máximo de iterações seja alcançado (JAIN; MAO, 1996).

É essencial mencionar que embora as RN sejam de extrema utilidade para configurações puramente preditivas, as RN são de aplicação complexa para fins de Data Mining devido ao efeito "caixa negra", que é criado com algumas das abordagens mais complexas. Este efeito ocorre quando apenas os valores de entrada e de saída são observáveis, e não os efeitos estimados, que são essenciais para a interpretação (mas não para a previsão).

### **3.2. Seleção de variáveis**

A seleção de variáveis reduz as variáveis independentes usadas para treinar o modelo, ao encontrar o menor conjunto de características mais relevantes. Essa seleção visa reduzir a complexidade, o que leva à diminuição a probabilidade de overfitting, e a um modelo preditivo mais preciso e um incremento na eficiência computacional. Neste estudo, usámos uma combinação da Recursive Feature Elimination (RFE), com as regressões Ridge e Lasso.

#### **Recursive Feature Elimination**

O RFE com a aplicação de uma medida de importância, procura melhorar o desempenho da generalização através da classificação recursiva das variáveis. A cada iteração, a importância das características é medida, e as variáveis com os menores efeitos no erro de treino são excluídas (GUYON *et al.*, 2002).

## Regressão Ridge

A regressão de Ridge impõe uma penalidade no tamanho dos coeficientes de regressão com a intenção de reduzi-los. Os coeficientes são reduzidos para zero (e entre si), reduzindo a complexidade e a multicolinearidade. As variáveis que retêm o maior coeficiente são escolhidas para fazer parte do modelo (HASTIE; TIBSHIRANI; FRIEDMAN, 2009).

## Regressão Lasso

Lasso é uma derivação de regressão linear que usa um encolhimento e minimiza o erro de previsão. Esse método começa por criar um limite superior. Uma restrição é aplicada na soma dos valores absolutos dos parâmetros do modelo. De seguida, o processo de encolhimento penaliza os coeficientes de regressão, reduzindo alguns dos valores a zero. Quanto maior a penalidade, mais as estimativas são reduzidas para zero. Após o processo de encolhimento, as variáveis que permanecem com coeficiente diferente de zero são selecionadas para fazer parte do modelo (FONTI, 2017).

### 3.3. Dados

Foi utilizado um conjunto de dados anónimos do ano letivo 2018-2019 disponibilizado pela Direção-Geral de Estatística da Educação e Ciência (DGEEC) do Ministério da Educação de Portugal. É composto por praticamente todos os alunos do ensino secundário público que realizaram exames nacionais de português, matemática ou ambos. O conjunto de dados contém 19.445 alunos e 35.780 exames, dos quais 14.207 são os exames nacionais de matemática e 21.573 os de português, permitindo a criação de três alvos: as notas dos exames nacionais de português e matemática e uma nota agregada que combina os dois conjuntos de dados. Os dados foram extraídos via Microsoft SQL Server Management Studio e a modelação através Python e SAS. É importante mencionar que este conjunto de dados fornece uma grande variedade de potenciais determinantes de DA que, até agora, não eram possíveis de serem analisados, especialmente nesta escala (nacional). Entre elas estão características do encarregado de educação (EE), como anos de escolaridade; a escola, como o ambiente socioeconómico geral; e os professores, como as licenciaturas.

### 3.4. Metodologia Aplicada

Inicialmente, acedemos a duas bases de dados diferentes, uma com dados dos alunos e outra com dados dos professores. Após a ligação dessas bases de dados, foi necessário limpar valores omissos e transformar variáveis categóricas em numéricas. No final do processo de limpeza de dados, tínhamos mais de 50 variáveis independentes.

Para combater os efeitos negativos daquilo a que chamamos “a maldição das dimensões”, procedemos a um processo de seleção de variáveis. Primeiro, a eliminação recursiva de recursos (RFE) é usada para selecionar o número ideal de recursos usados em cada conjunto de dados (N). Posteriormente, as regressões RFE, Lasso e Ridge foram usadas para determinar os N recursos mais essenciais para prever cada alvo. A combinação desses três métodos permitiu a redução de mais de 50 variáveis para cerca de 15 em cada conjunto de dados.

Para a fase de aprendizagem, foi aplicado às Redes Neurais uma pesquisa em grade de validação cruzada fornecida pelo *scikit-learn*, para que se ajustassem os parâmetros do algoritmo, de forma a otimizar os resultados. Os parâmetros otimizados foram: o tamanho das camadas ocultas, função de ativação, solucionador e taxa de aprendizagem.

Finalmente, a validação cruzada de dez vezes foi usada para avaliar os modelos objetivamente. A validação cruzada é uma ferramenta vital para entender o poder da generalização do modelo e contra o *overfitting*. Funciona dividindo os dados em K partes de tamanhos iguais (neste caso, dez amostras) e, em um *loop*, a função aprende usando K-1, enquanto o restante é usado para teste. Foram feitas três repetições desse processo, resultando na avaliação de 30 modelos. O resultado é o R-quadrado médio de todas as iterações para cada modelo. O R-Quadrado é a medida empregada, pois representa o quanto o modelo explica a variação da variável dependente (resumo na Figura 1).

Figura 1 - Metodologia Aplicada



Fonte: elaborado pelos autores

## 4. Resultados

O resultado é o R-Quadrado médio de todas as iterações modelo. As RN apresentaram um R-Quadrado médio 0,106 (0,026) no modelo com as notas agregadas, 0,118 (0,016) no modelo do exame de português e 0,121 (0,009) no modelo de matemática.

### 4.1 Importância das variáveis

Para compreender o valor das variáveis utilizadas, primeiramente, avaliou-se a significância das principais variáveis. A importância das variáveis foi medida através do scikit-learn, que embaralha aleatoriamente cada recurso e calcula as mudanças de desempenho dos modelos, onde avalia a diminuição média da impureza. As características que têm um impacto mais significativo no desempenho do modelo são consideradas as mais importantes (ver Tabela 1).

Embora a importância da variável forneça uma ponderação sobre a importância de cada recurso para os classificar, ainda não está claro como é que cada variável afeta a nota. Como mencionado anteriormente, este projeto visa fornecer uma compreensão de quais os impactos no DA. O problema com modelos mais complexos é o efeito “caixa negra” criado no processo de treino, sendo um trade-off entre a qualidade da previsão e a compreensão do processo. Uma nova abordagem foi testada para evitar esse dilema, com a criação de um conjunto de dados feito com “protótipos” para entender a importância das variáveis nas RN. Um conjunto de dados protótipo foi construído para cada alvo. Primeiramente, as linhas foram calculadas usando a média dessa variável no conjunto de dados original, exceto para uma variável em que o desvio padrão foi subtraído ou adicionado, para variáveis contínuas, ou considerado 0 ou 1, para variáveis binárias. O objetivo foi prever o impacto do aumento ou diminuição do valor da variável, comparando-o com o valor previsto da linha calculado apenas com médias. Em palavras simples, as RN pontuaram um pequeno conjunto de “protótipos”, ou seja, alunos imaginários, com todas as variáveis independentes com valores médios, mas uma de cada vez com um desvio padrão ou unidade, acima ou abaixo. No contexto educacional, essa metodologia inovadora ajuda-nos a alcançar uma contribuição crucial: ter os melhores métodos de ML (ou seja, poder preditivo) com o melhor dos tradicionais (quantificar o impacto de cada determinante). Após uma análise detalhada à literatura que utiliza métodos de machine learning na educação, no nosso melhor conhecimento, esta abordagem nunca foi utilizada no contexto educacional e permite-nos lançar luz sobre a escuridão que constitui as caixas negras das RN.

Tabela 1: Resultados previstos ao alterar uma variável, ceteris paribus, através de protótipos

Ordem	Nota Agregada			Português			Matemática		
	Var	$\Delta$	$\beta$	Var	$\Delta$	$\beta$	Var	$\Delta$	$\beta$
1 <sup>a</sup>	Est_Idade	0.58	-1.17	Est_Idade	0.61	-1.12	Est_Idade	0.54	-1.49
2 <sup>a</sup>	EE_Hab	3.30	0.90	EE_Hab	3.29	0.87	EE_Hab	3.27	1.54
3 <sup>a</sup>	Esc_Rep rov	0.08	-0.66	Est_Fem (Sim)	1	0.59	Esc_BS (Sim)	1	-0.88
4 <sup>a</sup>	Est_Fem (Sim)	1	0.48	Esc_Repr ov	0.08	-0.55	Esc_Repr ov	0.08	-0.80
5 <sup>a</sup>	Est_Num Matric	0.19	-0.33	Esc_Escal aoA	0.06	-0.18	Prof_MscP hD	0 (Não)	-0.47
6 <sup>a</sup>	Esc_Escal aoA	0.06	-0.12	Est_Apoio Soc	1 (Sim)	-0.24	EE_Mae (Sim)	1 (Sim)	0.32
7 <sup>a</sup>	Esc_MSs PhD	0.20	-0.09	EE_Mae (Sim)	1 (Sim)	0.20	Est_Fem (Sim)	1 (Sim)	0.17
8 <sup>a</sup>	Esc_BS (Sim)	1 (Sim)	-0.15	Esc_BS (Sim)	1 (Sim)	-0.17	Esc_Escal aoA	0.06	-0.12
9 <sup>a</sup>	EE_Mae (Sim)	1 (Sim)	0.03	Esc_MScP hD	0.20	0.11	EE_Proprio (Sim)	1 (Sim)	0.08
10 <sup>a</sup>	EE_Pai (Sim)	1 (Sim)	0.03	Est_Net (Sim)	1 (Sim)	0.10	Est_Net (Sim)	1 (Sim)	0.08

Fonte: elaborado pelos autores

Nota: idade do aluno (Est\_Idade), número de anos de escolaridade do Encarregado de Educação (EE\_Hab), taxa de alunos reprovados naquela escola (Esc\_Reprov), número de matrículas do aluno (Est\_NumMatric), taxa de alunos com apoio social escolar (Est\_ApoioSoc), taxa de professores com mestrado ou doutorado na escola (Esc\_MScPhD), sendo o aluno do sexo feminino (Est\_Fem), taxa de alunos com maior nível de apoio social na escola (Esc\_EscalaoA), sendo o responsável legal a mãe (EE\_Mae), o pai (EE\_Pai) ou sendo o próprio aluno (EE\_Proprio), professores com mestrado ou doutorado (Prof\_MscPhD), escolas que oferecem ensino fundamental e médio (Esc\_BS) e acesso dos alunos à Internet (Est\_Net).

## 5. Discussão e implicações

### 5.1. Discussão

No nosso melhor conhecimento, este trabalho trouxe uma abordagem inovadora, pelo menos nesta área, através do uso de protótipos para quantificar o impacto que cada variável tem no DA, medida que era, pelo menos até agora, quase exclusivo dos métodos tradicionais, que são, por sua vez, geralmente menos eficientes, de acordo com a literatura já conhecida. Além disso, o uso de dados institucionais e o método de ciência de dados empregado são vantagens deste trabalho, pois fornecem a oportunidade de usar uma grande amostra de dados, em comparação com a pesquisa de DA tradicional. Adicionalmente, implementamos métodos de ciência de dados mais sofisticados e comprovadamente superiores (por exemplo, Costa-Mendes *et al.* (2020) e Cruz-Jesus *et al.* (2020)) e comparamos a importância relativa de diversos determinantes no indivíduo (aluno) e nível escolar, no mesmo estudo. De acordo com nossos resultados, os cinco fatores mais importantes são a idade do aluno, a taxa de reprovação na escola, o nível de escolaridade dos EE, a taxa de alunos com apoio social na escola e a taxa de professores com pós licenciatura na escola.

Concluimos que a idade dos alunos é um proxy para um fator de extrema importância na procura dos impulsionadores do DA, ou seja, a retenção anterior, pois alunos mais velhos que a média de uma turma são, por norma, sinónimos de retenção anterior. A nossa análise destacou essa variável como uma das mais impactantes (negativamente) no DA. As políticas de retenção estão entre os temas mais debatidos na educação, e os nossos resultados ajudam a lançar alguma luz sobre esse debate. O nosso estudo mostra que quanto mais velhos os alunos, piores são suas notas nos exames, principalmente a matemática. Embora a retenção iminente possa persuadir e motivar alguns alunos a levar a escola mais a sério, a verdade é que encontramos algumas evidências de que a retenção também não está a impedir retenções futuras, sem mencionar que a maioria dos alunos que sofrem retenção sofre efeitos académicos e psicossociais adversos contínuos. O efeito negativo da retenção não é perceptível exclusivamente ao nível do aluno. Encontramos evidências de que isso também é aparente no nível escolar: nas escolas com maiores taxas de reprovação, os resultados dos exames tendiam a ser piores. Esses resultados destacam a necessidade de rever as políticas de reprovação em todos os anos escolares. Parece evidente que a retenção tem efeitos destrutivos no DA, não apenas para o progresso do aluno reprovado, mas também para toda a comunidade escolar, impactando na qualidade geral da escola. A retenção leva a um aumento do comportamento negativo em

todos os alunos nas escolas que têm níveis mais altos de retenção e que têm colegas de turma mais velhos. As políticas que ajudam os alunos a permanecerem no caminho acadêmico podem beneficiar os alunos em risco de fracasso acadêmico e melhorar o comportamento dos outros (MUSCHKIN; GLENNIE; BECK, 2014). Fornecer um programa de retenção personalizado com base nos melhores modelos preditivos pode melhorar as práticas atuais em relação à retenção (por exemplo, Coussement *et al.* (2020), Maldonado *et al.* (2021) e Olaya *et al.* (2020)).

Outro dos principais determinantes do DA destacado nestes resultados é o nível de escolaridade dos EE, que tem um efeito positivo nos resultados dos exames, pois os alunos com EE com mais de 12 anos de escolaridade superam seus colegas em até 1,5 valores. Esse efeito positivo da educação superior dos pais no DA do ensino secundário é especialmente significativo em matemática, um resultado importante a ser considerado nos dias de hoje, pois há uma procura alta e crescente por habilidades STEM (ciência, tecnologia, engenharia e matemática) no mercado de (EUROPEAN CENTRE FOR THE DEVELOPMENT OF VOCATIONAL TRAINING, 2018). As habilidades matemáticas são muito relevantes, pois os alunos alfabetizados matematicamente reconhecem o papel que a matemática desempenha no mundo para fazer julgamentos e decisões bem fundamentados necessárias para cidadãos construtivos e reflexivos (OECD, 2021). À luz de nossos resultados, conhecer o grau de escolaridade dos EE dentro de uma turma pode ser uma ferramenta essencial para ajudar a sinalizar os alunos que podem se beneficiar mais do apoio extra do professor para obter melhores notas. Determinantes relacionados à escola também estão entre os mais relevantes para explicar DA, de acordo com nossos resultados. Nas escolas com maiores taxas de alunos economicamente desfavorecidos que recebem bolsa escolar, o DA é prejudicado, e escolas com maior taxa de professores com mestrado ou doutorado têm um efeito cumulativo sobre o DA, mesmo que leve. Em Portugal, e noutros países onde as matrículas nas escolas públicas são determinadas pela morada da sua família (ou seja, os alunos frequentam as escolas mais próximas de casa), isto pode constituir uma desvantagem para os alunos que vivem em bairros desfavorecidos ou de pauperados, o que leva a um efeito de bola de neve do qual é quase impossível escapar. De acordo com estes resultados, surge a questão se faria sentido implementar programas de aprendizagem específicos nesses bairros para evitar ter escolas que são aglomerados de alunos desprivilegiados onde, como os resultados mostram, será mais desafiador prosperar academicamente e quebrar o ciclo de desvantagem.

Conclusões adicionais deste estudo mostram outros determinantes do DA relevantes, nomeadamente o género, uma vez que as mulheres superam os homens nas

notas dos, mas esta tendência é menos evidente em matemática, que é consistente com a investigação anterior (por exemplo, Brunner, Krauss e Kunter (2008)). O uso da Internet também surgiu como um fator positivo do DA, embora tenha um impacto insignificante. Estudos anteriores relatam descobertas conflituosas sobre o uso da Internet como determinante do DA, com efeitos positivos (BOWERS; BERLAND, 2013) e adversos (ROZGONJUK; TAHT; VASSIL, 2021). Em relação ao tamanho da escola, os nossos resultados mostram que em escolas maiores, onde os alunos do ensino básico e secundário estão juntos, as notas dos exames do ensino secundário tendem a ser mais baixas, em linha com os estudos anteriores (EGALITE; KISIDA, 2016).

Estes resultados destacam a importância de promover o DA nas gerações futuras. O DA aparentemente crescente funcionará como um efeito de contágio, uma vez que os pais ou EE com formação universitária e os professores com graus superiores com formação universitária e os professores com graus superiores (Mestrado ou Doutorado) são determinantes positivos do DA. Tendo em conta que níveis mais altos de educação levam a mais desenvolvimento pessoal e melhores condições de vida (ou seja, melhor emprego, SES mais alto e capital cultural) e que os EE nessas condições também ajudam seus filhos a ter sucesso na escola, é compreensivelmente o caminho a seguir pelos decisores na área da educação.

## 5.2. Implicações teóricas

Esta análise explorou os dados institucionais das escolas secundárias públicas em Portugal para compreender melhor o DA. Este estudo forneceu resultado inovador: a caracterização quantitativa do impacto de cada impulsionador do DA nas notas finais do exame nacional, o que permite uma comparação entre a importância de diversos determinantes relacionados a alunos, EE, professores e escolas no mesmo estudo, utilizando um grande conjunto de dados e técnicas avançadas de ciência de dados. Os nossos resultados contribuem para a literatura, pois, reconhecem a relevância de cada um dos determinantes identificados para o conhecimento científico nesta área. Os nossos resultados também convalidam a literatura ao destacar os efeitos adversos no DA na reprovação de um aluno, tanto no nível individual quanto no escolar. Existem extensos estudos sobre esse tópico (JIMERSON, 2001), e quase todas as evidências estão de acordo com nossos resultados, pois apontam para os impactos negativos da reprovação de um aluno nos anos do ensino primário, básico e secundário, diminuindo até mesmo a probabilidade de conclusão do ensino secundário (ANDREW, 2014). Adicionalmente, mostramos que quanto mais velhos os alunos são (o que significa que eles provavelmente foram retidos anteriormente), piores são suas notas nos exames e que estudar numa escola

com taxas de retenção mais altas também resulta em piores notas nos exames. Anderson, Whipple e Jimerson (2011) constataram que a retenção é um preditor significativo de abandono escolar entre estudantes do ensino secundário devido à ausência de estratégias eficazes para aumentar as competências e estigmatização por parte dos colegas que podem agravar problemas de ajuste emocional e comportamental. Quaisquer resultados acadêmicos positivos de curto prazo da retenção tendem a desaparecer (TOLEN; QUINLIN, 2017). A retenção também está relacionada com uma menor motivação acadêmica e maior absenteísmo (MARTIN, 2011).

Este estudo também demonstra que o DA do ensino secundário é influenciado pela escolaridade dos pais, pois filhos de pais com ensino terciário (pós-secundário ou superior) apresentam melhores resultados nos exames finais, principalmente a matemática. Embora a literatura evidencie o impacto positivo do nível de escolaridade dos pais (STEINMAYR; DINGER; SPINATH, 2010), apontar a universidade ou o ensino superior como um dos aspetos mais significativos para melhorar o DA do aluno é um resultado muito relevante. Tanto quanto sabemos, não tinha sido reportado na investigação do DA ao nível do ensino secundário, especialmente considerando que incluímos praticamente todos os alunos do ano 2018/2019 das escolas secundárias públicas portuguesas. De acordo com os estudos anteriores, a educação dos pais parece ser um impulsionador do DA em todos os anos escolares, o que demonstra efeitos positivos desde o jardim de infância (TAN; PENG; LYU, 2019), onde se destacam que as influências familiares (nível de educação traduzido em capital cultural, envolvimento dos pais e expectativas) desempenham um papel significativo nas experiências e resultados escolares. Essa dinâmica tem sido estudada tanto em nível de sala de aula, onde o nível de ensino superior dos pais dos colegas melhora o DA (WANG, 2021), frequentar escolar com maior proporção de estudantes de famílias em desvantagem educacional diminui o DA (CHESTERS; DALY, 2017).

### **5.3. Implicações práticas**

Este estudo ajuda a moldar o que parece ser um novo ciclo na investigação e na prática em educação. Ao permitir quantificar com precisão o impacto de cada impulsionador do DA nas notas do exame final dos alunos do ensino secundário, este fornece fundamentos significativos para os tomadores de decisão implementarem mudanças substanciais e aumentarem o sucesso. Os nossos resultados permitem classificar os determinantes mais críticos do DA e, assim, refletir sobre medidas para melhorar os resultados do ensino secundário. Primeiro, os critérios de retenção poderiam ser revistos, pois reprovar um ano no progresso escolar tem efeitos deletérios sobre o DA nos níveis individual e escolar. As políticas de retenção têm estado em debate em Portugal nos últimos anos. De facto, as

taxas de retenção têm vindo a diminuir, mas recomenda-se mais trabalho ao nível do aluno (acompanhamento especializado na escola para alunos em risco de rejeição, por exemplo) e ao nível institucional, nomeadamente a alteração dos critérios de retenção, de acordo com este estudo. O estudo da Eurydice (2011) encontrou uma cultura predominante de retenção em alguns países europeus. Nos países onde existe essa cultura, a crença dominante é que a repetição de um ano beneficia o estudante, embora estudos e resultados recentes mostrem o contrário. Alguns professores, comunidade escolar e EE, ainda compartilham essa ideia, embora tenha sido repetidamente questionada. Portanto, o desafio pode residir mais em questionar tais pressupostos antes de implementar mudanças legislativas.

Em relação à escolaridade dos EE, principal fator positivo do DA encontrado neste estudo; algumas medidas podem ser sugeridas, como considerar a escolaridade do EE na atribuição dos alunos às turmas. Esta medida irá criar grupos heterogêneos que integrem alunos com pais de baixa e alta escolaridade ou que forneçam aos professores informação sobre o nível de escolaridade dos pais no início de cada ano letivo, o que ajuda a sinalizar alunos que possam necessitar de apoio extra nas aulas, tendo em conta o elevado impacto da educação dos pais nas notas dos exames.

O nível de formação dos professores (licenciatura versus mestrado ou doutorado) também é um impulsionador positivo promissor de DA, de modo que os formuladores de políticas podem facilitar oportunidades de carreira para professores que decidem seguir a pós-licenciatura e motivar as escolas a contratar professores de pós-licenciatura. As políticas ao nível da escola também estão contempladas: os nossos resultados mostram que o uso da Internet pode ajudar a melhorar o DA e, portanto, as escolas podem implementar um reforço digital eficaz, permitindo que todos os alunos tenham acesso à Internet, de forma a garantir a igualdade de acesso e minimizar qualquer desvantagem (por exemplo, alunos que não têm computador ou acesso à Internet em casa). Da mesma forma, como nossos resultados também mostram notas mais baixas nos exames em escolas com taxas mais altas de alunos que recebem apoio social (ou seja, alunos de famílias com renda mais baixa), os decisores políticos podem querer implementar medidas para minimizar os efeitos de viver em territórios social e economicamente desfavorecidos. Apesar de 136 escolas portuguesas (que abrangem 10% do total de alunos do ensino público) nestas condições já realizarem um programa específico (TEIP - Programa Territórios Educativos de Intervenção Prioritária) para promover o sucesso escolar e prevenir o abandono escolar, parece ser necessário mais trabalho a este nível para aumentar o DA.

## 6. Conclusões

Como primeira limitação identificada, temos o facto de poder haver modelos mais eficazes a prever o DA que as RN, apesar do apoio da literatura. Como tal, no futuro, deveria ser feita uma comparação de diversos métodos para garantir os melhores resultados possíveis. Adicionalmente, para garantir a veracidade destas conclusões a longo prazo e para que medidas possam ser tomadas com base nos mesmos, sugerimos que a metodologia seja repetida para diversos anos letivos. De especial importância, após a pandemia e todas as alterações que vieram com a mesma. Por último, sugerimos que fossem analisados mais termos que se possam referir aquilo a que nos referimos como protótipos, para garantir que não existe nenhuma metodologia semelhante em qualquer possível denominação.

## 7. Referências

AARONSON, D.; BARROW, L.; SANDER, W. Teachers and student achievement in the Chicago public high schools. **Journal of Labor Economics**, v. 25, n. 1, p. 95–135, 2007. DOI: <https://doi.org/10.1086/508733>.

ABAD, F. M.; LÓPEZ, A. A. C. C. Data-mining techniques in detecting factors linked to academic achievement. **School Effectiveness and School Improvement**, v. 28, n. 1, p. 39–55, 2017. DOI: <https://doi.org/10.1080/09243453.2016.1235591>.

ANDERSON, B. G. E.; WHIPPLE, A. D.; JIMERSON, S. R. Mental health outcomes. **Learning Disability Practice**, v. 14, n. 3, p. 11–11, 2011. DOI: <https://doi.org/10.7748/ldp.14.3.11.s7>.

ANDREW, M. The Scarring Effects of Primary-Grade Retention? A Study of Cumulative Advantage in the Educational Career. **Social Forces**, v. 93, n. 2, p. 653–685, 2014. DOI: <https://doi.org/10.1093/sf/sou074>.

ARCHIBALD, S. Narrowing in on educational resources that do affect student achievement. **Peabody Journal of Education**, v. 81, n. 4, p. 23–42, 2006. DOI: [https://doi.org/10.1207/s15327930pje8104\\_2](https://doi.org/10.1207/s15327930pje8104_2).

ASIF, R.; MERCERON, A.; ALI, S. A.; HAIDER, N. G. Analyzing undergraduate students' performance using educational data mining. **Computers and Education**, v. 113, p. 177–194, Oct. 2017. DOI: <https://doi.org/10.1016/j.compedu.2017.05.007>.

BERTHELON, M.; BETTINGER, E.; KRUGER, D. I.; MONTECINOS-PEARCE, A. The Structure of Peers: the Impact of Peer Networks on Academic Achievement. **Research in Higher Education**, v. 60, n. 7, p. 931–959, 2019. DOI: <https://doi.org/10.1007/s11162-018-09543-7>.

BOSWORTH, R. Class size, class composition, and the distribution of student achievement. **Education Economics**, v. 22, n. 2, p. 141–165, 2014. DOI: <https://doi.org/10.1080/09645292.2011.568698>.

BOWERS, A. J.; BERLAND, M. Does recreational computer use affect high school achievement? **Educational Technology Research and Development**, v. 61, n. 1, p. 51–69, 2013. DOI: <https://doi.org/10.1007/s11423-012-9274-1>.

BRUNNER, M.; KRAUSS, S.; KUNTER, M. Gender differences in mathematics: Does the story need to be rewritten? **Intelligence**, v. 36, n. 5, p. 403–421 2008. DOI: <https://doi.org/10.1016/j.intell.2007.11.002>.

CHESTERS, J.; DALY, A. Do peer effects mediate the association between family socio-economic status and educational achievement? **Australian Journal of Social Issues**, v. 52, n. 1, p. 65–77, 2021. DOI: <https://doi.org/10.1002/ajs4.3>.

COHN, E.; COHN, S.; BALCH, D. C.; BRADLEY, J. Determinants of undergraduate GPAs: SAT scores, high-school GPA and high-school rank. **Economics of Education Review**, v. 23, n. 6, p. 577–586. DOI: <https://doi.org/10.1016/j.econedurev.2004.01.001>.

COLEMAN, J. S. Equality of Educational Opportunity. **Equity and Excellence in Education**, v. 6, n. 5, p. 19–28, 1968. DOI: <https://doi.org/10.1080/0020486680060504>.

COLEMAN, J. S. Equality of educational opportunity, reexamined. **Socio-Economic Planning Sciences**, v. 2, n. 2–4, p. 347–354, 1969. DOI: [https://doi.org/10.1016/0038-0121\(69\)90029-9](https://doi.org/10.1016/0038-0121(69)90029-9).

COSTA-MENDES, R.; OLIVEIRA, T.; CASTELLI, M.; CRUZ-JESUS, F. A machine learning approximation of the 2015 Portuguese high school student grades: a hybrid approach.

**Education and Information Technologies**, v. 26, p. 1527-15-47, Sept. 2020. DOI: <https://doi.org/10.1007/s10639-020-10316-y>.

COUNCIL OF EUROPE. **EUROPE AND NORTH AMERICA EDUCATION 2030**

**CONSULTATION**. Background documents. Oct. 2018. Disponível em:

<https://www.sdg4education2030.org/europe-and-north-america>. Acesso em: 19 Oct. 2022.

COUSSEMENT, K.; PHAN, M.; DE CAIGNY, A.; BENOIT, D. F.; RAES, A. Predicting student dropout in subscription-based online learning environments: the beneficial impact of the logit leaf model. **Decision Support Systems**, v. 135, n. 113325, May 20-20. DOI:

<https://doi.org/10.1016/j.dss.2020.113325>.

CRUZ-JESUS, F.; CASTELLI, M.; OLIVEIRA, T.; MENDES, R.; NUNES, C.; SA-VELHO, M.; ROSA-LOURO, A. Using artificial intelligence methods to assess academic achievement in public high schools of a European Union country. **Heliyon**, v. 6, n. 6, e04081, 2020. DOI:

<https://doi.org/10.1016/j.heliyon.2020.e04081>.

DELEN, D. A comparative analysis of machine learning techniques for student retention management. **Decision Support Systems**, v. 49, n. 4, p. 498–506, 2010. DOI:

<https://doi.org/10.1016/j.dss.2010.06.003>.

EGALITE, A. J.; KISIDA, B. School size and student achievement: a longitudinal analysis.

**School Effectiveness and School Improvement**, v. 27, n. 3, p. 406–417, 2016. DOI:

<https://doi.org/10.1080/09243453.2016.1190385>.

EUROPEAN CENTRE FOR THE DEVELOPMENT OF VOCATIONAL TRAINING. **Insights into skill shortages and skill mismatch**: learning from Cedefop's European skills and jobs survey. Luxemburgo: Publications Office, 2018. Cedefop reference series; n. 106. DOI:

<https://data.europa.eu/doi/10.2801/645011>.

EUROPEAN COMMISSION. **European Pillar of Social Rights**. Mar. 2017. DOI:

<https://doi.org/10.2792/95934>.

EURYDICE. Grade Retention during Compulsory Education in Europe: regulations and statistics. **European Education and Culture Executive Agency**. 2011.

FAN, X.; CHEN, M. Parental Involvement and Students' Academic Achievement: a Meta-Analysis. **Educational Psychology Review**, v. 13, n. 1, p. 1–22, 2001. DOI: <https://doi.org/10.1023/A:1009048817385>.

FISCHBEIN, S. Biosocial influences on sex differences for ability and achievement test results as well as marks at school. **Intelligence**, v. 14, n. 1, p. 127–139, 1990. DOI: [https://doi.org/10.1016/0160-2896\(90\)90018-O](https://doi.org/10.1016/0160-2896(90)90018-O).

FONTI, V. **Feature Selection using LASSO**. 2017. Disponível em: [https://www.researchgate.net/profile/David-Booth-7/post/Regression-of-pairwise-trait-similarity-on-similarity-in-personal-attributes/attachment/5b18368d4cde260d15e3a4e3/AS%3A634606906785793%401528313485788/download/werkstuk-fonti\\_tcm235-836234.pdf](https://www.researchgate.net/profile/David-Booth-7/post/Regression-of-pairwise-trait-similarity-on-similarity-in-personal-attributes/attachment/5b18368d4cde260d15e3a4e3/AS%3A634606906785793%401528313485788/download/werkstuk-fonti_tcm235-836234.pdf). Acesso em: 19 Oct. 2022.

GENESEE, F.; LINDHOLM-LEARY, K.; SAUNDERS, W.; CHRISTIAN, D. **Educating English Language Learners**. Cambridge: Cambridge University Press, 2006.

GUYON, I.; WESTON, J.; BARNHILL, S.; VAPNIK, V. Gene selection for cancer classification using support vector machines. **Machine Learning**, v. 46, n. 1–3, p. 389–422, 2002. DOI: <https://doi.org/10.1023/A:1012487302797>.

HASTIE, T.; TIBSHIRANI, R.; FRIEDMAN, J. The Elements of Statistical Learning: data mining, inference and prediction. **Springer**, v. 26, n. 4. Springer New York, 2009. DOI: <https://doi.org/10.1007/978-0-387-84858-7>.

HILL, N. E.; TAYLOR, L. C. Parental school involvement and children's academic achievement pragmatics and issues. **Current Directions in Psychological Science**, v. 13, n. 4, p. 161–164, 2004. DOI: <https://doi.org/10.1111/j.0963-7214.2004.00298.x>.

JAIN, A. K.; MAO, J. Artificial neural networks - a Tutorial. **Computer**, v. 29, n. 3, p. 31-44, March 1996. DOI: <https://doi.org/10.1109/2.485891>.

JAYANTHI, S. V.; BALAKRISHNAN, S.; CHING, A. L. S.; LATIFF, N. A. A.; NASIRUDEEN, A. M. A. Factors Contributing to Academic Performance of Students in a Tertiary Institution in Singapore. **American Journal of Educational Research**, v. 2, n. 9, p. 752–758, 2014. DOI: <https://doi.org/10.12691/education-2-9-8>.

JIMERSON, S. R. Meta-analysis of grade retention research: implications for practice in the 21st century. **School Psychology Review**, v. 30, n. 3, p. 420–437, 2001. DOI: <https://doi.org/10.1080/02796015.2001.12086124>.

KOTSIANTIS, S. B. Use of machine learning techniques for educational proposes: a decision support system for forecasting students' grades. **Artificial Intelligence Review**, v. 37, n. 4, p. 331–344, 2012. DOI: <https://doi.org/10.1007/s10462-011-9234-x>.

LIPTON, Z. C. The mythos of model interpretability. **Communications of the ACM**, v. 61, n. 10, p. 35–43, 2018. DOI: <https://doi.org/10.1145/3233231>.

MALDONADO, S.; MIRANDA, J.; OLAYA, D.; VÁSQUEZ, J.; VERBEKE, W. Redefining profit metrics for boosting student retention in higher education. **Decision Support Systems**, v. 143, Nov. 2020). DOI: <https://doi.org/10.1016/j.dss.2021.113493>.

MARKER, C.; GNAMBS, T.; APPEL, M. Active on Facebook and Failing at School? Meta-Analytic Findings on the Relationship Between Online Social Networking Activities and Academic Achievement. **Educational Psychology Review**, v. 30, n. 3, p. 651–677, 2018. DOI: <https://doi.org/10.1007/s10648-017-9430-6>.

MARTIN, A. J. Holding back and holding behind: grade retention and students' non-academic and academic outcomes. **British Educational Research Journal**, v. 37, n. 5, p. 739–763, 2011. DOI: <https://doi.org/10.1080/01411926.2010.490874>.

MENSAH, F. K.; KIERNAN, K. E. Gender differences in educational attainment: influences of the family environment. **British Educational Research Journal**, v. 36, n. 2, p. 239–260, 2010. DOI: <https://doi.org/10.1080/01411920902802198>.

MIGUÉIS, V. L.; FREITAS, A.; GARCIA, P. J. V.; SILVA, A. Early segmentation of students according to their academic performance: A predictive modelling approach. **Decision Support Systems**, v. 115, p. 36–51, Mar. 2018. DOI: <https://doi.org/10.1016/j.dss.2018.09.001>.

MÜNICH, D.; PSACHAROPOULOS, G. Education externalities: what they are and what we know. **European Expert Network on Economics of Education (EENEE)**, Issue 34, 2018. <https://doi.org/10.2766/309796>.

MUSCHKIN, C. G.; GLENNIE, E.; BECK, A. N. Peer contexts: do old for grade and retained peers influence student behavior in middle school? **Teachers College Record**, v. 116, n. 4, p. 1–37, 2014.

MUSSO, M. F.; HERNÁNDEZ, C. F. R.; CASCALLAR, E. C. Predicting key educational outcomes in academic trajectories: a machine-learning approach. **Higher Education**, v. 80, n. 5, p. 875–894, 2020. DOI: <https://doi.org/10.1007/s10734-020-00520-7>.

OECD. **Mathematics performance (PISA)** (indicator). 2021. DOI: <https://doi.org/10.1787/04711c74-en>.

OLAYA, D.; VÁSQUEZ, J.; MALDONADO, S.; MIRANDA, J.; VERBEKE, W. Uplift Modeling for preventing student dropout in higher education. **Decision Support Systems**, v. 134, n. 113320, Jan. 2021. DOI: <https://doi.org/10.1016/j.dss.2020.113320>.

OPDENAKKER, M.-C.; VAN DAMME, J. Relationship between School Composition and Characteristics of School Process and their Effect on Mathematics Achievement. **British Educational Research Journal**, v. 27, n. 4, p. 407–432, 2001. DOI: <https://doi.org/10.1080/01411920125573>.

POKROPEK, A.; BORGONOV, F.; JAKUBOWSKI, M. Socio-economic disparities in academic achievement: A comparative analysis of mechanisms and pathways. **Learning and Individual Differences**, v. 42, p. 10–18, Aug. 2015. DOI: <https://doi.org/10.1016/j.lindif.2015.07.011>.

RIVKIN, S. G.; HANUSHEK, E. A.; KAIN, J. F. Teachers, schools, and academic achievement. **Econometrica**, v. 2, n. 1, p. 101–104, 2005. DOI: <https://doi.org/10.5585/eccos.v2i1.211>.

ROZGONJUK, D.; TÄHT, K.; VASSIL, K. Internet use at and outside of school in relation to low-and high-stakes mathematics test scores across 3 years. **International Journal of STEM Education**, v. 8, n. 1, 2021. DOI: <https://doi.org/10.1186/s40594-021-00287-y>.

SALMELA-ARO, K.; TYNKKYNEN, L. Gendered pathways in school burnout among adolescents. **Journal of Adolescence**, v. 35, n. 4, p. 929–939, 2012. DOI: <https://doi.org/10.1016/j.adolescence.2012.01.001>.

ŞEN, B.; UÇAR, E.; DELEN, D. Predicting and analyzing secondary education placement-test scores: a data mining approach. **Expert Systems with Applications**, v. 39, n. 10, p. 9468–9476, 2012. DOI: <https://doi.org/10.1016/j.eswa.2012.02.112>.

STEINMAYR, R.; DINGER, F. C.; SPINATH, B. Parents' education and children's achievement: the role of personality. **European Journal of Personality**, v. 24, n. 6, p. 535–550, 2010. DOI: <https://doi.org/10.1002/per.755>.

TAN, C. Y.; PENG, B.; LYU, M. What types of cultural capital benefit students' academic achievement at different educational stages? Interrogating the meta-analytic evidence. **Educational Research Review**, v. 28, n. 100289, September). DOI: <https://doi.org/10.1016/j.edurev.2019.100289>.

TESFAGIORGIS, M.; TSEGAI, S.; MENGESHA, T.; CRAFT, J.; TESSEMA, M. The correlation between parental socioeconomic status (SES) and children's academic achievement: the case of Eritrea. **Children and Youth Services Review**, v. 116, p. 105242, July 2020. DOI: <https://doi.org/10.1016/j.childyouth.2020.105242>.

TOLEN, A.; QUINLIN, L. **The Efficacy of Student Retention: A Review of Research & Literature**. 2017. Disponível em: <https://www.sjsd.k12.mo>. Acesso em: 30 mar. 2022.

UNITED NATIONS. **Transforming our world: the 2030 Agenda for Sustainable Development**. p. 1–41. 2015. DOI: <https://doi.org/10.1201/b20466-7>. Disponível em:

<https://sdgs.un.org/sites/default/files/publications/21252030%20Agenda%20for%20Sustainable%20Development%20web.pdf>.

VANDELANNOTE, I.; DEMANET, J. “What’s High School Got to do With It?” Secondary School Composition, School-Wide Social Capital and Higher Education Enrollment.

**Research in Higher Education**, v. 62, p. 680-708, Nov. 2020. DOI:

<https://doi.org/10.1007/s11162-020-09617-5>.

WANG, T. Classroom Composition and Student Academic Achievement: the impact of peers’ parental education. **B. E. Journal of Economic Analysis and Policy**, v. 21, n. 1, p. 273–305, 2021. DOI: <https://doi.org/10.1515/bejeap-2020-01>.