

UNIVERZA NA PRIMORSKEM
FAKULTETA ZA MATEMATIKO, NARAVOSLOVJE IN
INFORMACIJSKE TEHNOLOGIJE

Diploma Thesis
**Predicting Secondary School Student Performance Using
Machine Learning**

()

Ime in priimek: Lucas Lorenzo Jakin

Študijski program: Računalništvo in informatika EN

Mentor: izr. prof. dr. Branko Kavšek

Koper, 2024

Ključna dokumentacijska informacija

Ime in PRIIMEK: Lucas Lorenzo JAKIN

Naslov zaključne naloge: Predicting Student Performance Using Machine Learning

Kraj: Koper

Leto: 2024

Število listov:

Število slik:

Število tabel:

Število referenc:

Mentor: izr. prof. dr. Branko Kavšek

Ključne besede:

Izvleček:

Key words documentation

Name and SURNAME: Lucas Lorenzo JAKIN

Title of final project paper: **Predicting Student Performance Using Machine Learning**

Place: Koper

Year: 2024

Number of pages:

Number of figures:

Number of tables:

Number of references:

Mentor: Assoc. Prof. Branko Kavšek, PhD

Keywords:

Abstract: The application of machine learning in educational settings has gathered significant attention, particularly in model predictions allowing businesses to make accurate guesses to the likely outcomes of a question based on historical data. My diploma thesis aims to analyze a dataset from Kaggle related to secondary school students in Portugal. By following the CRISP-DM methodology (Cross Industry Standard Process for Data Mining), I will systematically explore each phase, from data understanding and preparation to modeling and evaluation. In the modeling phase I will implement various ML algorithms, using different parameters every time, to predict students' academic success. The performance of these models will be then compared and tested to determine the most effective approach.

Zahvala

Kazalo vsebine

1	INTRODUCTION	1
2	Problem description	5
3	Data understanding	6
4	Data preparation	7
5	Methodology	8
5.1	K najbližjih sosedov	8
5.2	Metoda podpornih vektorjev	8
5.3	Umetna Nevronska mreža	9
5.4	Konvolucijska nevronska mreža	9
6	Modeling / Classification	10
7	Results	11
8	...	12
8.1	12
9	CONCLUSION / FINAL THOUGHTS	13
10	LITERATURA IN VIRI	14

Kazalo tabel

Kazalo slik

1	CRISPM-DM methodology	3
---	---------------------------------	---

Seznam kratic

...

1 INTRODUCTION

Machine Learning (ML) has become a revolutionary technology in recent years, transforming industries and changing the way data-based decisions are conducted. Machine Learning is primarily about creating algorithms that help computers learn from data and make predictions. This ability has enabled companies to discover insights, improve processes, and enhance decision-making in different areas, such as finance, healthcare, marketing, and manufacturing.

The CRISPM-DM (Cross-Industry Standard Process for Data Mining) methodology is crucial for the successful implementation of machine learning models. The structured framework provided by CRISP-DM helps guide data scientists and analysts in the complex process of developing data mining and machine learning models. This approach guarantees that projects are carried out in a structured manner, resulting in dependable and understandable results. Published in 1999 to standardize data mining processes across industries, it has since become the **most common methodology** for data mining, analytics, and data science projects.

The CRISP-DM methodology outlines six critical phases in a data mining project:

- **Business Understanding:** This initial phase focuses on understanding the objectives and requirements of the project. Establishing a strong business understanding is absolutely essential. Most of the tasks performed in this phase are foundational project management activities that are universal to most projects:
 1. **Determine business objectives:** It is crucial to first understand, from a business perspective, what the customer really wants to accomplish.
 2. **Asses situation:** Determine resources availability, project requirements, assess risks and contingencies, and conduct a cost-benefit analysis
 3. **Determine data mining goals:** Moreover, in addition to defining bussiness objectives, it is also important to define what success looks like from a technical data mining perspective.
 4. **Produce project plan:** Select technologies and tools and define detailed plans for each project phase.

- **Data Understanding:** In this phase, data is collected and explored to gain insights into its properties. It drives the focus to identify, collect, and analyze the data sets that can help you accomplish the project goals. This phase has four tasks:
 1. **Collect initial data:** Acquire the necessary data and load it into the analysis tool.
 2. **Describe data:** Examine the data and document its surface properties like data format, number of fields or field identities.
 3. **Explore data:** Query the data, visualize it, and identify relationships among the data.
 4. **Verify data quality:** Document any quality issues. How clean is the data?
- **Data Preparation:** Almost 80% of the project is data preparation. This phase, which is often referred to as "data munging", prepares the final data set for modeling. It has five tasks:
 1. **Select data:** Determine which data sets will be used and document reasons for inclusion/exclusion.
 2. **Clean data:** A common practice during this task is to correct, impute, or remove erroneous values. Without it, the quality of the data be innacurate.
 3. **Construct data:** Derive new attributes that will be helpful.
 4. **Integrate data:** Create new data sets by combining data from multiple sources.
 5. **Format data:** Re-format data as necessary, for example converting string values that store numbers to numeric for better mathematical operations.
- **Modeling** Even if this might be the most exciting work is also often the shortest phase of the project. Here various models will likely be built and assessed based on several different modeling techniques. This phase has four tasks:
 1. **Select modeling techniques:** Determine which algoritms to try (e.g. regression, neural network, KNN).
 2. **Generate test design:** Depending on the approach you took, there might be the need to split the data into training set, test, and validation sets.
 3. **Build model:** The selected modeling techniques are put into action to create predictive models. Beside running lines of code, this process requires careful consideration of various factors to ensure the models are robust and effective.

-
- ```

graph TD
 BU(Business Understanding) <--> DU(Data Understanding)
 DU --> DP(Data Preparation)
 DP <--> M(Modeling)
 M --> E(Evaluation)
 E --> D(Deployment)
 D --> BU
 D((Data)) --- DU
 D --- DP

```

Slika 1: CRISPM-DM methodology

For this diploma thesis, a dataset will be selected to demonstrate the application of the CRISPM-DM methodology. The focus of the thesis will be to systematically apply each phase of CRISPM-DM to the chosen dataset. This will involve selecting an appropriate dataset, understanding its structure and contents, preparing it for analysis, and implementing various machine learning models to predict target values. The models will be then evaluated to determine their accuracy and effectiveness. This final project aims to provide a thorough and structured approach to data mining and predictive modeling, showcasing the practical application of these techniques on real-world data.

The dataset that I have chosen is about student performance and it is available on both Kaggle and UCI Machine Learning Repository. The dataset approaches student achievement in secondary education in two Portuguese schools.

## 2 Problem description

Education is a key factor for achieving long-term economic progress. During the last two decades, the portuguese educational level has improved. However, the statistics keep Portugal at Europe's tail end due to its high student failure and dropping rates. In particular, failure in the core classes of Mathematics and Portuguese (native language) is extremely serious, since they provide fundamental knowledge for the success in other remaining subjects.

The chosen data set approach student achievements in secondary education of two Portuguese schools. It includes a wide range of attributes that capture various aspects of student life and academic achievement, consisting of student grades, demographic information, social factors, and school-related features. The data was collected through school reports and questionnaires, ensuring a detailed and multifaceted view of the factors influencing student performance.

**Objectives:** The main objective of this thesis is to build and evaluate predictive models that can accurately forecast the final grades of students in Mathematics. The problem will be approached from two angles:

- **Classification:** Categorizing students' final grades into predefined levels or classes. In this case including multi-class classification, for example grading on a scale from A to F.
- **Regression:** Predicting the exact numerical value of a student's final grade based on the input features.

Predicting student performance is of great importance to educators and students themselves. Accurate predictions can help in identifying students at risk of underperforming, allowing for timely interventions and support. By sistematically performing a thorough analysis of student performance data and develop reliable models we can achieve the above mentioned objectives.

## 3 Data understanding

## 4 Data preparation

MNIST (Modified NIST) je nabor podatkov, ki vsebuje 70000 slik ročno napisanih števk. Podatke je od srednješolcev in zaposlenih na Ameriškem uradu za statisko sprva zbral NIST (National institute of standards and technology). Podatki so ločeni na učno množico in testno množico. Učna množica vsebuje 60000 slik, testna pa 10000. Pisci števk, katerih podatki so v učni množici, nimajo slik v testni množici. Vsaka slika je označena s primernim številom med 0 in 9.

Slike so velikosti 28x28, vsak piksel zavzame vrednost med 0 in 255. Vse slike so bile normalizirane glede na velikost in centrirane glede na težišče. [1]

## 5 Methodology

### 5.1 K najbližjih sosedov

kNN je primer neparametričnega, nadzorovanega strojnega učenja. Dano podatkovno točko  $x$  neznanega razreda klasificiramo tako, da najdemo  $k$  točk v naboru podatkov za učenje, ki so tej točki najbližje, nato dodelimo  $x$  tisti razred, ki je večinski v množici sosednjih točk.

(Kot neparametričen algoritem) kNN ne zgradi modela iz nabora podatkov za učenje (training data, učilni podatki? ne vem), vendar shrani celoten učilni nabor, ki ga uporabi za klasifikacijo. Posledično je učenje hitrejše, sklepanje pa počasnejše v primerjavi z ostalimi algoritmi. [citat?]

Učinkovitost algoritma je močno odvisna od izbire  $k$ . Optimalna izbira je drugačna za vsak nabor podatkov. V splošnem velja, da vodijo premajhne vrednosti do prekomernega prilaganja (overfit) in prevelike vrednosti do nezadostnega (?) prilaganja (underfit).[citat]

Za uspešnost algoritma kNN je pomembna tudi izbira razdalje/metrike (), še posebej za podatke z visokim številom atributov. [citat] (je/ni pomembno za naš problem? curse of dimensionality?)

KNN se lahko izboljša z tako, da točke, ki so bližje  $x$  bolj upoštevamo pri klasifikaciji, torej utežimo bližnje točke glede na razdaljo. (weighted knn)

Konkretna implementacija(?)... [slika]

### 5.2 Metoda podpornih vektorjev

(M.Bishop, sci-kit.org) - konkretna implementacija (scikit)? - kako se izogne prekomernem prilaganju, regularizacija

Metoda podpornih vektorjev (support vector machine - SVM) je ena od metod strojnega učenja. Osnovni algoritem lahko loči dva razreda, če so podatki linearno ločljivi. V fazi učenja algoritem najde hiperravnino, ki maksimizira razdaljo med razredoma. Pri določanju ravnine so pomembni le določeni vektorji (podporni vektorji), ki so najbližje ravnini (?).

Metoda se lahko prilagodi tudi za večrazredno klasifikacijo. Za to obstajata dva



načina, *ena na ena* in *eden proti vsem*.

Pri *ena na ena* ...

Pri *eden proti vsem* ...

Osnovni algoritem se odreže dobro na linearno ločljivih podatkih, vendar ne zmore ločiti nelinearnih podatkov, tudi ko obstaja hiperploskev, ki jih loči. Ta problem rešujemo tako, da podatke najprej preslikamo v višjedimenzijski prostor, kjer postanejo linearno ločljivi. Preslikava v višjo dimenzijo poveča računsko zahtevnost problema. kernel trick...

Iskanje hiperravnine je konveksen optimizacijski problem, najdena rešitev pa je globalni minimum. V tem pogledu se metoda podpornih vektorjev razlikuje od umetnih nevronskih mrež.

[slika]

## 5.3 Umetna Nevronska mreža

## 5.4 Konvolucijska nevronska mreža

## 6 Modeling / Classification

## 7 Results

## 8 ...

### 8.1 ...

## 9 CONCLUSION / FINAL THOUGHTS

## 10 LITERATURA IN VIRI

- [1] MNIST handwritten digit database, Yann LeCun, Corinna Cortes and Chris Burges. (*Citirano na strani 7.*)