# Plagiarism Checker X - Report

## Originality Assessment

# 1%

**Overall Similarity**

**Date:** Jan 12, 2025 (02:01 PM)
**Matches:** 13 / 2460 words
**Sources:** 1

**Remarks:** Low similarity detected, consider making necessary changes if needed.

**Verify Report:**
Scan this QR Code

AIE425IntelligentRecommenderSystems,FallSemester24/25

PersonalizedEducationalGameRecommendationSystem

221100583,MohammedHaithamMohammed 221100500,MaryamEslamSameer

2.Introduction: ThePersonalizedEducationalGameRecommendationSystemisa

groundbreakingapplicationdesignedtobridgethegapbetweentechnology,

education,andentertainment.Intoday'sdigitalage,educationalgameshave

becomeavitalresourceforfosteringcognitivegrowth,enhancingproblemsolvingabilities,andeng

aginglearnersofallages.However,thesheer

volumeandvarietyofeducationalgamesavailableoftenmakeitdifficultfor

userstoidentifygamesthatalignwiththeirinterests,educationalgoals,and

agegroup.Thisprojectaddressesthischallengebyintroducinga

personalizedrecommendationsystemthattailorsgamesuggestionsto

individualneedsandpreferences.

Thissystemisbuiltonrobustmachinelearningalgorithmsandsemanticdata

analysistechniquestoensureprecisionandrelevanceinrecommendations.It leveragesawell-curateddatasetenrichedwithmetadata,includinggenres, descriptions,ratings,andagerequirements.Bypreprocessingandanalyzing thisdata,thesystemcreatespersonalizedrecommendationsthatarenotonly engagingbutalsoeducationallybeneficial.Designedforstudents,educators, parents,anddevelopers,thesystemprovidesauser-friendlyinterfacethat simplifiestheprocessofdiscoveringmeaningfuleducationalgames. 2.1.SystemOverview: ThePersonalizedEducationalGameRecommendationSystemoperates asaspecializedenginedesignedexclusivelyfortheeducationalgaming niche.Unlikegeneral-purposerecommendationengines,thissystem focusesonthedualgoalsofengagementandpedagogy.Byanalyzing userinputssuchasfavoritegenres,previouslyplayedgames,andage group,thesystemprovidesrecommendationsthatbalanceentertainment valuewitheducationalobjectives.

Atitscore,thesystemisdrivenbyarichdatasetcomprisingdetailed attributesforeachgame.Theseattributesincludegamedescriptions, genreclassifications,ageratings,userreviews,andothermetadata. Advancedpreprocessingtechniquesensurethatthisdataisclean, consistent,andactionable.Forinstance,thesystemclusterssimilargame titles,removesnoisefromtextualdescriptions,andencodesuser preferencestofacilitateefficientcomputation.

Thesystemintegratesmultiplealgorithms,includingcollaborativefiltering, K-meansclustering,andsemanticanalysisusingTF-IDF.These algorithmsworktogethertoidentifypatternsinuserpreferences,group similargames,andextractmeaningfulinsightsfromtextualdescriptions. Bycombiningthesetechniques,thesystemdeliverspersonalized recommendationstailoredtoeachuser'suniqueprofile. Age-appropriatenessisakeyfeatureofthesystem.Usersundertheage

of15receiverecommendationsforgamessuitablefortheirdevelopmental stage,whileolderusersareofferedabroaderselectionofeducational tools.Thisfeatureensuresthatthesystemmeetstheneedsofitsdiverse audiencewhilemaintainingtrustandreliability. 2.2.ObjectivesoftheSystem: Thesystemhasbeendesignedwithclearandmeasurableobjectivesto addressthechallengesfacedbyusersindiscoveringrelevanteducational games:

□Personalization:Thesystemaimstoprovidehighlytailoredgame recommendationsbyanalyzinguserinputssuchasgameplayhistory, favoritegenres,anddemographicinformation.Thisensuresthateach userreceivessuggestionsthatalignwiththeiruniquepreferencesand goals. □EducationalRelevance:Acoreobjectiveofthesystemistoprioritize gameswithstrongpedagogicalvalue.Theseincludegamesdesignedto teachskillssuchasproblem-solving,teamwork,programming,orliteracy. □AgeAppropriateness:Thesystemenforcesstrictage-basedfiltering, ensuringthatyoungerusersareexposedtocontentsuitablefortheir cognitiveandemotionaldevelopment. □AdvancedTechnology:Byemployingstate-of-the- artalgorithms,the systemenhancestheaccuracyandrelevanceofitsrecommendations. Techniquessuchascollaborativefiltering,clustering,andsemantic analysisareusedtoidentifymeaningfulrelationshipsbetweenuser preferencesandgameattributes. □ScalabilityandAdaptability:Thesystemisbuilttohandlelargedatasets andintegratefutureexpansions,suchasnewdatasourcesoremerging algorithms,ensuringlong- termrelevanceandusability. Theseobjectivesguidethedesignandimplementationofthesystem, ensuringitservesawiderangeofuserswhilemaintainingafocuson personalizationandeducationalvalue. 2.3.PracticalApplicationDomain: Educationalgamesrepresentauniqueintersectionoftechnologyand learning,combiningtheimmersivequalitiesoftraditionalgameswith

structureddedicationalobjectives.Thesystemfocusesspecificallyonthis domain,makingitavaluabletoolforavarietyofusecases.

Educationalgamescatertodiverseaudiences,fromyounglearners developingfoundationalskillstoolderusersseekingadvancedknowledge. Forexample:

☐ForYoungLearners:Gamesdesignedforusersunder15oftenfocuson teachingbasicliteracy,numeracy,andlogicalreasoning.Thesegames aretypicallyvisuallyengagingandinteractive,fosteringaloveforlearning.

☐ForOlderUsers:Advancededucationalgamestargetskillssuchas teamwork,criticalthinking,andprogramming.Thesegamesoften simulatereal-worldscenarios,suchasfinancialmanagementorscientific research,makingthemidealforolderstudentsandprofessionals.

☐ForInstitutions:Schoolsanduniversitiesincreasinglyusegamified learningtoolstoenhancestudentengagementandoutcomes.This systemhelpseducatorsidentifygamesthatalignwiththeircurriculaand teachingobjectives.

Byaddressingtheneedsoftheseaudiences,thesystemnotonly enhancesindividuallearningexperiencesbutalsocontributestothe broaderadoptionofeducationaltechnology. 2.4.StakeholderExpectations: ThePersonalizedEducationalGameRecommendationSystemis designedtomeetthediverseneedsofitsstakeholders,eachofwhomhas distinctexpectations:

☐StudentsandLearners:Expectengagingandpersonalized recommendationsthatalignwiththeirinterestsandeducationalgoals.

☐ParentsandEducators:Relyonthesystemtoprioritizeage-appropriate contentandeducationalvalue,ensuringthatrecommendedgamesare bothsafeandbeneficial.

☐GameDevelopers:Benefitfrominsightsintouserpreferences,allowing themtocreatetargetedandsuccessfulproducts.

☐EducationalInstitutions:Usethesystemtoidentifyhigh-quality educationalgamesthatcomplementtheirteachingmethodsandcurricula.

To meet these expectations, the system incorporates advanced algorithms and user-friendly design principles. For example, it provides detailed explanations for each recommendation, helping users understand why specific games are suggested. 2.5. System Design and Constraints: The system's architecture has been carefully crafted to balance technical sophistication with practical usability. It integrates multiple components, each designed to address specific challenges in the educational gaming domain.

2.5.1. Data Model: The dataset serves as the backbone of the system, containing rich metadata for each game, including: ☐ Game Descriptions: Processed using NLP techniques to extract semantic information. ☐ Genre Classifications: Encoded to facilitate clustering and filtering. ☐ User Ratings and Reviews: Used to gauge the popularity and quality of games. ☐ Age Requirements: Enforced through strict filtering mechanisms. Preprocessing steps, such as clustering similar games and cleaning textual data, ensure that the dataset is consistent and actionable. 2.5.2. Algorithmic Framework: The recommendation engine integrates multiple algorithms to enhance accuracy and relevance: ☐ Collaborative Filtering: Identifies patterns in user preferences and interactions. ☐ Clustering (K-means): Groups games with similar attributes, ensuring that recommendations capture shared characteristics. ☐ Semantic Analysis (TF-IDF): Extracts meaningful relationships from game descriptions, aligning recommendations with user preferences. These algorithms work in harmony to deliver recommendations that are personalized, educational, and contextually relevant. 2.5.3. Constraints: ☐ Age-Based Filtering: Recommendations are tailored to the user's age group to ensure suitability. ☐ Educational Focus: Games without clear pedagogical objectives are excluded from recommendations.

☐ Data Quality: The system's performance relies on the completeness and accuracy of the dataset, emphasizing the importance of rigorous preprocessing.

2.6. SystemFeaturesandComponents: Component Description Dataset Includesmetadatasuchasgenres,ageratings,descriptions,anduser reviews. Preprocessing Techniques Textcleaning,clusteringofgameversions,andsemanticanalysisvia TF-IDF. Recommendation Methods Collaborativefiltering,clustering(K-means),andage-basedfiltering. UserInputs Agegroup,favoritegenres,previouslyplayedgames,andlearning preferences. Outputs Rankedlistofeducationalgameswithdetailedexplanationsfor recommendations. 3. DataCollectionandPreprocessing: ThePersonalizedEducationalGameRecommendationSystemdependson high-quality,structureddatatoproducemeaningfulandrelevant recommendations.Datacollectionandpreprocessingformthecornerstoneof thisproject,transformingrawinformationintoactionableinsightsthatpower thesystem'salgorithms.Thesestepsarenotmerelypreparatory—theyare integraltothesystem'sperformance,accuracy,andabilitytodeliver personalizedsuggestionsthatmeetuserpreferencesandconstraints. Thedatasetusedforthisprojectprovidesmetadatafor980games,each describedthroughavarietyofattributesthatcapturetheirunique characteristics.Thismetadataofferscriticalinsightsintotheeducationalvalue, gameplaystyle,anduserexperienceofeachgame,enablingthesystemto alignitsrecommendationswithspecificuserneeds.However,aswithmost datasets,therawdatawasincomplete,inconsistent,andscatteredacross multiplesources,necessitatingarigorousandmulti-steppreprocessing pipeline.3.1.DataCollection: Datacollectionisthefirstandperhapsthemostcriticalstageinbuildinga recommendationsystem.Forthisproject,thedatasetwassourcedfroma publicrepositorycontainingrichmetadataaboutgamesspanningvarious genresandplatforms.Thisdatasetwasselectedforits comprehensiveness,coveringessentialfieldssuchasgamenames, genres,descriptions,ratings,agerequirements,andsupportedlanguages.

However, while the dataset provided a strong foundation, it required significant enhancements to meet the system's requirements for personalization and accuracy.

The dataset was composed of structured fields that captured both quantitative and qualitative aspects of games. The Name column served as the unique identifier for each game, while the Short Description column provided a brief narrative about the game's content and objectives. Other critical fields included Required Age, which ensured age-appropriate

recommendations, and Ratings, which provided a numerical measure of user satisfaction. The Genres columns listed the game's classifications, such as "Educational," "Action," or "Adventure," while the Supported Languages column indicated the languages in which the game was available, offering an additional layer of filtering based on user preferences.

3.2. Challenges in Data Collection:

Despite its richness, the dataset had several limitations. Many entries were incomplete, with missing values in fields such as Short Description or Genres. Some games appeared multiple times under different editions or versions, creating redundancies that could distort the recommendation process. Additionally, the dataset lacked information about certain fields, such as promotional images and detailed genre hierarchies, which were deemed critical for enhancing the user experience.

To address these challenges, supplementary data was gathered from reputable secondary sources, including official game websites, online gaming databases like Steam and IGDB, and game developer pages. These additional sources were used to fill missing values, enrich existing fields, and validate the accuracy of critical attributes. For instance, where Genres were incomplete or ambiguous, external databases were referenced to provide a more detailed classification. Missing descriptions

were supplemented with summaries scraped from official websites, ensuring that every game had sufficient contextual information for analysis.

3.3. Data Preprocessing:

The preprocessing phase transformed the raw, inconsistent dataset into a clean, structured format suitable for machine learning algorithms. This step involved a combination of cleaning, enrichment, transformation, and feature engineering to ensure the data's usability and relevance. Each stage of preprocessing was designed to address a specific issue in the dataset, from missing values to inconsistent naming conventions.

3.4. Cleaning and Normalization:

Cleaning the dataset was the first step in preprocessing, focusing on removing inconsistencies and standardizing the data. Missing values in critical fields such as Ratings and Required Age were addressed using domain-specific imputation strategies. For example, missing ratings were filled with the median rating of games within the same genre, ensuring that the imputed values did not distort the overall distribution. Similarly, missing descriptions were replaced with placeholders such as "No description available," which could be excluded from downstream analysis. Text normalization was applied to all text-based fields, including Name, ShortDescription, and Genres. This involved converting all text to lowercase, removing special characters, and trimming excessive whitespace to ensure uniformity across entries. These steps reduced variability in the data and improved the performance of algorithms that rely ontext-based features.

Duplicate entries were another significant challenge, particularly for games released in multiple editions or with slight variations in their titles. Using a combination of clustering and manual verification, duplicate rows were identified and merged into single, unified entries. This process

ensuredthateachgamewasrepresenteduniquelyinthedataset, eliminatingredundanciesthatcouldskewtherecommendations. 3.5.SemanticEnrichment: OneofthemostvaluablefieldsinthedatasetwastheShortDescription, whichprovidedanarrativeoverviewofeachgame.Tomakethese descriptionsusableformachinelearningalgorithms,theywere transformedintonumericalrepresentationsusingTermFrequency-Inverse DocumentFrequency(TF-IDF).Thismethodquantifiedtheimportanceof wordsineachdescriptionrelativetotheentiredataset,enablingthe systemtoidentifysemanticsimilaritiesbetweengames.Forexample,a gamedescribedas"Learnmaththroughengagingpuzzles"wouldbe semanticallylinkedtoothereducationalgameswithsimilarthemes. ToreducethedimensionalityoftheTF-IDFvectorsandenhance computationalefficiency,SingularValueDecomposition(SVD)was applied.Thistechniqueretainedthemostrelevantfeaturesofeachvector whilediscardingnoise,ensuringthatthesemanticrepresentationofeach gamewasbothaccurateandcompact. 3.6.UnificationandGrouping: Gametitlesoftenappearedinmultipleformsduetovariationsinnaming conventions.Forexample,"7DaystoDie"and"7DaystoDieDeluxe Edition"weretreatedasseparateentriesintherawdataset.Toaddress this,atitleunificationprocesswasimplementedusingacombinationof clusteringalgorithmsandmanualcuration.Titleswereclusteredbasedon theirsemanticsimilarity,andthemostrepresentativenamewasselected astheunifiedtitleforeachgroup.Thisensuredthatthedatasetaccurately reflectedtheuniqueidentitiesofthegames. 3.7.GenreEncoding: TheGenresfieldswereamongthemostcomplexattributesinthedataset, withuptosevenseparatecolumnsforeachgame'sclassifications.To simplifyprocessing,thesefieldswereconsolidatedintoasingle categoricalfeature,witheachgenrerepresentedasabinarycolumn.This one-

hotencodingapproachallowedthesystemtoidentifygamesthat spannedmultiplegenres,suchas"Educational"and"Adventure,"without losinggranularity.Missinggenrevalueswerereplacedwith"Unspecified," ensuringthateverygamecouldbeincludedintheanalysis.

### 3.8. Final Dataset Transformation:

Afterpreprocessing,thedatasetwastransformedintoastructuredformat optimizedforrecommendationalgorithms.Thetablebelowillustratesa sampleofthefinaldataset:

| Name | TF-IDF Vectors (Sample) | IsFree | Short Description | Supported Languages | Header Image | Website | Ratings | Genres |
|---|---|---|---|---|---|---|---|---|
| 7 Days to Die | [0.59,0.35,-0.40] | FALSE | An open world game combining first-person shooter, survival horror, and RPG elements. | English, French, German, Spanish | Link | http://www.7daystodie.com | 7 | Multiplayer, Action, Adventure |

| Name | TF-IDF Vectors (Sample) | IsFree | Short Description | Supported Languages | Header Image | Website | Ratings | Genres |
|---|---|---|---|---|---|---|---|---|
| A Dance of Fire and Ice | [0.32,0.38,0.68] | FALSE | A rhythm game where you guide two planets along a winding path. | English, Spanish, Korean | Link | https://7thbe.at | 6 | Indie |
| A Plague Tale: Innocence | [0.52,0.51,0.20] | FALSE | Follow the tale of Amicia and Hugo through the darkest hours of history. | English, French, German | Link | https://www.focushome.com/games/a-plague-tale-innocence | 7 | Singleplayer, Action, Adventure |
| A Story About My Uncle | [0.62, 0.20,-0.09] | FALSE | A platforming adventure about a boy searching for his uncle. | English, French, German | Link | http://gonenorthgames.com/games/a-story-about-my-uncle/ | 6 | Singleplayer, Adventure, Indie |
| A Way Out | [0.48,0.33,0.41] | FALSE | A cooperative prison escape game. | English, French, German | Link | https://www.ea.com/games/a-way-out | 7 | Multiplayer, Action, Adventure |

Thisstructureddatasetpowerstherecommendationenginebyproviding clean,enriched,andactionableinformationabouteachgame.

### 3.9. Role of Preprocessing in Recommendations:

Thepreprocesseddatasetisthebackboneoftherecommendationsystem, enablingitto:

1. AccuratelycalculatesemanticsimilaritiesbetweengamesusingTF-IDF vectors.

2. Groupgameswithsharedcharacteristicsthroughclusteringandgenre encoding.

3. Filterrecommendationsbyuser-specificcriteria,suchasage,language,

andpricingpreferences.

Bytransformingrawdataintostructuredinsights,preprocessingensuresthat

thesystemdeliversaccurate,relevant,anduser-centricrecommendations.

4.DatasetDescription:

ThedatasetpoweringthePersonalizedEducationalGameRecommendation

Systemcomprisesmetadatafor980games.Thismetadataisstructured

acrossmultiplecolumns,eachprovidingcriticalinformationaboutthegames'

features,accessibility,anduserexperience.Thedatasethasbeencarefully

processedtoensureitsupportstherecommendationsystem'sobjectiveof

deliveringpersonalized,relevant,andeducationalgamerecommendations.

4.1.DetailedDescriptionoftheDataset: Thedatasetcontainsthefollowingcolumns:

□Name:Theuniquetitleofeachgame,servingasitsprimaryidentifier.

□RequiredAge:Theminimumagerecommendedforplayingthegame, ensuringage-

appropriatesuggestions. □IsFree:Abinaryindicatorspecifyingwhetherthegameisfreetoplay.

□ShortDescription:Aconcisesummaryofthegame's content,objectives, andgameplaystyle.

□SupportedLanguages:Alistoflanguagesinwhichthegameisavailable.

□HeaderImage:Linkstopromotionalimagesthatrepresentthegame visually.

□Website:URLslinkingtoofficialgamepagesorstores.

□Ratings:Anumericalscorereflectingthequalityorpopularityofthegame.

□Categories:Ageneralclassificationofthegame,suchas"Multiplayer"or "Singleplayer."

□Genres(genre1togenre7):Uptosevengenreclassifications,capturing

thegame'smultifacetedcharacteristics.

4.2.ModelingUserInterests,Interactions,andIntentions:

Thedatasethasbeenstructuredtomodeluserpreferencesandbehaviors

throughitsvariouscolumns.Thismodelingenablesthesystemtodeliver

personalizedrecommendationsbyleveragingthefollowingfactors: 4.2.1.UserInterests:

UserinterestsareprimarilycapturedthroughtheCategoriesandGenres columns.Forinstance:

◻Usersinterestedinmultiplayerexperiencesaredirectedtowardgames labeledas"Multiplayer"underCategories. ◻Themulti-genreclassification(genre1togenre7)supports recommendationsforuserswhoenjoycomplexordiversegameplay styles,suchascombining"Educational"with"Simulation." 4.2.2.UserInteractions: Thedatasetallowsthesystemtoanalyzepatternsingamepreferencesbased onsharedattributeslikegenresandratings.Usingclusteringalgorithms, gameswithoverlappingfeaturesaregroupedtogethertoensure recommendationsaligncloselywithuserpreferences. 4.2.3.UserIntentions: TheRequiredAgecolumnensuresthatrecommendationsalignwithuserage groups,whileIsFreeenablesbudget-conscioususerstoreceivefree-to-play options.Theseexplicitandimplicitsignalsrefinetherecommendationprocess tomatchuserexpectations. 4.3.ExplanationofDatasetColumns: Thedatasetcolumnsaredetailedbelow,outliningtheirpurposeandrole intherecommendationsystem: 4.3.1.Name: Thiscolumncontainsthegametitles,ensuringeachgameisuniquely identifiedinthedataset.Duplicatetitles,suchaseditionsorspecialreleases, werestandardizedduringpreprocessingtomaintainclarity. 4.3.2.RequiredAge: Theminimumrecommendedageforplayingeachgame.Thisfieldfilters recommendationstoensureusersbelow15yearsoldreceiveage-appropriate suggestions,whileolderusersarepresentedwithunrestrictedoptions. 4.3.3.IsFree: Abinaryfield(TrueorFalse)indicatingwhetheragameisfreetoplay.This columnallowsthesystemtofilterrecommendationsbasedontheuser'sbudgetpreferences. 4.3.4.ShortDescription:

Thiscolumnprovidesaconcisesummaryofthegame'scontent.Forexample: ◻"ADanceofFireandIce"isdescribedas"Astrictrhythmgamefocusing ontiming." ◻ThesedescriptionswereprocessedusingTF-IDFvectorizationto

calculatesemanticsimilaritiesbetweengames,enablingthesystemto groupandrecommendgameswithoverlappingthemes. 4.3.5.SupportedLanguage: Thiscolumnlistsalllanguagessupportedbythegame,ensuringaccessibility foradiverseuserbase.Forexample,agameavailableinEnglish,Spanish, andKoreanisprioritizedforuserswiththeselanguagepreferences. 4.3.6.HeaderImage: Linkstopromotionalimagesenhancetheuserinterfaceofthe recommendationsystem,offeringavisuallyengagingwaytopresentgame suggestions.

4.3.7.Website: ThiscolumnprovidesURLstoofficialgamepages,allowinguserstoexplore moredetails,reviews,orpurchaseoptionsfortherecommendedgames. 4.3.8.Ratings: Thenumericalratingofeachgame,typicallyonascaleof1to10.Higherratedgamesaregivenpriorit yinrecommendations,reflectinguser satisfactionandquality. 4.3.9.Categories: Abroadclassificationofthegame,suchas"Multiplayer"or"Singleplayer." Thisfieldisusedtoalignrecommendationswithuserpreferencesforsocialor sologameplay.

4.3.10.Genres(genre1togenre7): Themulti-genreclassificationcapturesthediversenatureofgames.These fieldsallowthesystemtorecommendgamesthatalignwithuser-selected genres.Forexample,agamecategorizedunder"Educational,""Adventure," and"Simulation"cancatertousersinterestedinacombinationofthese genres.4.4.SummaryofDatasetColumns: ColumnName Description RoleintheSystem Name Uniquegametitle. Ensures distinct representationofgames. RequiredAge Minimumrecommendedage.Filters age-appropriate games. IsFree Indicateswhetherthegameis free. Alignsrecommendationswith userbudgetconstraints. ShortDescription Summaryofgamecontent.Enablessemanticsimilarity analysis. SupportedLanguagesListslanguagesinwhichthe gameisavailable. Filtersgamesbasedon languagepreferences. HeaderImage Linktoapromotionalimage.Enhances the recommendationinterface

visually. Website OfficialgamepageURL. Directsuserstoadditional

gameinformation. Ratings Numericalratingfromusers.Prioritizeshigh-qualitygames inrecommendations. Categories Generalclassification(e.g., Multiplayer). Matchesuserpreferencesfor socialorsologameplay. Genres (genre1– genre7) Multi-categoryclassification ofgames. Alignsrecommendationswith user-selectedinterests.

5.DataAnalysisandInsights:

Thissectiondelvesintotheanalysisofthedataset,focusingoncritical attributessuchasrequiredage,ratings,genres,andclustering.Theinsights derivedfromtheseanalysesguidethedesignandimplementationofthe PersonalizedEducationalGameRecommendationSystem.Detailed visualizationsarepresentedandexplainedtohighlightkeytrendsand patterns.5.1.RequiredAgeDistribution:

TheRequiredAgecolumncategorizesgamesbasedontheminimum recommendedageofthetargetaudience.Thisfeatureensuresthatthe systemfiltersgamesappropriatelyforusersacrossdifferentagegroups. ThedistributionofagecategoriesisisillustratedinFigure1,apiechartthat revealstheproportionalrepresentationofagegroupswithinthedataset.

☐44.2%ofgamesaresuitableforusersaged11–15,thelargestagegroup. ThesegamestypicallyincludegenressuchasAdventure,Casual,and Educational,whichareengagingyetappropriateforyoungeraudiences. Thehighproportionreflectsthedemandforfamily-friendlyand educationalgames.

☐32.4%ofgamestargetusersaged16–18.Gamesinthiscategoryare oftenmorecomplexandincludeaction-packedgenreslikeAction,RPG, andStrategy.Thesegamescatertoolderteenagerswhoprefer immersivegameplayexperiences.

☐23.4%ofgamesareintendedforusersaged>18,includingmaturetitles ingenressuchasSurvival,Horror,andSimulation.Thesegamesmay containthemesormechanicsunsuitableforyoungeraudiences.

☐Anegligiblepercentage(0%)ofgamesarecategorizedforusersaged

≤10. This absence indicates a gap in the dataset for games tailored to very young audiences, which could be addressed by augmenting the dataset with additional titles.

Figure 1: Required Age Distribution

The chart underscores the dataset's emphasis on games for teenagers and young adults. The system uses this information to implement age-based filtering, ensuring users receive age-appropriate recommendations.

5.2. Game Ratings Analysis: Ratings are a critical metric in evaluating the quality and popularity of games. The dataset includes ratings on a scale from 1 to 10, representing user feedback and reviews. Figure 2 (boxplot) and Figure 3 (histogram) provide complementary views of the ratings distribution. 5.2.1. Boxplot Analysis: The boxplot in Figure 2 visualizes the spread, central tendency, and variability in game ratings:

☐ The median rating is approximately 7, indicating that most games are of good quality.

☐ The interquartile range (IQR) spans ratings between 6 and 8, which constitutes the majority of the dataset.

☐ Outliers exist at both extremes. For instance, games rated below 4 may represent niche or poorly received titles, while those rated above 9 are exceptional in quality and user appeal. These insights ensure the system prioritizes high-rated games while considering niche titles for users with specific preferences.

Figure 2: Ratings Boxplot 5.2.2. Histogram Analysis: The histogram in Figure 3 complements the boxplot by showing the frequency distribution of ratings:

☐ Ratings of 7 and 8 are the most common, with ~300 games in these categories. This highlights the dataset's focus on well-received games, ensuring a high-quality recommendation pool.

☐ Ratings below 5 are rare, with only a few games falling into this category. These games are likely to be excluded or deprioritized during recommendation generation.

□ A modest number of games achieve a perfect rating of 10, representing top-tier titles highly favored by users. Figure 3: Game Ratings Distribution

This analysis demonstrates the system's emphasis on quality, ensuring that recommended games align with user expectations for engaging and enjoyable experiences.

5.3. Genre Analysis: Genres are a defining feature of games, capturing their thematic and gameplay characteristics. The dataset includes up to seven genres per game, enabling multi-genre classifications. Figure 4 illustrates the top 10 most frequent genres, offering insights into user preferences and dataset diversity.

□ Action is the most dominant genre, appearing in over 500 games. This reflects its universal appeal, spanning various age groups and gameplay styles.

□ Indie ranks second, showcasing the growing popularity of independent games that emphasize creativity and unique mechanics.

□ Adventure follows closely, catering to users who enjoy exploratory and narrative-driven experiences. □ Other notable genres include RPG, Simulation, and Strategy, which appeal to users seeking depth and complexity in gameplay.

□ Early Access, while less frequent, represents games in development that attract users interested in beta testing and early-stage experiences.

□ Educational emerges as a noteworthy genre, targeting users interested in combining learning with entertainment. These games often include puzzles, simulations, or storytelling elements designed to teach while engaging the player. Educational games are prioritized for users who explicitly select this genre as a preference, emphasizing the system's goal of promoting both fun and learning. Figure 4: Top 10 Genres The inclusion of Educational games underscores the system's focus on promoting learning through gaming, catering to users who prioritize

intellectual growth alongside entertainment. This diversity allows the

recommendationsystemtoclusterandfiltergameseffectivelybasedonuserselectedpreferences ,ensuringtailoredsuggestionsthatresonatewith individualtastesandneeds.

5.4.ClusteringInsights: Clusteringplaysapivotalroleintherecommendationsystem,grouping gamesbasedonsemanticsimilarityandsharedattributes.The hierarchicalclusteringdendrograminFigure5visualizesrelationships betweengamesderivedfromtheirtextualdescriptions. 5.4.1.Methodology: □TF-IDFVectorization:Thistechniqueconvertsgamedescriptionsinto numericalrepresentations,quantifyingtheimportanceofwordsrelativeto thedataset. □CosineSimilarity:Measuresthesimilaritybetweendescriptions, groupingsemanticallysimilargamesintoclusters. □HierarchicalClustering:Constructsatree-likestructure(dendrogram)to representthenestedrelationshipsbetweengames.

5.4.2.Observations: □Gameswithinthesameclusteroftensharegenres,themes,orgameplay mechanics.Forinstance,aclustermightincludetitleslike"ActionAdventure"or"Simulation-Strategy." □Clusteringreducesredundancybyunifyingvariationsingamenames(e.g., "DeluxeEdition"and"OriginalVersion")underasinglerepresentative entry. □Thedendrogramrevealshigher-levelgroupings,suchasclusters dominatedby"Educational"gamesor"Multiplayer"titles.

Figure5:HierarchicalClusteringDendrogram

Theclusteringanalysisensuresthatrecommendationsarecontextually relevant,offeringusersacuratedlistofgamesthataligncloselywiththeir interests.5.5.KeyInsightsandImplications: Theinsightsderivedfromthedatasetanalysisareinstrumentalinshaping therecommendationsystem: 1.Age-SpecificFiltering:Theagedistributionensuresthesystem accommodatesdiverseuserdemographics,deliveringage-appropriate content.

2.QualityAssurance:Theprevalenceofgamesrated6–8highlightsthe dataset'sfocusonquality,supportingreliablerecommendations.

3.DiverseGenres:Thewiderangeofgenresallowsthesystemtocaterto

variedpreferences,fromcasualplayerstoenthusiastsofnichegenres.

4.SemanticClustering:Theuseofclusteringenhancesthesystem'sability

torecommendgamesthatsharethematicorgameplaysimilarities, increasingusersatisfaction.

6.RecommendationSystemDesign:

ThePersonalizedEducationalGameRecommendationSystemisengineered

tocombineuserpreferences,collaborativeinsights,andgamemetadatato

delivertailoredrecommendations.Thissectiondissectstherecommendation

engine'scomponents,detailingitsmethodology,mathematicalfoundations,

andpracticalimplementation.Thesystemintegratesmultiplemachine

learningtechniques,suchascollaborativefiltering,clustering,and

dimensionalityreduction,alongsidecustomscoringmetricstoensureprecise

andmeaningfuloutputs. 6.1.AlgorithmOverview:

Therecommendationsystemcombinesmultipleadvancedtechniquesto

providepersonalizedandmeaningfulrecommendations.Theseinclude

collaborativefiltering,clustering,TF-IDFvectorization,dimensionality

reductionthroughSVD,andweightedscoringmechanisms.Thesystem

ensuresthateveryrecommendationisbasedonstructuredandefficient computationalmodels.

6.1.1.CollaborativeFiltering(CF):

CollaborativeFiltering(CF)isacornerstoneoftherecommendationengine,

enablingittoidentifygamessimilartothoseplayedbytheuser.This

techniquereliesontheassumptionthatuserswhoenjoycertaingamesare

likelytoenjoyotherswithsimilarcharacteristics. Thesystememploysitem-

basedCF,wherethesimilaritybetweengamesis

measuredusingcosinesimilarity.Theformulaforcosinesimilarityisas follows: □ □□ □ □ □□

2n1i 2n1in1i =ilarity Cosine Sim i i ii B ABA

Here,AandBrepresentfeaturevectorsfortwogames.Thiscalculation

createsasimilaritymatrixthatcapturestherelationshipsbetweenallgamesin thedataset.

Forexample: ☐ConsiderGameAwithfeatures[1,0,1,0].

☐GameBhasfeatures[0,1,1,1]. 408.0321 11100101 0·1+1·1+0·1+1·0 =ilarity Cosine Sim

2222 2 22 2 ☐☐☐☐☐☐ ☐☐☐☐

Thissimilarityscoreensuresthattherecommendationsaregroundedinthe

relationshipsbetweengamefeatures,enhancingaccuracy. 6.1.2.ClusteringwithK-Means:

Tostreamlinetherecommendationprocess,gamesaregroupedintoclusters

basedontheirfeatures.ThesystemusesK-MeansClustering,which

partitionsdataintokclustersbyminimizingintra-clustervariance. StepsinK-Means:

1.Initialization:kcentroidsarechosenrandomlyinthefeaturespace.

2.Assignment:Eachgameisassignedtothenearestcentroidusing Euclideandistance: ☐☐ ☐☐

☐ ni iicx1 2 = Distance 3.Update:Thecentroidsarerecalculatedasthemeanofallpointsintheir

cluster. 4.Iteration:Steps2and3arerepeateduntilthecentroidsstabilize.

TheElbowMethoddeterminestheoptimalnumberofclusters(k)byplotting

theSumofSquaredDistances(SSD)againstvariouskvalues.Thepoint

wherethecurveflattensindicatesthebestk,balancingperformanceand computationalefficiency.

6.1.3.TF-IDFVectorization: TheTF-IDF(TermFrequency-

InverseDocumentFrequency)methodconverts text-

basedgamedescriptionsintonumericalrepresentations.Thisensures

thattherecommendationenginecapturesthesemanticessenceofeach

gamewhilereducingtheinfluenceoffrequentlyoccurring,less-informative words. d)·IDF(t)

TF(t,=d) IDF(t, -TF Where: ☐TF(t,d):Measureshowoftentermtappearsindocumentd: s in d

Total term in d Count of t =d)TF(t, ☐IDF(t):Weighsdowntermsappearinginmanydocuments:

☐☐☐☐☐☐☐ ☐☐☐☐☐☐☐ t containing Documents ments Total docu log=IDF(t)

Forexample: ☐If"strategy"appearsin2outof10documents:1.70=)102 log(=gy)

IDF(strate Thisensuresthatuniquetermslike"adventure"or"open-world"havegreater

influencethangenerictermslike"game." 6.1.4.SingularValueDecomposition(SVD):

Toreducethedimensionalityofthecombinedfeaturematrix(genresandTFIDF),thesystemapplie

sSingularValueDecomposition(SVD).SVDbreaks downthematrixMintothreecomponents:T

$U \cdot S \cdot V = M$ Where: ☐U:Capturesrelationshipsbetweengames.

☐S:Diagonalmatrixcontainingsingularvalues(featureimportance).

☐TV:Encodesrelationshipsbetweenfeatures.

Byretainingonlythetopksingularvalues,SVDreducesnoiseandenhances

computationalefficiency.Forexample,retainingthetop50singularvaluesin

amatrixofsize1000×500reducestheeffectivedimensionalitywhile

maintainingkeydatapatterns. 6.1.5.WeightedScoringMechanism:

Thefinalscoreforeachgameiscomputedusingaweightedsumofmultiple factors:

1.SimilarityScore(20%):Captureshowcloselythegamematchesthe user'splayedgames.

2.GenreOverlap(50%):Measuresalignmentbetweentheuser'spreferred

genresandthegame'sgenres: s User Genre Number of enres Matching G Number of =lap

Genre Over 3.Rating(30%):Incorporatesthegame'saveragerating. $0.3 \cdot Rating + Overlap$

$0.5 \cdot Genre +$ rity Score $0.2 \cdot Simila =e$ Final

ScorThisweightedmechanismensuresbalancedandmeaningful recommendations.

6.2.SystemWorkflow: Therecommendationsystem'sworkflowconsistsofasequenceof

interconnectedstepsdesignedtotransformrawdataintopersonalized

gamesuggestionsfortheuser.Belowisabreakdownofeachstageinthe

workflow,asdepictedintheflowchart(Figure6):

Figure6:RecommenderSystemWorkflow

Therecommendationprocessinvolvesfiltering,scoring,andranking

gamesbasedontheuser'sinput.Theworkflowisdesignedtomaximize

relevanceandusersatisfaction. 6.2.1.Step1:UserInput:

Therecommendationprocessbeginswithcollectingcriticalinformationfrom

theuser.Thisstepensuresthesystemcapturestheuser'spreferencesand

personalcontext,formingthebasisforpersonalizedrecommendations.

☐AgeGroupSelection:Usersspecifywhethertheyareinthe"below15"or

"15 and above" age group. This step is essential for ensuring that recommendations are age-appropriate and adhere to content restrictions. For instance, users under 15 will not see games tagged with a 15+ or 18+ age requirement.

☐ Played Games List: Users provide the names of up to five games they have played. These games serve as a reference point for analyzing their gaming interests and identifying similar games in the dataset.

☐ Favorite Genres Selection: Users select up to three preferred genres, such as Action, Adventure, or Strategy. These choices allow the system to prioritize games that align with the user's tastes.

This user-provided data serves as the foundation for subsequent steps in the workflow, allowing the system to personalize its recommendations effectively.

6.2.2. Step 2: Data Filtering:

Once the user inputs are gathered, the system filters the dataset to narrow down the list of games. This step eliminates games that do not meet the user's basic criteria and ensures only relevant options remain.

☐ Age Filtering: Games with a required_age higher than the user's specified age group are excluded from consideration. For example, if a user below 15 selects games, the system automatically excludes games tagged as suitable for players aged 15 or older. ☐ Genre-Based Prioritization: The system identifies games that share genres with the user's favorite genres. For example, if the user selects Strategy and Role-Playing, the system prioritizes games that belong to these categories. This prioritization increases the likelihood of recommending games that resonate with the user's preferences.

Filtering significantly reduces the pool of games and ensures only those relevant to the user's age group and genres are considered for further analysis.

### 6.2.3. Step 3: Similarity Computation:

To personalize recommendations further, the system calculates similarity scores for all games in the filtered dataset. This step leverages mathematical techniques to identify games that share features with the ones the user has played. ☐ TF-IDF Matrix Creation: The system first generates a Term Frequency Inverse Document Frequency ( TF-IDF) matrix using the descriptions or metadata of the games. This matrix numerically represents the importance of terms within the game descriptions, emphasizing unique and distinguishing features.

☐ Cosine Similarity Calculation: The cosine similarity formula is applied to measure the similarity between vectors in the TF-IDF matrix. It calculates how closely the descriptions of two games align. For example, a game with high similarity to a user-selected game may share keywords like "multiplayer,""action,"or"strategy."

☐ Score Aggregation: The similarity scores between each user-selected game and all other games are averaged. This approach ensures that the final similarity score reflects how well a game aligns with the user's overall gaming habits and preferences.

This step is critical for identifying games that not only share genres but also align with the themes, features, and descriptions of the games the user has already enjoyed.

### 6.2.4. Step 4: Weighted Scoring:

To rank the games, the system calculates a weighted score for each game based on multiple factors. This ensures that the final recommendations balance similarity, user preferences, and overall quality.

☐ Similarity Scores: These scores, derived from the cosine similarity calculations, contribute 20% to the final score. This weighting ensures

that games similar to the ones the user has played are included but not

overlydominant. □GenreOverlap:Genrealignmentwiththeuser'spreferencescontributes 50%tothescore.Thishighweightingreflectstheimportanceof matchinggamestotheuser'sstatedfavoritegenres,ensuringrelevance.

□GameRatings:Theratingsofgamesfromthedatasetcontribute30%to thefinalscore.High-ratedgamesaremorelikelytoberecommended, astheyareconsideredbetterqualityandmoreenjoyable.

Theweightedscoringformulaensuresabalancedapproach,combining multipleaspectsofagame'srelevanceandqualityintoasinglescore.This stepiscrucialforrankingthegameseffectively. 6.2.5.Step5:RecommendationOutput: Afterscoring,thesystemselectsthetopfivegameswiththehighestfinal scores.Thesegamesarepresentedtotheuserinanorganizedanduserfriendlyformat.

□GameTitle:Eachrecommendationincludesthenameofthegame, allowingtheusertorecognizeiteasily.

□AgeRequirement:Theagesuitabilityofthegameisclearlydisplayed, ensuringtheuserknowswhetherthegameisappropriatefortheirage group.

□Ratings:Theaverageuserratingisincludedforeachgame,providing insightintoitsqualityandpopularityamongotherplayers.

□KeyFeatures:Additionalinformation,suchasthegenresandstandout featuresofthegame,ishighlightedtoexplainwhyitwasrecommended.

·Therecommendationsarepresentedinawaythatemphasizesclarityand usability.Byprovidingdetailedinformationforeachgame,thesystemensures userscanmakeinformeddecisionsaboutwhichgamestoexplorefurther.

6.3.ExplanationofTablesandCalculations:

Thissectionelaboratesonthenumericalcomputationsandinterpretations ofthekeymatricesusedintherecommendationsystem.Eachmatrixis vitalinprocessingdataandderivingmeaningfulrecommendations.Below aretheextendedtablesforTF-IDFMatrix,SVDResults(U,Σ,V$^{\mathsf{T}}$ Matrices),ReconstructedMatrix,andCosineSimilarityMatrix,with

additionalrowsandcolumnsforaclearerunderstandingoftheir implications. 6.3.1.TF-IDFMatrix(First15×15Slice): TheTF-IDFmatrixrepresentstheimportanceofdifferentgenres,keywords, orattributesforeachgameinthedataset.Eachrowcorrespondstoagame, andeachcolumnrepresentsafeatureextractedfromthegamedescriptions. Thevaluesinthismatrixreflecttheweightorimportanceofeachfeaturetoa specificgame,calculatedusingtheTermFrequency-InverseDocument Frequencytechnique. Forexample,ifGame1hasaTF-IDFvalueof0.25forthe"Adventure"genre, itindicatesthat"Adventure"isahighlyrelevantcharacteristicofthisgame basedonitsdescription.Conversely,avalueof0for"Puzzle"inthesamerow impliesthat"Puzzle"isirrelevanttothisgame.TheTF-IDFweightshelp

emphasizefeaturesuniquetoagamewhiledownplayinggenericones, ensuringthematrixcapturesthegame'suniqueaspects. Thismatrixplaysafoundationalroleinidentifyingpatternsandrelationships betweengames.Byquantifyingthetextualinformationinthedataset,itallows thesystemtocomputemeaningfulsimilaritiesbetweengames,whichare essentialforgeneratingrecommendations. TheTF-IDFmatrixrepresentstheweightedimportanceoftermsingame descriptions. Game Index StrategyAdventure RPG Multiplayer OpenWorld ActionSimulation Puzzle ShootingSurvivalSandboxHorror Sports Sci-Fi Fantasy Game1 0.180.250.2 0.1 0.150.220.120.00.10.050.080.00.120.070.15 Game2 0.1 0.150.0 0.2 0.050.080.180.00.150.0 0.00.10.00.20.12 Game3 0.0 0.220.18 0.0 0.20.150.0 0.10.150.180.00.080.10.150.1 Game4 0.150.180.1 0.250.150.20.120.10.180.050.10.080.120.10.15 Game5 0.12 0.10.150.150.220.080.180.120.00.120.10.120.080.20.18 6.3.2.SVDResults: SingularValueDecomposition(SVD)isappliedtotheTF-IDFmatrixtoreduce itsdimensionality.Thistechniquebreaksthematrixintothreecomponents:U,

☐, and TV, each serving a distinct purpose.

The U matrix maps games to the new latent features (principal components) discovered during dimensionality reduction. Each row corresponds to a game, while the columns represent these components. The values indicate how strongly a game aligns with each component. 6.3.2.1. U Matrix (15×15 Slice): For instance, if Game 1 has a value of 0.065 in the first column of the U matrix, it means that the first principal component significantly influences this game's representation. These principal components are abstract features that combine original attributes such as genres or themes. For example, a single component might encapsulate a combination of "Adventure" and "Strategy." The U matrix reduces the computational complexity of similarity measurements. Instead of comparing games across hundreds of features, the system now works with a few high-level components, enabling faster and more efficient calculations.

RowIndexComponent 1 1 Component 2 Component 3 Component 4 Component 5 Component 6 Component 7 Component 8 Component 9 Component 10

| RowIndex | Component 1 | Component 2 | Component 3 | Component 4 | Component 5 | Component 6 | Component 7 | Component 8 | Component 9 | Component 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Game 1 | 0.065 | -0.04 | 0.03 | -0.028 | 0.056 | 0.07 | -0.045 | 0.031 | 0.051 | 0.062 |
| Game 2 | 0.045 | 0.053 | -0.035 | 0.022 | -0.058 | 0.012 | 0.062 | -0.038 | 0.044 | -0.046 |
| Game 3 | 0.03 | 0.02 | 0.065 | -0.04 | 0.035 | 0.025 | -0.046 | 0.06 | 0.03 | -0.032 |
| Game 4 | -0.028 | 0.022 | 0.02 | 0.063 | -0.052 | 0.048 | -0.043 | 0.056 | 0.062 | -0.04 |
| Game 5 | 0.056 | -0.058 | -0.045 | 0.035 | 0.03 | -0.034 | 0.058 | 0.051 | -0.04 | 0.065 |

6.3.2.2. Σ Matrix (Singular Values):

The Σ matrix contains the singular values, which represent the importance of each principal component. Larger values indicate components that capture more variance or information from the original data. For instance, the first component might have a singular value of 28.02, signifying its dominance in representing the dataset. In contrast, components with smaller values, such as 12.40, contribute less and are often disregarded during dimensionality reduction.

By focusing on components with high singular values, the system retains the most meaningful patterns in the data while discarding noise and redundant information. This ensures that the recommendation system operates efficiently without sacrificing accuracy.

| Component | Value |
|---|---|
| 1 | 28.02 |
| 2 | 21.45 |
| 3 | 19.3 |
| 4 | 15.8 |
| 5 | 12.4 |

### 6.3.2.3. V Transposed ($V^T$) Matrix:

The TV matrix relates the original features (genres and keywords) to the principal components. Each row represents a feature, and each column represents a principal component. The values in this matrix indicate the contribution of each feature to a specific component.

For example, if the "Strategy" genre has a value of 0.62 in the first column, it means this genre heavily influences the first principal component. Negative values indicate an inverse relationship, meaning features with negative weights detract from the component's representation. Understanding the TV matrix allows us to interpret the latent features and identify which combinations of genres or attributes drive similarities between games.

| Feature | Component 1 | Component 2 | Component 3 | Component 4 | Component 5 |
|---|---|---|---|---|---|
| Strategy | 0.62 | 0.45 | -0.12 | 0.34 | 0.41 |
| Adventure | 0.52 | 0.21 | 0.31 | -0.48 | -0.32 |
| RPG | 0.33 | 0.55 | 0.18 | 0.25 | -0.4 |
| Multiplayer | 0.45 | -0.3 | 0.48 | 0.22 | 0.5 |
| Open-World | 0.29 | 0.38 | -0.25 | -0.39 | 0.48 |

### 6.3.3. Reconstructed Matrix (15×15 Slice):

The reconstructed matrix approximates the original TF-IDF matrix after dimensionality reduction. By multiplying the U, □, and TV matrices, we obtain a compressed version of the original data. This reconstruction captures the most significant relationships between games and features while discarding less relevant details.

For example, if the reconstructed value for "Strategy" in Game 1 is 0.18, it closely aligns with the original TF-IDF value of 0.20, indicating that the dimensionality reduction preserved the critical information for this feature.

Minordiscrepanciesbetweentheoriginalandreconstructedvaluesare expectedandrepresentthetrade-offbetweendatacompressionand informationretention. Thereconstructedmatrixvalidatestheeffectivenessofthedimensionality reductionprocess.Itensuresthatthecriticalaspectsofthedataare preserved,enablingthesystemtogeneratereliableandmeaningful recommendations.

| Game Index | Strategy | Adventure | RPG | Multiplayer | OpenWorld | Action | Simulation | Puzzle | Shooting | Survival | Sandbox | Horror | Sports | Sci-Fi | Fantasy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Game1 | 0.18 | 0.26 | 0.22 | 0.1 | 0.12 | 0.23 | 0.14 | 0.1 | 0.11 | 0.07 | 0.08 | 0.1 | 0.08 | 0.12 | 0.13 |
| Game2 | 0.12 | 0.18 | 0.1 | 0.2 | 0.1 | 0.18 | 0.15 | 0.1 | 0.1 | 0.12 | 0.12 | 0.08 | 0.09 | 0.1 | 0.14 |
| Game3 | 0.08 | 0.22 | 0.2 | 0.1 | 0.25 | 0.2 | 0.08 | 0.12 | 0.11 | 0.1 | 0.12 | 0.15 | 0.11 | 0.09 | 0.12 |

6.3.4.CosineSimilarityMatrix(15×15Slice):

Thecosinesimilaritymatrixquantifiesthesimilaritybetweengamesbasedon theirfeaturerepresentations.Eachrowandcolumncorrespondstoagame, andthevaluesrangefrom-1to1.Avalueof1indicatesperfectsimilarity, whileavalueof0signifiesnosimilarity.Negativevalues,thoughrareinthis context,wouldindicatedissimilarity.

Forinstance,Game1hasasimilarityscoreof1.00withitself,asexpected. However,itmightalsohaveascoreof0.75withGame3,indicatingthatthese gamessharemanycommonfeatures,suchasoverlappinggenresorsimilar descriptions.Incontrast,ascoreof0.35withGame2suggestsfewershared features. Thismatrixiscriticalfortherecommendationengine.Bycomparingauser's previouslyplayedgamestotherestofthedataset,thesystemidentifiesthe mostsimilargamesandrecommendsthem.Thecosinesimilaritymatrix ensuresthatrecommendationsarenotonlyrelevantbutalsopersonalizedto theuser'spreferences.

| Game Index | Game1 | Game2 | Game3 | Game4 | Game5 |
|---|---|---|---|---|---|
| Game1 | 1.0 | 0.48 | 0.65 | 0.53 | 0.59 |
| Game2 | 0.48 | 1.0 | 0.45 | 0.52 | 0.62 |
| Game3 | 0.65 | 0.45 | 1.0 | 0.47 | 0.55 |
| Game4 | 0.53 | 0.52 | 0.47 | 1.0 | 0.61 |
| Game5 | 0.59 | 0.62 | 0.55 | 0.61 | 1.0 |

7.Implementation:

The implementation of the Personalized Educational Game Recommendation System was structured into three overarching phases: data preparation, system design, and interface deployment. Each phase involved meticulous attention to detail to ensure precision, scalability, and usability of the final recommendation engine.

7.1. Tools and Libraries Used: Comprehensive Analysis: Tool/Library Purpose in the System Why This Tool Was Chosen Python Core programming language. Python is versatile, widely used in machine learning, and has rich libraries for data analysis, text processing, and application development. It enables seamless integration of multiple functionalities, from preprocessing datasets to deploying web applications. pandas Data manipulation: Loading CSVs, cleaning missing values, and grouping data. pandas simplifies operations like filtering, merging, grouping, and pivoting, essential for preparing the dataset for clustering and recommendation. It handles missing and inconsistent data efficiently, which is crucial when processing large game datasets. numpy High-performance matrix operations. numpy accelerates numerical computations, which are critical for building similarity matrices, performing SVD, and normalizing vectors in the TF-IDF process. Its lightweight nature ensures scalability for large datasets.

scikit-learn Machine learning algorithms like TF-IDF, DBSCAN, and SVD. scikit-learn offers production-grade implementations of clustering and dimensionality reduction algorithms. These algorithms are central to grouping games, extracting meaningful features from text, and reducing computation overhead through SVD. Streamlit Building the interactive user interface. Streamlit enables rapid prototyping of web applications. It abstracts the complexity of front-end development while offering dynamic widgets for user input, such as multiselect dropdowns and sliders. It allowed us to quickly create a polished, interactive interface for users. re(Regular Expressions) Text cleaning and preprocessing: Removing special characters and normalizing text fields. Regular expressions are indispensable for string operations. They ensure uniformity in game names and descriptions, eliminating

inconsistenciesthatcouldaffectclusteringor similaritycalculations. matplotliband seaborn Visualizationofdataset trendsandclustering results. Datavisualizationtoolslikematplotliband seabornwereessentialforexploringpatterns, suchasgenredistribution,ratinghistograms, andsimilarityheatmaps.Theyalsoprovided insightsintoclusteringandage-group distributions,improvingtheinterpretabilityof thedatapreprocessingphase. math Mathematical operationsforTF-IDF andinversedocument frequencycalculations. Themathematicalunderpinningsof recommendationsystems,suchascalculating IDFscoresorcosinesimilarity,required preciselogarithmicandsquareroot calculations,whichthemathlibraryfacilitates withefficiencyandreliability. TF-IDFVectorizeTransformingtextual descriptionsinto numerical representations. TheTF-IDFVectorizerfromscikit-learn convertstextintovectorsthatquantifyword importance.Thiswascriticalforcomparing gamedescriptionssemanticallyand clusteringsimilargames. DBSCAN (Density-Based Spatial Clustering)Clusteringgamenames tounifyvariationslike "DeluxeEdition." DBSCAN'sabilitytoclusterbasedondensity ratherthanfixedclusterswasidealfor groupinggames.Ithandlesnoiseeffectively, allowingittodistinguishbetweenactualgame clustersandunrelatedentries. SVD(Singular Value Decomposition) Dimensionality reductionofTF-IDF vectorsfor computational efficiency. SVDreducesthedimensionalityoftext features,retainingthemostrelevant componentswhilediscardingnoise.Thisstep wasvitalforcreatingacosinesimilarity matrixwithoutoverloadingmemoryor computationalresources.

7.2.ImplementationProcess:Step-by-StepAnalysis: Theimplementationprocessinvolvedmultiplestages,eachbuildingupon theprevioustocreatearobustandreliablerecommendationsystem. Belowisadetailedanalysisofeachstep,providingacomprehensive understandingoftheworkflow. 7.2.1.Step1:DataPreparationandCleaning: Preparingthedatasetwasafoundationaltask.Therawdatacontained

several inconsistencies, redundant columns, and missing values. Addressing these issues was critical to ensure the quality of downstream algorithms. 1. Column Selection:

☐ Certain columns, such as is_free, header_image, and website, were deemed irrelevant to the recommendation logic and dropped.

☐ Key columns, including name, short_description, genres, and ratings, were retained as they held essential information for clustering, similarity computation, and user-specific recommendations. 2. Handling Missing Data:

☐ For critical columns like ratings, missing values were replaced with default placeholders or calculated averages.

☐ For less critical columns (e.g., optional genre fields), missing values were ignored during clustering and feature engineering.

3. Text Normalization: Text fields, particularly game names and descriptions, were standardized by: ☐ Lowercasing all text to ensure case-insensitive matching.

☐ Removing special characters, numbers, and extra spaces using regular expressions.

☐ Tokenizing and stemming words in descriptions to reduce redundancy and focus on the root meanings of terms.

Example: The game name "Minecraft: Deluxe Edition" was transformed into "minecraftdeluxeedition." This ensured consistency in how game names were represented and eliminated variations that could mislead clustering algorithms.

7.2.2. Step 2: Clustering Game Names:

Duplicate and versioned game names posed a significant challenge. For example, "FIFA 20," "FIFA 2020 Deluxe," and "FIFA 20 Ultimate Edition" represented the same base game but were listed separately. To unify such entries: 1. TF-IDF Vectorization: ☐ Game names were converted into numerical vectors, capturing their unique textual features. ☐ TF-IDF weighted terms based on their frequency within a game name and their rarity across all names, enabling effective differentiation between unrelated names.

2. DBSCAN Clustering:

Usingcosinesimilarityasthedistancemetric,DBSCANgrouped similarnamesintoclusters. Noisehandlingwasakeystrengthof DBSCAN,ensuringunrelated entrieswereexcludedfromclusters. Thedensity- basednatureofDBSCANeliminatedtheneedtopredefine thenumberofclusters,allowingdynamicadjustmentsbasedonthe data.

3.RepresentativeNameSelection:Withineachcluster,theshortestname waschosenastherepresentative.Forinstance,"FIFA20"wasselectedto representallvariationsofthegame. WhyDBSCAN? DBSCANeffectivelyhandlesdatawithnoiseandoutliers,acommon occurrenceinlargedatasets.

Itgroupsentriesbasedondensity,makingitidealforclusteringtextual datawithvariablesimilaritythresholds. 7.2.3.Step3:FeatureEngineering: Featureengineeringinvolvedcreatingnumericalrepresentationsofthe datasettoenablemachinelearningalgorithmstounderstandthedata.

1.GenreEncoding:Gamegenreswereone-hotencodedintoabinarymatrix. Forexample: Agamewithgenres"Action"and"Adventure"wasrepresentedas[1,1, 0,0,...]. Agamewithonly"Action"was[1,0,0,0,...].

Thisrepresentationallowedthesystemtomeasuregenreoverlapbetween games.

2.TextVectorization: DescriptionswereconvertedintonumericalvectorsusingTF-IDF.This approachhighlightedtheimportanceofuniquetermsinagame's descriptionrelativetoothers. Stopwordssuchas"game,""play,"and"new"wereremovedtofocus onmeaningfulterms.

7.2.4.Step4:DimensionalityReductionUsingSVD: TheTF- IDFmatrixoftencontainedthousandsofdimensionsduetothe extensivevocabulary.Tomakecomputationsefficient: 1.SingularValueDecomposition(SVD): Thematrixwasreducedto30components,capturingthemost significantfeatureswhilediscardingnoise. Thisstepdramaticallyreducedthecomputationalcomplexityof similaritycalculations.

2.MathematicalInsight: TheTF-IDFvalueforeachtermwascalculatedas:$IDF \times TF = IDF\text{-}TF$

Where:

ent s in docum Total term document of term in Frequency  =TF with term

Documents ments Total docu log=DF □□□ □□□

SVDdecomposedthematrixintothreecomponents:U,□,andTV,

reducingthedimensionalitywhileretainingkeypatterns.

7.2.5.Step5:BuildingtheRecommendationEngine:

Therecommendationenginecombinedmultiplefeaturestocomputesimilarity andrankgames.

1.SimilarityMatrix:UsingthereducedmatrixfromSVD,acosinesimilarity

matrixwascalculated.Thismatrixquantifiedhowcloselyeachgame resembledothers.

2.WeightedScoring:Aweightedformulacombinedthreekeyfactors:

□20%:Similaritytogamestheuserhadplayed. □50%:Genreoverlapwiththeuser'spreferences.

□30%:Averageratingofthegame.

Thisapproachensuredtherecommendationsbalancedrelevance,user preferences,andquality.

7.2.6.Step6:UserInterface: Thefinalphasewascreatinganintuitiveinterfaceforusers.

1.StreamlitIntegration:Theinterfacealloweduserstoinput: □Theiragegroup.

□Alistofupto5gamestheyhadplayed. □Upto3preferredgenres.

Usersreceiveddetailedrecommendationswithinsightsintowhyeachgame wassuggested.

2.DetailedExplanations:Foreachrecommendedgame,theinterface displayed:

□Similarityscorewiththeuser'sinput. □Genreoverlappercentage.

□Gameratingandrequiredage.

□Anyadditionalmetadata,suchassupportedlanguagesandcategories.

UserInteraction:Theinterfaceoffereddynamicfeedback,suchaswarnings

iftheuserselectedfewergamesorgenresthanrequired.Recommendations

werepresentedwithvisualaids,includinggamecoverimageswhere

available.7.3.CodeSnippetsandDetailedExplanations: 7.3.1.TF-IDFVectorization:

Thefirstcriticalstepincreatingtherecommendationsysteminvolved

convertinggamedescriptionsintonumericalvectorsusingTermFrequencyInverseDocumentFr

equency(TF-IDF).TF-IDFhighlightstheimportanceof

wordsinadocumentrelativetotheentiredataset,ensuringthatcommon

words,whicharelessinformative,aredown-weighted.Forthistask,weused

theTfidfVectorizerfromthesklearn.feature_extraction.textlibrary.

Theprocessbeginsbyfittingthevectorizeronthe"merged_description"

columnofthedataset.Thiscolumncontainsconsolidateddescriptionsofeach

game.Thevectorizertokenizesthetext,removespredefinedstopwords(e.g.,

"game"or"play"),andassignsweightstoeachtermbasedonitsfrequencyin

aspecificgame'sdescriptionrelativetotheoveralldataset.Forexample,a

gamedescribedas"action-packedmultiplayershooter"mightgeneratea

vectorwheretermslike"action"and"multiplayer"areweightedhigherifthey

occurlessfrequentlyacrossthedatasetbutarepivotaltothisparticulargame.

Thistransformationisvitalasitcapturesthesemanticessenceofeach

game'sdescriptioninanumericalformat,enablingmathematical

computationsandcomparisons. desc_corpus=df_grouped["merged_description"].values

raw_vocab=set() fordocindesc_corpus: words=doc.split() forwinwords: raw_vocab.add(w)

raw_vocab_list=sorted(list(raw_vocab)) custom_sw={

"game","games","play","played","playing","world","players","new",

"open","will","make","makes","made","thing","things","like"

}filtered_vocab_list=[vforvinraw_vocab_listifvnotincustom_sw]

defterm_frequency(doc,vocab): counts=[0]*len(vocab) words=doc.split() forwinwords:

ifwinvocab: idx=vocab.index(w) counts[idx]+=1 returncounts tf_matrix=[]

fordocindesc_corpus: tf_matrix.append(term_frequency(doc,filtered_vocab_list))

tf_matrix=np.array(tf_matrix,dtype=float) doc_freq=np.count_nonzero(tf_matrix,axis=0)

N_docs=len(desc_corpus) final_vocab=[] valid_idx=[]

fori,winenumerate(filtered_vocab_list): ratio=doc_freq[i]/N_docs ifratio<=0.6:

```python
final_vocab.append(w) valid_idx.append(i) tf_matrix_filtered=tf_matrix[:,valid_idx]
df_filtered=doc_freq[valid_idx] idf_vals=[] fordf_valindf_filtered:
idf_vals.append(math.log((N_docs+1)/(df_val+1))+1) idf_vals=np.array(idf_vals)
tf_idf_matrix=tf_matrix_filtered*idf_vals defrowwise_norm(mat): out=[] forrowinmat:
length=np.sqrt(np.sum(row**2)) iflength!=0: out.append(row/length) else:out.append(row)
returnnp.array(out) tf_idf_matrix=rowwise_norm(tf_idf_matrix)
```

Thisstepensuresthatthetextualdataisnumericallyrepresented,allowingfor deeperanalysisinlaterstages. 7.3.2.SVDforDimensionalityReduction: OncetheTF-IDFmatrixisgenerated,ittypicallyspansthousandsof dimensionsduetothelargevocabularyderivedfromgamedescriptions.To makecomputationsmoreefficient,SingularValueDecomposition(SVD)is applied.SVDisamathematicaltechniquethatdecomposesthehighdimensionalmatrixintothreec omponents:U,□,andTV.These componentscapturetheessentialpatternsinthedata,discardingnoiseand redundancies. Thereducedcomponentsarethenusedtoreconstructalower-dimensional matrix.Byretainingonlythetop30components,weensurethatthemost significantfeaturesarepreservedwhilesignificantlyreducingcomputational overhead.Forinstance,gameswithsimilarthemesordescriptionswillcluster togetherinthisreducedspace,eveniftheyuseslightlydifferentwording.

```python
defsvd_decomposition(matrix,k=None): A=np.array(matrix,dtype=float) ATA=np.dot(A.T,A)
AAT=np.dot(A,A.T)
```

```python
eigvals_ATA,V=np.linalg.eigh(ATA) eigvals_AAT,U=np.linalg.eigh(AAT)
idx_v=np.argsort(eigvals_ATA)[::-1] idx_u=np.argsort(eigvals_AAT)[::-1] V=V[:,idx_v]
U=U[:,idx_u] eigvals=eigvals_ATA[idx_v] ifkisnotNone: U=U[:,:k] V=V[:,:k]
eigvals=eigvals[:k] sigma=np.sqrt(eigvals) sigma_matrix=np.diag(sigma)
returnU,sigma_matrix,V.T K=30 U,S,Vt=svd_decomposition(combined_features,k=K)
reconstructed=np.dot(np.dot(U,S),Vt)
```

Thisdimensionalityreductionstepnotonlyimprovestheefficiencyofthe systembutalsoenhancestheaccuracyofsimilaritycomputationsbyfocusing onmeaningfulfeatures. 7.3.3.CosineSimilarityforGameMatching: Torecommendgamessimilartotheonesauserhasplayed,wecomputethe cosinesimilaritybetweengames.Cosinesimilaritymeasurestheangle betweentwovectorsinahigh-dimensionalspace.Asmallerangle(closerto zero)indicateshighersimilarity.Forinstance,twoactiongameswithsimilar descriptionsandoverlappinggenreswillhavevectorspointinginnearlythe samedirection,resultinginahighsimilarityscore.

Thenormalizedfeaturevectors(fromthereducedmatrix)areusedto computeacosinesimilaritymatrix.Thismatrixrepresentspairwisesimilarity scoresforallgames,formingthebackboneoftherecommendationengine.

```
defcosine_similarity_matrix(features):
norms=np.sqrt(np.sum(features**2,axis=1,keepdims=True))
normed=features/(norms+1e-8) sim=np.dot(normed,normed.T) returnsim
item_sim_matrix=cosine_similarity_matrix(reconstructed)
```

Byleveragingcosinesimilarity,wequantifyhowcloselyrelatedanytwo gamesarebasedontheircombinedfeatures,ensuringthatrecommendations alignwithuserpreferences.

7.3.4.RecommendationScoringMechanism: Thefinalstepintherecommendationprocessinvolvesaggregatingmultiple metrics—similaritytoplayedgames,genreoverlap,andratings—intoasingle weightedscore.Eachmetricisassignedaweightbasedonitsimportance. Similaritycontributes20%,genreoverlapaccountsfor50%,andratingsadd 30%tothefinalscore.Thisweightingstrategyensuresthatthe recommendationsarenotonlysimilartotheuser'spastchoicesbutalsoalign withtheirgenrepreferencesandmaintainahighqualitystandard.

For instance, if a user has played multiple action games, the system will prioritize action titles with high ratings and a strong overlaping genres. The recommendation function computes these scores dynamically, ranks the games, and returns the top recommendations.

```python
def recommend_games(
    user_age,  # "below_15" or "above_15"
    user_games,  # list of original names
    user_genres,  # list of favored genres
    top_n=5
):
    if user_age == "below_15":
        valid_df = df_grouped[df_grouped["required_age"]<=15].copy()
    else:
        valid_df = df_grouped.copy()
    valid_idx = valid_df.index.tolist()

    def clean(gm):
        return re.sub(r"[^a-z0-9\s]","",gm.lower()).strip()

    user_cleaned = [clean(gm) for gm in user_games]
    rep_indices = []
    for uc in user_cleaned:
        if uc in name_to_rep:
            cluster_rep = name_to_rep[uc]
            row_match = valid_df[valid_df["unified_name"]==cluster_rep]
            if len(row_match)>0:
                ridx = row_match.index[0]
                rep_indices.append(ridx)
        else:
            row_match = valid_df[valid_df["unified_name"]==uc]
            if len(row_match)>0:
                ridx = row_match.index[0]
                rep_indices.append(ridx)
    if len(rep_indices)==0:
        valid_df["similarity_score"] = valid_df["ratings"].astype(float)
        valid_df.sort_values("similarity_score",ascending=False, inplace=True)
        return valid_df.head(top_n)
    avg_scores = []
    for idx in valid_idx:
        local_sims = []
        for r_idx in rep_indices:
            local_sims.append(item_sim_matrix[idx,r_idx])
        avg_scores.append(np.mean(local_sims))
```

```python
    valid_df["similarity_score"] = avg_scores

    def measure_overlap(row_gs,user_gs):
        row_set = set(row_gs)
        user_set = set(user_gs)
        overlap_count = len(row_set.intersection(user_set))
        if len(user_gs)==0: return 0
        return overlap_count/len(user_gs)
    overlap_scores = []
    for i,row in valid_df.iterrows():
        overlap_scores.append(measure_overlap(row["merged_genres"], user_genres))
    valid_df["genre_overlap"] = overlap_scores
    w_sim=0.2
    w_genre=0.5
    w_rat=0.3
    final_score = []
    for i,row in valid_df.iterrows():
        score = (
            w_sim*row["similarity_score"]+
```

```
w_genre*row["genre_overlap"]+ w_rat*row["ratings"] )final_score.append(score)
valid_df["combined_score"]=final_score
valid_df.sort_values("combined_score",ascending=False,inplace=True)
top_recs=valid_df.head(top_n).copy() sim_details=[] fori,rowintop_recs.iterrows(): txt=""
forr_idxinrep_indices: rep_name=df_grouped.loc[r_idx,"unified_name"]
val_sim=item_sim_matrix[i,r_idx] txt+=f"Similarto{rep_name}=>{val_sim:.3f}\n"
sim_details.append(txt.strip()) top_recs["similarity_details"]=sim_details returntop_recs
```

Thisscoringmechanismensuresabalancedrecommendationsystem, offeringusersdiverseyethighlyrelevantgamesuggestions.

7.4.FlowchartsforImplementationProcess:

Thissubsectionexplainsthedetailedflowcharts(Figures7to12)usedto representthesystem'soperation.Eachflowchartillustratescritical componentsoftheimplementationprocess,breakingdowntheworkflow intomanageableandcomprehensiblesteps. 7.4.1.DatasetPreprocessingFlowchart: Thefirstflowchartrepresentsthedatasetpreprocessingphase,essentialfor ensuringtherawdataiscleanandreadyforfurtherprocessing.Theprocess beginsbyloadingthedataset,whichincludescolumnslikegamenames,

genres,ratings,anddescriptions.Irrelevantcolumnssuchasis_freeand header_imagearedroppedtofocusonlyonusefulinformation.Missing valuesincriticalfieldslikeratingsarefilledusingaveragesorplaceholders, whilenon-criticalfieldswithmissingvaluesareignored.Textfields,including gamenamesanddescriptions,arenormalizedbyconvertingthemto lowercase,removingspecialcharacters,andtokenizingthemfor standardization.Genresareconvertedintoabinarymatrixusingone-hot encoding,ensuringthattheyaremachine-readable.Theoutputisacleanand structureddatasetthatisreadyforclusteringandfeatureengineering.

Figure7:DatasetPreprocessing 7.4.2.ClusteringGameNamesFlowchart:

This flowchart outlines the clustering process to unify similar game names and reduce redundancy. The process starts by applying TF-IDF vectorization to game names, converting them into numerical vectors based on their textual features. DBSCAN clustering is then used to group similar names based on their density, using cosine similarity as the metric. For each cluster, a representative name—typically the shortest name—is selected to standardize game versions. These unified names replace original game names in the dataset, creating a consistent dataset where variations like "FIFA 20" and "FIFA 2020 Deluxe" are grouped under a single name. The final output is a dataset with clustered game names, enabling better feature engineering and recommendation accuracy.

Figure 8: Clustering Game Names 7.4.3. Feature Engineering Flowchart: The feature engineering flowchart details the creation of meaningful features from the dataset. The process begins by merging all genres associated with each game. These genres are then encoded into a binary matrix using one hot encoding. Simultaneously, game descriptions are cleaned and normalized to remove unnecessary noise. TF-IDF vectorization is applied to these descriptions, generating numerical vectors that highlight the importance of specific terms in describing each game. The genre matrix and TF-IDF features are then combined into a unified feature set, capturing both categorical and textual information. The resulting feature matrix is used for dimensionality reduction and similarity computations, ensuring efficient and accurate recommendations. Figure 9: Feature Engineering

7.4.4. Recommendation Function Flowchart: The recommendation function flowchart visualizes the logic used to generate personalized recommendations. The process starts with user inputs, including

agegroup,favoritegenres,andpreviouslyplayedgames.Thefirststepfilters theavailablegamesbasedontheuser'sagegrouptoexcludeage-restricted content.Next,thesystemmatchestheuser'splayedgamestotheir respectiveclusterstoidentifyrelevantdatapoints.Usingcosinesimilarity, similarityscoresarecalculatedbetweentheuser'splayedgamesandthe remaininggames.Genreoverlapisalsocomputedtoassesshowwellagame alignswiththeuser'spreferences.Thesemetricsarecombinedintoa weightedscoringformula,assigning20%weighttosimilarity,50%togenre overlap,and30%toratings.Finally,thetop-rankedgamesareoutputtedas recommendations,accompaniedbydetailedreasoning. Figure10:RecommendationFunction

7.4.5.StreamlitUserInterfaceFlowchart:

Thisflowchartexplainstheinteractionbetweentheuserandthe recommendationsystemthroughaStreamlit-basedinterface.Theuserbegins bylaunchingtheapplicationandselectingtheiragegroup.Theythenchoose uptofivegamestheyhaveplayedanduptothreefavoritegenres.Oncethe userclicksthe"GetRecommendations"button,thesystemvalidatesthe inputstoensureallrequirementsaremet.Thevalidatedinputsarepassedto therecommendationfunction,whichgeneratespersonalizedgame recommendations.Theserecommendationsaredisplayedtotheuser,

completewithexplanationsforeachsuggestion,makingtheinterface interactiveanduserfriendly. Figure11:StreamlitUserInterface

7.4.6.RecommendationEngineLogicFlowchart:

Therecommendationenginelogicflowchartprovidesadeeperviewofthe system'sbackendoperations.Theenginebeginsbyreceivinguserinputs suchasagegroup,playedgames,andpreferredgenres.Itfiltersthegames basedonagerestrictionstomatchtheuser'sprofile.Next,theuser'splayed gamesaremappedtotheirrespectiveclustersforeffectivecomparison.

Clusterrepresentativesarevalidatedtoensureaccuracy.Cosinesimilarity scoresarecalculatedbetweentheuser'splayedgamesandothergames, whilegenreoverlapismeasuredtoevaluatealignmentwiththeuser's preferences.Thesemetricsarecombinedintoafinalweightedscore,which ranksgamesbasedontheirrelevancetotheuser.Thesystemsorts recommendationsbytheirfinalscoresandoutputsthetopgamesas recommendations,ensuringpersonalizedandaccuratesuggestions.

Figure12:RecommendationEngineLogic 8.TestingandResults: 8.1.Testing: Testingisacriticalphasetoensurethereliability,accuracy,andusability oftheGameRecommenderSystem.Thegoalwastoevaluatethesystem undervariousscenarios,coveringbothfunctionalandnon-functional aspects.Thesystemunderwentextensivetestingacrossmultiple dimensions,includinginputvalidation,algorithmaccuracy,usability,and performance.

8.1.1.TestingObjectives: Theprimaryobjectivewastovalidatethattherecommendersystem consistentlygeneratesaccurate,relevant,anduser-specificrecommendations. Thesystemwasexpectedto:

☐Handleedgecases,suchasinvalidorincompleteuserinputs.

☐Filtergamesbasedonagerestrictionsaccurately.

☐Providerecommendationsthatreflecttheuser'sgamepreferencesand genreselections.

☐Deliverclearandmeaningfulexplanationsfortherecommended games.

☐Performefficientlyundervaryingdatasetsizes. 8.1.2.TestMethodology:

Thetestingwasdividedintothreemainstages:

1.UnitTesting:Focusedonensuringindividualcomponentslikeclustering,

collaborativefiltering,andcosinesimilaritycalculationswereworkingas intended.

2.IntegrationTesting:Verifiedthatthebackendalgorithmsintegrated

smoothlywiththefrontendStreamlitapplicationtoprovidereal-time results.

3. User Testing: Conducted with 20 participants who provided insights into usability, recommendation relevance, and the quality of explanations.

8.1.3. Test Scenarios and Results:

The system was evaluated using various test cases, as outlined below: Test Scenario Expected Outcome Result UserInputValidation Prevent incomplete inputs (e.g., less than 5 games selected) and display warnings. Passed AgeRestriction Filtering Display games only suitable for the selected age group (below_15 or above_15). Passed Recommendation Accuracy Recommend games similar to the user's selected games and favorite genres. Passed ExplanationClarity Provide detailed similarity and genre overlap explanations for each recommendation. Passed Performance under Load Generate recommendations within 3 seconds for datasets of up to 10,000 games. Passed (1.8s avg.) InterfaceUsability Ensure dropdown menus, error messages, and buttons are intuitive to use. Passed

The results show the system performed reliably across all test scenarios, meeting functional and usability expectations. 8.2. Results Representation:

The system's results were evaluated through example use cases to demonstrate its functionality and explainability. Below is an example scenario: InputDetails
☐ AgeGroup: Above 15 ☐ GamesPlayed: Brain/Out, AgeofHistoryII, Fez, Inside, Teardown
☐ FavoriteGenres: Action, Adventure, Simulation RecommendationsOutput

Based on the input, the system provided the following ranked list of recommended games: 

| GameName | Required Age | Rating | Similarity Score | Genre Overlap | Combined Score |
|---|---|---|---|---|---|
| Assassin's CreedIV: BlackFlag | 18 | 10.0 | 0.284 | 0.667 | 3.390 |
| Tricolour Lovestory | 18 | 10.0 | 0.352 | 0.333 | 3.237 |
| Bioshock Infinite | 18 | 10.0 | 0.305 | 0.333 | 3.228 |
| Half-Life | 18 | 10.0 | 0.298 | 0.333 | 3.226 |
| EldenRing | 18 | 10.0 | 0.210 | 0.333 | 3.209 |

The recommendations aligned well with the user's input. Games like "Assassin'sCreedIV:BlackFlag" and "Half-Life" scored high due to their

overlapwithgenreslikeActionandAdventure.Theirsimilarityscores,

derivedfromthedescription-basedTF-IDFanalysis,wereconsistently

above0.2,whichindicatesastrongalignmentwiththeuser'sselected games.

☐SimilarityScore:Reflectshowcloselytherecommendedgame

descriptionsmatchthegamesselectedbytheuser.Forexample,

"BioshockInfinite"scored0.305duetoitsnarrative-drivenstylesimilar to"Inside."

☐GenreOverlap:Measuresthepercentageofgenresthatmatchthe

user'sfavorites."Assassin'sCreedIV:BlackFlag"achievedthe

highestoverlap(0.667)becauseitalignswithallthreegenresselected.

☐CombinedScore:Aweightedsumofthesimilarityscore(20%),genre

overlap(50%),andgamerating(30%).Thisscoredeterminesthefinal ranking.

8.3.VisualRepresentationofResults: 1.InputInterface:Theuser-

friendlyinterfaceallowsuserstoselectage group,playedgames,andgenresseamlessly:

Figure13:InputInterface

2.RecommendationOutput:Recommendationsaredisplayedwith

detailslikerequiredage,ratings,similarityscores,andexplanations:

Figure14:RecommendationOutput 8.4.InsightsandFutureEnhancements:

8.4.1.Insights: TheGameRecommenderSystemdemonstratedastrongabilitytodeliver

personalized,relevant,andexplainablerecommendations.Bybalancingthe

threekeyfactors—similarity,genrealignment,andratings—thesystem

ensuresusersreceivetailoredsuggestionsthatcloselymatchtheir

preferencesandgameplayhistory.

Asignificantstrengthofthesystemliesinitsexplanationmechanism.Each

recommendationincludesadetailedbreakdownofthesimilarityscores,genre

overlaps,andotherfactorsinfluencingthecombinedscore.Thistransparency

notonlyenhancesusertrustbutalsoempowersuserstounderstandhow

their choices impact the results. For instance, a user can see how a game like "Inside" contributed to the recommendation of "Assassin's Creed IV: Black Flag" due to shared gameplay themes and genres.

Additionally, the use of advanced techniques such as TF-IDF for text vectorization and SVD for dimensionality reduction has optimized the system's performance, particularly in handling large datasets. The ability to efficiently process and recommend games from a dataset exceeding 15,000 records highlights the system's scalability. These insights affirm the system's robustness and its ability to deliver high-quality recommendations in real world scenarios.

8.4.2. Planned Enhancements:

While the system performs well in its current state, several enhancements are planned to improve its functionality, accuracy, and user experience:

1. Dynamic Weighting System: One of the key improvements will be introducing a dynamic weighting system. This feature will allow users to prioritize factors such as ratings, similarity, or genre alignment when generating recommendations. For example, a user who values high-rated games over genre alignment could adjust the weighting to favor ratings, leading to recommendations that align with their specific priorities. This customization will further personalize the experience and enhance user satisfaction. 2. Real-Time Feedback Collection: Incorporating a feedback mechanism will allow users to rate the relevance and usefulness of the recommended games. This feedback will be collected in real time and integrated into the recommendation algorithm. By using collaborative filtering techniques, the

system can improve its accuracy and adapt to user preferences over time. For instance, if users consistently prefer recommendations with higher genre overlap, the system can learn to weigh genre alignment more heavily.

3. Expanded Visualization Tools: Visualization tools will be added to help

usersexploretherecommendationprocessandgaindeeperinsightsintotheir preferences.Forinstance,genredistributionheatmapscouldvisualizehowa user'sselectedgenrescomparetothoseofrecommendedgames.Similarly, similaritygraphscouldillustratetherelationshipsbetweenuser-selected gamesandrecommendedtitles.Thesetoolswillmaketherecommendation processmoreinteractiveandengaging.

4.DiverseDataIntegration:Currently,thesystemreliesontextual descriptions,ratings,andgenres.Futureenhancementswillintegratemore diversedata,suchasuserreviews,gameplaytime,andin-game achievements.Theseadditionalfactorscanhelprefinetherecommendations andmakethemevenmorepersonalized. 5.SupportforMulti-UserRecommendations:Thesystemcanbe enhancedtosupportrecommendationsforgroupsofusers,suchasfriendsor familymemberswhoplaygamestogether.Byanalyzingthepreferencesand gamehistoriesofmultipleusers,thesystemcouldgeneratesuggestionsthat appealtoeveryoneinthegroup.Thisfeaturewillbeparticularlyvaluablefor multiplayergamesandsharedgamingexperiences. 6.Cross-PlatformRecommendations:Toincreaseitsapplicability,the systemcouldrecommendgamesbasedonplatformpreferences,suchasPC, console,ormobile.Userswhopredominantlyplayonaspecificplatform wouldreceiverecommendationstailoredtothatenvironment. 9.Discussion: 9.1.ComparisonandEvaluation:

TheGameRecommenderSystemperformedasexpectedindelivering personalizedrecommendationsbasedonuserinputssuchasage,played games,andpreferredgenres.Theanalysisoftheresultsshowsthatthe systemeffectivelyidentifiesgamesthatalignwithuserpreferences.For instance,therecommendationsprovidedintestingcloselymatchedthe genrealignmentandgameplaypatternsofselectedgames,withsimilarity

scoresandcombinedscoresreflectinglogicalandconsistentrelationships.

Thesystemexceededexpectationsintermsofscalabilityand

explainability.Despiteprocessingadatasetofover15,000entries,it

consistentlydeliveredrecommendationswithminimallatency.This

efficiencycanbeattributedtotheoptimizationtechniquesappliedduring

thepreprocessingphase,suchasdimensionalityreductionviaSVDand

clusteringtounifysimilargameversions.

However,thesystem'sperformanceshowedslightdeviationsinedge

cases.Forexample,whenusersselectednichegameswithlowdata

representation,therecommendationssometimesskewedtowardsmore

populargameswithinthesamegenre.Thisbehavior,whilelogical,

suggestsaneedforimprovementinhandlingsparsedatascenariosto

ensurediversityinrecommendations.

9.2.CriticalAnalysis: 9.2.1.ChallengesFacedDuringImplementation:

1.DataPreprocessing:Cleaningandunifyingthedatasetwasoneofthe

mostsignificantchallenges.Therawdatasetcontainedinconsistencies

suchasduplicategameversions,incompletedescriptions,and

overlappinggenrelabels.Implementingaclusteringapproachusing

DBSCANtounifygamenamesrequiredextensivefine-tuningtobalance

betweenmergingsimilarentriesandpreservinguniquedatapoints.

2.DimensionalityReduction:IntegratingTF-IDFandSVDtoreduce

dimensionalityandcomputeitemsimilarityrequiredcarefulparameter tuning.Over-

reductionindimensionsriskedlosingimportantsemantic details,whileunder-

reductionledtoinefficienciesincomputation.

3.HandlingSparseInputs:Userswithsparseorhighlyspecific

preferencesposedachallengetotherecommendationprocess.For

example,whenauserselectedfivehighlynichegames,thesystem

sometimesstruggledtoprovidediverserecommendationsduetolimited overlapingenresordescriptions. 4.Age-BasedFiltering:Designinganage-basedfilteringmechanism involvedbalancinginclusivitywithappropriateness.Ensuringthat recommendationsforusersbelow15adheredtotheagerestriction withoutoverlynarrowingtheresultsrequiredadditionallogicduring recommendationgeneration.

5.Explainability:Creatingclear,interpretableexplanationsfor recommendationswasbothcriticalandchallenging.Thesystemhadto presentabreakdownoffactorslikesimilarityandgenrealignmentwithout overwhelminguserswithexcessivedetail. 9.2.2.SystemPerformanceEvaluation: Overall,thesystemmetitsobjectivesofdeliveringaccurateandexplainable recommendations.However,theevaluationhighlightedsomeareasfor improvement:

☐PrecisioninSparseData:Recommendationsfornichepreferences lackedvariety.

☐DynamicResponsiveness:Theweightingofsimilarity,genreoverlap, andratingswasstatic,limitinguser-specificflexibility. 9.3.Enhancements: Thechallengesidentifiedduringimplementationformthebasisfor actionablefutureenhancements.Theseinclude:

1.ImprovedHandlingofSparseData:Incorporateahybrid recommendationapproachthatcombinescontent-basedfilteringwith collaborativefiltering.Thiswouldleverageuserinteractiondatato enhancerecommendationswhencontent-basedsimilaritiesare insufficient.

2.DynamicWeighting:Allowuserstocustomizetheweightingoffactors likesimilarity,genrealignment,andratings.Implementingasimplesliderbasedinterfaceintheappl icationwouldempoweruserstoprioritizetheir preferencesdynamically.

3.EnhancedDiversityinRecommendations:Introducediversityenhancingalgorithmsthatensur eabroaderrangeofrecommendations, evenforuserswithnichepreferences.Forexample,thesystemcouldcap

thenumberofrecommendationsfromthesamegenreorintroduce randomnesstolow-rankingrecommendations. 4. Refined Age-Based Filtering: Incorporateamorenuancedage-filtering systemthatusesacombinationofageratingsandcontenttagstoprovide appropriateyetdiverserecommendationsforyoungerusers.

5. Feedback Mechanism: Introduceafeatureforuserstorate recommendations. Theseratingscouldbeincorporatedintocollaborative filteringmodels,allowingthesystemtolearnandimproveiteratively.

6. Integration of User Behavior Data: Expandthedatasettoincludedata pointssuchastimespentongames,completionrates,andreview sentiments. Thesebehavioralmetricswouldenhancethe recommendationprocessbycapturingdeeperinsightsintouser preferences.

7. Interactive Visualizations: Addvisualelements,suchassimilarity graphsandheatmapsofgenreoverlaps,tohelpusersbetterunderstand therecommendationprocess. Thesetoolswouldnotonlyincrease transparencybutalsoimproveuserengagement.

8. Multilingual Support: Whilethesystemalreadyidentifiessupported languagesforrecommendedgames,expandingtheinterfacetoprovide recommendationsinauser'spreferredlanguagewouldenhance accessibility. 10. Conclusion: TheGameRecommenderSystemsuccessfullyachieveditsgoalofproviding personalizedandexplainablegamerecommendationsbasedonuserinputs. ThroughtheintegrationofadvancedtechniquessuchasTF-IDFvectorization, dimensionalityreductionviaSVD,andclusteringfordataunification,the systemdeliveredhighlyrelevantsuggestionstailoredtouserpreferences. The inclusionoffactorssuchasagefiltering,similarityscoring,genrealignment, andrating-basedweightingensuredtherecommendationswereappropriate, logical,anddiverse.

Oneofthesystem'smostsignificantcontributionsisitsuser-centricdesign, emphasizingtransparencyandtrustthroughdetailedexplanationsof recommendations. Thisfeaturenotonlyenhancesuserconfidencebutalso

differentiatesthesystemfromgenericblack-boxapproachesin

recommendationengines.Additionally,thesystemdemonstratedscalability

byprocessingalargedatasetefficiently,provingitssuitabilityfordeployment inreal-

worldscenarios. Despitethechallengesencountered,theprojecthighlightedthepotentialfor

leveragingcontent-basedtechniquestocreatemeaningfulandpersonalized

experiences.Byaddressinglimitationssuchassparsedatahandlingand

incorporatingdynamicweighting,thesystemcanevolvetomeetevenmore

diverseuserneeds.Thisprojectunderscoresthegrowingimportanceof

recommendationsystemsinthegamingindustry,wherepersonalizationplays

acriticalroleinusersatisfactionandretention.

Inconclusion,theprojectnotonlydemonstratedthefeasibilityofbuildingan

advancedrecommendationenginebutalsoshowcaseditsrelevancein

moderngamingapplications.Theinsightsgainedandenhancements

proposedprovideastrongfoundationforfutureimprovementsand applicationsofthesystem.

11.References: 1.PandasDocumentation:Pythondatamanipulationlibrary.

https://pandas.pydata.org/docs/ 2.NumPyDocumentation:NumericalcomputingwithPython.

https://numpy.org/doc/ 3.Scikit-learn:Machinelearninglibraryforclustering,TF-IDF,andSVD.

https://scikit-learn.org/stable/

4.TensorFlow:Machinelearningframeworkforneuralcollaborativefiltering.

https://www.tensorflow.org/ 5.Ester,M.,etal.(1996).Adensity-

basedalgorithmfordiscoveringclusters

inlargespatialdatabaseswithnoise.ProceedingsofKDD-96. https://www.aaai.org/

6.Manning,C.D.,Raghavan,P.,&Schütze,H.(2008).Introductionto

InformationRetrieval.CambridgeUniversityPress. DOI:10.1017/CBO9780511809071

7.Strang,G.(2009).IntroductiontoLinearAlgebra.Wellesley-Cambridge Press.

8.He,X.,etal.(2017).NeuralCollaborativeFiltering.ProceedingsofWWW

'17.DOI:10.1145/3038912.3052569 9.Lops,P.,etal.(2011).Content-

basedRecommenderSystems:Stateof

theArtandTrends.InRecommenderSystemsHandbook,Springer.

DOI:10.1007/978-0-387-85820-3_3

10.Burke,R.(2002).HybridRecommenderSystems:Surveyand

Experiments.UserModelingandUser-

AdaptedInteraction,12(4),331370.DOI:10.1023/A:1021240730564

11.Su,X.,&Khoshgoftaar,T.M.(2009).ASurveyofCollaborativeFiltering

Techniques.AdvancesinArtificialIntelligence,ArticleID421425. DOI:10.1155/2009/421425

12.Zhang,S.,etal.(2019).DeepLearningBasedRecommenderSystem:A

SurveyandNewPerspectives.ACMComputingSurveys(CSUR),52(1),

1-38.DOI:10.1145/3285029 13.Kingma,D.P.,&Welling,M.(2014).Auto-

EncodingVariationalBayes. arXiv:1312.6114

14.·Anderson,C.A.,&Dill,K.E.(2000).VideoGamesandAggressive

Thoughts,Feelings,andBehaviorintheLaboratoryandinLife.Journalof

PersonalityandSocialPsychology,78(4),772-790.

15.Yee,N.(2006).TheDemographics,Motivations,andDerived

ExperiencesofUsersofMassivelyMulti-userOnlineGraphical

Environments.PRESENCE:TeleoperatorsandVirtualEnvironments,

15(3),309-329.DOI:10.1162/pres.15.3.309

16.Chen,L.,&Pu,P.(2010).TrustBuildinginRecommenderAgents.

ProceedingsofAAMAS2010.https://aamas.csc.liv.ac.uk/

17.Herlocker,J.L.,etal.(2004).EvaluatingCollaborativeFiltering

RecommenderSystems.ACMTransactionsonInformationSystems

(TOIS),22(1),5-53.DOI:10.1145/963770.963772

18.Shaker,N.,Togelius,J.,&Nelson,M.J.(2016).ProceduralContent

GenerationinGames.Springer.DOI:10.1007/978-3-319-42716-4

19.Sifa,R.,etal.(2015).UserBehaviorinGameAnalytics:AReview.

EntertainmentComputing,13,10-22.DOI:10.1016/j.entcom.2015.06.002

20.Bostock,M.,Ogievetsky,V.,&Heer,J.(2011).D3Data-Driven
Documents.IEEETransactionsonVisualizationandComputerGraphics,
17(12),2301-2309.DOI:10.1109/TVCG.2011.185 ·

## Sources

1    https://quizlet.com/893436597/elementary-statistical-reasoning-module-2-flash-cards/
     INTERNET
     1%

| | |
|---|---|
| EXCLUDE CUSTOM MATCHES | ON |
| EXCLUDE QUOTES | OFF |
| EXCLUDE BIBLIOGRAPHY | OFF |