

Instituto Tecnológico de Costa Rica

Unidad de Computación

Proyecto 1

Carlos Akion Garro Campos

Sede San Carlos

14/06/2022

Tabla de contenidos

Contenido

Resumen Ejecutivo	3
Objetivo General	4
Objetivos Específicos	4
Introducción	5
Descripción del problema	6
Desarrollo	8
Análisis de resultados	12
Conclusiones	13
Recomendaciones	13
Bibliografía	14

Resumen Ejecutivo

Home Credit es una entidad financiera encargada de brindar créditos a clientes que no cuentan con un historial, para hacer esto debe de hacer uso de una gran cantidad de datos para determinar si este cuenta con capacidad de pago.

Por lo que para indicar cuales clientes cumplen es necesario la visualización de la información presente en la base de datos, por lo que en el presente proyecto se desarrolla un sistema gestor de bases de datos que permita visualizar los datos de los usuarios para ver si pueden optar por un crédito.

Objetivo General

- Crear un sistema gestor de bases de datos robusto para mostrar la información de Home Credit.

Objetivos Específicos

- Normalizar los datos de Home Credit para hacer el modelo relacional del sistema.
- Hacer el proceso de migración de Home Credit hacia la base de datos.
- Utilizar técnicas de optimización para la ejecución de consultas de manera rápida.
- Crear un sistema para la visualización de los datos financieros de Home Credit.

Introducción

Los créditos son muy importante a la hora de optar por alguna mejora o adquisición de bienes, pero las entidades financieras tienen que asegurarse que los clientes tengan la capacidad económica para pagar, por lo que Home Credit no es la excepción, pero muchas veces por historiales de crédito insuficientes o inexistentes los clientes no pueden optar por uno.

Home Credit utiliza varios métodos estadísticos y de aprendizaje automático para hacer predicciones sobre si es factible o no brindar un crédito a una persona por lo que, en el presente proyecto se va a hacer uso de los datos brindados para desarrollar un sistema capaz de mostrar información importante acerca de Home Credit.

Por otra parte la información brindada por el dataset de Home Credit es información desnormalizada en formato CSV por lo que en el presente trabajo se hará énfasis en el proceso de normalización, migración, seguridad y visualización de datos de manera eficiente haciendo uso de las buenas prácticas de SQL Server.

Descripción del problema

Home Credit es una institución financiera que se encarga de realizar estudios para ampliar los créditos a personas que no se encuentran asociadas a historiales crediticios o personas que cuentan con pocos créditos en una institución financiera. Por lo que esta empresa va a tratar de brindar la mejor experiencia a las personas que van a optar por realizar un crédito, pero a su vez también determinar por medio de una gran cantidad de datos alternativos brindados por los usuarios para predecir las capacidades de pago de sus clientes.

De esta manera, Home Credit utiliza todos los datos mencionados anteriormente para ingresarlos en varios métodos estadísticos y de aprendizaje automático para hacer estas predicciones, ya asegurará que los clientes capaces de pagar no sean rechazados y que los préstamos se otorguen con un calendario de capital, vencimiento y reembolso que permita a sus clientes tener éxito.

Para la utilización del dataset de Home Credit, primeramente se analiza que se encuentran almacenados los datos en archivos CSV, por lo que se debe de hacer la obtención de los títulos de cada columna que estos se van a convertir en los atributos de cada entidad de la base de datos. Una vez obtenidos estos títulos se tiene que empezar a analizar de manera muy cuidadosa, tratando de identificar posibles anomalías para proceder a aplicar la normalización de estas y empezar con la creación del modelo relacional y E/R.

Después una vez normalizados se empieza la creación de todas las tablas, restricciones y procedimientos almacenados que se van a encargar de ingresar los datos provenientes del CSV hacia la base de datos.

Posteriormente se presenta el problema de que los datos se encuentran en archivos CSV, por lo que se debe de crear toda la lógica del migrador, donde un diagrama que nos puede facilitar este proceso se visualiza en la figura 1.

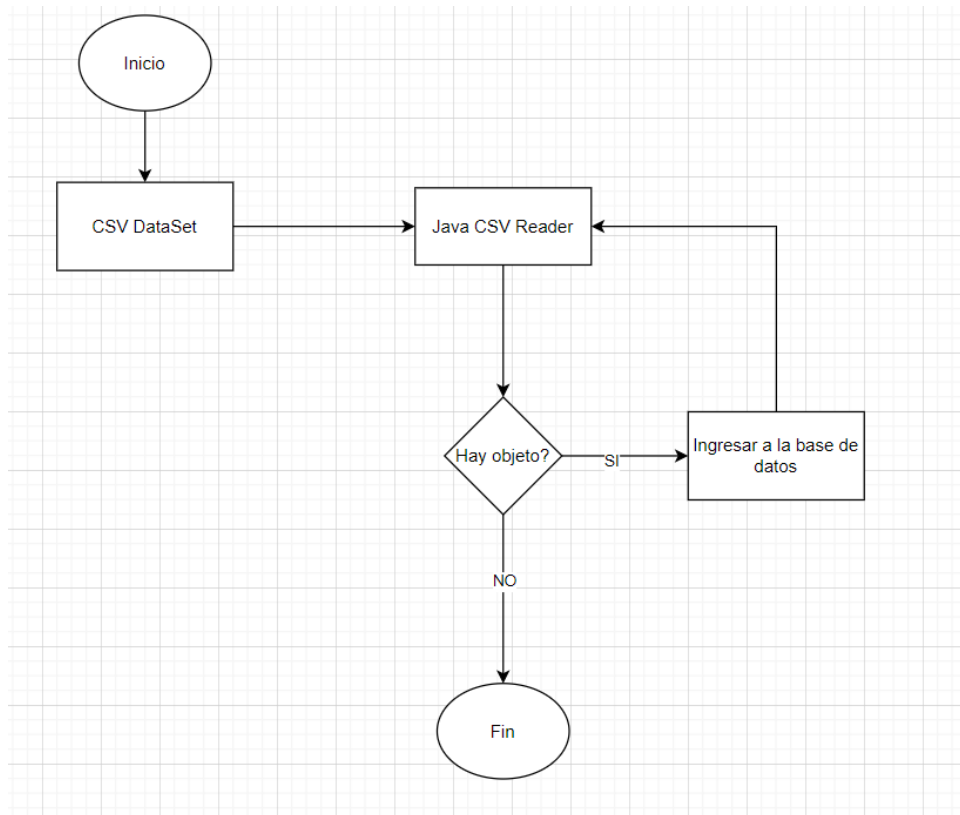


Figura 1. Diagrama que nos ayuda a comprender el ingreso de los datos a la base.

Por último se debe de crear una aplicación para mostrar todos los datos ingresados en la base de datos, esta debe de contener filtros y componentes que nos demuestran la información de los clientes de Home Credit.

Desarrollo

Primeramente para la creación del proyecto se planteó y se dividieron varias secciones de importancia para poder modular el problema, esta división está enfocada en el método de divide y vencerás donde se resuelve tareas más pequeñas y luego se agrupan estas soluciones para la creación de la solución general. En la figura 2 se puede ejemplificar la manera en que se dividieron los módulos para posteriormente empezar a trabajarlos de manera individual.

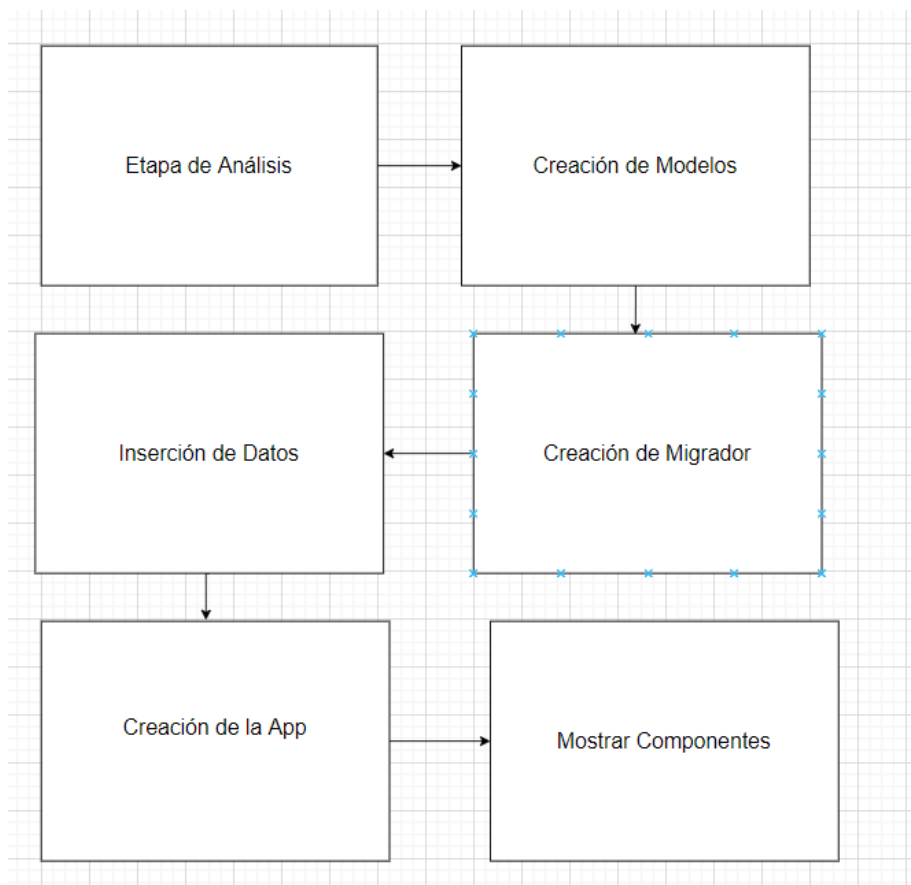


Figura 2. División inicial del problema a resolver

En la etapa de análisis se obtuvieron todos los títulos de las columnas de los archivos CSV, ya que estos son los atributos de las entidades en la base de datos. En este proceso se verificaron que se encuentra con varias anomalías ya que algunos de estos datos no se encontraban ni en la primera, segunda ni tercera forma normal qué es lo que se recomienda. En las siguientes figuras se muestra de manera inicial el proceso de análisis y especificación de algunas entidades.

<p>0: "SK_ID_CURR" PK</p> <p>1: "TARGET"</p> <p>1: "NAME_CONTRACT_TYPE"</p> <p>2: "CODE_GENDER"</p> <p>5: "CNT_CHILDREN"</p> <p>6: "AMT_INCOME_TOTAL"</p> <p>7: "AMT_CREDIT"</p> <p>8: "AMT_ANNUITY"</p> <p>9: "AMT_GOODS_PRICE"</p> <p>10: "NAME_TYPE_SUITE"</p> <p>11: "NAME_INCOME_TYPE"</p> <p>12: "NAME_EDUCATION_TYPE"</p> <p>13: "NAME_FAMILY_STATUS"</p> <p>14: "NAME_HOUSING_TYPE"</p> <p>15: "REGION_POPULATION_RELATIVE"</p> <p>16: "DAYS_BIRTH"</p> <p>17: "DAYS_EMPLOYED"</p> <p>18: "DAYS_REGISTRATION"</p> <p>19: "DAYS_ID_PUBLISH"</p> <p>20: "OWN_CAR_AGE"</p> <p>27: "OCCUPATION_TYPE"</p> <p>28: "CNT_FAM_MEMBERS"</p> <p>39: "ORGANIZATION_TYPE"</p> <p>90: "OBS_30_CNT_SOCIAL_CIRCLE"</p> <p>91: "DEF_30_CNT_SOCIAL_CIRCLE"</p> <p>92: "OBS_60_CNT_SOCIAL_CIRCLE"</p> <p>93: "DEF_60_CNT_SOCIAL_CIRCLE"</p> <p>94: "DAYS_LAST_PHONE_CHANGE"</p>	<p>ClientApartmentsDetails (<u>Toda la información acerca del apartamento en que vive el cliente</u>)</p> <p>0: "SK_ID_CURR" FK</p> <p>43: "APARTMENTS_AVG"</p> <p>44: "BASEMENTAREA_AVG"</p> <p>45: "YEARS_BEGINEXPLUATATION_AVG"</p> <p>46: "YEARS_BUILD_AVG"</p> <p>47: "COMMONAREA_AVG"</p> <p>48: "ELEVATORS_AVG"</p> <p>49: "ENTRANCES_AVG"</p> <p>50: "FLOORSMAX_AVG"</p> <p>51: "FLOORSMIN_AVG"</p> <p>52: "LANDAREA_AVG"</p> <p>53: "LIVINGAPARTMENTS_AVG"</p> <p>54: "LIVINGAREA_AVG"</p> <p>55: "NONLIVINGAPARTMENTS_AVG"</p> <p>56: "NONLIVINGAREA_AVG"</p> <p>57: "APARTMENTS_MODE"</p> <p>58: "BASEMENTAREA_MODE"</p> <p>59: "YEARS_BEGINEXPLUATATION_MODE"</p> <p>60: "YEARS_BUILD_MODE"</p> <p>61: "COMMONAREA_MODE"</p> <p>62: "ELEVATORS_MODE"</p> <p>63: "ENTRANCES_MODE"</p> <p>64: "FLOORSMAX_MODE"</p>	<p>Flag Document</p> <p>0: "SK_ID_CURR" FK</p> <p>1: "Document"</p> <p>AMT_REQ_Credit</p> <p>0: "SK_ID_CURR" FK</p> <p>BUREAU_HOUR"</p> <p>BUREAU_DAY"</p> <p>BUREAU_WEEK"</p> <p>BUREAU_MON"</p> <p>BUREAU_QRT"</p> <p>BUREAU_YEAR"</p>
<p>FlagsProvidedInformation (<u>Si el cliente proporcionó la siguiente información</u>)</p> <p>0: "SK_ID_CURR" FK</p> <p>FLAG_OWN_CAR</p> <p>FLAG_OWN_REALTY</p> <p>FLAG_MOBIL</p> <p>FLAG_EMP_PHONE"</p> <p>FLAG_WORK_PHONE"</p> <p>FLAG_CONT_MOBILE"</p> <p>FLAG_PHONE"</p> <p>FLAG_EMAIL"</p> <p>AddresClientInformation</p> <p>0: "SK_ID_CURR" FK</p> <p>29: "REGION_RATING_CLIENT"</p> <p>30: "REGION_RATING_CLIENT_W_CITY"</p> <p>31: "WEEKDAY_APPR_PROCESS_START"</p> <p>32: "HOUR_APPR_PROCESS_START"</p> <p>33: "REG_REGION_NOT_LIVE_REGION"</p> <p>34: "REG_REGION_NOT_WORK_REGION"</p> <p>35: "LIVE_REGION_NOT_WORK_REGION"</p> <p>36: "REG_CITY_NOT_LIVE_CITY"</p> <p>37: "REG_CITY_NOT_WORK_CITY"</p> <p>38: "LIVE_CITY_NOT_WORK_CITY"</p> <p>40: "EXT_SOURCE_1"</p> <p>41: "EXT_SOURCE_2"</p> <p>42: "EXT_SOURCE_3"</p>	<p>65: "FLOORSMIN_MODE"</p> <p>66: "LANDAREA_MODE"</p> <p>67: "LIVINGAPARTMENTS_MODE"</p> <p>68: "LIVINGAREA_MODE"</p> <p>69: "NONLIVINGAPARTMENTS_MODE"</p> <p>70: "NONLIVINGAREA_MODE"</p> <p>71: "APARTMENTS_MEDI"</p> <p>72: "BASEMENTAREA_MEDI"</p> <p>73: "YEARS_BEGINEXPLUATATION_MEDI"</p> <p>74: "YEARS_BUILD_MEDI"</p> <p>75: "COMMONAREA_MEDI"</p> <p>76: "ELEVATORS_MEDI"</p> <p>77: "ENTRANCES_MEDI"</p> <p>78: "FLOORSMAX_MEDI"</p> <p>79: "FLOORSMIN_MEDI"</p> <p>80: "LANDAREA_MEDI"</p> <p>81: "LIVINGAPARTMENTS_MEDI"</p> <p>82: "LIVINGAREA_MEDI"</p> <p>83: "NONLIVINGAPARTMENTS_MEDI"</p> <p>84: "NONLIVINGAREA_MEDI"</p> <p>85: "FONDKAPREMONT_MODE"</p> <p>86: "HOUSETYPE_MODE"</p> <p>87: "TOTALAREA_MODE"</p> <p>88: "WALLSMATERIAL_MODE"</p> <p>89: " WALLSMATERIAL_MODE "</p>	

Bureau

0: "SK_ID_CURR" FK PK
1: "SK_ID_BUREAU" PK
2: "CREDIT_ACTIVE"
3: "CREDIT_CURRENCY"
4: "DAYS_CREDIT"
5: "CREDIT_DAY_OVERDUE"
6: "DAYS_CREDIT_ENDDATE"
7: "DAYS_ENDDATE_FACT"
9: "CNT_CREDIT_PROLONG"
14: "CREDIT_TYPE"
15: "DAYS_CREDIT_UPDATE"

AMT_CREDIT

1: "SK_ID_BUREAU" PK
8: "MAX_OVERDUE"
10: "SUM"
11: "SUM_DEBT"
12: "SUM_LIMIT"
13: "SUM_OVERDUE"
16: "ANNUITY"

installments_payments

SK_ID_PREV FK
SK_ID_CURR FK
NUM_INSTALLMENT_VERSION
NUM_INSTALLMENT_NUMBER
DAYS_INSTALLMENT
DAYS_ENTRY_PAYMENT
AMT_INSTALLMENT
AMT_PAYMENT

POS CASH balance

0: "SK_ID_PREV" FK
1: "SK_ID_CURR" FK
2: "MONTHS_BALANCE"
3: "CNT_INSTALLMENT"
4: "CNT_INSTALLMENT_FUTURE"
5: "NAME_CONTRACT_STATUS"
6: "SK_DPD"
7: "SK_DPD_DEF"

credit_card_balance

0: "SK_ID_PREV" FK
1: "SK_ID_CURR" FK
1: ID_CCB PK
2: "MONTHS_BALANCE"
20: "NAME_CONTRACT_STATUS"
21: "SK_DPD"
22: "SK_DPD_DEF"

AMT_BALANCE_DETAILS

1: ID_CCB FK
3: BALANCE
4: CREDIT_LIMIT_ACTUAL
5: DRAWINGS_ATM_CURRENT
6: DRAWINGS_CURRENT
7: DRAWINGS_OTHER_CURRENT
8: DRAWINGS_POS_CURRENT
9: INST_MIN_REGULARITY
10: PAYMENT_CURRENT
11: PAYMENT_TOTAL_CURRENT
12: RECEIVABLE_PRINCIPAL
13: RECEIVABLE
14: TOTAL_RECEIVABLE

CNT_BALANCE_DETAILS

1: ID_CCB FK
15: DRAWINGS_ATM_CURRENT
16: DRAWINGS_CURRENT
17: DRAWINGS_OTHER_CURRENT
18: DRAWINGS_POS_CURRENT
19: INSTALLMENT_MATURE_CUM

AMT_PREVIOUS_DETAILS

0: "SK_ID_PREV" FK
3: AMT_ANNUITY
4: AMT_APPLICATION
5: AMT_CREDIT
6: AMT_DOWN_PAYMENT
7: AMT_GOODS_PRICE

bureau_balance

0: "SK_ID_BUREAU" PK FK
1: "MONTHS_BALANCE"
2: "STATUS"

previous_application

0: "SK_ID_PREV" PK
1: SK_ID_CURR FK
2: NAME_CONTRACT_TYPE
8: WEEKDAY_APPR_PROCESS_START
9: HOUR_APPR_PROCESS_START
10: FLAG_LAST_APPL_PER_CONTRACT
11: NFLAG_LAST_APPL_IN_DAY
12: RATE_DOWN_PAYMENT
13: RATE_INTEREST_PRIMARY
14: RATE_INTEREST_PRIVILEGED
15: NAME_CASH_LOAN_PURPOSE
16: NAME_CONTRACT_STATUS
17: DAYS_DECISION
18: NAME_PAYMENT_TYPE
19: CODE_REJECT_REASON
20: NAME_TYPE_SUITE
21: NAME_CLIENT_TYPE
22: NAME_GOODS_CATEGORY
23: NAME_PORTFOLIO
24: NAME_PRODUCT_TYPE
25: CHANNEL_TYPE
26: SELLERPLACE_AREA
27: NAME_SELLER_INDUSTRY
28: CNT_PAYMENT
29: NAME_YIELD_GROUP
30: PRODUCT_COMBINATION
31: DAYS_FIRST_DRAWING
32: DAYS_FIRST_DUE
33: DAYS_LAST_DUE_1ST_VERSION
34: DAYS_LAST_DUE
35: DAYS_TERMINATION
36: NFLAG_INSURED_ON_APPROVAL

Figura 3. División de las entidades

Una vez realizado el proceso de análisis y normalización se inicia con la creación de los modelos de E/R y modelo relacional, que son de suma importancia para la creación de la base de datos, donde en estos se termina de arreglar los errores del proceso de anterior y se define de la manera en que va a quedar la base de datos. En la siguiente figura se muestra un ejemplo de algunas entidades del modelo relacional de Home Credit.

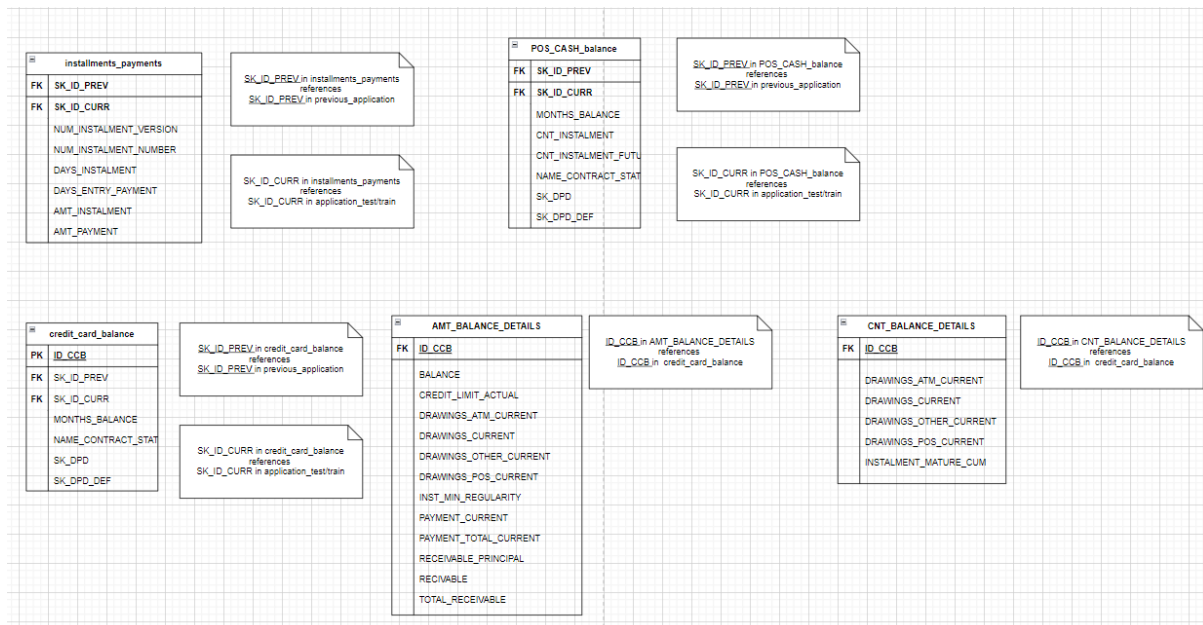


Figura 4: Entidades del modelo relacional de Home Credit

En la creación del migrador se hizo uso del lenguaje Java junto a una librería llamada Open CSV, que proporciona todos los métodos necesario la para lectura de los archivos CSV, a estos se le hizo un casteo para transformarlos en datos aceptados por los procedimientos almacenados creados con anterioridad. Se empieza a ejecutar el migrador de forma en que las entidades que son independientes se ingresen y posteriormente las que tienen llaves foráneas a estas otras entidades. En la figura 5 se muestra la manera en que se agregaron.

```

m1.migrationAplicationTrain(conex);
m1.migrationClientApartments(conex);
m1.migrationAddressClientInformation(conex);
m1.migrationFlagDocument(conex);
m1.migrationFlagsProvidedInfo(conex);
m1.migrationAMTREQCredit(conex);
m1.migrationBureau(conex);
m1.migrationBureauBalance(conex);
m1.migrationAMTCreditDetails(conex);
m1.migrationPreviousApplication(conex);
m1.migrationPosCashBalance(conex);
m1.migrationInstallmentsPayments(conex);
m1.migrationAMT_PREVIOUS_DETAILS(conex);
m1.migrationCreditCardBalance(conex);
m1.migrationAMT_BALANCE_DETAILS(conex);
m1.migrationCNTBalanceDetails(conex);

```

Figura 5. Ingreso de los datos a la base Home Credit.

Por último se utilizó Java con la parte gráfica de Swing, la cuál permitió crear una interfaz sencilla pero funcional para mostrar los datos de la base de datos Home Credit.

Análisis de resultados

Tarea/Requerimiento	Estado	Observaciones
Modelo E/R	Completo	
Modelo Relacional	Completo	
Esquemas	Completo	Se crearon 3 esquemas
Migrador	Completo	Se utilizó Open CSV de Java
Funciones para insertar los datos	Completo	
Consultas	Completo	Se crearon 3 consultas
Índices No Clúster	Completo	
Crear usuarios	Completo	
Componentes	Completo	Se crearon dos componentes con +1000 000 de datos. Y se muestran con JFreeChart

Conclusiones

Al trabajar con datasets tan grandes es importante darle énfasis en la optimización y utilizar las buenas prácticas de las bases de datos, ya que estas están fuertemente vinculadas al tiempo de espera de un usuario final, por lo que va a ser muy importante que cuando el usuario solicita una información tratar de brindarla en el menor tiempo posible.

El proceso de migración es sumamente importante ya que es un proceso que se realiza con frecuencia en el ámbito profesional ya que los sistemas se actualizan a tecnologías más eficientes y por lo que los datos deben de migrarse.

Recomendaciones

Para la creación de estos sistemas gestores de bases de datos se recomienda que se utilice Web en vez de una aplicación de escritorio, ya que hay muchas más información y además hay muchas más opciones que permiten dibujar los gráficos y mostrar los datos de una manera más profesional, ya que al adentrarse un poco en la investigación de herramientas para crear gráficos en java fueron pocas las opciones que se presentaron.

Bibliografía

How to add row of data to Jtable from values received from jtextfield and comboboxes. (2014, 15 enero). Stack Overflow. <https://stackoverflow.com/questions/21135452/how-to-add-row-of-data-to-jtable-from-values-received-from-jtextfield-and-combob>

Melgoza, J. (2019, 1 abril). Como Hacer Graficos con Java – Pastel, Barras, Lineas, 3D . . . Jonathan Melgoza. <https://jonathanmelgoza.com/blog/como-hacer-graficos-con-java/>

Gigena, M., & Perfil, V. T. M. (s. f.). Todo Java. Todo Java. <http://labojava.blogspot.com/2012/06/graficos-estadisticos-jfreechart>

Elaboración de cuadros estadísticos en netbeans. (2018, 15 febrero). YouTube. https://www.youtube.com/watch?v=MBWZMQ-IneI&ab_channel=TICUTMACH

Home Credit Default Risk | Kaggle. (s. f.). Kaggle. <https://www.kaggle.com/competitions/home-credit-default-risk/overview>

M. (2020, 26 diciembre). How to read and parse CSV file in Java. Mkyong.Com. <https://mkyong.com/java/how-to-read-and-parse-csv-file-in-java/>