



طراحی DSL تحلیل داده های آزمایشگاهی

مستندات جامع فنی و پیاده سازی و تست سناریو

پروژه پایانی درس کامپایلر

اعضای تیم پروژه:

محمدعلی احمدیان
علی جباری پور
سید احمد موسوی مالواجردی

تاریخ ارائه:

۱۰ بهمن ۱۴۰۴

فهرست مطالب

۳	۱	مقدمه و نمای کلی سیستم
۳	۱.۱	معماری کلان سیستم
۳	۲	کتابخانه‌های مورد استفاده و نقش هر کدام
۴	۳	موتور ترجمه فارسی به DSL میانی: کلاس PersianToDSLMapper
۴	۱.۳	انگیزه و ضرورت طراحی این لایه
۵	۲.۳	معماری و ساختار کلاس PersianToDSLMapper
۵	۳.۳	تحلیل قواعد نگاشت (Mapping Rules)
۶	۴.۳	مزایای استفاده از زبان میانی
۶	۴	تعریف گرامر DSL در Lark
۶	۱.۴	نمای کلی گرامر
۷	۲.۴	Non-Terminals (قواعد غیربنیادی)
۷	۳.۴	دستور بارگذاری داده (load_stmt)
۷	۴.۴	دستور ذخیره‌سازی (save_stmt)
۷	۵.۴	دستور کپی داده (duplicate_stmt)
۷	۶.۴	دستورات پاکسازی (clean_stmt)
۸	۷.۴	دستورات محاسبات آماری (calc_stmt)
۸	۸.۴	دستور رسم نمودار (plot_stmt)
۸	۹.۴	دستورات فیلتر (filter_range_stmt و filter_stmt)
۸	۱۰.۴	فیلتر شرطی پیچیده (filter_complex_stmt)
۹	۱۱.۴	دستور جستجو (search_stmt)
۹	۱۲.۴	سطح‌بندی داده (level_stmt)
۹	۱۳.۴	مرتب‌سازی (sort_stmt)
۹	۱۴.۴	گروه‌بندی و تجمیع (groupby_stmt)
۹	۱۵.۴	سایر دستورات
۱۰	۵	مصورسازی AST: تابع ast_to_dot
۱۱	۶	کلاس CodeGenerator: تبدیل AST به کد پایتون
۱۱	۱.۶	پیکربندی اولیه و سیستم گزارش‌دهی
۱۱	۱.۱.۶	سیستم لاگ‌گیری
۱۱	۲.۶	مدیریت ساختار کلی AST
۱۲	۳.۶	دستور LOAD: بارگذاری داده
۱۲	۴.۶	دستورات اطلاعاتی: HEAD و DESCRIBE
۱۲	۵.۶	دستور DUPLICATE و SAVE
۱۲	۶.۶	دستور CLEAN: پاکسازی داده
۱۲	۱.۶.۶	تشخیص نوع عملیات
۱۲	۲.۶.۶	حذف داده‌های پرت (IQR)
۱۳	۳.۶.۶	جایگزینی مقادیر خالی
۱۳	۷.۶	دستورات CALC: محاسبات آماری
۱۳	۸.۶	دستورات PLOT: تولید نمودار
۱۳	۹.۶	دستورات FILTER
۱۳	۱۰.۶	LEVELING: دسته‌بندی بازه‌ای
۱۳	۱۱.۶	دستورات CRUD روی ستون‌ها
۱۳	۱۲.۶	NORMALIZE: نرمال‌سازی Min-Max
۱۴	۱۳.۶	CORRELATE: محاسبه هم‌بستگی
۱۴	۷	نحوه استفاده از CodeGenerator و تولید کد نهایی
۱۴	۱.۷	جریان کامل کامپایل
۱۴	۲.۷	نحوه استفاده از CodeGenerator
۱۴	۳.۷	خروجی متد transform
۱۵	۴.۷	تولید فایل generated_code.py
۱۵	۵.۷	اجرای کد تولیدشده
۱۵	۶.۷	تولید گزارش نهایی
۱۵	۷.۷	خلاصه خروجی‌های CodeGenerator
۱۶	۸.۷	جمع‌بندی
۱۶	۸	پایپ‌لاین run_compiler و ارتباط با کد تولیدی

۱۸	۱۰ توضیح سناریوی نمونه و نحوه اجرای آن
۱۸	۱.۱۰ معرفی سناریو
۱۹	۲.۱۰ جریان اجرای سناریو در کامپایلر
۲۰	۳.۱۰ توضیح گام به گام عملیات سناریو
۲۰	۴.۱۰ تحلیل نتایج اجرا بر اساس داده‌ی ورودی نمونه
۲۰	۱.۴.۱۰ بررسی داده‌ی ورودی
۲۱	۲.۴.۱۰ تأثیر پاکسازی و فیلتر داده‌ها
۲۱	۳.۴.۱۰ محاسبه شاخص سلامت و نرمال‌سازی
۲۱	۴.۴.۱۰ سطح‌بندی معنایی بیماران
۲۲	۵.۴.۱۰ تحلیل نقشه حرارتی همبستگی
۲۲	۶.۴.۱۰ نمودار پراکندگی سن و شاخص سلامت
۲۲	۵.۱۰ فایل رپورت
۲۴	۶.۱۰ کد تولید شده
۲۵	۷.۱۰ تصویر ast نهایی از دستورات
۲۶	۸.۱۰ دیتاست نهایی ذخیره شده
۲۶	۹.۱۰ جمع‌بندی نهایی سناریو

۱ مقدمه و نمای کلی سیستم

این پروژه یک `Compiler` برای یک زبان خاص منظوره^۱ (DSL) فارسی در حوزه تحلیل داده های آزمایشگاهی است. هدف اصلی سیستم این است که:

- کاربر بتواند دستورات تحلیل داده را به زبان فارسی (نیمه طبیعی) و یا حتی انگلیسی، بنویسد.
- این دستورات به یک DSL میانی ساخت یافته و سپس به `Python + Pandas + Matplotlib + Seaborn` تبدیل شوند.
- در نهایت علاوه بر اجرای کد، موارد زیر تولید شود:
 - کد پایتون تولید شده که اگر همین کد را به تنهایی اجرا کنیم، کار هایی که از DSL خواستیم را برایمان انجام میدهد ولی نتایج `cpmiller` و ریزگزارشات و درخت را دیگر نمیسازد!
 - گزارش متنی آماری (`report.txt`)
 - نمودارها (`.png`) در پوشه `plots` و همچنین یک تصویر از `AST` ساخته شده توسط دستورات ما
 - فایل `ast.dot` برای مصورسازی درخت نحو (`AST`)

۱.۱ معماری کلان سیستم

معماری پروژه به طور خلاصه شامل این مراحل است:

۱. نگاشت فارسی به **DSL میانی**: تبدیل الگوهای فارسی به دستورات DSL انگلیسی مانند (مثلاً " بگیر از ... به نام ... " به " `LOAD ... INTO ...` ")
۲. تحلیل نحوی (**Parsing**): استفاده از گرامر EBNF در `Lark` برای تولید `AST`.
۳. تولید کد (**Code Generation**): پیمایش `AST` توسط `Transformer` و تولید کد پایتون، گزارش آماری و دستورات رسم نمودار.
۴. اجرا و خروجی: اجرای کد تولید شده در محیطی کنترل شده و ذخیره کد، گزارش، نمودارها و شمای `AST`.
این ساختار عملاً یک **pipeline کامپایلر کامل** را نشان می دهد از زبان سطح بالا (DSL فارسی) تا کد اجرایی پایتون.

۲ کتابخانه های مورد استفاده و نقش هر کدام

در ابتدای فایل پایتون، کتابخانه های زیر فراخوانی شده اند:

```
import re
import os
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from lark import Lark, Transformer
from datetime import datetime
```

تحلیل نقش هر ماژول

- **re**: برای کار با عبارات باقاعده (`Regex`) هسته ی نگاشت فارسی به DSL از این ماژول استفاده می کند.
- **os**: برای کار با فایل ها و پوشه ها:
 - ساخت پوشه ی `output` و `plots`
 - بررسی وجود فایل ورودی
 - ساخت مسیرهای نسبی برای نمودارها و کد تولیدی

¹Domain Specific Language

- **pandas**: کتابخانه‌ی اصلی تحلیل داده. هر `NORMALIZE`، `MERGE`، `GROUPBY`، `FILTER`، `LOAD` و غیره در نهایت به توابع `pandas` نگاشت می‌شود.
- **matplotlib.pyplot (plt) & seaborn (sns)**: برای رسم نمودارهای
 - نمودارهای هیستوگرام، خطی، پراکندگی، جعبه‌ای با `matplotlib`
 - نقشه‌ی حرارتی (Heatmap) با `seaborn`
- **lark.Lark, lark.Transformer**:
 - `Lark`: تحلیل نحوی بر اساس گرامر DSL
 - `Transformer`: تبدیل AST به کد پایتون (Code Generation)
- **datetime**: برای ثبت زمان عملیات در گزارش آماری (timestamp در لاگ‌ها).

۳ موتور ترجمه فارسی به DSL میانی: کلاس `PersianToDSLMapper`

۱.۳ انگیزه و ضرورت طراحی این لایه

پردازش مستقیم زبان طبیعی فارسی با استفاده از گرامرهای رسمی (نظیر Context-Free Grammar) به دلایل متعددی بسیار دشوار، شکننده و پرهزینه است. از جمله این چالش‌ها می‌توان به موارد زیر اشاره کرد:

- تنوع بسیار زیاد ساختارهای نحوی در زبان فارسی،
 - آزادی بالای ترتیب کلمات در جملات،
 - وجود مترادف‌ها و بیان‌های مختلف برای یک مفهوم واحد،
 - دشواری نگهداری و توسعه گرامرهای بزرگ زبان طبیعی.
- به همین دلیل، در این پروژه به جای تعریف مستقیم گرامر برای جملات فارسی، یک معماری دومارحله‌ای اتخاذ شده است. در این معماری:

۱. ابتدا جملات فارسی کاربر با استفاده از عبارات باقاعده (Regular Expressions) به یک زبان میانی (Inter-DSL `mediate`) ساده، محدود و ساخت‌یافته تبدیل می‌شوند؛

۲. سپس روی این زبان میانی، یک گرامر رسمی و دقیق با استفاده از ابزار `Lark` تعریف می‌شود.

کلاس `PersianToDSLMapper` مسئول پیاده‌سازی مرحله‌ی اول این فرآیند است و نقش یک «مترجم سطح بالا» از زبان طبیعی فارسی به DSL میانی را ایفا می‌کند.

۲.۳ معماری و ساختار کلاس PersianToDSLMapper

شکل ۱ نمای کلی ساختار این کلاس را نشان می‌دهد.

```
class PersianToDSLMapper:
    def __init__(self):
        self.rules = [
            (r'([^\s]+)' "به نام (\w+)", r'LOAD "\1" INTO \2'),
            (r'(\w+)' "خلاصه (\w+)", r'DESCRIBE \1'),
            (r'(\w+) : (\d+)' "سطر اول (\w+)", r'HEAD \1 \2'),
            (r'(\w+) : (\w+)' "حذف تکراری (\w+)", r'CLEAN \1 DROP_DUPLICATES'),
            (r'(\w+) : (\w+)' "حذف کلی (پرت|خالی)", r'CLEAN \1 DROP_ALL \2'),
            (r'(\w+) : (\w+)' "حذف (\w+) (پرت|خالی)", r'CLEAN \1 DROP_SPECIFIC \2 \3'),
            ...
        ]

    def translate(self, text):
        for p, r in self.rules:
            text = re.sub(p, r, text)
        return text.strip()
```

شکل ۱: ساختار کلی کلاس PersianToDSLMapper

در این کلاس، مجموعه‌ای از قواعد نگاشت به صورت یک لیست از زوج‌های (pattern, replacement) نگهداری می‌شود. هر قاعده بیان می‌کند که:

- یک الگوی خاص در زبان فارسی چگونه شناسایی شود،
- و به چه دستور استاندارد در DSL میانی تبدیل گردد.

فرآیند ترجمه به صورت ترتیبی انجام می‌شود؛ به این معنا که ورودی کاربر به ترتیب از روی تمام قواعد عبور داده شده و هر قاعده‌ای که منطبق باشد، جایگزینی متناظر خود را اعمال می‌کند.

۳.۳ تحلیل قواعد نگاشت (Mapping Rules)

هر عنصر از `self.rules` شامل دو بخش اصلی است:

- **Pattern (الگو):** یک عبارت باقاعده که ساختار یک دستور فارسی خاص را شناسایی می‌کند. این الگو می‌تواند شامل گروه‌های کپیچر باشد تا اجزای متغیر جمله (مانند نام دیتافریم یا ستون) استخراج شوند.
- **Replacement (جایگزین):** یک رشته‌ی متنی در قالب DSL میانی که با استفاده از گروه‌های استخراج‌شده، دستور معادل استاندارد را تولید می‌کند.

به عنوان مثال، جملات فارسی زیر:

- « پاکسازی df حذف همه خالی‌ها »
- « پاکسازی df حذف ستون age که خالی است »
- « پر کردن مقادیر خالی با میانگین در df »

می‌توانند به صورت دستورات DSL زیر نگاشت شوند:

- " CLEAN df DROP_ALL خالی "
- " CLEAN df DROP_SPECIFIC age خالی "
- " CLEAN df FILL_ALL میانگین خالی "

این استانداردسازی باعث می‌شود که تنوع بالای زبان طبیعی به مجموعه‌ای کوچک و کنترل‌شده از دستورات تقلیل یابد.

۴.۳ مزایای استفاده از زبان میانی

استفاده از لایه‌ی DSL میانی مزایای مهمی به همراه دارد:

- گرامر Lark تنها با یک زبان ساده، صریح و بدون ابهام سروکار دارد؛
- افزودن قابلیت‌های جدید صرفاً با اضافه کردن یک قاعده‌ی نگاشت انجام می‌شود، بدون نیاز به تغییر گرامر اصلی؛
- خطاهای نحوی فارسی پیش از ورود به مرحله‌ی پارس کاهش می‌یابند؛
- امکان ترکیب دستورات فارسی و انگلیسی در یک اسکریپت فراهم می‌شود.

در نتیجه، کلاس `PersianToDSLMapper` نقش یک لایه‌ی حیاتی در افزایش پایداری، توسعه‌پذیری و خوانایی کل کامپایلر را ایفا می‌کند.

۴ تعریف گرامر DSL در Lark

۱.۴ نمای کلی گرامر

گرامر DSL پروژه با استفاده از ابزار Lark و در قالب یک رشته‌ی چندخطی (raw string) تعریف شده است. این گرامر، زبان میانی استاندارد را توصیف می‌کند که پس از ترجمه‌ی جملات فارسی توسط لایه‌ی `PersianToDSLMapper`، به عنوان ورودی مرحله‌ی پارس و تولید AST استفاده می‌شود.

```
grammar = r"""
start: statement+

statement: load_stmt
         | duplicate_stmt
         | save_stmt
         | clean_stmt
         | calc_stmt
         | plot_stmt
         | filter_stmt
         | filter_range_stmt
         | level_stmt
         | sort_stmt
         | groupby_stmt
         | drop_col_stmt
         | describe_stmt
         | head_stmt
         | filter_complex_stmt
         | search_stmt
         | merge_stmt
         | create_col_stmt
         | rename_stmt
         | convert_time_stmt
         | normalize_stmt
         | corr_stmt

load_stmt: "LOAD" STRING "INTO" ID

...
"""
```

گرامر ریشه‌ی

```
start: statement+
```

این قاعده بیان می‌کند که یک برنامه‌ی DSL شامل یک یا چند دستور است که به صورت ترتیبی اجرا می‌شوند.

دستورها (statement) هر دستور DSL می‌تواند یکی از حالت‌های زیر باشد:

- بارگذاری داده
- پاکسازی و پیش‌پردازش
- محاسبات آماری
- فیلتر و جستجو
- رسم نمودار
- تبدیل و مهندسی ویژگی
- تحلیل آماری و همبستگی

در ادامه، تمام های statement گرامر به تفصیل توضیح داده می‌شوند.

۳.۴ دستور بارگذاری داده (load_stmt)

```
load_stmt: "LOAD" STRING "INTO" ID
```

این دستور برای بارگذاری یک فایل داده (مانند CSV یا Excel) و ذخیره‌ی آن در یک متغیر (دیتافریم) استفاده می‌شود.
ساختار کلی:

```
LOAD "file.csv" INTO df
```

مفهوم: فایل مشخص‌شده خوانده شده و به صورت یک دیتافریم با نام داده‌شده در حافظه قرار می‌گیرد.

۴.۴ دستور ذخیره‌سازی (save_stmt)

```
save_stmt: "SAVE" ID "TO" STRING
```

این دستور برای ذخیره‌ی یک دیتافریم در قالب فایل خروجی استفاده می‌شود.
مثال:

```
SAVE df TO "output.csv"
```

۵.۴ دستور کپی داده (duplicate_stmt)

```
duplicate_stmt: "DUPLICATE" ID "TO" ID
```

یک کپی مستقل از یک دیتافریم ایجاد می‌کند.
مثال:

```
DUPLICATE df TO df_backup
```

۶.۴ دستورات پاکسازی (clean_stmt)

```
clean_stmt: "CLEAN" ID clean_op
```

این دستور مجموعه‌ای از عملیات پاکسازی داده را پوشش می‌دهد.

داده‌ها حذف

```
۱ drop_all: "DROP_ALL" TARGET
۲ drop_specific: "DROP_SPECIFIC" ID TARGET
```

مثال:

CLEAN df DROP_ALL خالی

داده‌ها کردن پر

```
۱ fill_all: "FILL_ALL" TARGET METHOD
۲ fill_specific: "FILL_SPECIFIC" ID TARGET METHOD
```

مثال:

CLEAN df FILL_ALL می‌انگین خالی

تکراری داده‌های حذف

```
۱ drop_duplicates: "DROP_DUPLICATES"
```

—

۷.۴ دستورات محاسبات آماری (calc_stmt)

```
۱ calc_stmt: "CALC" ID calc_op+
۲ calc_op: MEAN "OF" ID
۳          | STD "OF" ID
```

این دستور برای محاسبه‌ی شاخص‌های آماری مانند میانگین و انحراف معیار استفاده می‌شود.
مثال:

CALC df MEAN OF age STD OF salary

—

۸.۴ دستور رسم نمودار (plot_stmt)

```
۱ plot_stmt: "PLOT" ID PLOT_TYPE "OF" (ID|ALL) "IN" (ID|ALL)
```

این دستور انواع نمودارهای تحلیلی را پشتیبانی می‌کند.
مثال‌ها:

PLOT df OF age IN ALL •

PLOT df OF salary IN age •

PLOT df HEATMAP OF ALL IN ALL •

—

۹.۴ دستورات فیلتر (filter_stmt و filter_range_stmt)

```
۱ filter_stmt: "FILTER" ID "WHERE" ID COMP (STRING|NUMBER)
۲ filter_range_stmt: "FILTER_RANGE" ID ID NUMBER NUMBER
```

مثال:

FILTER df WHERE age >= 18

—

۱۰.۴ فیلتر شرطی پیچیده (filter_complex_stmt)

```
۱ filter_complex_stmt: "FILTER_COMPLEX" ID CONDITION_EXPR
```

این دستور امکان استفاده از شروط ترکیبی و پیچیده را فراهم می‌کند.

—

۱۱.۴ دستور جستجو (search_stmt)

```
۱ search_stmt: "SEARCH" ID "IN" ID "CONTAINS" STRING
```

مثال:

```
SEARCH df IN name CONTAINS "Ali"
```

—

۱۲.۴ سطح‌بندی داده (level_stmt)

```
۱ level_stmt: "LEVELING" ID ID level_op+  
۲ level_op: NUMBER ":" ID
```

برای تبدیل مقادیر عددی به دسته‌های معنایی استفاده می‌شود.
مثال:

```
LEVELING df score 0:low 0.5:medium 0.8:high
```

—

۱۳.۴ مرتب‌سازی (sort_stmt)

```
۱ sort_stmt: "SORT" ID "BY" ID ORDER
```

—

۱۴.۴ گروه‌بندی و تجمیع (groupby_stmt)

```
۱ groupby_stmt: "GROUPBY" ID "BY" ID "OP" AGG_FUNC "OF" ID
```

مثال:

```
GROUPBY df BY gender OP mean OF salary
```

—

۱۵.۴ سایر دستورات

- DROP_COL: حذف ستون
- RENAME: تغییر نام ستون
- CREATE_COL: ایجاد ستون جدید
- CONVERT_TIME: تبدیل ستون زمانی
- NORMALIZE: نرمال‌سازی داده
- DESCRIBE: خلاصه آماری
- HEAD: نمایش سطرهای ابتدایی
- CORRELATE: محاسبه‌ی همبستگی

این مجموعه دستورات، یک زبان DSL قدرتمند، منعطف و قابل توسعه برای تحلیل داده فراهم می‌کند.

۵ مصورسازی AST: تابع ast_to_dot

برای مشاهده‌ی ساختار AST و استفاده در مستندات/دییگ، تابع ast_to_dot در نظر گرفته شده است:

```
1 def ast_to_dot(tree, output_name="ast", output_dir="outputs"):
2     """
3     Generates AST image (PNG) using Graphviz and saves it in outputs directory.
4     """
5
6     os.makedirs(output_dir, exist_ok=True)
7     dot_path = os.path.join(output_dir, f"{output_name}.dot")
8     img_path = os.path.join(output_dir, f"{output_name}.png")
9     node_id = 0
10    lines = [
11        "digraph AST {",
12        "node [shape=box, fontname=\"Helvetica\"];",
13    ]
14
15    def visit(node, parent=None):
16        nonlocal node_id
17        my_id = node_id
18        label = node.data if hasattr(node, "data") else str(node)
19        # Escape quotes for Graphviz
20        label = label.replace('"', '\\')
21        lines.append(f'node{my_id} [label="{label}"];')
22        if parent is not None:
23            lines.append(f'node{parent} -> node{my_id};')
24        node_id += 1
25        if hasattr(node, "children"):
26            for child in node.children:
27                visit(child, my_id)
28    visit(tree)
29    lines.append("}")
30
31    # Write DOT file
32    with open(dot_path, "w", encoding="utf-8") as f:
33        f.write("\n".join(lines))
34
35    # Convert DOT to PNG
36    subprocess.run(
37        ["dot", "-Tpng", dot_path, "-o", img_path],
38        check=True
39    )
```

خلاصه عملکرد

- این تابع یک درخت AST تولیدشده توسط Lark را دریافت می‌کند.
- در ابتدای اجرا، پوشه‌ی خروجی (به‌طور پیش‌فرض outputs/) در صورت عدم وجود ساخته می‌شود.
- متغیر node_id به‌عنوان یک شمارنده‌ی سراسری برای اختصاص شناسه‌ی یکتا به هر گره استفاده می‌شود.
- تابع بازگشتی visit کل درخت را پیمایش کرده و:
 - برای هر گره یک نود Graphviz با قالب nodeX [label="..."] تولید می‌کند.
 - در صورت وجود گره‌ی والد، یال nodeParent -> nodeChild ایجاد می‌نماید.
 - تمامی فرزندان گره (از طریق node.children) را به‌صورت بازگشتی پردازش می‌کند.
- برای جلوگیری از خطا در Graphviz، کاراکترهای نقل‌قول (") در برچسب گره‌ها Escape می‌شوند.
- پس از تکمیل پیمایش، فایل DOT در مسیر مشخص‌شده ذخیره می‌شود.
- در مرحله‌ی نهایی، با استفاده از دستور dot از ابزار Graphviz، فایل DOT به تصویر PNG تبدیل شده و در همان پوشه ذخیره می‌گردد.

۶ کلاس CodeGenerator: تبدیل AST به کد پایتون

کلاس CodeGenerator وظیفه‌ی تبدیل درخت نحوی (AST) حاصل از پارس گرامر DSL به کد اجرایی پایتون را بر عهده دارد. این کلاس از Transformer در کتابخانه‌ی Lark ارث‌بری می‌کند و برای هر statement در گرامر، یک متد متناظر دارد.

۱.۶ پیکربندی اولیه و سیستم گزارش‌دهی

```
1 class CodeGenerator(Transformer):
2     def __init__(self, output_dir):
3         self.output_dir = output_dir
4         self.plots_dir = os.path.join(output_dir, "plots")
5         os.makedirs(self.plots_dir, exist_ok=True)
6
7         self.current_var = None
8         self.report_lines = []
9         self.log_counter = 1
```

- output_dir: پوشه‌ی اصلی خروجی که شامل گزارش‌ها و فایل‌های تولیدی است.
- plots_dir: زیرپوشه‌ای برای ذخیره‌ی نمودارهای تولیدشده.
- current_var: آخرین دیتافریم فعال در فرآیند کامپایل.
- report_lines: لیستی از متن‌های گزارش که در نهایت در فایل report.txt ذخیره می‌شوند.
- log_counter: شماره‌گذاری ترتیبی عملیات‌ها.

۱.۱.۶ سیستم لاگ‌گیری

```
1 def add_log(self, title, body=""):
2     ts = datetime.now().strftime("%Y-%m-%d %H:%M:%S")
3     block = f"[{self.log_counter}] {title}\n{body}\n\nTimestamp: {ts}\n"
4     self.report_lines.append(block)
5     self.log_counter += 1
```

این متد بعد از هر عملیات:

- عنوان عملیات (مثلاً LOAD یا CLEAN)
- جزئیات آن عملیات
- زمان اجرا

را در گزارش نهایی ذخیره می‌کند.

۲.۶ مدیریت ساختار کلی AST

```
1 def statement(self, items):
2     return items[0]
3
4 def start(self, items):
5     return "\n".join(items)
```

- هر statement به یک قطعه کد پایتون تبدیل می‌شود.
- در متد start تمام دستورات به صورت متوالی کنار هم قرار می‌گیرند.

۳.۶ دستور LOAD: بارگذاری داده

```
1 def load_stmt(self, items):
2     file_path, var = items
3     self.current_var = var
```

در مرحله‌ی **compile-time** اگر فایل وجود داشته باشد، اطلاعات آماری آن استخراج می‌شود:

```
1 if os.path.exists(file_name):
2     df_tmp = pd.read_csv(file_name)
3     body += f"Rows: {len(df_tmp)}\nCols: {df_tmp.shape[1]}"
```

کد تولیدشده برای زمان اجرا:

```
1 if not os.path.exists(file_path):
2     raise FileNotFoundError("File not found")
3
4 if file_path.endswith(".csv"):
5     df = pd.read_csv(file_path)
6 elif file_path.endswith(".xlsx"):
7     df = pd.read_excel(file_path)
```

۴.۶ دستورات اطلاعاتی: HEAD و DESCRIBE

```
1 print(df.describe())
2 print(df.head(n))
```

این دستورات فقط خروجی تحلیلی تولید می‌کنند و دیتافریم را تغییر نمی‌دهند.

۵.۶ دستور SAVE و DUPLICATE

```
1 dest = source.copy()
```

برای ذخیره:

```
1 df.to_csv(path, index=False)
2 df.to_excel(path, index=False)
3 df.to_json(path, orient="records")
```

۶.۶ دستور CLEAN: پاکسازی داده

۱.۶.۶ تشخیص نوع عملیات

```
1 op_type = actual_op.data.upper()
2 params = [str(c) for c in actual_op.children]
```

۲.۶.۶ حذف داده‌های پرت (IQR)

$$IQR = Q_3 - Q_1$$

$$x < Q_1 - 1.5 \times IQR \quad \text{یا} \quad x > Q_3 + 1.5 \times IQR$$

پیاده‌سازی:

```
1 q1, q3 = df[col].quantile([0.25, 0.75])
2 iqr = q3 - q1
3 df = df[(df[col] >= q1-1.5*iqr) & (df[col] <= q3+1.5*iqr)]
```

۳.۶.۶ جایگزینی مقادیر خالی

```
۱ df[col] = df[col].fillna(df[col].mean())
۲ df[col] = df[col].fillna(df[col].mode()[0])
```

۷.۶ دستورات CALC: محاسبات آماری

```
۱ mean_col = df["col"].mean()
۲ std_col = df["col"].std()
```

مقادیر محاسبه شده در محیط اجرا ذخیره شده و بعداً در گزارش استفاده می شوند.

۸.۶ دستورات PLOT: تولید نمودار

```
۱ plt.figure(figsize=(10,6))
۲ df["col"].hist(bins=20)
۳ plt.savefig("plots/hist_col.png")
۴ plt.close()
```

ویژگی ها:

- پشتیبانی از نمودارهای Histogram، Line، Box Scatter و Heatmap
- ذخیره خودکار در پوشه plots
- عدم استفاده از plt.show() برای جلوگیری از بلاک شدن برنامه

۹.۶ دستورات FILTER

```
۱ df = df[df["age"] >= 18]
۲ df = df.query("age > 18 & salary < 5000")
```

۱۰.۶ LEVELING: دسته بندی بازه ای

```
۱ df["age_level"] = pd.cut(
۲     df["age"],
۳     bins=[18,30,60,inf],
۴     labels=["Teen","Adult","Senior"],
۵     right=False
۶ )
```

۱۱.۶ دستورات CRUD روی ستون ها

- pd.merge(df1, df2, on=key):MERGE
- df.eval(expr):CREATE_COL
- DROP_COL: حذف ستون
- RENAME: تغییر نام ستون
- CONVERT_TIME: تبدیل به datetime

۱۲.۶ NORMALIZE: نرمال سازی Min-Max

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

```
۱ df[col] = (df[col] - df[col].min()) / (df[col].max() - df[col].min())
```

```
corr_val = df["col1"].corr(df["col2"])
```

این مقدار به صورت عددی در خروجی چاپ و در گزارش ثبت می‌شود.

۷ نحوه استفاده از CodeGenerator و تولید کد نهایی

پس از تعریف گرامر، DSL، نگاشت دستورات و ساخت، کلاس CodeGenerator در مرحله‌ی نهایی کامپایل استفاده می‌شود تا درخت نحوی به کد اجرایی پایتون تبدیل گردد.

۱.۷ جریان کامل کامپایل

فرآیند کلی کامپایل به صورت زیر است:

۱. دریافت کد DSL (فارسی یا انگلیسی)
۲. تبدیل به DSL میانی (در صورت وجود Mapper)
۳. پارس کد با Lark و تولید AST
۴. عبور از AST به CodeGenerator
۵. تولید کد پایتون قابل اجرا
۶. اجرای کد تولیدشده و تولید خروجی‌ها

۲.۷ نحوه استفاده از CodeGenerator

در تابع اصلی کامپایل (مثلاً run_compiler):

```
۱ parser = Lark(grammar, parser="lalr")
۲ ast_tree = parser.parse(dsl_code)
۳
۴ generator = CodeGenerator(OUTPUT_DIR)
۵ python_code = generator.transform(ast_tree)
```

در این مرحله:

- ast_tree: درخت نحوی برنامه DSL
- generator: نمونه‌ای از کلاس CodeGenerator
- transform: متد اصلی Transformer که:
 - روی تمام نودهای AST پیمایش می‌کند
 - برای هر دستور DSL متد متناظر را صدا می‌زند
 - خروجی هر دستور را به صورت رشته کد پایتون جمع‌آوری می‌کند

۳.۷ خروجی متد transform

```
python_code = generator.transform(ast_tree)
```

نوع خروجی:

- یک string شامل کد کامل پایتون
- این کد شامل:
 - دستورات بارگذاری داده
 - پاکسازی
 - محاسبات آماری
 - فیلترها
 - تولید نمودار
 - عملیات CRUD

نکته مهم: کلاس CodeGenerator فایل تولید نمی‌کند، بلکه فقط کد را به صورت رشته بازمی‌گرداند.

۴.۷ تولید فایل generated_code.py

کد بازگشتی در فایل ذخیره می‌شود:

```
1 with open("generated_code.py", "w", encoding="utf-8") as f:
2     f.write("import pandas as pd\n")
3     f.write("import matplotlib.pyplot as plt\n")
4     f.write("import seaborn as sns\n")
5     f.write("import os\n\n")
6
7     f.write(f'OUTPUT_DIR = r"{OUTPUT_DIR}"\n')
8     f.write(f'PLOTS_DIR = "plots"\n')
9     f.write(f'PLOT_PATH = os.path.join(OUTPUT_DIR, PLOTS_DIR)\n\n')
10
11     f.write(f'os.makedirs(PLOT_PATH, exist_ok=True)\n\n')
12     f.write(python_code)
```

در نتیجه فایل generated_code.py:

- کاملاً مستقل است
- بدون نیاز به کامپایلر DSL قابل اجراست
- تمام خروجی‌ها را در پوشه‌ی output/ تولید می‌کند

۵.۷ اجرای کد تولیدشده

کد تولیدشده بلافاصله اجرا می‌شود:

```
1 env = {"pd": pd, "plt": plt, "os": os, "sns": sns}
2 exec(python_code, env)
```

- env: محیط اجرای کنترل‌شده
- تمام متغیرهای محاسبه‌شده (میانگین‌ها، انحراف معیارها و ...) در این دیکشنری ذخیره می‌شوند

۶.۷ تولید گزارش نهایی

پس از اجرای کد:

```
1 for line in generator.report_lines:
2     f.write(line.format(**env))
```

- report_lines: خروجی side-effect کلاس CodeGenerator
- شامل:

- لاگ تمام عملیات‌ها
- نتایج محاسبات آماری
- زمان اجرای هر دستور

۷.۷ خلاصه خروجی‌های CodeGenerator

مورد	نوع	توضیح
transform()	string	کد پایتون تولیدشده
report_lines	list[string]	گزارش عملیات‌ها
نمودارها	فایل	plots.png در پوشه plots
گزارش نهایی	فایل	report.txt
generated_code.py	فایل	اسکرپت مستقل پایتون

۸.۷ جمع‌بندی

کلاس `CodeGenerator` قلب کامپایلر DSL است که:

- AST را به کد اجرایی تبدیل می‌کند
- گزارش تحلیلی تولید می‌کند
- نمودارها و فایل‌های خروجی را مدیریت می‌کند
- خروجی نهایی قابل اجرا و مستقل تحویل می‌دهد

۸ پایپ‌لاین `run_compiler` و ارتباط با کد تولیدی

تابع `run_compiler` قلب اصلی کامپایلر فارسی-انگلیسی است و کل فرآیند `Translation` → `Parsing` → `Code Generation` → `Execution` → `Reporting` را به صورت یکپارچه مدیریت می‌کند.

```
1 def run_compiler(persian_code, capture_output=True):
2     import sys
3     import io
4
5     OUTPUT_DIR = "output"
6     PLOTS_DIR = "plots"
7     os.makedirs(OUTPUT_DIR, exist_ok=True)
8
9     if capture_output:
10         old_stdout = sys.stdout
11         old_stderr = sys.stderr
12         captured = io.StringIO()
13         sys.stdout = sys.stderr = captured
14
15     try:
16         mapper = PersianToDSLMapper()
17         dsl_code = mapper.translate(persian_code)
18
19         parser = Lark(grammar, parser="lalr")
20         ast_tree = parser.parse(dsl_code)
21
22         ast_to_dot(ast_tree, output_name="ast", output_dir=OUTPUT_DIR)
23
24         generator = CodeGenerator(OUTPUT_DIR)
25         python_code = generator.transform(ast_tree)
26
27         gen_path = os.path.join("./", "generated_code.py")
28         with open(gen_path, "w", encoding="utf-8") as f:
29             f.write("import pandas as pd\n"
30                    "import matplotlib.pyplot as plt\n"
31                    "import os\n"
32                    "import seaborn as sns\n"
33                    "import warnings\n"
34                    "warnings.filterwarnings('ignore')\n\n")
35
36             f.write(f'OUTPUT_DIR = r"{OUTPUT_DIR}"\n')
37             f.write(f'PLOTS_DIR = r"{PLOTS_DIR}"\n')
38             f.write(f'PLOT_PATH = os.path.join(OUTPUT_DIR, PLOTS_DIR)\n\n')
39
40             f.write('os.makedirs(OUTPUT_DIR, exist_ok=True)\n')
41             f.write('os.makedirs(PLOT_PATH, exist_ok=True)\n\n')
42
43             f.write(python_code)
44
45     env = {
46         "pd": pd,
47         "plt": plt,
48         "os": os,
49         "sns": sns,
```

```

۵۰     "OUTPUT_DIR": OUTPUT_DIR,
۵۱     "PLOTS_DIR": PLOTS_DIR,
۵۲     "PLOT_PATH": os.path.join(OUTPUT_DIR, PLOTS_DIR)
۵۳ }
۵۴
۵۵ exec(python_code, env)
۵۶
۵۷ report_path = os.path.join(OUTPUT_DIR, "report.txt")
۵۸ with open(report_path, "w", encoding="utf-8") as f:
۵۹     for line in generator.report_lines:
۶۰         try:
۶۱             f.write(line.format(**env) + "\n\n")
۶۲         except Exception:
۶۳             f.write(line + "\n\n")
۶۴
۶۵     return (True, captured.getvalue() if capture_output else "", "")
۶۶
۶۷ except Exception as e:
۶۸     import traceback
۶۹     return (False, captured.getvalue(), str(e))
۷۰
۷۱ finally:
۷۲     if capture_output:
۷۳         sys.stdout = old_stdout
۷۴         sys.stderr = old_stderr

```

تحلیل گام به گام پایپ لاین

۱. ایجاد پوشه های خروجی

پوشه ی output و زیرپوشه ی plots ساخته می شوند تا تمام خروجی ها (کد، نمودار، گزارش و (AST) در یک ساختار منظم ذخیره شوند.

۲. ترجمه DSL فارسی به DSL میانی

`:re.sub` با استفاده از `PersianToDSLMapper.translate`

- کلمات فارسی را به دستورات DSL انگلیسی تبدیل می کند
- خروجی آن یک رشته ی DSL استاندارد است

۳. Parsing و تولید AST

کد DSL توسط Lark تجزیه شده و یک درخت نحوی انتزاعی (AST) ساخته می شود که نمایش ساختاری برنامه است.

۴. تولید فایل AST

درخت AST به صورت گراف (dot/png) ذخیره می شود تا تحلیل ساختار برنامه ممکن باشد.

۵. تولید کد پایتون

کلاس `CodeGenerator`:

- AST را پیمایش می کند
- برای هر دستور، DSL کد پایتون معادل می سازد
- خروجی نهایی، رشته ی `python_code` است

۶. نوشتن فایل `generated_code.py`

این فایل شامل:

- import کتابخانه ها
- تعریف مسیرهای خروجی
- کد پایتون تولید شده از AST

۷. اجرای کد تولیدی

کد با دستور `exec(python_code, env)` اجرا می شود. دیکشنری `env` نقش محیط اجرای ایمن را دارد و شامل:

- `DataFrame` ها

- متغیرهای محاسبه‌شده
- مسیر نمودارها

۸. تولید گزارش نهایی

لیست `generator.report_lines` شامل قالب‌های متنی گزارش است. این خطوط با استفاده از:

```
line.format(**env)
```

به مقادیر واقعی محاسبه‌شده در زمان اجرا متصل می‌شوند.

۹ رابط کاربر (get_user_input) و تابع __main__

تابع `get_user_input` مسئول دریافت کد DSL از کاربر است و دو حالت اصلی برای ورودی ارائه می‌دهد:

- **ورودی دستی:** کاربر می‌تواند کد DSL را خط به خط در کنسول وارد کند. برای این منظور، تابع یک حلقه `while` باز می‌کند و تا زمانی که کاربر خط خالی وارد نکند، هر خط را به یک لیست اضافه می‌کند.
 - **بارگذاری از فایل متنی:** کاربر می‌تواند مسیر یک فایل `.txt` را وارد کند. تابع بررسی می‌کند که فایل وجود داشته باشد و سپس محتوای آن را با `encoding="utf-8"` می‌خواند.
- پس از دریافت کد، تابع آن را به صورت یک رشته واحد (`string`) بازمی‌گرداند.

در فایل اصلی، بخش زیر وظیفه اجرای برنامه را بر عهده دارد:

```
1 if __name__ == "__main__":
2     try:
3         user_code = get_user_input()
4         run_compiler(user_code)
5     except Exception as e:
6         print("Error:", e)
```

توضیح گام به گام:

۱. شرط: `if __name__ == "__main__":` تضمین می‌کند که این بخش تنها زمانی اجرا شود که اسکریپت به صورت مستقیم اجرا شده باشد و نه وقتی که به عنوان ماژول وارد (`import`) شود.
 ۲. `get_user_input()` فراخوانی می‌شود تا کد DSL از کاربر دریافت گردد.
 ۳. سپس `run_compiler(user_code)` برای پردازش و اجرای کد وارد شده استفاده می‌شود.
 ۴. در صورت بروز هرگونه خطا، با بلوک `try-except` آن خطا گرفته شده و پیغام مناسب به کاربر نمایش داده می‌شود.
- این ساختار ساده، کامپایلر را به یک ابزار CLI قدرتمند تبدیل می‌کند که هم انعطاف‌پذیری دریافت ورودی از کاربر را دارد و هم مدیریت خطاها را به شکل مناسبی انجام می‌دهد.

۱۰ توضیح سناریوی نمونه و نحوه اجرای آن

در این بخش، یک سناریوی کامل نوشته‌شده با زبان خاص منظوره فارسی (Persian DSL) معرفی و تحلیل می‌شود. این سناریو با هدف پردازش، تحلیل آماری و مصورسازی داده‌های آزمایشگاهی طراحی شده و توسط کامپایلر DSL فارسی پیاده‌سازی شده در این پروژه اجرا می‌گردد.

۱.۱۰ معرفی سناریو

کد زیر یک اسکریپت DSL فارسی است که عملیات مختلفی از جمله بارگذاری داده، پاکسازی، فیلتر، ایجاد ستون‌های محاسباتی، نرمال‌سازی، سطح‌بندی، رسم نمودار و ذخیره خروجی را انجام می‌دهد:

۲.۱۰ جریان اجرای سناریو در کامپایلر

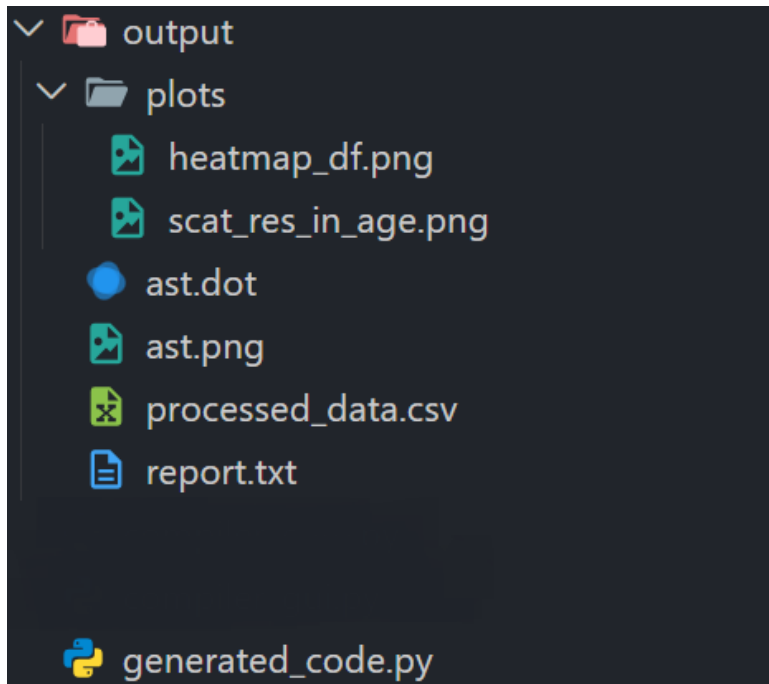
در ابتدا برنامه را اجرا میکنیم :

```
1 PS S:\FivethSemester\Compiler\lastProject> python LabDsl.py
2 =====
3 Persian / English DSL Compiler
4 =====
5 HINT: You may mix Persian and English DSL commands in the same script.
6 -----
7
8 Select input method:
9 1) Enter code manually
10 2) Load code from text file
11 Your choice (1/2): 2
12
```

سپس اجرا به شکل خواندن دستورات از روی فایل رو اجرا میکنیم :

```
1 Your choice (1/2): 2
2 Enter .txt file path: path.txt
3
```

نگاهی به فایل های ایجاد شده پس از اجرای دستورات :



فایل دستورات با نام path.txt حاوی دستور های زیر است به ترتیب :

- بگیر از "lab_data.csv" به نام df
- پاکسازی df : حذف تکراری
- پاکسازی df : حذف کلی خالی
- فیلتر df : وقتی $age < 15$
- ایجاد ستون df : $((cholesterol * 9.0) + (glucose * 1.0)) = result$
- تغییر نام df : result به sick
- نرمال سازی df : ستون sick
- نمودار df : نقشه حرارتی
- مرتب سازی df : sick صعودی

- ایجاد ستون `sick = res : df`
 - سطح بندی `df` : در `res` به `low:۳.۰۰ normal:۶.۰۰ heigh:۸.۰۰ danger`
 - نمودار `df` : پراکندگی `res` در `age`
 - ذخیره `df` : در `"processed_data.csv"`
- اجرای این سناریو در کامپایلر DSL فارسی شامل مراحل زیر است:

۱. نگاشت دستورات فارسی به DSL میانی در این مرحله، کلاس `PersianToDSLMapper` با استفاده از عبارات باقاعده (Regex)، دستورات فارسی را به معادل انگلیسی DSL تبدیل می کند (مانند `LOAD, CLEAN, FILTER, PLOT`).
۲. تحلیل نحوی (Parsing) کد DSL میانی با استفاده از گرامر تعریف شده در `Lark` تحلیل شده و یک درخت نحو انتزاعی (AST) تولید می شود.
۳. تولید AST و مصورسازی آن AST تولید شده به کمک تابع `ast_to_dot` به فایل `ast.dot` تبدیل می شود که امکان مشاهده ساختار برنامه را فراهم می کند.
۴. تولید کد پایتون کلاس `CodeGenerator` که از `Transformer` ارث بری می کند، AST را پیمایش کرده و کد معادل پایتون را تولید می نماید.
۵. اجرای کد و تولید خروجی ها کد پایتون تولید شده به صورت پویا اجرا شده و فایل ها، نمودارها و گزارش آماری ایجاد می شوند.

۳.۱۰ توضیح گام به گام عملیات سناریو

- ابتدا فایل `lab_data.csv` بارگذاری شده و در `DataFrame` با نام `df` ذخیره می شود.
- سطرهای تکراری حذف می شوند تا داده ها یکتا گردند.
- تمام سطرهایی که شامل مقادیر خالی هستند حذف می شوند.
- فقط رکوردهایی که مقدار ستون `age` آن ها بزرگ تر از ۱۵ است نگه داشته می شوند.
- یک ستون محاسباتی جدید با نام `result` بر اساس ترکیب وزنی `glucose` و `cholesterol` ایجاد می شود.
- نام ستون `result` به `sick` تغییر می یابد.
- ستون `sick` با استفاده از روش `Min-Max` نرمال سازی می شود.
- یک نقشه حرارتی (Heatmap) از همبستگی ویژگی های عددی رسم می گردد.
- داده ها بر اساس ستون `sick` به صورت صعودی مرتب می شوند.
- ستون جدیدی به نام `res` به عنوان کپی از `sick` ایجاد می شود.
- ستون `res` به سطوح معنایی `low, normal, heigh, danger` تقسیم بندی می گردد.
- نمودار پراکندگی بین `res` و `age` رسم می شود.
- در نهایت، داده های پردازش شده در فایل `processed_data.csv` ذخیره می شوند.

۴.۱۰ تحلیل نتایج اجرا بر اساس داده ی ورودی نمونه

۱.۴.۱۰ بررسی داده ی ورودی

داده ی اولیه از فایل `lab_data.csv` بارگذاری شده است. این فایل شامل ۱۰ رکورد و ۵ ستون به شرح زیر است:

- `id` : شناسه بیمار
- `gender` : جنسیت
- `age` : سن
- `glucose` : میزان قند خون
- `cholesterol` : میزان کلسترول

بررسی آماری اولیه نشان می‌دهد که:

- ستون‌های glucose و cholesterol شامل مقادیر خالی هستند.
- حداقل سن برابر ۱۰ بوده که در مراحل بعدی فیلتر می‌شود.
- همبستگی اولیه بین قند خون و کلسترول نسبتاً مثبت است.

۲.۴.۱۰ تأثیر پاکسازی و فیلتر داده‌ها

در مراحل اولیه پردازش، عملیات زیر انجام شده است:

- حذف سطرهای تکراری
- حذف کامل سطرهای دارای مقدار خالی
- فیلتر داده‌ها بر اساس شرط $age > 15$
- پس از اعمال این مراحل:
- تعداد رکوردها از ۱۰ به ۶ کاهش یافته است.
- تمامی ستون‌های عددی بدون مقدار خالی هستند.
- داده‌ها برای انجام تحلیل‌های آماری آماده شده‌اند.

۳.۴.۱۰ محاسبه شاخص سلامت و نرمال‌سازی

در این سناریو، یک ستون محاسباتی جدید با رابطه زیر ایجاد شده است:

$$sick = 0.1 \times glucose + 0.9 \times cholesterol$$

این شاخص، وزن بیشتری به کلسترول داده و نمایانگر یک معیار ترکیبی از وضعیت سلامت بیمار است. سپس ستون sick با استفاده از روش Min-Max نرمال‌سازی شده است، به طوری که مقادیر آن در بازه $[0, 1]$ قرار می‌گیرند. نتیجه نهایی نشان می‌دهد:

- کمترین مقدار sick برابر ۰
- بیشترین مقدار برابر ۱
- میانگین شاخص سلامت برابر ۳۹۲.۰

۴.۴.۱۰ سطح‌بندی معنایی بیماران

برای تفسیر ساده‌تر شاخص سلامت، ستون res به سطوح معنایی زیر تقسیم شده است:

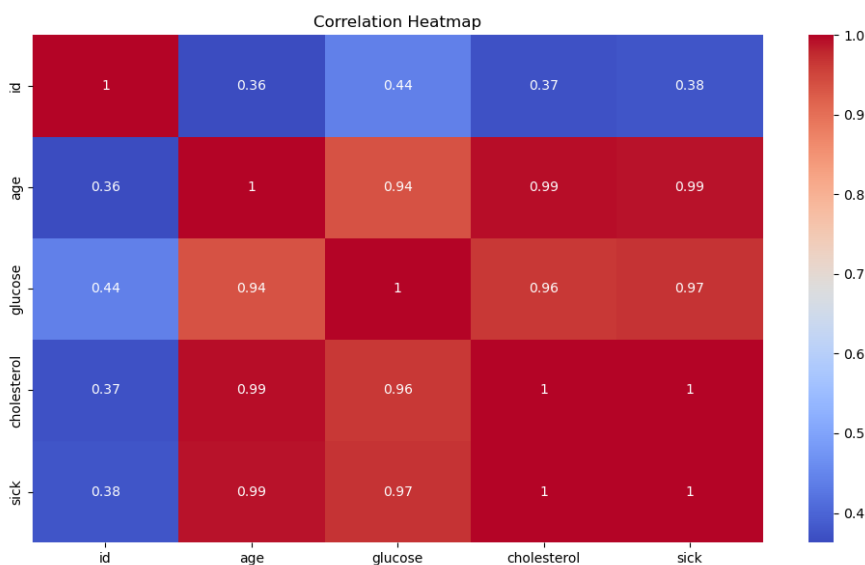
سطح سلامت	بازه مقدار
low	$[0, 0.3)$
normal	$[0.3, 0.6)$
heigh	$[0.6, 0.8)$
danger	$[0.8, 1.0]$

بر اساس این تقسیم‌بندی:

- بیماران جوان‌تر عمدتاً در سطح low قرار دارند.
- بیماران با سن بالاتر و کلسترول بیشتر در سطح danger مشاهده می‌شوند.

۵.۴.۱۰ تحلیل نقشه حرارتی همبستگی

شکل ۲ نقشه حرارتی همبستگی بین ویژگی‌های عددی را نمایش می‌دهد.

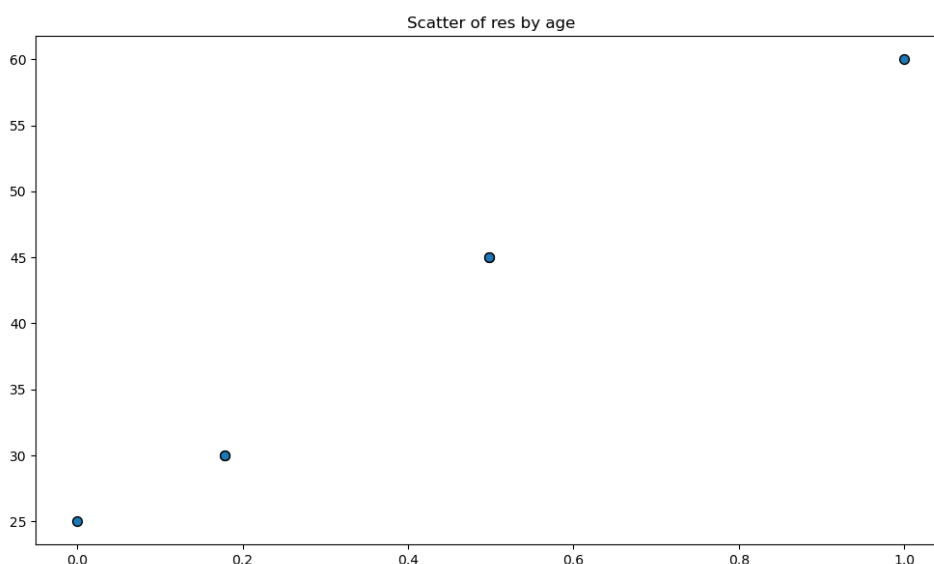


شکل ۲: نقشه حرارتی همبستگی ویژگی‌های عددی

این نمودار نشان می‌دهد که بین glucose و cholesterol همبستگی مثبتی وجود دارد که توجیه‌کننده استفاده از آن‌ها در شاخص سلامت است.

۶.۴.۱۰ نمودار پراکندگی سن و شاخص سلامت

نمودار پراکندگی شکل ۳ رابطه بین سن بیماران و شاخص سلامت نهایی را نمایش می‌دهد.



شکل ۳: نمودار پراکندگی سن و شاخص سلامت

بر اساس این نمودار، با افزایش سن، تمایل به افزایش شاخص بیماری مشاهده می‌شود، هرچند رابطه کاملاً خطی نیست.

۵.۱۰ فایل ریپورت

یکی از خروجی‌های ما فایل ریپورت است که در آن تمام کارهایی که انجام شده به ترتیب به شکل لاگ موجود است :

```

1 Laboratory Data Statistical Report
2 =====
3
4 [1] LOAD
5 Loaded file: lab_data.csv
6 Rows: 10
7 Cols: 5
8
9 Timestamp: 2026-01-31 23:22:39
10
11
12 [2] CLEAN - DROP_DUPLICATES
13 Params:
14
15 Timestamp: 2026-01-31 23:22:39
16
17
18 [3] CLEAN - DROP_ALL
19 Params:
20
21 Timestamp: 2026-01-31 23:22:39
22
23
24 [4] FILTER
25 Condition: age > 15
26
27 Timestamp: 2026-01-31 23:22:39
28
29
30 [5] CREATE_COL
31 Created column: result in df
32
33 Timestamp: 2026-01-31 23:22:39
34
35
36 [6] RENAME
37 Renamed column result to sick in df
38
39 Timestamp: 2026-01-31 23:22:39
40
41
42 [7] NORMALIZE
43 Normalized column: sick in df
44
45 Timestamp: 2026-01-31 23:22:39
46
47
48 [8] PLOT
49 Type: HEATMAP
50 Target: ALL
51 Group: ALL
52
53 Timestamp: 2026-01-31 23:22:39
54
55
56 [9] SORT
57 Sort by sick ))
58
59 Timestamp: 2026-01-31 23:22:39
60
61
62 [10] CREATE_COL
63 Created column: res in df
64
65 Timestamp: 2026-01-31 23:22:39

```



```

66
67
68 [11] LEVELING
69 Column: res
70 Levels: ['low', 'normal', 'heigh', 'danger']
71
72 Timestamp: 2026-01-31 23:22:39
73
74
75 [12] PLOT
76 Type: SCAT
77 Target: res
78 Group: age
79
80 Timestamp: 2026-01-31 23:22:39
81
82
83 [13] SAVE
84 Saved df to "processed_data.csv"
85
86 Timestamp: 2026-01-31 23:22:39
87

```

۶.۱۰ کد تولید شده

در این قسمت هم کد تولید شده که مد به زبان سطح پایین است آورده شده و این فایل با عنوان generated_code در فایل output موجود است :

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import os
4 import seaborn as sns
5
6
7 # --- Load Data ---
8 if not os.path.exists("lab_data.csv"):
9     raise FileNotFoundError("File not found")
10
11 if "lab_data.csv".endswith(".csv"):
12     df = pd.read_csv("lab_data.csv")
13 elif "lab_data.csv".endswith(".xlsx"):
14     df = pd.read_excel("lab_data.csv")
15 else:
16     raise ValueError("Only CSV or Excel files are supported")
17
18
19 # --- Cleaning: DROP_DUPLICATES ---
20 df = df.drop_duplicates()
21
22 # --- Cleaning: DROP_ALL ---
23 df = df.dropna()
24 df = df[df["age"] > 15]
25 df['result'] = df.eval('((0.1 * glucose) + (0.9 * cholesterol))')
26 df = df.rename(columns={"result": "sick"})
27
28 # --- Min-Max Normalization ---
29 col_min = df["sick"].min()
30 col_max = df["sick"].max()
31 df["sick"] = (df["sick"] - col_min) / (col_max - col_min)
32
33
34 # --- Plotting HEATMAP ---
35 plt.figure(figsize=(10, 6))
36 corr = df.select_dtypes(include=["number"]).corr()
37 sns.heatmap(corr, annot=True, cmap="coolwarm")

```

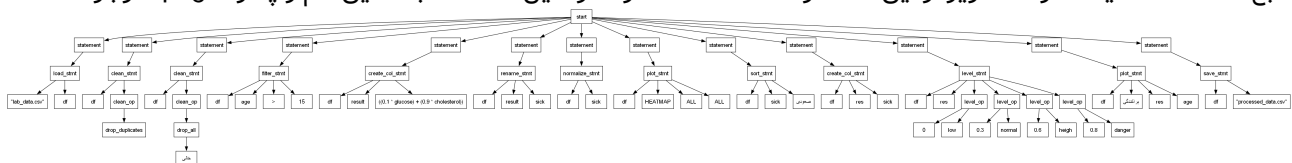
```

۳۸ plt.title("Correlation Heatmap")
۳۹ plt.tight_layout()
۴۰ plt.savefig(r"output\plots\heatmap_df.png")
۴۱ plt.show()
۴۲ plt.close()
۴۳
۴۴ df = df.sort_values(by="sick", ascending=True)
۴۵ df['res'] = df.eval('"'sick"'')
۴۶
۴۷ # --- Leveling Of res From df ---
۴۸ bins = [0.0, 0.3, 0.6, 0.8]
۴۹ bins.append(float('inf'))
۵۰ labels = ['low', 'normal', 'heigh', 'danger']
۵۱ df["res_level"] = pd.cut(df["res"], bins=bins, labels=labels, right=False)
۵۲
۵۳
۵۴ # --- Plotting SCAT ---
۵۵ plt.figure(figsize=(10, 6))
۵۶ plt.scatter(df["res"], df["age"], edgecolors="black", s=50)
۵۷ plt.title("Scatter of res by age")
۵۸ plt.tight_layout()
۵۹ plt.savefig(r"output\plots\scat_res_in_age.png")
۶۰ plt.show()
۶۱ plt.close()
۶۲
۶۳
۶۴ # --- Save DataFrame ---
۶۵ if "processed_data.csv".endswith(".csv"):
۶۶     df.to_csv(r"processed_data.csv", index=False)
۶۷ elif "processed_data.csv".endswith(".xlsx"):
۶۸     df.to_excel(r"processed_data.csv", index=False)
۶۹ elif "processed_data.csv".endswith(".json"):
۷۰     df.to_json(r"processed_data.csv", orient="records")
۷۱ else:
۷۲     raise ValueError("Supported formats: .csv, .xlsx, .json")
۷۳
۷۴

```

۷.۱۰ تصویر نهایی از دستورات

تابع ast_to_dot یک نمونه تصویر از این ast خواهد ساخت که در کنار فایل ast.dot با همین نام و پسوند png موجود است :



۸.۱۰ دیتاست نهایی ذخیره شده

دیتاست پروسس شده جدید به اضافه مقادیر محاسبه شده در فایل output به نام process_data.csv ذخیره خواهد شد. قبل پروسس در فایل lab_data :

	id	gender	age	glucose	cholesterol
1	1	Male	25	92	180
2	2	Female	30	105	190
3	3	Male	45	110	210
4	4	Female	50	98	
5	5	Male	10	115	220
6	6	Female	40		195
7	7	Male	60	130	240
8	8	Female	55	120	
9	9	Male	45	110	210
10	10	Female	30	105	190
11					

بعد از پروسس در فایل processed_data :

	id	gender	age	glucose	cholesterol	sick	res	res_level
1	1	Male	25	92.0	180.0	0.0	0.0	low
2	2	Female	30	105.0	190.0	0.17820069204152264	0.17820069204152264	low
3	10	Female	30	105.0	190.0	0.17820069204152264	0.17820069204152264	low
4	3	Male	45	110.0	210.0	0.4982698961937717	0.4982698961937717	normal
5	9	Male	45	110.0	210.0	0.4982698961937717	0.4982698961937717	normal
6	7	Male	60	130.0	240.0	1.0	1.0	danger
7								

۹.۱۰ جمع‌بندی نهایی سناریو

نتایج این سناریو نشان می‌دهد که DSL فارسی طراحی شده قادر است:

- داده‌های واقعی دارای نقص را پاکسازی کند
- شاخص‌های سلامت قابل تفسیر ایجاد نماید
- تحلیل‌های آماری و تصویری معنادار تولید کند
- خروجی‌های استاندارد قابل استفاده در تحلیل‌های بعدی ارائه دهد

این سطح از خودکارسازی، DSL را به ابزاری مناسب برای کاربران غیرمتخصص حوزه داده و علوم آزمایشگاهی تبدیل می‌کند.