
مستندات پروژه تحلیل مصرف آب

عنوان و نمای کلی پروژه 1.

- **عنوان پروژه:** داشبورد تحلیل مصرف آب و پیش‌پردازش داده
- **هدف:** این پروژه بر تحلیل جامع داده‌های مصرف آب برای سال‌های ۱۴۰۱، ۱۴۰۲ و ۱۴۰۳ متمرکز است. این پروژه شامل بارگذاری، پیش‌پردازش، پاکسازی داده‌ها، تحلیل داده‌های اکتشافی و توسعه یک داشبورد وب تعاملی برای بصری‌سازی معیارهای کلیدی و شناسایی مشکلات کیفیت داده و الگوهای مصرف می‌باشد.
- **اجزا:**
 - (Excel/CSV) فایل‌های داده خام
 - برای پیش‌پردازش داده و Jupyter (preprocessing and cleaning.ipynb) نوت‌بوک تحلیل اولیه
 - که یافته‌های کلیدی نوت‌بوک (preprocessing_summary.md) خلاصه تحلیل آماری را خلاصه می‌کند
 - که توابع قابل استفاده مجدد برای تحلیل و بصری‌سازی (analyzer.py) کد پایتون داده‌ها را فراهم می‌کند
 - که یک داشبورد وب تعاملی برای کاوش داده‌ها ایجاد (app.py) Streamlit اپلیکیشن می‌کند

توضیحات داده 2.

استفاده می‌کند که سال‌های ۱۴۰۱، Excel/CSV این پروژه از داده‌های مصرف آب ارائه شده در قالب حاوی (1401.xlsx، 1402.xlsx، 1403.xlsx) و ۱۴۰۳ را پوشش می‌دهد. فایل‌های داده خام در دو ردیف هستند MultiIndex اطلاعات مفصلی با ساختار هدر

داده‌های هر سال به نظر می‌رسد ترکیبی از موارد زیر است:

1. **اطلاعات ثابت مشتری/کنتور:** در ستون‌های ابتدایی قرار دارد. این شامل جزئیاتی مانند نام مشترک، کد اشتراک، نوع پروانه، موقعیت مکانی (استان، شهرستان، امور آب، محدوده مطالعاتی)، جزئیات کنتور (سریال، سایز، تاریخ نصب)، زمان آخرین اتصال و وضعیت قطع می‌باشد.

2. **داده‌های ماهانه مصرف:** پس از اطلاعات ثابت، ستون‌های ماهانه، داده‌های سری زمانی برای هر مشتری را ارائه می‌دهند. این ستون‌های ماهانه زیر عنوان ماه مربوطه (مثلاً '۱۴۰۱/۰۱'، '۱۴۰۱/۰۲') گروه‌بندی شده‌اند و چندین معیار را برای آن دوره زمانی مشخص ارائه می‌دهند.

ستون‌های کلیدی شناسایی شده (با نام‌های ترجمه شده استفاده شده در نوت‌بوک)

• **اطلاعات ثابت**

- Customer Name / نام مشترک
- Subscription Code / کد اشتراک
- License Type / نوع پروانه
 - Province / استان
 - County / شهرستان
- Water Authority / امور آب
- Study Area / محدوده مطالعاتی
- Meter Serial / سریال کنتور
 - Meter Size / سایز کنتور
- Installation Date / تاریخ نصب
- Last Connection Time / زمان آخرین اتصال
- Disconnect Status / قطع کن

- بر اساس موقعیت در بخش) Consumption in Period (m^3) / مصرف بازه (m^3) (اطلاعات ثابت، این به نظر می‌رسد مصرف کل در یک دوره بزرگتر باشد، نه ماهانه).
- مشابهاً، مجموع ساعت) Operating Hours in Period (h) / ساعت کارکرد بازه (h) (کارکرد در یک دوره بزرگتر).
- میانگین دبی در یک دوره) Average Flow Rate in Period (l/s) / میانگین دبی بازه (l/s) (بزرگتر).

• **YYYY/MM برای هر ماه مصرف ماهانه**

- (میانگین دبی برای ماه مشخص) Flow Rate (l/s) / دبی (l/s)
- (تعداد نمونه‌های دبی منفی در ماه) Number of Negative Flows / تعداد دبی منفی
- درصد نمونه‌های دبی منفی در) Percentage of Negative Flows / درصد دبی منفی (ماه).
- (مجموع ساعت کارکرد برای ماه) Operating Hours (h) / ساعت کارکرد
- (مصرف کل برای ماه) Consumption (m^3) / مصرف (m^3)
- تعداد نقاط داده معتبر) Number of Available Data Points / تعداد اطلاعات موجود (در ماه).
- تعداد مورد انتظار نقاط داده) Expected Number of Data Points / تعداد مورد انتظار (در ماه).
- درصد نقاط داده) Percentage of Available Data Points / درصد اطلاعات موجود (معتبر در ماه).

3. اهداف

اهداف این پروژه عبارتند از: * پاکسازی و پیش‌پردازش داده‌های خام از چند سال و قالب‌های مختلف. * شناسایی و مدیریت مشکلات کیفیت داده، به ویژه مقادیر اشتباه مصرف منفی و ساعت کارکرد. * تبدیل داده‌ها به قالبی مناسب برای تحلیل سری‌های زمانی. * انجام تحلیل داده‌های

اکتشافی برای درک الگوهای مصرف، توزیع‌ها و روابط بین متغیرها. * ارائه یک داشبورد تعاملی برای کاربران جهت کاوش آسان داده‌ها و بصری‌سازی‌ها. * آماده‌سازی زمینه برای کارهای آتی احتمالی، مانند تشخیص ناهنجاری یا مدل‌سازی پیش‌بینانه

4. فرآیند پیش‌پردازش و پاکسازی (preprocessing and cleaning.ipynb)

مراحل انجام شده برای آماده‌سازی داده‌های خام برای preprocessing and cleaning.ipynb نوت‌بوک. تحلیل را شرح می‌دهد.

1. بارگذاری داده‌های خام:

- `pd.read_excel(..., header=[0, 1])` با استفاده از (1401.xlsx، 1402.xlsx، 1403.xlsx) فایل‌های اکسل دو ردیفی به درستی مدیریت `MultiIndex` بارگذاری می‌شوند تا هدر `header=[0, 1]` شود.
- برای داده‌های کامل هر سال (df_1401, df_1402, df_1403) جداگانه `DataFrames` ایجاد می‌شود.
- ستون‌های اطلاعات ثابت و ستون‌های مصرف ماهانه به صورت مفهومی بر اساس محدوده‌های ایندکس ستون جدا می‌شوند.

2. بررسی و کاوش اولیه داده‌ها

- بارگذاری شده انجام می‌شود، از جمله نمایش `DataFrames` بررسی‌های اولیه روی و تولید (`.info()`) بررسی انواع داده و تعداد مقادیر غیرتهی (`.head()`)، چند ردیف اول برای ستون‌های عددی (`.describe()`) آمار توصیفی.
- بررسی می‌شود `DataFrame` شکل هر
- مانند) شناسایی و شمارش می‌شوند `.isnull().sum()` مقادیر گم‌شده با استفاده از این بلافاصله ستون‌هایی با مقادیر گم‌شده را مشخص (`df_info`) سلول ۳۵ برای می‌کند.

- بصری‌سازی‌ها (هیستوگرام، باکس‌پلات، ماتریس همبستگی) برای درک توزیع ویژگی‌های عددی و شناسایی مشکلات احتمالی مانند ناهنجاری‌ها استفاده می‌شوند. برای بصری‌سازی الگوهای داده گمشده استفاده می‌شود missingno کتابخانه

3. Long فرمت) بازسازی ساختار داده :

- است که برای Long به قالب Wide یک مرحله حیاتی، تبدیل داده‌های ماهانه از قالب تحلیل سری‌های زمانی و رسم نمودار در طول ماه‌ها مناسب‌تر است.
 - ستون کد اشتراک استخراج می‌شود.
 - شناسایی می‌شوند MultiIndex ستون‌های ماهانه با استفاده از سطوح
 - stack یک حلقه از طریق ماه‌ها تکرار می‌شود و معیارهای ماهانه برای هر مشتری می‌شود.
- دارای ردیف‌هایی است که هر ترکیب مشتری-ماه را (df_long) حاصل DataFrame و معیارهای 'month'، 'Subscription Code' نشان می‌دهند، با ستون‌هایی برای و 'Consumption (m³)'، 'Operating Hours (h)'، 'Flow Rate (l/s)' ماهانه این فرآیند برای ۱۴۰۱ (سلول‌های ۲۴، ۳۱)، ۱۴۰۳ (سلول ۲۵، ۲۶) و ۱۴۰۲ (سلول ۲۸). (غیره نشان داده شده است
- را تأیید می‌کند و حجم داده‌های گمشده در Long فرمت df_long روی info().
 - زیادی را برای NaNs مانند سلول ۳۱ (که) ستون‌های معیارهای ماهانه را نشان می‌دهد (نشان می‌دهد Flow Rate، Operating Hours، Consumption).

4. تغییر نام ستون‌ها :

- برای نگاشت نام ستون‌های فارسی به نام‌های column_translations یک دیکشنری انگلیسی توصیفی‌تر تعریف می‌شود (سلول ۳۵)
- با استفاده از این دیکشنری تغییر نام داده می‌شوند (سلول‌های ۳۸، df_long و df_info
 - (۳۹). توجه داشته باشید که یک عدم تطابق جزئی وجود دارد که در آن 'کد اشتراک' به stacking در دیکشنری نگاشت شده است، اما پس از فرآیند 'Subscription Code'

‘subscription code’ و تغییر نام در برخی سلول‌ها (مانند سلول ۲۳) به صورت ظاهر می‌شود. به نظر می‌رسد این یک اثر جانبی df_long در (حروف کوچک) و مرحله تغییر نام بعدی باشد. برای حفظ یکپارچگی در مستندات، از نام reset_index استفاده می‌شود مگر اینکه به صراحت به نام ‘Subscription Code’ استاندارد اشاره شود df_long موقت در خروجی

5. مدیریت داده‌های اشتباه (مقادیر منفی)

- بر اساس دانش دامنه، مصرف آب و ساعت کارکرد منفی نامعتبر هستند
- (df_info_1402, df_info_1403) اطلاعاتی برای ۱۴۰۲ و ۱۴۰۳ DataFrames ردیف‌های در منفی است، ‘Consumption in Period (m³)’ که در آن‌ها (df_info_1403) سلول‌های ۹۲ و ۹۳ خروجی‌هایی برای ردیف‌های مصرف منفی در) شناسایی می‌شوند (اطلاعاتی ۱۴۰۲ و ۱۴۰۳ نشان می‌دهند DataFrames).
- متناظر با این رکوردهای اشتباه در لیست‌هایی ‘Subscription Code’ مقادیر (customers_to_drop_1402, customers_to_drop_1403) جمع‌آوری می‌شوند
- (df_info_clean_1402, df_info_clean_1403) اطلاعاتی پاکسازی شده DataFrames آن‌ها در این لیست‌ها قرار دارد، ‘Subscription Code’ با فیلتر کردن ردیف‌هایی که ایجاد می‌شوند (سلول ۹۴). این مورد جدی‌ترین مشکل کیفیت داده در اطلاعات خلاصه را برطرف می‌کند.
- و ‘Consumption (m³)’ برای (df_long) مقادیر منفی در داده‌های ماهانه مصرف شناسایی می‌شوند (سلول‌های ۳۳، ۳۴). اگرچه این موارد به ‘Operating Hours (h)’ عنوان ناهنجاری لیست شده‌اند، اما نوت‌بوک آن‌ها را شناسایی می‌کند ولی در حذف نمی‌کند قبل از ذخیره داده‌ها df_long اسنیپت‌های ارائه شده صراحتاً آن‌ها را از ایده‌های تحلیل اشاره به بررسی/مدیریت این موارد دارند Long در قالب

6. ذخیره داده‌های پاکسازی شده/پردازش شده

- DataFrames (df_info_clean_1402, df_info_clean_1403) اطلاعاتی پاکسازی شده CSV به صورت فایل‌های (1402_clean_info.csv, 1403_clean_info.csv) (سلول ۹۶).
- DataFrames Long (df_long ۱۴۰۱، برای df_long_1402, df_long_1403) ذخیره می‌شوند CSV به صورت فایل‌های (long_usage1401.csv, long_usage1402.csv, long_usage1403.csv) (سلول‌های ۱۰۰، ۱۰۱، ۱۰۲).

و یافته‌های کلیدی (EDA) تحلیل داده‌های اکتشافی 5.

نوت‌بوک شامل مراحل مختلفی برای درک ویژگی‌های داده است:

1. بینشی در مورد تمایل df_long **آمار توصیفی کلی**: خلاصه آمار برای ستون‌های عددی در مرکزی، پراکندگی و دامنه معیارهای ماهانه ارائه می‌دهد (سلول ۳۰). به ویژه، این نشان می‌دهد که میانگین مصرف به دلیل ناهنجاری‌های منفی شدید، منفی است.
2. و محاسبه 'Subscription Code' بر اساس df_long **آمار در سطح مشتری**: گروه‌بندی میانگین، بینشی در مورد رفتار میانگین هر مشتری در طول سال ارائه می‌دهد (سلول ۱۱۷).
3. 'Percentage of' و 'Number of Negative Flows' **معیارهای کیفیت داده**: تحلیل به اندازه‌گیری کامل بودن داده و نرخ خطا در ترکیب‌های مشتری- 'Available Data Points' ماه کمک می‌کند (سلول‌های ۴۰، ۴۱). درصد قابل توجهی از رکوردها دسترسی به داده پایینی دارند (5۵).
4. روابط را نشان df_long **تحلیل همبستگی**: یک ماتریس همبستگی برای ستون‌های عددی در و 'Operating Hours (h)' می‌دهد (سلول ۴۲). یک همبستگی منفی قوی بین مشاهده می‌شود که غیرمعمول است و یک پرچم برای مشکل کیفیت 'Consumption (m³)' داده است که نیاز به بررسی بیشتر دارد.

5. **شناسایی موارد پرت:** مشتریان برتر بر اساس میانگین مصرف ماهانه و تعداد دبی منفی شناسایی می‌شوند (سلول‌های ۴۳، ۴۵)، که مناطق احتمالی برای بررسی عمیق‌تر را برجسته می‌کنند.

6. شناسایی ناهنجاری‌ها

- df_long ناهنجاری‌های واضح مانند مصرف/ساعت کارکرد منفی به صراحت در شناسایی می‌شوند (سلول ۴۸)
- با حذف 'Consumption (m³)' برای تشخیص موارد پرت احتمالی در IQR روش اعمال می‌شود (سلول ۴۹). تعداد df_long در (مقادیر منفی برای استحکام محاسبه قابل توجهی موارد پرت یافت می‌شود)
- برای تولید بصری‌سازی برای تشخیص InteractiveDataAnalyzer از متدهای کلاس parallel، همبستگی heatmap توزیع (هیستوگرام) و روابط، (Z امتیاز) ناهنجاری استفاده می‌شود (سلول‌های ۹۸، ۹۹، ۱۰۱). این df_info_clean روی coordinates Consumption تحلیل بصری، کجی توزیع‌ها و وجود موارد پرت در معیارهایی مانند را تأیید می‌کند 'Average Flow Rate in Period (l/s)' و 'in Period (m³)'

هستند که نیاز به MultiIndex یافته‌های کلیدی و نکات کیفیت داده: * داده‌ها دارای ساختار برای تحلیل سری‌های زمانی ضروری است. Long بارگذاری و پردازش خاصی دارد. * تبدیل به فرمت داده‌های گمشده در ستون‌های معیارهای ماهانه رایج هستند. * مقادیر منفی در مصرف و ساعت کارکرد وجود دارند و به عنوان داده‌های اشتباه در نظر گرفته می‌شوند، به ویژه در داده‌های اطلاعاتی که رکوردهای مشتری متناظر حذف می‌شوند. * یک همبستگی منفی قوی و غیرمعمول بین ساعت کارکرد ماهانه و مصرف وجود دارد، که نشان‌دهنده مشکلات احتمالی داده یا الگوهای عملیاتی غیرمعمول است. * تحلیل بصری، توزیع‌های بسیار کج و موارد پرت در معیارهای عددی کلیدی را تأیید می‌کند.

6. داشبورد تعاملی (app.py و analyzer.py)

پیاده‌سازی می‌کند و از Streamlit یک داشبورد تحلیل تعاملی را با استفاده از فریم‌ورک app.py فایل بهره می‌برد analyzer.py تعریف شده در InteractiveDataAnalyzer کلاس

را تعریف می‌کند. این کلاس توابع InteractiveDataAnalyzer این فایل پایتون کلاس **analyzer.py**: پانداس کپسوله می‌کند. متدهای آن شامل DataFrame مختلفی را برای تحلیل و بصری‌سازی یک موارد زیر هستند: * فیلتر کردن داده‌ها بر اساس مقدار، محدوده یا بازه تاریخی. * تولید نمودارهای هیستوگرام، باکس‌پلات، نمودار پراکندگی، نمودار خطی، تشخیص (Plotly) تعاملی با استفاده از همبستگی و نمودار خاص مصرف بر اساس نوع heatmap، ناهنجاری، نمودار دایره‌ای، نمودار میله‌ای * (و خلاصه تفصیلی شامل کجی و کشیدگی (describe(). معادل) انجام خلاصه‌های آماری * (پروانه ارائه * Z-score. نمایش تعداد مقادیر گم‌شده. * پیاده‌سازی تشخیص ناهنجاری با استفاده از روش با استفاده (scale_column, normalize_all_numeric) قابلیت‌های مقیاس‌بندی/نرمال‌سازی داده‌ها scikit-learn. از پیش‌پردازشگرهای

را برای انجام تحلیل و تولید Streamlit این کلاس منطق اصلی مورد استفاده توسط اپلیکیشن بصری‌سازی به صورت پویا بر اساس انتخاب کاربر فراهم می‌کند.

به عنوان رابط کاربری تعاملی برای پروژه عمل می‌کند. * به کاربران Streamlit اپلیکیشن **app.py**: را بارگذاری کنند. * پس از بارگذاری فایل‌ها، یک نوار CSV امکان می‌دهد یک یا چند فایل اکسل یا کناری ظاهر می‌شود که به کاربران امکان انتخاب یک مجموعه داده از فایل‌های بارگذاری شده را پانداس بارگذاری می‌شود. * یک شیء DataFrame می‌دهد. * فایل انتخاب شده در یک بارگذاری شده نمونه‌سازی می‌شود. * نوار کناری DataFrame با InteractiveDataAnalyzer گزینه‌های پیکربندی را ارائه می‌دهد: * انتخاب یک ستون عددی. * انتخاب یک ستون دسته‌ای. * 'None', 'Histogram', 'Boxplot', 'Line Chart', 'Outlier Detection', 'Pie Chart', 'Bar Chart', 'Correlation Heatmap'. * بر اساس نوع نمودار و ستون‌های انتخاب شده، اپلیکیشن متد مربوطه را از شیء

نمایش `st.plotly_chart` را با استفاده از `Plotly` فراخوانی کرده و نمودار `InteractiveDataAnalyzer` می‌دهد. * اپلیکیشن همچنین نمایش می‌دهد: * خلاصه آماری برای ستون عددی انتخاب شده جدول مقادیر گمشده در هر ستون * (`analyzer.statistical_summary`). بخشی برای پیش‌پردازش که به کاربران امکان انتخاب یک * (`analyzer.show_missing_values`). و اعمال آن بر روی یک ستون عددی انتخاب ('standard', 'minmax', 'robust') روش مقیاس‌بندی نمایش می‌دهد. * `DataFrame` را می‌دهد و نتیجه را در یک (`analyzer.scale_column`) شده مدیریت خطا برای بارگذاری فایل گنجانده شده است.

این ساختار به کاربران امکان می‌دهد به صورت پویا جنبه‌های مختلف داده را کاوش کرده، توزیع‌ها و روابط مختلف را بصری‌سازی کرده و به سرعت به آمار کلیدی و اطلاعات کیفیت داده از طریق یک رابط کاربری دوستانه دسترسی پیدا کنند. این برنامه از داده‌های پاکسازی شده تولید شده توسط پاکسازی شده را بارگذاری می‌کند یا خود برنامه داده‌های CSV با فرض اینکه کاربر فایل‌های) نوت‌بوک برای پشتیبانی از بصری‌سازی‌ها و ویژگی‌های خود (پاکسازی شده را به صورت داخلی بارگذاری می‌کند استفاده می‌کند.

اجرای پروژه 7.

و تعامل با داده‌ها: 1. اطمینان حاصل کنید که پایتون و کتابخانه‌های `Streamlit` برای اجرای اپلیکیشن (`streamlit`, `pandas`, `plotly`, `seaborn`, `matplotlib`, `scipy`, `scikit-learn`) لازم نصب شده‌اند `bash`: نصب کنید `pip` می‌توانید آن‌ها را با استفاده از. (در صورت نیاز برای نوت‌بوک `missingno` `pip install streamlit pandas plotly seaborn matplotlib scipy scikit-learn missingno` در یک دایرکتوری `app.py` و `analyzer.py` کدهای ارائه شده را در فایل‌هایی با نام‌های `openpyxl` 2. و `long_usage1401.csv`, `1402_clean_info.csv` مانند) ذخیره کنید. 3. فایل‌های داده پاکسازی شده را در مکانی (اصلی، بسته به اینکه قصد دارید برنامه چه فایل‌هایی را بارگذاری کند `.xlsx` - غیره - یا فایل‌های `Command` که توسط اسکریپت قابل دسترسی باشد (مثلاً همان دایرکتوری) قرار دهید. 4. ترمینال یا خود را باز کنید. 5. به دایرکتوری که فایل‌ها را در آن ذخیره کرده‌اید بروید. 6. اپلیکیشن `Prompt` مرورگر وب شما به طور 7. `bash streamlit run app.py` را با دستور زیر اجرا کنید `Streamlit`

را نمایش می‌دهد. از قسمت بارگذاری فایل برای آپلود فایل‌های Streamlit خودکار باز شده و داشبورد داده خود و از نوار کناری برای تعامل با گزینه‌های تحلیل استفاده کنید.

کارهای آتی احتمالی 8.

بر اساس وضعیت فعلی پروژه و یافته‌ها، کارهای آتی می‌تواند شامل موارد زیر باشد:

- **پاکسازی داده قوی‌تر**

- پیاده‌سازی مدیریت پیچیده‌تر برای مصرف/ساعت کارکرد منفی در داده‌های ماهانه (imputation) فراتر از صرفاً شناسایی. این می‌تواند شامل جایگزینی (df_long)، یا پرچم‌گذاری برای تحلیل جداگانه بر اساس تخصیص دامنه (capping) محدود کردن باشد.

- توسعه استراتژی‌هایی برای مدیریت مقادیر گم‌شده در معیارهای ماهانه (مثلاً روش‌های جایگزینی خاص سری‌های زمانی)

- **تشخیص ناهنجاری پیشرفته**

- پیاده‌سازی روش‌های تشخیص ناهنجاری خاص سری‌های زمانی (مثلاً بر اساس تجزیه فصلی، آمار متحرک یا الگوریتم‌های پیشرفته‌تر) برای شناسایی الگوهای مصرف غیرمعمول در طول زمان برای مشتریان منفرد.
- کاوش تکنیک‌های تشخیص ناهنجاری چند متغیره

- **(Feature Engineering) مهندسی ویژگی**

- استخراج ویژگی‌های جدید از داده‌های سری زمانی (مانند تغییر مصرف ماهانه، شاخص‌های فصلی، معیارهای بهره‌وری عملیاتی)
- Long اطلاعاتی و DataFrames ترکیب موثر ویژگی‌ها از

- **تحلیل و مدل‌سازی سری‌های زمانی**

- تحلیل روندها و فصلی بودن مصرف در سطح کلی و برای بخش‌های خاص مشتریان
- توسعه مدل‌هایی برای پیش‌بینی مصرف آب در آینده

- ساخت مدل‌هایی برای خوشه‌بندی مشتریان بر اساس الگوهای مصرف آن‌ها

- **توسعه داشبورد**

- اضافه کردن بصری‌سازی‌های تعاملی و گزینه‌های فیلتر بیشتر
- ادغام نتایج مدل‌های خوشه‌بندی یا پیش‌بینی در داشبورد
- بهبود مدیریت و نمایش مقادیر منفی یا ناهنجاری‌ها، شاید با اجازه دادن به کاربران برای تغییر وضعیت نمایش آن‌ها یا مشاهده خلاصه‌هایی به طور خاص برای رکوردهای دارای مشکل
- در داخل برنامه که معیارها و (tooltips) اضافه کردن مستندات واضح یا توضیحات بصری‌سازی‌ها را توضیح دهند

نتیجه‌گیری 9.

این پروژه پایه محکمی برای تحلیل داده‌های مصرف آب فراهم می‌کند. نوت‌بوک با موفقیت به پاکسازی اولیه داده‌ها و بازسازی ساختار پرداخته و مشکلات کلیدی کیفیت داده مانند مصرف منفی و InteractiveDataAnalyzer همبستگی‌های غیرمنتظره را شناسایی کرده است. توسعه کلاس ابزار ارزشمندی را برای کاوش و بصری‌سازی تعاملی ایجاد می‌کند و داده‌ها را Streamlit اپلیکیشن مناطق حیاتی (مقادیر منفی، EDA قابل دسترس‌تر و قابل فهم‌تر می‌سازد. یافته‌های تحلیل ناهنجاری‌های همبستگی، داده‌های گمشده) را برجسته می‌کنند که نیاز به بررسی بیشتر برای ساخت مدل‌های قابل اعتماد یا استنتاج قطعی در مورد الگوهای مصرف آب دارند
