

Absolutely—let’s zoom in on **Phase 1: Triage** and build it into a **thorough, repeatable system** that helps you quickly understand the shape of your dataset without getting stuck in the weeds.

Phase 1: Triage – “Map the Landscape”

Goal:

To quickly scan **every feature** in the dataset and assess:

- What type of data it is
 - Whether it’s clean or problematic
 - Whether it’s likely to be **useful, redundant, or ignorable**
 - Where to **dig deeper** in later phases
-

What You’re Looking For

Aspect	Why It Matters
Data type (numeric, categorical, date, etc.)	Guides what kind of analysis is appropriate
Cardinality (how many unique values)	Helps identify IDs, enums, flags, junk fields
Missing values	Tells you about data quality and where you’ll need imputation
Basic range or distribution info	Alerts you to outliers or scaling needs
Constant values or near-constant	Waste of model space—drop or deprioritize
Inconsistent or strange values	Highlights cleaning needs (e.g., typos, casing issues)

Triage Checklist (Per Feature)

For every column, answer the following:

1. **What is the column name?**
2. **What type of data is it?**
 - Numeric
 - Categorical
 - Date/Time
 - Boolean
 - Text
 - Identifier

3. **How many unique values?**
 - Constant (1)
 - Low (<10)
 - Medium (10–50)
 - High (>50)
 - Very High (>100 or matches row count)
4. **What % of values are missing?**
5. **Is the data usable at a glance?**
 - Clean / Dirty / Ambiguous
6. **Quick action label:**
 - Keep
 - Drop
 - Needs cleaning
 - Investigate
 - Derive features

Triage Tracker Table (Markdown or Spreadsheet)

Feature Name	Type	Unique	Missing (%)	Example Values	Notes / Red Flags	Quick Action
customer_id	Categorical	2500	0%	C123, C124, ...	Unique ID, no use in modeling	Drop
region	Categorical	6	2.5%	North, East, ...	Slight imbalance	Keep
purchase_date	Date	700	0%	2022-01-01 → 2023-01-01	Time-based, could extract features	Derive
units_sold	Numeric	1,345	0%	1, 2, 5, 10, 500	Right-skewed, possible outliers	Keep
email_opt_in	Boolean	2	0.1%	Yes, No	Clean binary	Keep
discount_rate	Numeric	1	100%		Always null	Drop
notes	Text	870	25%	Long-form text	Probably not usable in current model	Drop or NLP later

Tips for Efficiency

- You don't need to run detailed stats—use `df.describe(include='all')` and basic summaries.
 - For cardinality, try: `df[column].nunique()`
 - For missing %: `(df[column].isna().mean() * 100).round(1)`
 - For type: `df[column].dtype`
 - Set up a **Jupyter Notebook cell** that loops through each column and logs the triage table.
-

Optional Notebook Cell Template (Python)

```
for col in df.columns:
    print(f"Feature: {col}")
    print(f" - Type: {df[col].dtype}")
    print(f" - Unique Values: {df[col].nunique()}")
    print(f" - Missing %: {df[col].isna().mean() * 100:.2f}%")
    print(f" - Sample Values: {df[col].unique()[:5]}")
    print("-" * 50)
```

What You End Phase 1 With

- A **shortlist** of:
 - Columns to drop
 - Columns to clean
 - Columns that are priority targets
- A quick **qualitative map** of your data
- A **foundation** for deeper analysis in Phase 2

Absolutely. Phase 2 is where you **move beyond quick scans and start uncovering meaningful insights** in your data. This is where the EDA becomes thoughtful and strategic—but still manageable.

Phase 2: Prioritize — “Focus on What Matters”

Purpose:

Now that you’ve mapped out all features in Phase 1, the goal here is to **go deeper on high-priority features** to understand their behavior and potential relevance.

You’re answering: - What patterns are emerging? - What variables may drive outcomes? - What needs transformation or cleanup? - Where might business value hide?

Step-by-Step Workflow for Phase 2

1. Identify High-Priority Features

Choose variables to focus on using these criteria:

Priority Criteria	Why it Matters
Strong domain relevance	Business-critical fields
Strong statistical signal	Skewed, highly variable, high correlation
Data quality risk	Messy fields with missing or inconsistent values
Known to drive outcomes	From SMEs or prior models
Unexpected behavior	Surprises or anomalies worth digging into

Tip: Don’t try to analyze every column—just the ones that might affect insights, modeling, or decisions.

2. Explore in Detail

Here's what to look for by variable type:

A. Numeric Variables

Question	What to Check	Why It Matters
What's the distribution?	Histogram, boxplot	Skew = transformation?
Are there outliers?	Min/max vs IQR	May distort model
Central tendency?	Mean vs median	Skewed or symmetric?
Variance?	Std dev, CV	Useful or constant?
Correlation?	With target or other variables	Predictor signal or multicollinearity

Example Note: `> unit_price` is right-skewed with a long tail of premium products. Median is \$18, mean is \$27. Strongly correlated (0.62) with total revenue. Suggest log transform or segmentation by pricing tier.

B. Categorical Variables

Question	What to Check	Why It Matters
How many categories?	Cardinality	Too many? May need grouping
Are some categories rare?	Frequency table	Combine or drop
Is the distribution balanced?	Value counts	Model bias risk if not
Correlation with target?	Group target means	Driver of outcome?
Clean values?	Inconsistent labels	May need standardization

Example Note: `> payment_type` has 6 values, but 2 are rare. "CC" and "Credit Card" both used—needs consolidation. Credit Card users have 1.8x higher return rate. Keep and encode.

C. Date/Time Variables

Question	What to Check	Why It Matters
Time range?	Min/max dates	Any gaps or future dates?
Trends over time?	Line plots	Sales growing? Dropping?
Seasonality?	Monthly/weekly cycles	Business rhythm? Forecasting signal?
Feature potential?	Extract day, month, week, etc.	Create useful predictors

Example Note: `> signup_date` ranges from Jan 2021 to Mar 2023. Weekly cycles are visible (spikes on Mondays). Could derive weekday or signup cohort features.

D. Target Variable Relationship

Question	What to Check	Why It Matters
Grouped summary?	Mean target per category or bin	Segment insight
Scatter or trend?	Numeric vs target	Slope or clustering?
Is the relationship stable over time?	Group by date	Drift or concept change?

Example Note: `> customer_age` shows a U-shaped relationship with churn: younger (<25) and older (>60) users churn more. Possible segmentation by age group.

3. Use Structured Notetaking for Each Feature

Use a consistent format for medium-depth observations.

Example Markdown Format:

```
### Feature: product_rating
- Type: Numeric (1-5 scale)
- Observation: Left-skewed, most ratings are 5; 5% are 1-star
- Interpretation: Possible review bias, most customers happy or silent
- Business Insight: Need to monitor low-rated items for customer retention
- Next Step: Engineer binary flag for 1-star ratings; possible NLP on reviews
```

Common Transformations Identified in Phase 2

Issue Found	Recommended Fix
Right-skewed variable	Log transform
Rare categorical levels	Group or drop
Flat or constant column	Drop
Text field with tags or codes	Extract structured features
Seasonality in date field	Create lag/lead or cyclic features

Output of Phase 2: Prioritized Feature Notes

You now have: - A medium-detailed set of notes on 10-25 variables - Key features tagged for transformation, modeling, or feature engineering - A shortlist of EDA visuals and tables to include in reports - Clear rationale for what's going into your model or dashboard

Absolutely. **Phase 3: Synthesize** is where your EDA transitions from technical exploration to **insight communication**. This is where your findings become actionable and valuable to others—whether that's a stakeholder, a client, or your future self. It's also the hardest phase to master because it requires **connecting patterns to purpose**.

Phase 3: Synthesize – “Turn Observations into Insights”

Goal:

To translate statistical or structural findings from your EDA into:

- Business-relevant interpretations
- Hypotheses and explanations
- Actionable recommendations
- Key messages for reporting

This is where your technical observations get **framed in the language of decisions**.

What You're Doing in This Phase:

Task	Purpose
Interpreting patterns	Why does this happen? What could explain it?
Linking to outcomes	How does this affect the target or goal?
Translating into domain terms	What does this mean in the real world?
Recommending action	What should the business/data team do with this?
Structuring for communication	How will I summarize this in a report or slide?

Four-Part Synthesis Framework (O-I-I-R)

Use this framework per variable or theme:

Element	Explanation	Example
Observation	What did the data show?	"Customer churn rate is 18% and higher in Segment C."
Interpretation	Why might that be happening?	"Segment C includes customers with fewer than 3 support interactions."
Impact (Business Meaning)	Why does this matter to the org?	"These customers may be disengaged or underserved, leading to preventable churn."
Recommendation	What do we do about it?	"Launch a targeted outreach campaign or improve onboarding for Segment C."

Examples by Feature Type

Numeric Variable Example

Feature: ``delivery_delay_days``

- Observation: Right-skewed with a median of 3 days but a long tail up to 21 days.
 - Interpretation: Most orders ship on time, but a few cause extreme delays, possibly from specific suppliers.
 - Business Impact: Outliers are inflating average delivery KPIs and could harm customer trust.
 - Recommendation: Report median instead of mean. Investigate long-tail delays and consider SLA renegotiations.
-

Categorical Variable Example

Feature: ``product_category``

- Observation: Accessories category shows a 24% return rate, double the overall average.
 - Interpretation: Accessories may be impulse buys or poorly described online.
 - Business Impact: High returns increase logistics costs and may signal a UX issue.
 - Recommendation: Improve product detail pages and track conversions vs. returns by subcategory.
-

Time-Based Variable Example

Feature: ``signup_date``

- Observation: Signups peak on Mondays and drop on weekends.
 - Interpretation: Most users sign up during work hours-possibly B2B context.
 - Business Impact: Marketing spend might be misaligned if weekend ads are underperforming.
 - Recommendation: Shift campaigns to weekdays or A/B test timing windows.
-

When to Synthesize by Theme Instead of Feature

Sometimes it's better to write about **themes** rather than individual columns. For example:

Theme: High Churn in Low-Engagement Users

- Users with fewer than 2 logins in the first 7 days are 3x more likely to churn.
 - This suggests early activity is a leading indicator of retention.
 - A strong onboarding experience could help drive activation and reduce churn.
-

Tools to Use in Phase 3

Tool	Use
Visuals	Use plots as evidence, not the insight. Annotate trends.
Narrative templates	Use the 4-part format to standardize clarity.
Business questions	Frame findings in terms of cost, risk, opportunity, or experience.
Insight tables	For bulk reporting, use a structured format per feature.

Sample Insight Table

Feature	Observation	Interpretation	Impact	Recommendation
units_sold	Right-skewed	Bulk buyers inflate average	Misrepresents typical behavior	Use median or segment users
device_type	Mobile = 72%	Users prefer mobile access	UX must be mobile-first	Prioritize mobile testing
churn_flag	Higher in Region C	Possibly due to support delays	Lost revenue and trust	Add CS resources to Region C

Tip: Avoid the “Just Reporting” Trap

Don't just say: > “Hydraulic pressure is skewed.”

Instead say: > “Hydraulic pressure is skewed, which may indicate system stress in certain conditions. Since this variable correlates with failure events, it may be valuable for early fault detection.”

Final Deliverables from Phase 3

You now have: - 5–10 synthesized insights tied to business outcomes - Feature-level write-ups or a theme-based executive summary - A strong foundation for a written report or stakeholder presentation

Want a Template?

Here's a **Markdown Template** for this synthesis phase:

```
### Feature: `feature_name_here`
```

- **Observation**:
- **Interpretation**:
- **Business Impact**:
- **Recommendation**:

Or, if you'd prefer, I can create: - A **Google Doc** or **PDF version** of these templates - A **Notion board**, **Airtable**, or **Markdown dashboard** for organizing insights across a project

Let me know your preferred format and I'll deliver it!