# Critical Issues Batch - Comprehensive Analysis

## Executive Summary

This document analyzes and fixes 4 critical issues found in the soldcomp-analyser2 application:

1. **Target URL Preservation** - Target property URL being overwritten
2. **EPC Lookup Row Corruption** - EPC Lookup row contains property data
3. **Duplicate Detection Failures** - Duplicates with wrong data appearing in output
4. **UTF-8 Encoding** - Display issues with pound sign (£)

## Issue #1: Target URL Preservation

### Problem

**Input:** `data (4).csv` line 5

```
isTarget: 1
Address: 317 Wharf Road, Ealand
URL: https://www.rightmove.co.uk/properties/160516301#/?channel=RES_BUY
```

**Output:** `output (55).csv` line 4

```
isTarget: 1
Address: 317 Wharf Road, Ealand
URL: View
```

The full Rightmove URL is being replaced with just "View".

### Root Cause Analysis

The target property URL is being overwritten during one of these stages:
1. CSV parsing (unlikely - other properties retain URLs)
2. URL classification (unlikely - classification doesn't modify URLs)
3. Scraping/merging (LIKELY - merge logic may be overwriting with incomplete data)
4. Hyperlink generation (LIKELY - excelHelper may be replacing URLs)

### Investigation Points

- Check `mergeScrapedData()` in main.js
- Check `mergeProperties()` in duplicateDetector.js
- Check `addHyperlinks()` in excelHelper.js
- Check if target property is being treated as a duplicate

### Proposed Fix

1. Add explicit target URL protection in merge operations
2. Skip target property during URL overwrites

3. Add validation before final output to ensure target URL is preserved

---

# Issue #2: EPC Lookup Row Corruption

## Problem

**Input:** `data (4).csv` line 4

```
Address: EPC Lookup
Postcode: DN17 4JW
Sq. ft: (empty)
Distance: (empty)
Ranking: (empty)
```

**Output:** `output (55).csv` line 8

```
Address: EPC Lookup
Postcode: DN17 4JW
Sq. ft: 2551
Distance: 0.1mi
Ranking: 68
Latitude: 53.5906
Longitude: -0.817001
```

The EPC Lookup row has been enriched with property data, making it indistinguishable from a real property.

## Root Cause Analysis

**CRITICAL FINDING:** EPC Lookup rows exist in INPUT files!

This is the core problem - EPC Lookup rows should ONLY be created by the code, not exist in input CSV files. When output CSV from a previous run is used as input for a new run (iterative processing):

1. Previous run creates EPC Lookup row correctly (empty property data)
2. Output CSV is saved with EPC Lookup row
3. User downloads output CSV and uploads it as new input
4. New run parses input CSV, sees "EPC Lookup" row
5. Current code skips EPC Lookup during duplicate detection (line 65-70 in duplicateDetector.js)
6. BUT the row still gets processed through geocoding, ranking, and other enrichment
7. Result: EPC Lookup row gets filled with property data

## Why Current Fix (PR #12) Didn't Work

The previous fix added detection/skipping in duplicate detector, but the corruption happens AFTER duplicate detection:
- Step 8: Duplicate detection (EPC Lookup skipped ✓)
- Step 9: Geocoding (EPC Lookup gets geocoded ❌)
- Step 10: EPC enrichment (EPC Lookup processed ❌)
- Step 11: Ranking (EPC Lookup gets ranked ❌)
- Step 13: Output preparation (EPC Lookup has data ❌)

## Proposed Fix

### Solution 1: Filter out EPC Lookup rows from input

```javascript
// In csvParser.js or early in main.js
properties = properties.filter(p =>
    p.Address !== 'EPC Lookup' && !p._isEPCLookupRow
);
log.warning(`Removed ${beforeCount - properties.length} EPC Lookup rows from input`);
```

### Solution 2: Add _isEPCLookupRow marker preservation

```javascript
// Ensure _isEPCLookupRow marker survives CSV serialization
// Add it to output column list
// Skip marked rows in ALL processing stages
```

### Solution 3: Skip EPC Lookup in all enrichment stages

```javascript
// In geocodeAndCalculateDistances()
properties.filter(p => p.Address !== 'EPC Lookup')

// In enrichWithEPCData()
properties.filter(p => p.Address !== 'EPC Lookup')

// In rankProperties()
properties.filter(p => p.Address !== 'EPC Lookup')
```

### Recommended: Combination of Solutions 1 & 3

- Remove EPC Lookup rows from input early (Solution 1)
- Add defensive checks in enrichment stages (Solution 3)
- Log warnings when EPC Lookup rows are found in input

---

# Issue #3: Duplicate Detection Failures

## Problem

**Input:** `data (4).csv`

```
ONE entry for "The Vicarage":
- Price: £362,000
- Sq. ft: 2024
- Sqm: 188
- URL: https://propertydata.co.uk/transaction/3DCCB7CA-0094-5B9D-E063-4704A8C0331E
```

**Output:** `output (55).csv`

```
TWO entries for "The Vicarage":
1. Line 6:
   - Price: £605,000
   - Sq. ft: 2390
   - Sqm: 222
   - URL: View
   - needs_review: "Price conflict: 605000 vs 195000; Sq. ft conflict: 2390 vs 797"

2. Line 11:
   - Price: £362,000
   - Sq. ft: 2024
   - Sqm: 188
   - URL: https://propertydata.co.uk/transaction/3DCCB7CA-0094-5B9D-E063-4704A8C0331E
```

## Root Cause Analysis

The duplicate is being CREATED during processing, not from input:

**Processing Flow:**
1. Input has ONE "The Vicarage" (from PropertyData)
2. Duplicate detection runs (Step 8) - no duplicates found ✓
3. Scraping happens - a different source returns different data for same property
4. Merge happens - but matching fails because:
- Different URLs
- Slightly different address format
- Merge creates duplicate instead of merging
5. Output has TWO "The Vicarage" entries

## Why Duplicate Detection Fails

The code has duplicate detection at line 91 in main.js:

```
allProperties = detectAndMergeDuplicates(allProperties);
```

But this runs BEFORE scraping. Duplicates created DURING/AFTER scraping are not caught.

## Proposed Fix

**Solution 1: Run duplicate detection AFTER scraping**

```
// Current flow:
Step 8: Detect duplicates
Step 9: Geocoding
Step 10: EPC enrichment

// Proposed flow:
Step 8: Detect duplicates (pre-scraping)
Step 9: Geocoding
Step 10: EPC enrichment
Step 10.5: Detect duplicates AGAIN (post-enrichment)
```

**Solution 2: Improve address normalization**

```
// Current normalization may miss variations
// Need better handling of:
// - "The Vicarage, Church Street" vs "The Vicarage"
// - Case variations
// - Comma vs no comma
```

**Solution 3: Add URL-based matching**

Current code has URL matching, but it may not be catching all cases. Need to:

- Normalize URLs more aggressively

- Match on transaction IDs in PropertyData URLs

- Match on coordinates (lat/long within threshold)

**Recommended: Solution 1 + Solution 3**

- Run duplicate detection twice (before and after enrichment)

- Improve URL-based matching for PropertyData transaction URLs

- Add coordinate-based matching as fallback

---

# Issue #4: UTF-8 Encoding

## Problem

Files display "Â£" instead of "£" in markdown rendering.

## Investigation

```
# Check actual file encoding
hexdump -C "data (4).csv" | grep -A2 "c2 a3"
# Result: c2 a3 is present (correct UTF-8 for £)
```

**Finding:** Files are correctly UTF-8 encoded!
- Bytes: `c2 a3` = UTF-8 encoding of £ (U+00A3)
- The "Â£" display is a rendering issue, not a file encoding issue

## Root Cause

When a UTF-8 file is read but interpreted as ISO-8859-1 (Latin-1):
- Byte `c2` displays as "Â" (U+00C2)
- Byte `a3` displays as "£" (U+00A3)
- Combined: "Â£"

This can happen in:
1. CSV readers configured for wrong encoding
2. Markdown renderers expecting different encoding
3. Excel opening UTF-8 files without BOM

## Proposed Fix

**Solution 1: Ensure consistent UTF-8 handling**

```
// In kvsHandler.js writeCSVToKVS()
const csvContent = '\uFEFF' + csvString; // Add BOM
```

**Solution 2: Clean existing files**

```
// Add utility to clean £ → £ during parsing
text = text.replace(/£/g, '£');
```

**Solution 3: Validate encoding during CSV write**

```
// Ensure CSV is written as UTF-8 with explicit encoding parameter
fs.writeFileSync(path, content, { encoding: 'utf8' });
```

**Recommended: Solution 1 + Solution 2**

- Add UTF-8 BOM to output files (helps Excel)
- Add defensive cleaning during parsing (fixes existing files)

---

# Testing Strategy

## Test 1: Target URL Preservation

```
// Create test with target property
// Verify URL survives entire pipeline
// Ensure URL not overwritten during merge
```

## Test 2: EPC Lookup Isolation

```
// Test input with EPC Lookup row (simulating iterative processing)
// Verify row is removed early
// Verify no property data added to EPC Lookup
// Verify new EPC Lookup row created correctly
```

## Test 3: Duplicate Detection

```
// Create property that will be scraped from multiple sources
// Verify duplicates are merged correctly
// Verify correct data is kept
// Verify conflicts are flagged in needs_review
```

## Test 4: UTF-8 Encoding

```
// Read file with £ symbols
// Verify no £ corruption
// Verify output has correct encoding
```

---

# Implementation Plan

## Phase 1: Critical Fixes (This PR)

1. ✅ Filter EPC Lookup rows from input

2. ✅ Skip EPC Lookup in all enrichment stages
3. ✅ Preserve target URL in all operations
4. ✅ Add post-enrichment duplicate detection
5. ✅ Add UTF-8 cleaning utility

### Phase 2: Testing

1. ✅ Create comprehensive test suite
2. ✅ Test with data.csv (LN11 9WB area)
3. ✅ Test with data (4).csv (DN17 4JW area)
4. ✅ Verify all issues resolved

### Phase 3: Documentation

1. ✅ Document fixes in this file
2. ✅ Create test results summary
3. ✅ Update README if needed

---

## Expected Outcomes

After implementing all fixes:

1. **Target URL preserved** - Full Rightmove URL maintained throughout processing
2. **No EPC Lookup corruption** - EPC Lookup rows filtered from input, new ones created correctly
3. **No duplicates** - Post-enrichment deduplication catches all cases
4. **Clean encoding** - £ displays correctly, no Â£ corruption

---

## Files Modified

1. `src/main.js` - Add EPC Lookup filtering, post-enrichment deduplication, target URL protection
2. `src/utils/duplicateDetector.js` - Improve URL matching, add coordinate-based matching
3. `src/utils/csvParser.js` - Add UTF-8 cleaning utility
4. `src/utils/excelHelper.js` - Protect target URL during hyperlink generation
5. `test-critical-issues-batch.js` - Comprehensive test suite

---

## Notes

- Previous fix (PR #12) was incomplete - it only fixed duplicate detection, not enrichment stages
- EPC Lookup rows in input files are a sign of iterative processing (output → input)
- Need to educate users NOT to use output CSV as new input, OR make code robust to this pattern
- Coordinate-based duplicate detection is a powerful fallback for address variations