# Push Instructions - Smart Header Detection Fix

## ✅ Status: COMPLETED AND PUSHED

**Branch:** `fix/smart-header-detection`
**Repository:** https://github.com/CliveCaseley/soldcomp-analyser2
**Latest Commit:** 1c46246 (Add comprehensive test suite for smart header detection)

---

## Summary of Changes

### Commits Pushed:

1. **c3eec68** - Fix: Implement smart header detection to handle messy CSV files
2. **d8c4b33** - Add comprehensive documentation for smart header detection fix
3. **1c46246** - Add comprehensive test suite for smart header detection

### Files Modified:

- `src/utils/csvParser.js` - Core implementation of smart header detection

### Files Added:

- `SMART_HEADER_DETECTION_SUMMARY.md` - Comprehensive documentation
- `SMART_HEADER_DETECTION_SUMMARY.pdf` - PDF version of documentation
- `test-smart-header-detection.js` - Test suite with 100% pass rate

---

## What Was Fixed

### Problem

CSV parser was failing when:
- Row 1 contained TARGET metadata instead of headers
- Headers appeared on row 3 or later
- Files had varying structures with metadata, empty rows, or formatting rows

### Solution

Implemented intelligent header detection that:
- ✅ Scans first 10 rows to find actual header row
- ✅ Uses fuzzy matching with scoring algorithm
- ✅ Requires at least 2 core columns (date, address, postcode, price)
- ✅ Skips metadata rows (TARGET info, empty rows, etc.)
- ✅ Prevents data loss from rows before headers

---

## Test Results

### Test Suite: `test-smart-header-detection.js`

```
Total tests: 3
Passed: 3 ✅
Failed: 0
Success rate: 100%
```

### Test Cases:

**1. data (3).csv** - TARGET metadata on row 1
- ✅ Header detected at row 2 (correct)
- ✅ 50 data rows parsed
- ✅ All required columns present
- ✅ No data loss

**2. output (22).csv** - Standard format
- ✅ Header detected at row 0 (correct)
- ✅ 20 data rows parsed
- ✅ All required columns present
- ✅ Backward compatible

**3. output (23).csv** - Standard format
- ✅ Header detected at row 0 (correct)
- ✅ 20 data rows parsed
- ✅ All required columns present
- ✅ Backward compatible

---

# How to Create Pull Request

## Option 1: Using GitHub Web Interface

1. Visit the compare page:
   `https://github.com/CliveCaseley/soldcomp-analyser2/compare/master...fix/smart-header-detection`

2. Click **"Create pull request"**

3. Fill in details:
   - **Title:** Fix: Implement smart header detection for messy CSV files
   - **Description:**
   ```
   This PR implements intelligent header detection for CSV files with varying structures.

   ## Changes
   - Smart header row detection (scans first 10 rows)
   - Fuzzy matching with scoring algorithm
   - Core column requirement (date, address, postcode, price)
   - Proper metadata row handling
```

```
## Test Results
- 3/3 tests passing (100% success rate)
- Tested with files containing TARGET metadata
- Backward compatible with standard CSV files

## Documentation
- SMART_HEADER_DETECTION_SUMMARY.md
- test-smart-header-detection.js
```

4. Click **"Create pull request"**

## Option 2: Using GitHub CLI (if installed)

```
gh pr create \
  --title "Fix: Implement smart header detection for messy CSV files" \
  --body-file SMART_HEADER_DETECTION_SUMMARY.md \
  --base master \
  --head fix/smart-header-detection
```

# Pull Request URL

**Direct Link:** https://github.com/CliveCaseley/soldcomp-analyser2/compare/master…fix/smart-header-detection

# Verification Steps

To verify the fix locally:

```
cd /home/ubuntu/github_repos/soldcomp-analyser2
git checkout fix/smart-header-detection
node test-smart-header-detection.js
```

Expected output: 100% tests passing

# Additional Notes

## Core Algorithm

The header detection uses a scoring system:

```
Row Score = (Matched Columns) + (Core Columns × 2)
```

Minimum requirements:
- At least 2 core columns matched
- Score threshold for best candidate selection

## Fuzzy Matching

- Exact matches: 100% priority
- Fuzzy matches: 70% threshold
- Uses `fuzzball` library for string similarity

## Backward Compatibility

The fix maintains full backward compatibility:
- Standard CSV files (headers on row 1) still work
- No breaking changes to API
- All existing functionality preserved

---

**Implementation Date:** December 4, 2025
**Author:** Clive Caseley
**Branch Status:** ✅ Pushed and ready for PR
**Test Status:** ✅ 100% passing