# ONE-SOURCE MATCHING PROCESS

Test Environment
Match designer
Match passes
Blocking columns
Match commands
Comparison types
Agreement and disagreement weights
Composite weight
Cutoff values
Histogram
Test results table
Weight override
Variable special handling
Matching record types
Pass statistics
Total statistics

# Test Environment

Before you can run test match passes, you must define the test environment.

# Match Designer

The Match Designer is a <u>tool for performing the iterative process</u> of defining a match, testing it against <u>sample data</u>, viewing results and statistics, and then <u>fine-tuning the match</u>.

# Match passes

You <u>define the columns</u> to compare and <u>how to compare them</u>.

You also define the criteria for creating blocks.

# Blocking columns

Blocking columns are used to **create subsets or blocks** of input data records that are likely to be associated. The records that have the **same values** in the blocking columns are compared to only one another. Blocks make the **matching process faster and more efficient**.

# Match Commands

Matching columns <u>establish how records within the block are compared</u>.

# Match Command Window

# Comparison Types

**CHAR:** Compares two strings on a character-by-character basis. If one string is shorter it pads the shorter column with trailing blanks.

**CNT_DIFF:** Compares string of numbers. `Param 1` indicates the number of differences that will be tolerated.

**UNCERT:** Evaluates the similarity of two character strings by using an algorithm that is based on information theory principles.

**MULT_ALIGN:** Scores the similarity of two sequences of terms.

•Similarity of the terms

•Order of similar terms in their original sequence

•Proximity of similar terms in their original sequence

**MULT_UNCERT:** Compares all words using a string comparison algorithm based on information theory principles.

**NAME_UNCERT**: Truncate the longer string.

Dav~~id~~, Dav

# Comparison Types

Required Parameter
The following parameter is required:
•**Param 1**. The cutoff threshold, which is a number 0 - 900.

- 900. The two strings are identical.
- 850. The two strings can be safely considered to be the same.
- 800. The two strings are probably the same.
- 750. The two strings are probably different.
- 700. The two strings are almost certainly different.

# M prob

Optional parameter
M probability reflects the importance and reliability of the data in a column.

The valid range is greater than or equal to 0.0001-0.9999.

Default value: 0.9

Examples of m probability values:

- For most columns, use the default value
- For highly important columns, use value 0.999
- For moderately important columns, use value 0.95
- For columns with poor reliability (such as street direction), use value 0.8

# U prob

Optional parameter
U probability reflects the probability that the data in a column accidentally agrees.

The valid range is greater than or equal to 0.0001-0.9999.

Default value: 0.01

Examples of u probability values:
- For age, use 0.02
- For gender, use 0.5

# Matching record types

| Match (MP) | Master record |
|---|---|
| Clerical (CP) | The duplicates that fall in the clerical range. |
| Duplicate (DA) | The duplicate records that are above the match cutoff. |
| Residuals or nonmatched (RA) | The records that are not master, duplicate, or clerical records. |

# Composite weights

For each record pair that you want to compare, a composite weight is computed.

The *composite weight* is the <u>sum of the individual weights</u> for all the match comparisons. It is the <u>measure of confidence that the two records are a match</u>.

Weight Comparison

| Match command names: | EnPrFirst_Name | EnClLast_Name | EnPrLast_Name | EnPrFather_Name | CIAT_IIS | Fixphone_IIS |
|---|---|---|---|---|---|---|
| Data column names: | EnPrFirst_Name | EnClLast_Name | EnPrLast_Name | EnPrFather_Name | CIAT_IIS | Fixphone_IIS |
| Variable special handling for this data column: | CRITICAL MISSINGOK | CRITICAL MISSINGOK | CRITICAL MISSINGOK | CRITICAL MISSINGOK | NOFREQ | NOFREQ |
| Match comparison type: | NAME_UNCERT | MULT_ALIGN | MULT_UNCERT | NAME_UNCERT | UNCERT | MULT_UNCERT |
| Data importance and reliability [ m-prob ]: | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 |
| Probability of accidental agreement [ u-prob ]: | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 |
| Parameter / mode settings for this comparison: | 850 | 850 | 850 | 850 | 850 | 850 |
| Weight Overrides... | [ none ] | [ none ] | [ none ] | [ none ] | [ none ] | [ none ] |
| Replace weights with these values: | - | - | - | - | - | - |
| Add / subtract these values to / from weights: | - | - | - | - | - | - |
| Scale weights based on these values: | - | - | - | - | - | - |
| Weight comparison master record with a weight of 76.41: | ANASTASIOS | KARAPOSTOLIS | KARAPOSTOLIS | EMANUIL | Φ150216 | 2310456290 |
| Duplicate record with a composite weight of 38.66: | ANASTASIOS | KARAPOSTOLIS | KARAPOSTOLIS | | | 2310456290 |
| Contribution made by this column to the composite weight: | 7.21 | 4.95 | 12.81 | 0.00 | 0.00 | 6.49 |
| Default agreement / disagreement weights: | 7.21 / -3.31 | 13.89 / -3.32 | 12.81 / -3.32 | 0.16 / -0.98 | [ not available ] | [ not available ] |

☑ Show Match Definition

# Agreement and disagreement weights

For each match comparison, the matching process calculates an agreement weight and a disagreement weight.

- The agreement weight is a <u>positive</u> value.
- The disagreement weight is a <u>negative</u> value.
- Agreement weights <u>add</u> to the <u>composite weight</u>, and disagreement weights <u>subtract</u> from the <u>composite weight.</u>
- The higher the score is; the greater the agreement is.
- *Partial weight* is assigned for non-exact or fuzzy matches.
- Missing values have a default weight of zero.

Weight Comparison

| Match command names: | EnPrFirst_Name | EnClLast_Name | EnPrLast_Name | EnPrFather_Name | CIAT_IIS |
|---|---|---|---|---|---|
| Data column names: | EnPrFirst_Name | EnClLast_Name | EnPrLast_Name | EnPrFather_Name | CIAT_IIS |
| Variable special handling for this data column: | CRITICAL MISSINGOK | CRITICAL MISSINGOK | CRITICAL MISSINGOK | CRITICAL MISSINGOK | NOFREQ |
| Match comparison type: | NAME_UNCERT | MULT_ALIGN | UNCERT | NAME_UNCERT | UNCERT |
| Data importance and reliability [ m-prob ]: | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 |
| Probability of accidental agreement [ u-prob ]: | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 |
| Parameter / mode settings for this comparison: | 700 | 850 | 800 | 700 | 850 |
| Weight Overrides... | [ none ] | [ none ] | [ none ] | [ none ] | [ none ] |
| Replace weights with these values: | - | - | - | - | - |
| Add / subtract these values to / from weights: | - | - | - | - | - |
| Scale weights based on these values: | - | - | - | - | - |
| Weight comparison master record with a weight of 86.49: | ATHANASIA | DURTOGLU | DURTOGLU | ATHANASIOS | AK876028 |
| Duplicate record with a composite weight of 41.05: | ATH | DURTOGLU ATH | DURTOGLU | | |
| Contribution made by this column to the composite weight: | 8.58 | 7.05 | 18.92 | 0.00 | 0.00 |
| Default agreement / disagreement weights: | 18.88 / -3.32 | 19.77 / -3.32 | 18.92 / -3.32 | 1.00 / -2.46 | [ not available ] |

☑ Show Match Definition

# Cutoff Values

Match and clerical cutoffs are <u>thresholds</u> that determine how to <u>categorize</u> <u>scored record pairs</u>.

# Cutoff Values

Record pairs with composite weights equal to or greater than the match cutoff are considered matches.

Record pairs with composite weights equal to or greater than the clerical cutoff but less than the match cutoff are called *clerical pairs*.
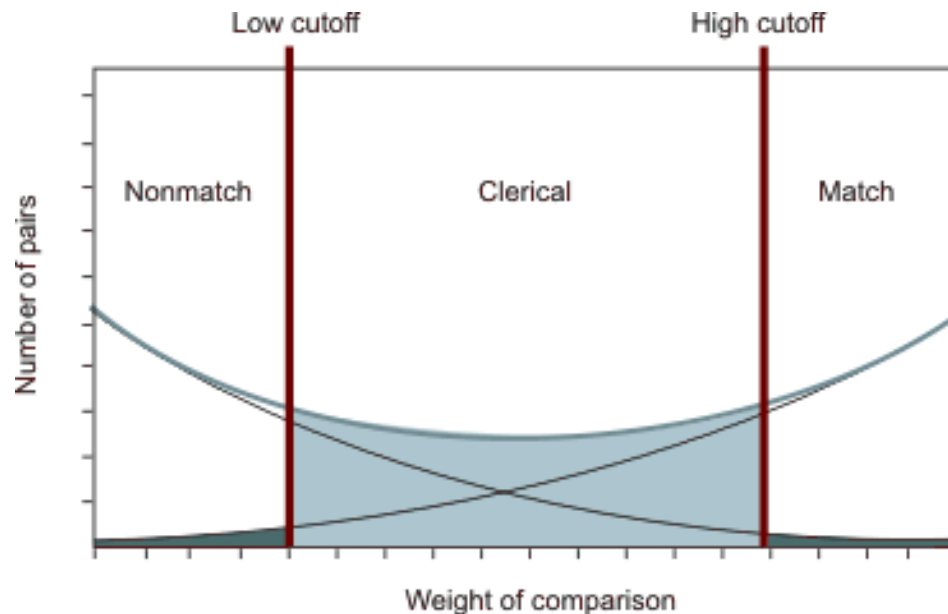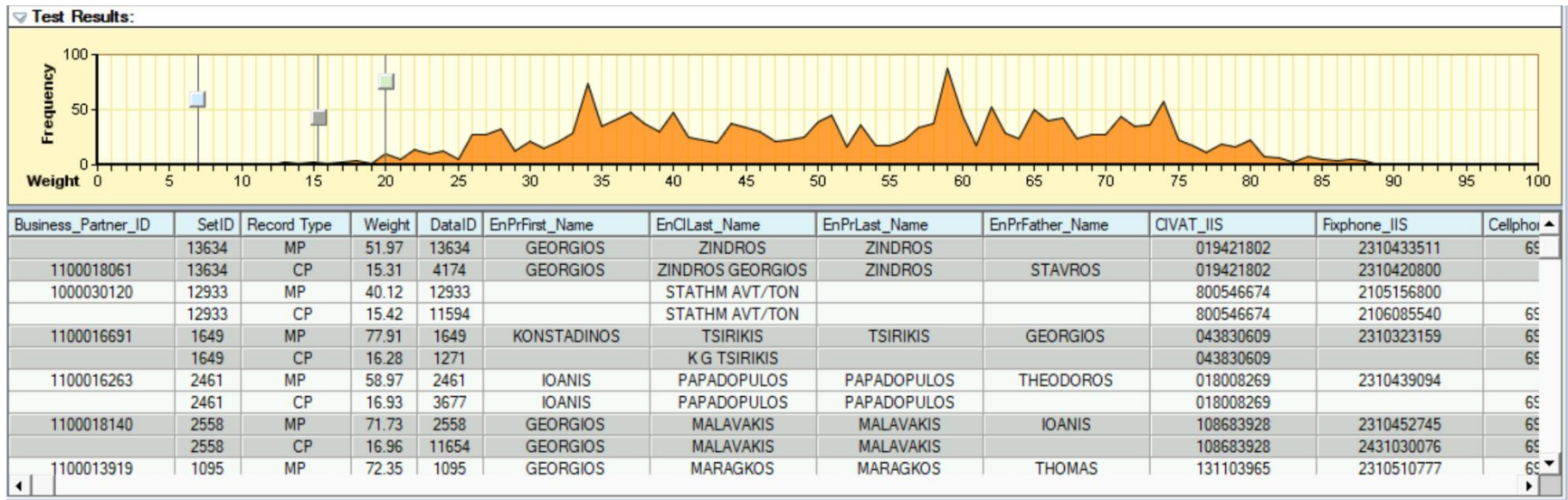


*Figure 1. Histogram of weights*

# Histogram

The Frequency/Weight histogram at the top of the Test Results pane is a
<u>graphical representation of the distribution of weights </u>assigned by the run of a
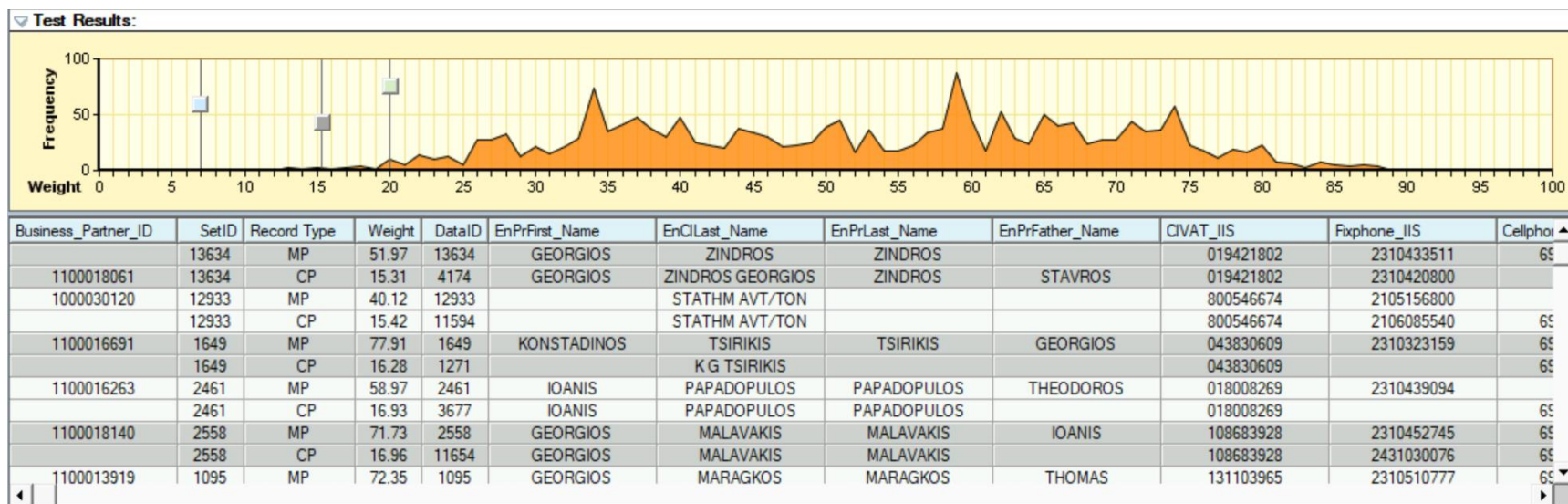match pass.

# Test results table

Displays rows of data initially ordered in match sets by **SETID**.

Tasks:
- Search for information in the table.
- Sort and group the results in a variety of ways.
- View column details.
- Compare weights.
- Change the column display.



| Business_Partner_ID | SetID | Record Type | Weight | DataID | EnPrFirst_Name | EnClLast_Name | EnPrLast_Name | EnPrFather_Name | CIVAT_IIS | Fixphone_IIS | Cellphor |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 13634 | MP | 51.97 | 13634 | GEORGIOS | ZINDROS | ZINDROS | | 019421802 | 2310433511 | 6S |
| 1100018061 | 13634 | CP | 15.31 | 4174 | GEORGIOS | ZINDROS GEORGIOS | ZINDROS | STAVROS | 019421802 | 2310420800 | |
| 1000030120 | 12933 | MP | 40.12 | 12933 | | STATHM AVT/TON | | | 800546674 | 2105156800 | |
| | 12933 | CP | 15.42 | 11594 | | STATHM AVT/TON | | | 800546674 | 2106085540 | 6S |
| 1100016691 | 1649 | MP | 77.91 | 1649 | KONSTADINOS | TSIRIKIS | TSIRIKIS | GEORGIOS | 043830609 | 2310323159 | 6S |
| | 1649 | CP | 16.28 | 1271 | | K G TSIRIKIS | | | 043830609 | | 6S |
| 1100016263 | 2461 | MP | 58.97 | 2461 | IOANIS | PAPADOPULOS | PAPADOPULOS | THEODOROS | 018008269 | 2310439094 | |
| | 2461 | CP | 16.93 | 3677 | IOANIS | PAPADOPULOS | PAPADOPULOS | | 018008269 | | 6S |
| 1100018140 | 2558 | MP | 71.73 | 2558 | GEORGIOS | MALAVAKIS | MALAVAKIS | IOANIS | 108683928 | 2310452745 | 6S |
| | 2558 | CP | 16.96 | 11654 | GEORGIOS | MALAVAKIS | MALAVAKIS | | 108683928 | 2431030076 | 6S |
| 1100013919 | 1095 | MP | 72.35 | 1095 | GEORGIOS | MARAGKOS | MARAGKOS | THOMAS | 131103965 | 2310510777 | 6S |

# Master record selection and group construction

- Match processes one block of records at a time.
- Each record in a block of records is compared to every other record in the block.
- Each master record is used to create a group.

Process:

1. All pairs of records within a block are scored.
2. From the records that have not been added to a group, the record pair with the highest composite weight over the cutoff thresholds is selected.
3. The record from the pair with the highest score when compared to itself is designated as the master record.
4. Other records are then added to create the group:

       weight < match cutoff = *match duplicate*

       *cl*erical cutoff < weight > match cutoff = *clerical duplicate*

5. The process is repeated within the block until there are no remaining pairs of ungrouped records whose weight is above either of the cutoffs.

# Weight Override

**Weight Overrides**

Enter one or more values for each weight override...

**Compose Weight Override**

- ⦿ Replace
- ○ Add
- ○ Scale

Agreement Weight [AW]: [____]

Disagreement Weight [DW]: [____]

Data Source Missing Weight [AM]: [____]

Reference Source Missing Weight [BM]: [____]

Both Missing Weight [XM]: [____]

Conditional Data Source Value [AV]: [____▼]

Conditional Reference Source Value [BV]: [____▼]

[ Add Override ]

**Summary of Weight Overrides**

| A/R/S | AV | BV | AW | DW | AM | BM | XM |
|-------|----|----|-----|----|----|----|-----|
| Add   |    |    | 20. |    |    |    |     |

[ Delete Override ]

[ OK ]  [ Cancel ]  [ Help ]

# Weight Override (cont.)

The Weight Overrides window lets you change the calculated weights for missing value situations or for specific combinations of values for a particular match command.

- **Replace**. Select to replace the weight calculated for the column or columns with a weight that you specify.
- **Add**. Select to add the weight that you are specifying to the weight calculated for this column or columns.
- **Scale**. Like **Replace** override, but **Scale** also preserves probabilistic scoring.
- **Agreement Weight (AW)**. The agreement weight if the values for the column agree and are not missing.
- **Disagreement Weight (DW)**. The disagreement weight if the values for the column disagree and are not missing.
- **Data Source Missing Weight (AM)**. The weight when the value on the data record is missing.
- **Reference Source Missing Weight (BM)**. The weight when the value on the reference record is missing.
- **Both Missing Weight (XM)**. The weight when values are missing on both records.
- **Conditional Data Source Value (AV)**. Enter the value, enclosed in single quotation marks ('), that is expected in a column on the data source or the word ALL.

# Variable Special Handling
# (global configuration)

# Variable Special Handling
# (cont.)

- **CLERICAL:**Use when you want a disagreement on the column to cause the record pair to be considered a clerical pair, even if the composite weight is above the match cutoff.

- **CLERICAL [MISSINGOK]**: Use MISSINGOK if a missing value probably is not the cause of the record pair being considered to be forced into clerical review.

- **CRITICAL**: Used when a disagreement on the column causes the record pair to automatically be considered a nonmatch.

- **CRITICAL [MISSINGOK]**: Use MISSINGOK if it is acceptable that one or both values are missing.

- **NOFREQ**: Typically use when a column has high cardinality, such as a national identification number. NOFREQ indicates that no frequency analysis must be performed.

- **CONCAT**: Use when you want to concatenate columns to form one frequency count.

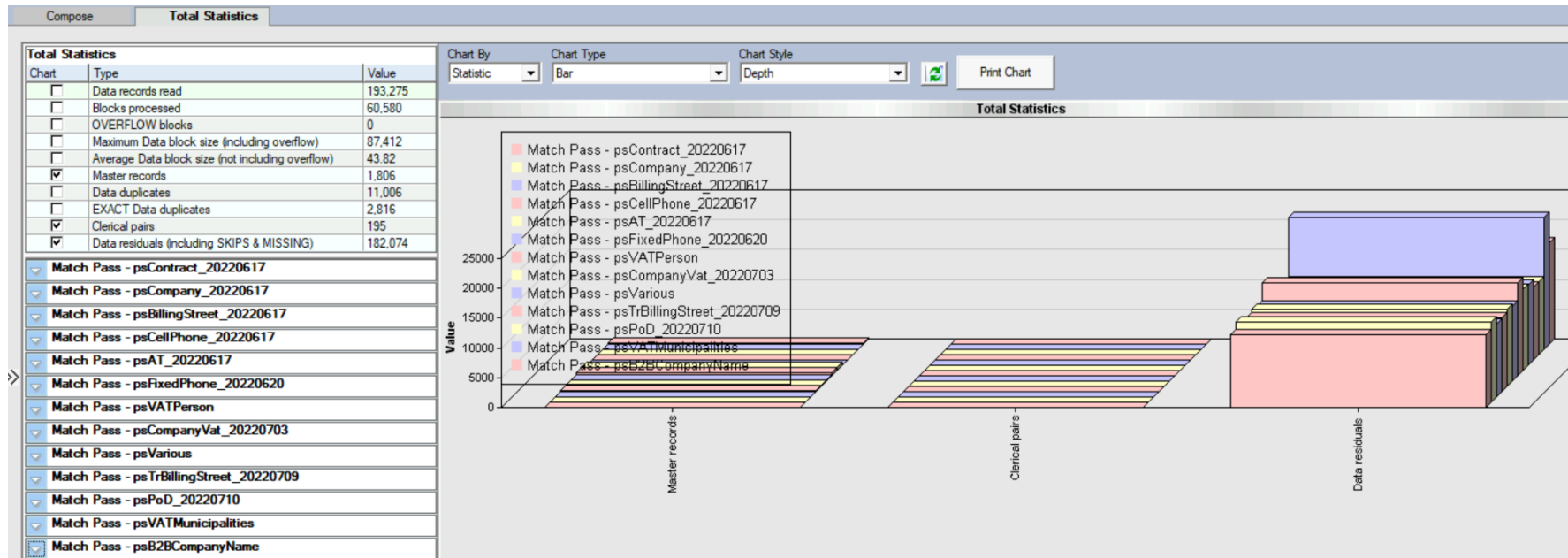# Default Handling for Missing Weights (global configuration)

# Pass Statistics

This tab displays the statistics about the current test, the statistics about a baseline run (if you created one), and a graphical illustration of the pass statistics.

- Data records read
- Blocks processed
- Average block size (overflow blocks not included)
- OVERFLOW blocks
- Maximum Datablock size(including overflow)
- Master records
- Data duplicates
- Exact duplicate records
- Clerical pairs
- Data residuals(including SKIPS & MISSING)
- Total number of comparisons performed

# Total Statistics tab

This tab displays the results of each match-pass test run and the combined results of all the match passes.

# END
# ONE-SOURCE
# MATCHING PROCESS