

PANDAS MERGE

SHARP SIGHT

WHAT YOU'LL LEARN

- How to use the Pandas `merge()` function
- How to "merge" data
 - i.e., combine by looking for a match on a variable
- How to perform different types of merges
 - left merge
 - right merge
 - inner merge

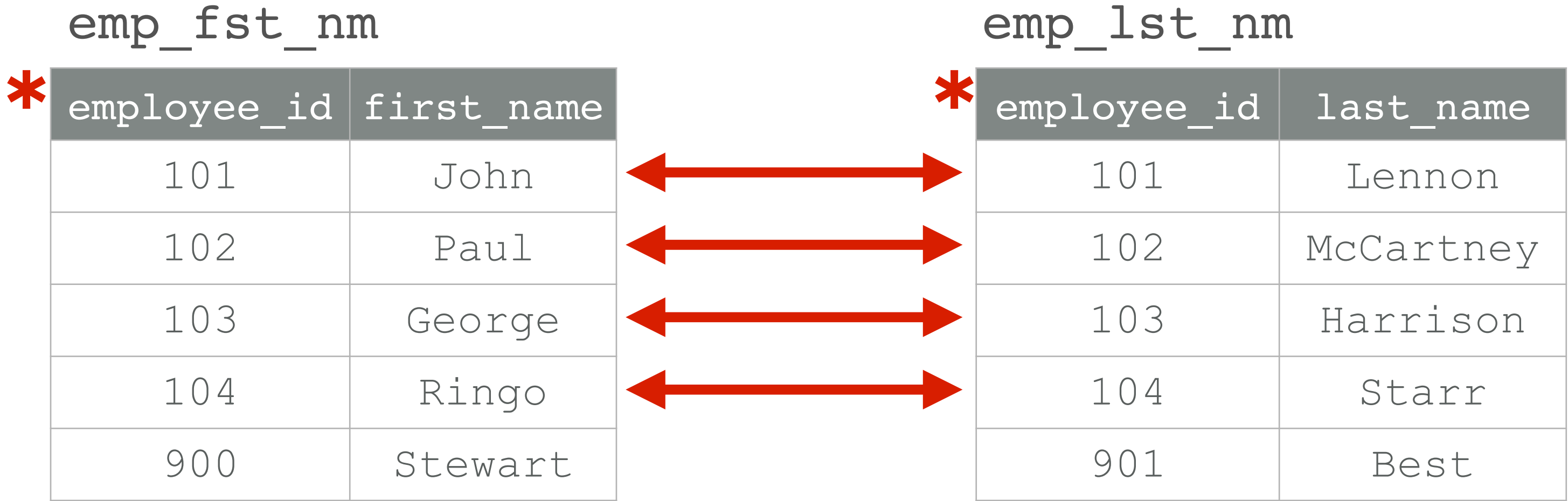
PANDAS MERGE OVERVIEW

WHAT IS A MERGE?

- A merge combines records from 2 datasets
- Merges are essentially the same as:
 - merge (SAS)
 - join (R's dplyr)
 - vlookup (Excel)
 - SQL (join)

HIGH LEVEL EXAMPLE OF A MERGE

```
pd.merge(emp_fst_nm, emp_lst_nm, on = 'employee_id')
```



HIGH LEVEL EXAMPLE OF A MERGE

```
pd.merge(emp_fst_nm, emp_lst_nm, on = 'employee_id')
```

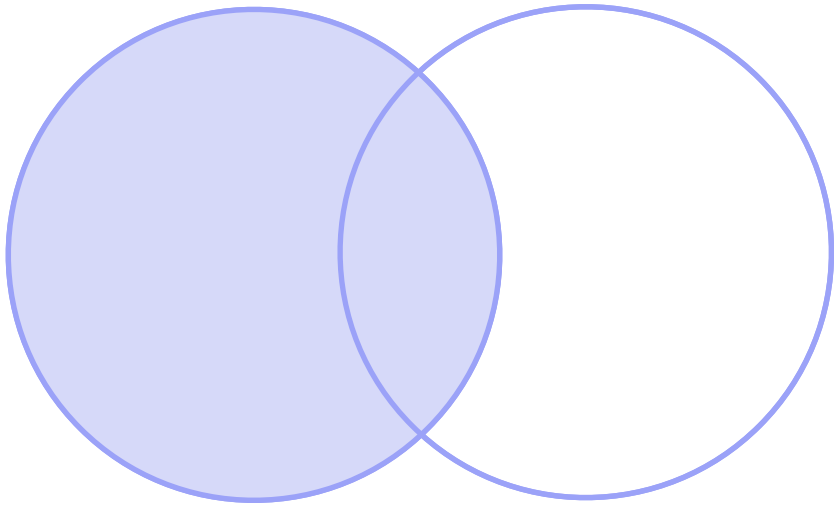
OUT :

employee_id	first_name	last_name
101	John	Lennon
102	Paul	McCartney
103	George	Harrison
104	Ringo	Starr

2 MERGE TYPES YOU'LL LEARN

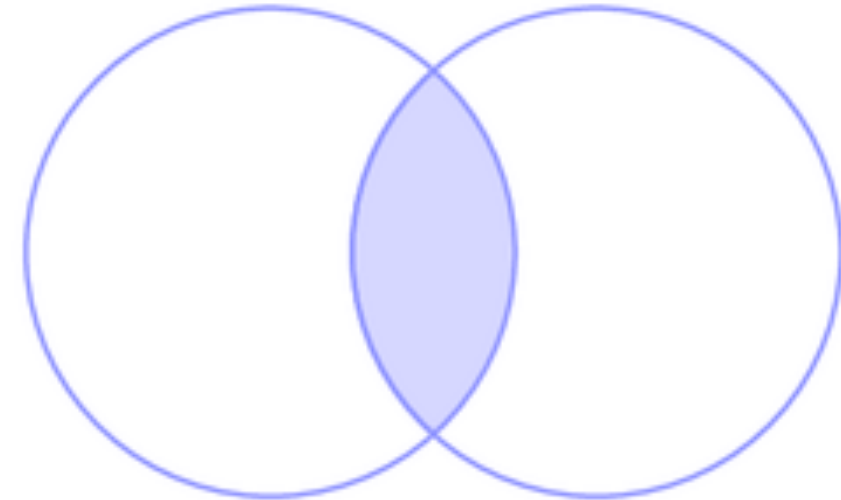
Left merge

(keep all records from
“left” dataset)



Inner merge

(only keep records that
are in BOTH datasets)



WHY INNER MERGE & LEFT MERGE?

- These two are most common
- There are two other merge types
 - you'll probably never use them
 - so, focus on most common
- Right merge is “inverted” version of left merge
 - once you know left, you can do right

SYNTAX: PANDAS MERGE

SYNTAX: PANDAS MERGE FUNCTION

The name of the function



```
pd.merge(left-data, right-data, on=, how=)
```

A DataFrame that you want to merge



Another DataFrame that you want to merge



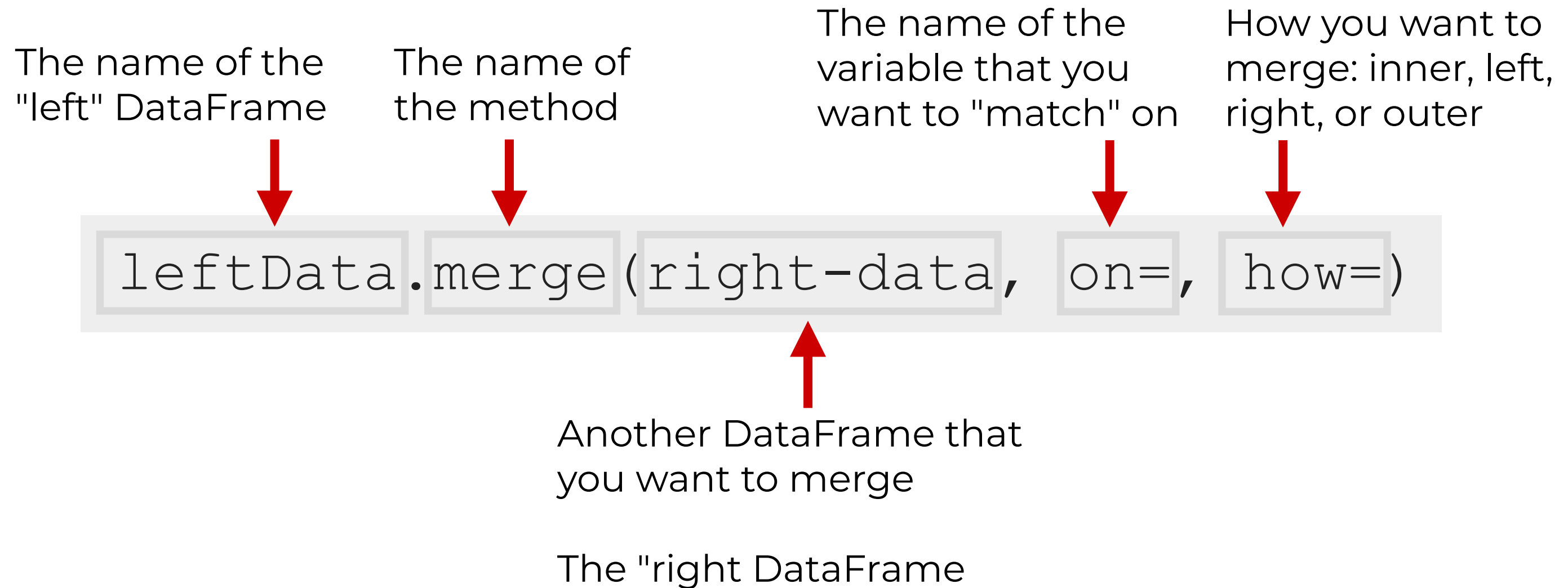
The name of the variable that you want to "match" on



How you want to merge: inner, left, right, or outer



NOTE: THERE'S ALSO A METHOD VERSION



PARAMETERS OF PANDAS MERGE

THE PARAMETERS OF THE PANDAS MERGE FUNCTION

Parameter	What it does	Format	Default	Required?
<code>left-data</code>	The the dataset on the left side of the syntax that you want to merge	DataFrame		Yes
<code>right-data</code>	The the dataset on the right side of the syntax that you want to merge	DataFrame		Yes
<code>on=</code>	Which variable upon which you want to look for "matches" for the merge	A column or list of columns	Intersection of the columns in both DataFrames. (You should always provide an argument to <code>on</code>)	No
<code>how=</code>	How you want to merge the data ... inner, left, right, or outer	One of the following: <code>inner</code> , <code>left</code> , <code>right</code> , <code>outer</code>	<code>inner</code>	No

THE OUTPUT OF PANDAS MERGE

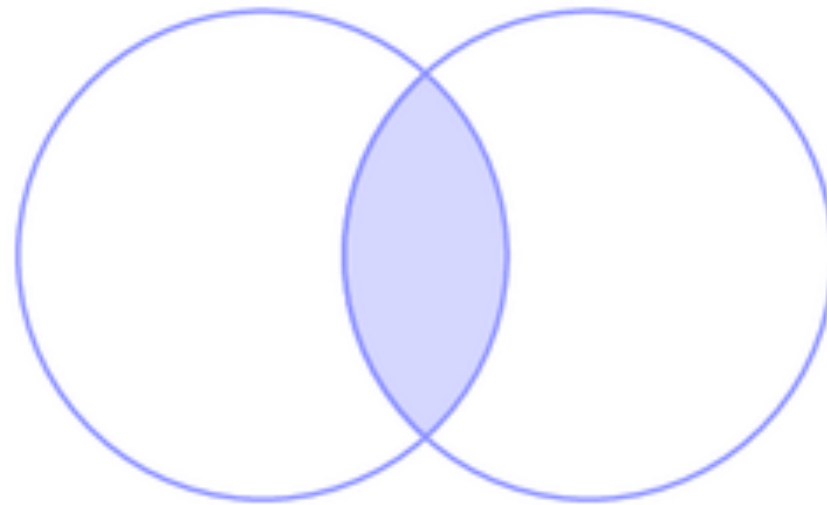
- The output of `pd.merge()` is a DataFrame
 - Contains the merged data
 - Merged as specified in your syntax

EXAMPLE: INNER MERGE

INNER MERGE

- Only keep records that are in BOTH datasets

Inner merge



SYNTAX: PANDAS INNER MERGE

The name of the
variable that you
want to "match" on

Execute an "inner"
merge

```
pd.merge(left-data, right-data, on=, how= 'inner')
```

A DataFrame
that you want to
merge

Another
DataFrame that
you want to merge

EXAMPLE: WE HAVE TWO SEPARATE DATASETS

- Want to combine them
 - keep only records that are in both

emp_fst_nm

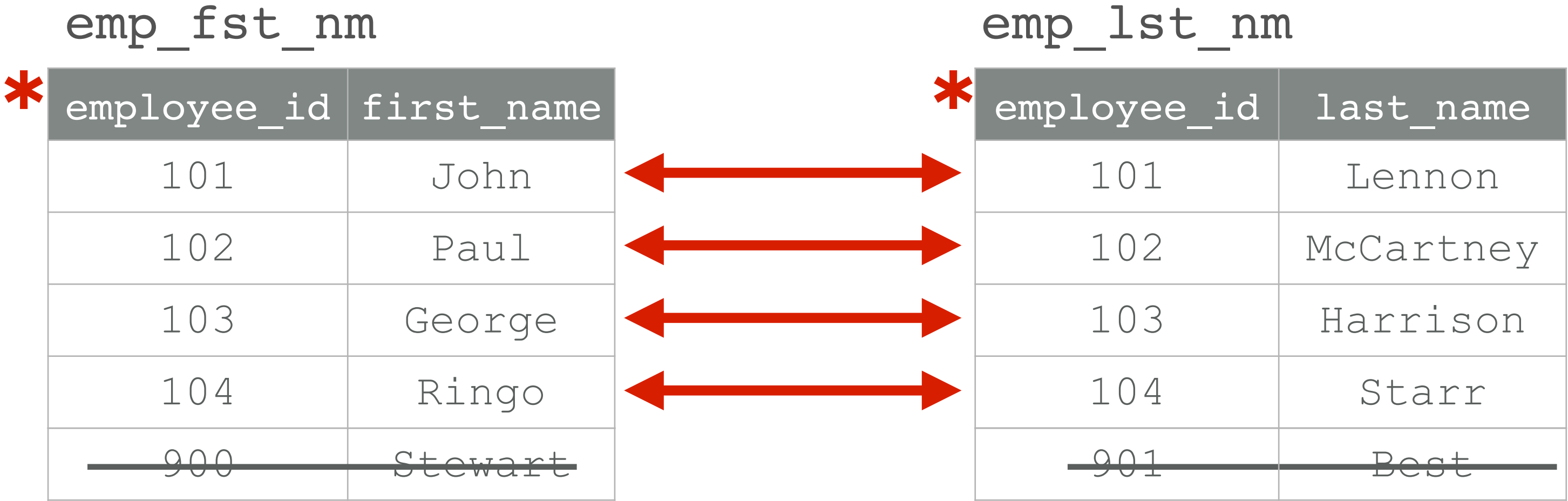
employee_id	first_name
101	John
102	Paul
103	George
104	Ringo
900	Stewart

emp_lst_nm

employee_id	last_name
101	Lennon
102	McCartney
103	Harrison
104	Starr
901	Best

EXAMPLE: INNER MERGE

```
pd.merge(emp_fst_nm, emp_lst_nm, on = 'employee_id', how = 'inner')
```



INNER MERGE RESULTS

```
pd.merge(emp_fst_nm, emp_lst_nm, on = 'employee_id', how = 'inner')
```

- Result is a combined dataset
 - Columns from both data sets
 - Excludes rows that were not in both

employee_id	first_name	last_name
101	John	Lennon
102	Paul	McCartney
103	George	Harrison
104	Ringo	Starr

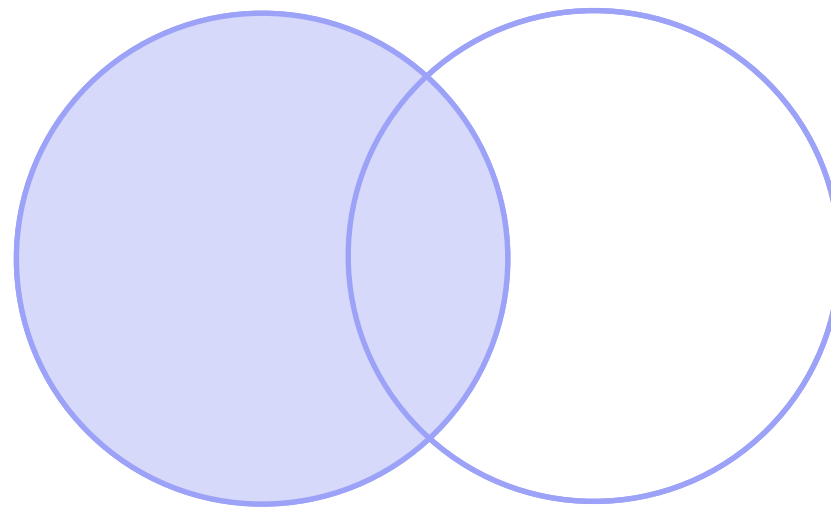
EXAMPLE: LEFT MERGE

LEFT MERGE

- Keep all records in the LEFT dataset
 - drop records in the right dataset that do not match

Left merge

(keep all records from
“left” dataset)



SYNTAX: PANDAS LEFT MERGE

The name of the
variable that you
want to "match" on

Execute a "left"
merge

```
pd.merge(left-data, right-data, on=, how= 'left')
```



A DataFrame
that you want to
merge

Another
DataFrame that
you want to merge

EXAMPLE: WE HAVE TWO SEPARATE DATASETS

- Keep all records in the LEFT dataset
 - drop records in the right dataset that do not match

emp_fst_nm

employee_id	first_name
101	John
102	Paul
103	George
104	Ringo
900	Stewart

emp_lst_nm

employee_id	last_name
101	Lennon
102	McCartney
103	Harrison
104	Starr
901	Best

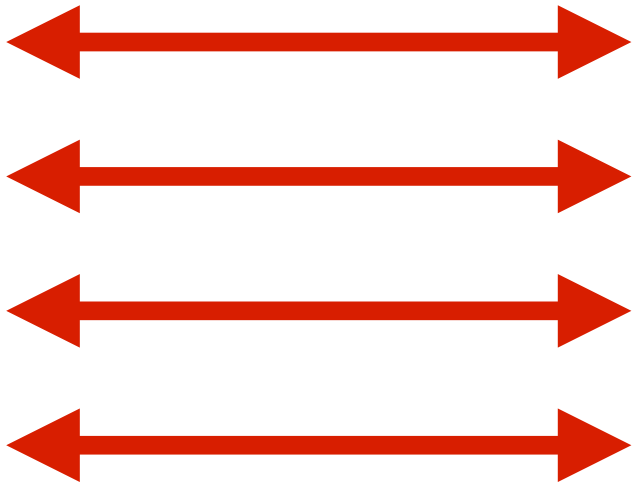
EXAMPLE: LEFT MERGE

```
pd.merge(emp_fst_nm, emp_lst_nm, on = 'employee_id', how = 'left')
```

emp_fst_nm

*

employee_id	first_name
101	John
102	Paul
103	George
104	Ringo
900	Stewart



emp_lst_nm

*

employee_id	last_name
101	Lennon
102	McCartney
103	Harrison
104	Starr
901	Best

LEFT MERGE RESULTS

```
pd.merge(emp_fst_nm, emp_lst_nm, on = 'employee_id', how = 'left')
```

- Result is a combined dataset
 - Include all rows in left dataset
 - Excludes rows from right dataset where no match was found
 - Replace missing data with missing value

employee_id	first_name	last_name
101	John	Lennon
102	Paul	McCartney
103	George	Harrison
104	Ringo	Starr
900	Stewart	NaN

RECAP

RECAP OF WHAT WE LEARNED

- You can "merge" DataFrames together with the `pd.merge()` function
- How to perform different types of merges
 - inner
 - left
- **Next Steps:** Watch the code walkthrough video for step-by-step examples of how to merge data