

# Merging Data

---

1405

Instructor: Ruiqing (Sam) Cao

# Merging Data from Different Sources

---

Suppose we need to merge two DataFrames `df1` and `df2` by matching the foreign key in `df1` to the primary key in `df2`

There are three possible match outcomes

1. primary key value exists in `df2` but no foreign key in `df1` is matched to it
2. foreign key value exists in `df2` but finds no match in `df1`
3. foreign key value exists in `df1` that is matched to a primary key value in `df2`

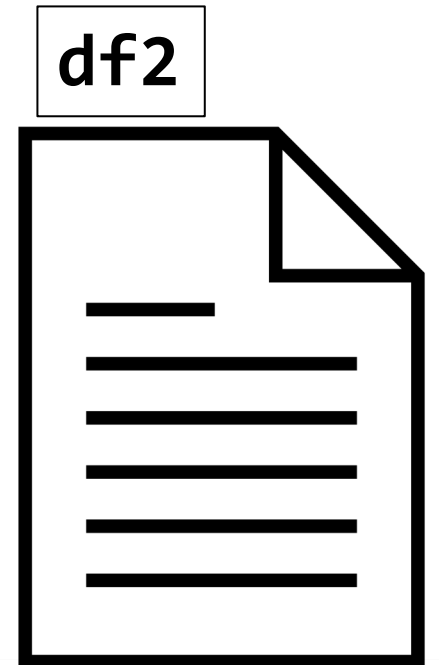
# Merging Two DataFrames

---

Suppose we need to merge two DataFrames `df1` and `df2` by matching the foreign key in `df1` to the primary key in `df2`

**df1** : The DataFrame on the *left*

**df2** : The DataFrame on the *right*



# Three Types of Joins

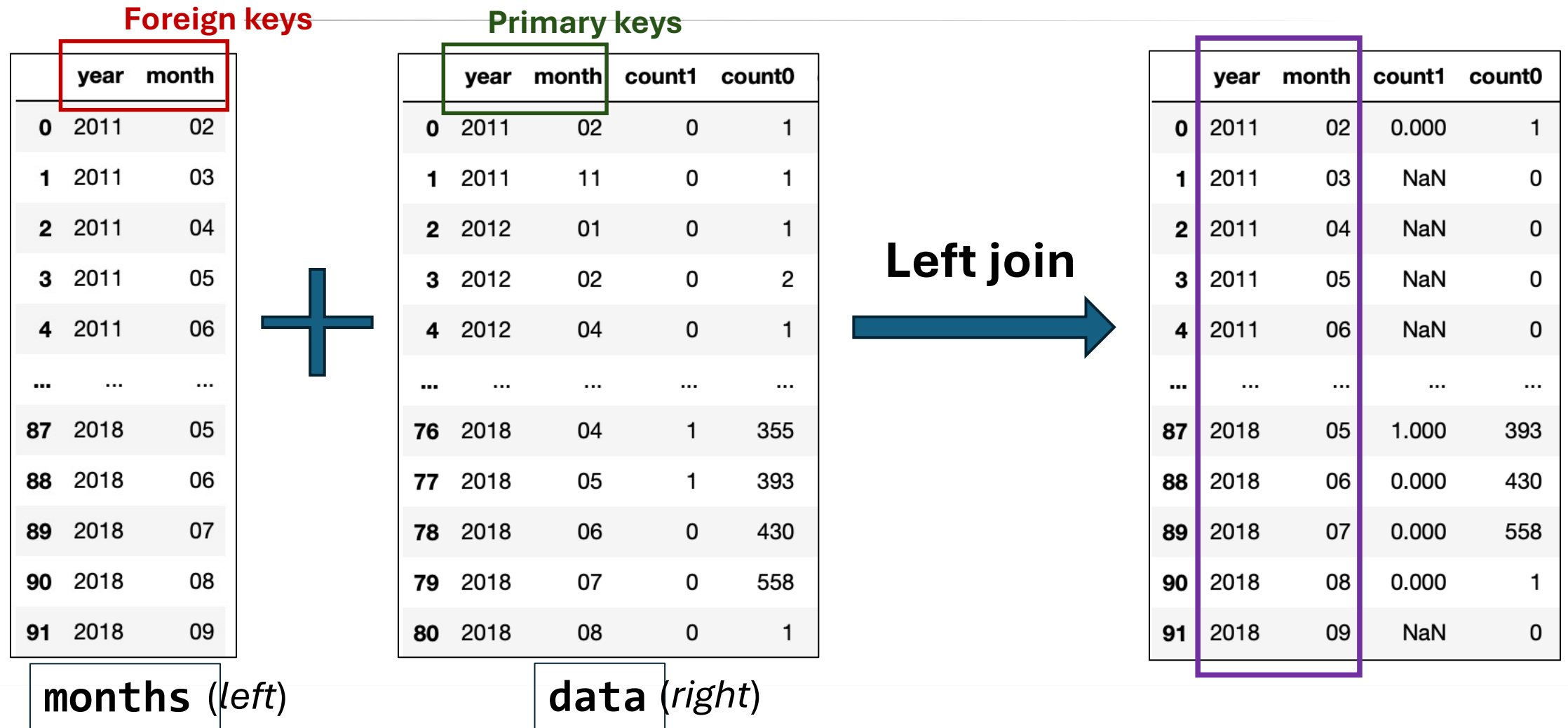
Join type	Match outcomes to keep	Match outcomes to drop
Inner	3 (matched)	1 (left only) & 2 (right only)
Left	1 (left only) and 3 (matched)	2 (right only)
Right	2 (right only) and 3 (matched)	1 (left only)
Outer	1 (left only), 2(right only), and 3 (matched)	N/A

- **Left** join preserves all information in the left-hand-side data (**df1**)
- **Right** join preserves all information in the right-hand-side data (**df2**)
- **Outer** join preserves all information both sides (**df1** and **df2**)

# Merging Two DataFrames

	<code>df1.merge(df2, how, left_on, right_on)</code>
Arguments	<code>how</code> : one of “inner”, “left”, “right”, or “outer” <code>left_on</code> : list of column names in df1 as foreign key <code>right_on</code> : list of column names in df2 as primary key
Returns	The resulting DataFrame from merging df1 ( <i>left</i> ) with df2 ( <i>right</i> )

# Merging Two DataFrames: an Example



# Notes on *Right* Join

---

```
df1.merge(df2, how='right', left_on, right_on)
```

- If you think that you need to perform right join, you should often switch the order of *df1* and *df2* and perform left join instead
- It's easier to think about merging a *more complex* data object (*df2*) onto a *simpler* one (*df1*), as implied by `df1.merge(df2)`