

PYSPARK TIP

# CLEAN EVERY COLUMN NAME IN ONE LINE

```
○○○  
1 # Before  
2 ["CustID", "trnxAmt", "ACCT-bal", "pmt$Status"]  
3  
4 # After  
5 ["cust_id", "trnx_amt", "acct_bal", "pmt_status"]
```

# THE MESS

[Sign Up to the newsletter](#)

 Databricks.news

Find more databricks tips @

[dailydatabricks.tips](#)

## EVERY SOURCE SYSTEM IS DIFFERENT

Vendor feeds, CSVs, APIs, legacy databases. Column names arrive in every format imaginable.

○ ○ ○

```
1 # Source system column names
2 df.columns = [
3 "CustID", # PascalCase
4 "custName", # camelCase
5 "trnx_date", # abbreviated
6 "ACCT-bal", # UPPER + hyphen
7 "pmt$Status", # dollar sign
8 "addr.Line1", # dot notation
9 ]
```



# ONE LINE FIX

[Sign Up to the newsletter](#)

 Databricks.news

Find more databricks tips @

[dailydatabricks.tips](#)

## DICT COMPREHENSION + .WITHCOLUMNSRENAME()

Loop over df.columns, apply a regex, pass the dict. One plan node, every column renamed.

```
○ ○ ○  
1 import re  
2  
3 # Strip special chars in one dict comprehension  
4 df_clean = df.withColumnsRenamed(  
5 {  
6 c: re.sub(r'[$\-@#!]', '_', c).lower()  
7 for c in df.columns  
8 }  
9 )
```

"ACCT-bal" becomes "acct\_bal", "pmt\$Status" becomes "pmt\_status"



# MAKE IT GENERIC

[Sign Up to the newsletter](#)

 Databricks.news

Find more databricks tips @

[dailydatabricks.tips](#)

## A REUSABLE CLEAN\_COLUMNS() FUNCTION

Strips special characters, splits camelCase, lowercases everything. Works on any DataFrame.

```
○ ○ ○  
1 def clean_columns(df):  
2     return df.withColumnsRenamed({  
3         c: re.sub(r'[$\-@#!]', '_',  
4         re.sub(r'([a-z])([A-Z])',  
5             r'\1_\2', c)).lower()  
6     for c in df.columns  
7 })  
8  
9     # Use it anywhere  
10    df_clean = clean_columns(df)
```



# EXPAND ABBREV -IATIONS

[Sign Up to the newsletter](#)

 Databricks.news

Find more databricks tips @

[dailydatabricks.tips](#)

## CHAIN IT WITH ABBREVIATION EXPANSION

"trnx\_amt" becomes "transaction\_amount", "acct\_bal" becomes "account\_balance"

```
○ ○ ○  
1 ABBR = {"trnx": "transaction",  
2 "cust": "customer",  
3 "acct": "account",  
4 "pmt": "payment",  
5 "amt": "amount",  
6 "bal": "balance"}  
7  
8 def expand(name):  
9     return "_".join(ABBR.get(p, p)  
10    for p in name.split("_"))  
11  
12 # Chain: clean then expand  
13 df.withColumnsRenamed(  
14 {c: expand(clean(c)) for c in df.columns})
```



# WANT MORE TIPS?

FOLLOW FOR DAILY  
DATABRICKS TIPS

SIGN UP TO THE NEWSLETTER



[dailydatabricks.tips](https://dailydatabricks.tips)