

Conquering confounds and covariates in machine learning and neuroscience

Pradeep Reddy Raamana



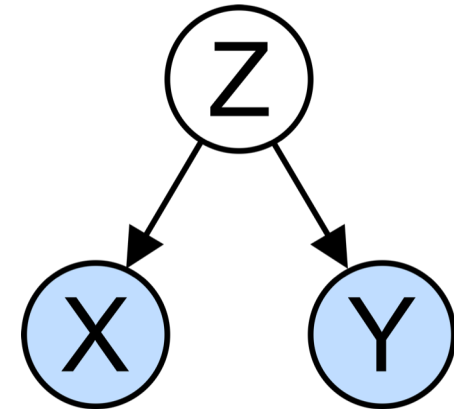
crossinvalidation.com

github.com/raamana



What are confounds?

- Causal definition[^] from statistics: A confound* is a variable (z) that influences both the dependent variable (y) and independent variable (X), causing a *spurious* association between X and y
 - X and Y are confounded by z, if z causally influence both X and y
 - To estimate the effect of X on Y, we must suppress or remove the effects of z, that influence both X and y
- Neuroimaging example: X is often imaging features, y is patient diagnosis, and z are age + gender.



[^] From Wikipedia *Also called a covariate, confounding factor, nuisance variable, lurking variable etc

Typical confounds/covariates

- Age
- Sex, Gender
- ICV, GMV, TBV
- Site, scanner
- Education, Intelligence
- Income, SES
- BMI, blood pressure, hypertension
- Lifestyle: exercise, smoking, alcohol
- Morphometrics
 - mean cortical thickness

And they are quite common in most of the studies and datasets now!

Open Challenges

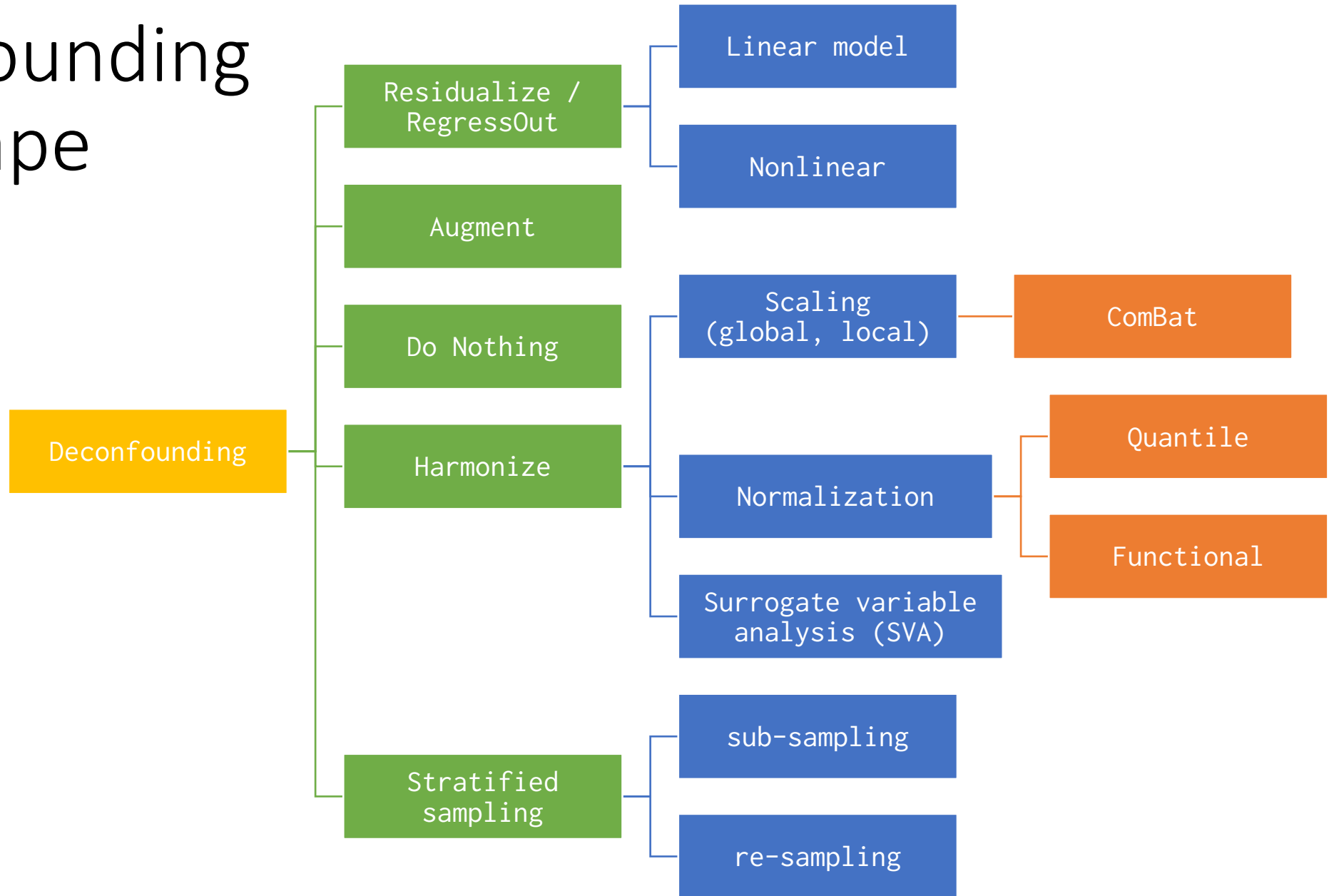
- Common approach: *pick some potential confounding variables, regress them out, and pretend everything is fine!*
- BUT life and science aren't that simple!
 - there is no clear, accepted definition and method to establish whether a given variable is indeed a confound!
 - It is often assumed they are, but their level of confounding is not properly quantified
- when should we try to de-confound it?
 - how do we know we did it right?
- how do we properly assess their impact?
 - how are the confounds / deconfounding methods affecting a given analysis?

confounds library aims to help solve these challenges

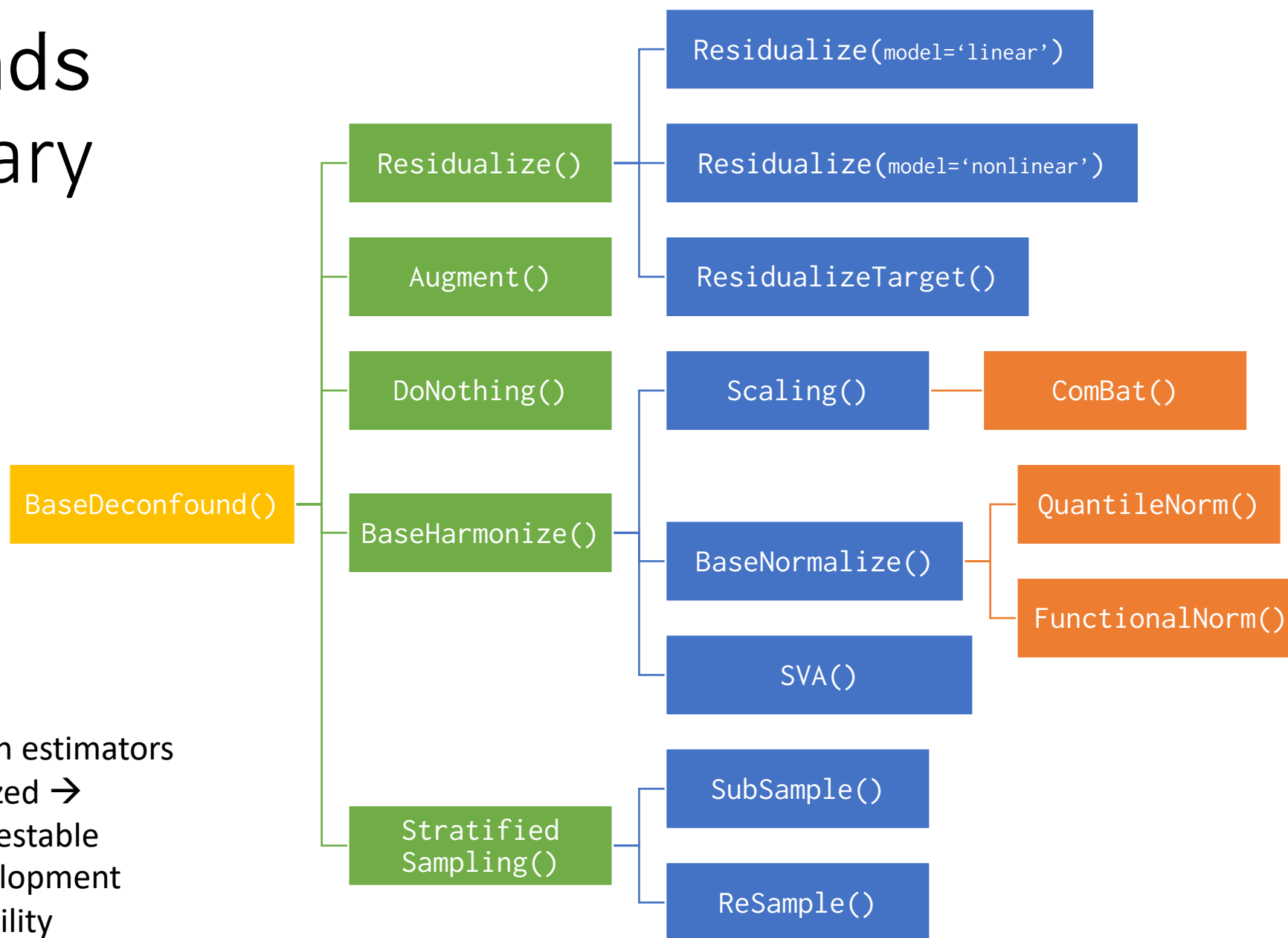
- with an open-source high-quality well-tested library
- to visualize and establish the presence of confounds
 - e.g. quantifying confound-to-target relationships
- to offer solutions to handle them appropriately via correction or removal etc
- analyze the effect of the deconfounding methods in the processed data e.g.
 - ability to check if deconfounding methods worked at all
 - **or if they introduced new or unwanted biases etc**



Deconfounding Landscape



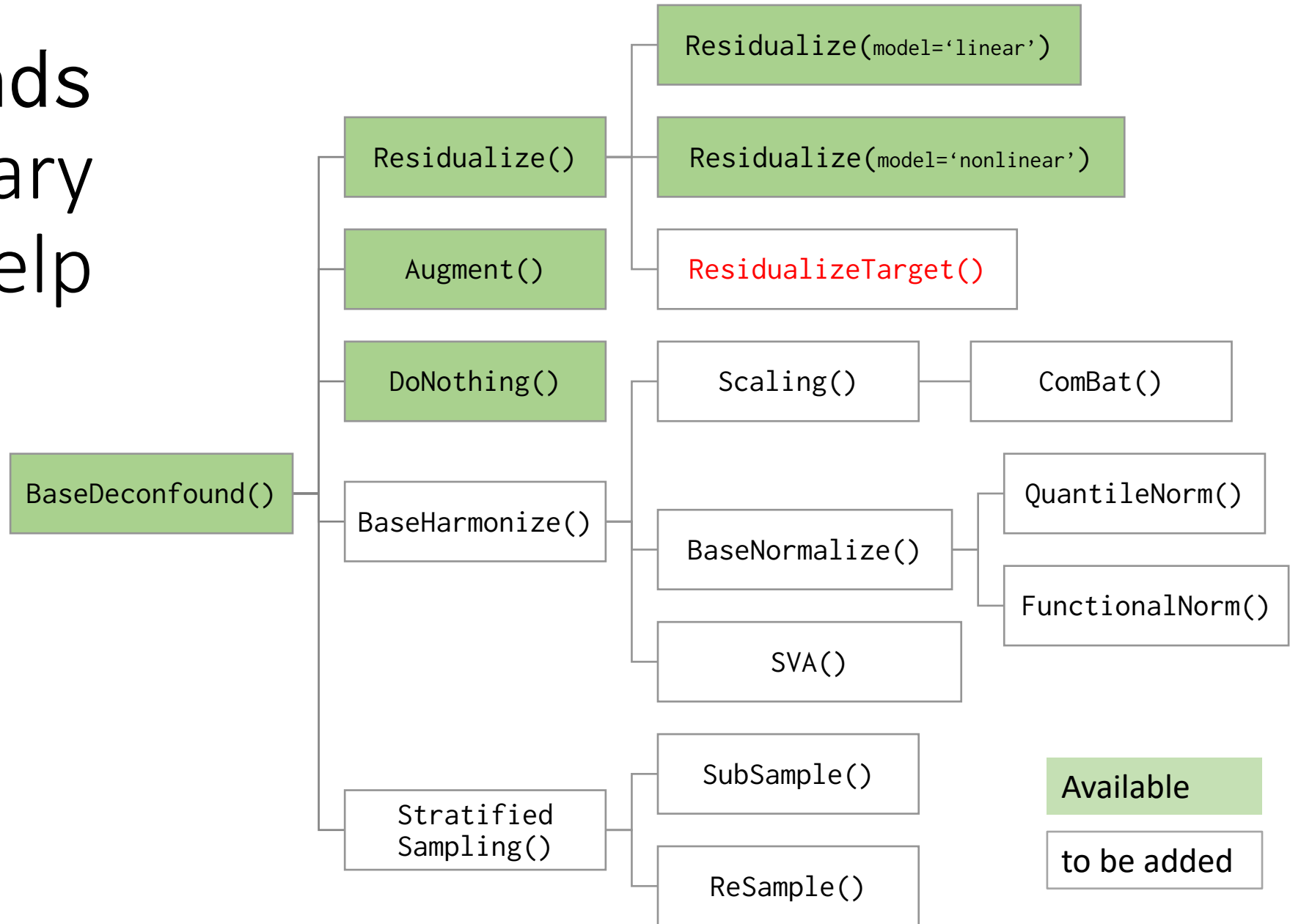
confounds library



Great features:

- act as scikit-learn estimators
- highly modularized → extensible and testable
- test-driven development
- focused on usability

confounds
library
needs help



Contribute to a great cause!

- github.com/raamana/confounds
- Focus of my OHBM'20 hackathon
 - Implement ComBat
 - from the ground up in Python
 - properly tested
 - individual components
 - as well as the whole
 - Related helpers and utils
 - viz tools
 - metrics etc

