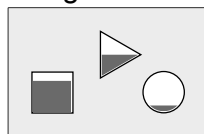


Original data



Resampling

Clustering

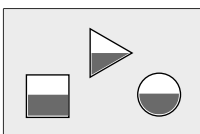
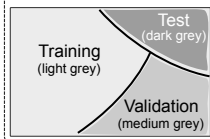
Imbalance legend

Mostly actives
Balanced
Mostly inactives

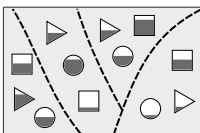
Protein A
Protein B
Protein C

Splitting sets legend

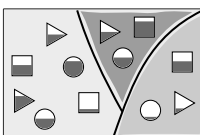
----- Clusters
----- Final splits



Clustering

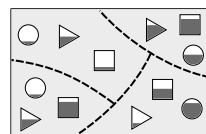


Splitting

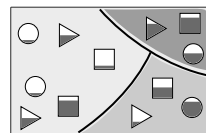


Model fit

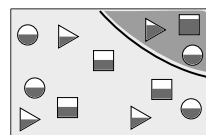
Resampling
before clustering



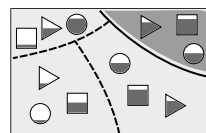
Splitting



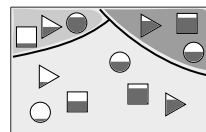
Resampling
(train/val only)



Clustering



Splitting

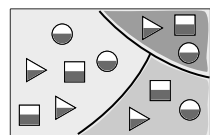


Model fit

Semi resampling

Model fit

Resampling



Model fit

Resampling
after clustering

No resampling