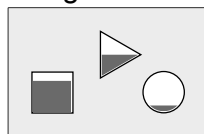
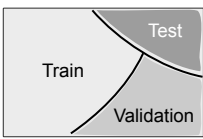


Original data



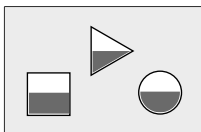
Legend

	Mostly actives	Balanced	Mostly inactives
Protein A			
Protein B			
Protein C			

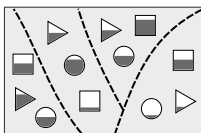


Resampling

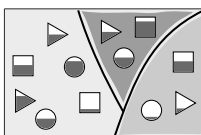
Clustering



Clustering

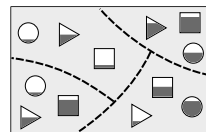


Splitting

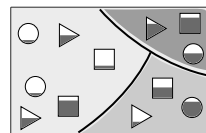


Model fit

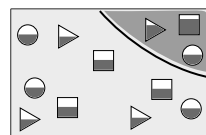
Resampling
before clustering



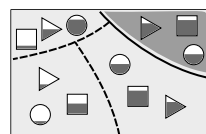
Splitting



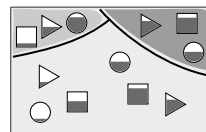
Resampling
(train/val only)



Clustering



Splitting



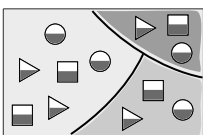
Model fit

Semi resampling

Model fit

No resampling

Resampling



Model fit

Resampling
after clustering