

Jchemo.jl

Chemometrics and machine learning on high-dimensional data with Julia

docs stable docs dev CI passing repo status Active

Lesnoff, M. 2024. UMR SELMET, Montpellier, France
<https://github.com/mlesnoff/Jchemo>



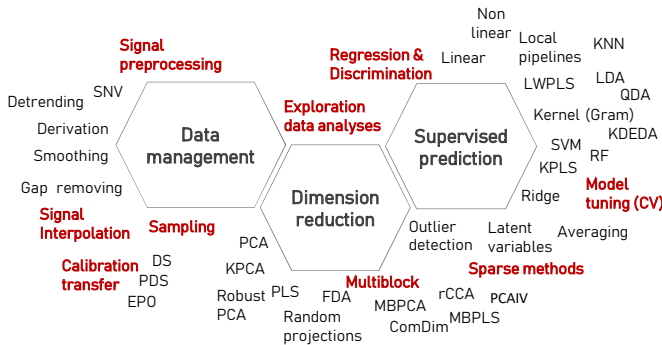
julia.org

Discourse



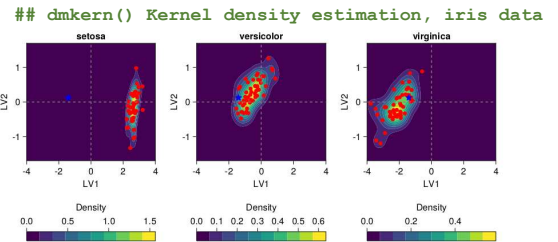
- Fast computing
 - Easy to learn/read
 - Open source
- Multiple dispatch
 Inplace functions
 Simple multi-threading, GPU
 High-quality graphics

Packages manager
 Project environments



See: <https://mlesnoff.github.io/Jchemo.jl/dev/domains>

```
## SNV preprocessing transformation
function transf!(object::Srv, X::Matrix)
    n, p = size(X)
    centr = object.par.centri
    scal = object.par.scal
    centr ? mu = rowmean(X) : mu = zeros(n)
    scal ? s = rowstd(X) : s = ones(n)
    @inbounds for j = 1:p
        X[:, j] .= (vcol(X, j) .- mu) ./ s
    end
end
```



Easy to use

Simple installation procedure

- Jchemo ∈ official Julia packages repository (equivalent to R Cran)

Consistent syntax

Generic functions of model tuning

- Grid-search cross-validation, validation-set
- Sampling designs (Kennard-Stone, Duplex, Systematic, etc.)

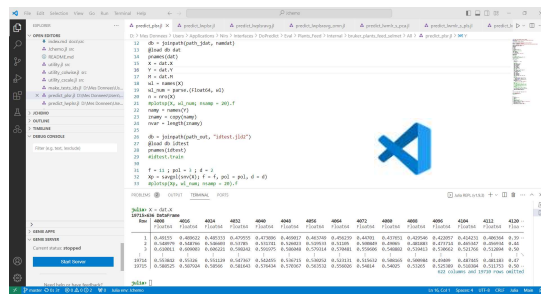
Pipelines

- mod = pip(mod1, mod2, mod3)

Each function is documented with examples

Training materials are available (Github)

```
## Pipelines
## Ex: SNV -> Savitzky-Golay preprocessing
mod1 = model(snv; centr = true, scal = true)
mod2 = model(savgol; npoint = 11, deriv = 2, degree = 3)
mod = pip(mod1, mod2)
fit!(mod, X)
Xp = transf(mod, X)
plotsp(Xp; xlabel = "Wavelength (nm)", ylabel = "Absorbance").f
```



IDE: Visual Studio Code

```
## PLS2 with 1e6 observations
Platform Info:
OS: Windows (x86_64-mingw32)
CPU: 16 x Intel(R) Core(TM) i9-10885H CPU @ 2.40GHz
n = 10^6 # nb. observations (samples)
p = 500 # nb. X-variables (features)
q = 10 # nb. Y-variables to predict
X = rand(n, p)
Y = rand(n, q)
nlv = 25 # nb. PLS latent variables

mod = model(plskern; nlv)
@time fit!(mod, X, Y)
7.532 seconds (299 allocations: 4.130 GiB, 6.58% gc time)
```

```
## Replicated K-fold CV
K = 3 # nb. folds (segments)
segm = segmkf(ntrain, K; rep = 30)
nlv = 0:20
mod = model(plskern)
res = gridcv(mod, Xtrain, ytrain; segm, score = rmsep, nlv)
plotf(res.nlv, res.yl; step = 2, xlabel = "Nb. LVs", ylabel = "RMSEP").f
```

```
## PLS-QDA
mod = model(plsqda; nlv = 25)
fit!(mod, Xtrain, ytrain)
res = predict(mod, Xtest)
errp(res.pred, ytest) > 0.093
cf = confusion(pred, ytest)
plotconf(cf).f
```

