

# CAS and PubChem Chemical Descriptors (Draft)

David Hall

08/09/2020

## Packages needed for this lab

We'll be using the `webchem` and `tidyverse` packages. If you haven't already installed them you can run the following code in the console below or install them using the "Install" button under the "Packages" tab in the bottom right window of RStudio.

```
install.packages("tidyverse")
install.packages("webchem")
```

After you've installed the packages, load them into your current R session by running the code below.

```
library(webchem)
library(tidyverse)
```

## Obtaining your assigned chemicals

## Obtaining pre-calculated chemical descriptors from PubChem

This needs to be written but I'll probably talk about:

- What are CAS numbers
- What is PubChem & pubchem CIDs
- Point by Point description of each molecular descriptor

Depending on the how we assign their chemicals, I was planning on having them construct a **data.frame from scratch and build it up column by column**. What I did below was done on a list of all 24 chemicals, and we can use it for later reference.

Note: Can't use `cir_query` because some CASRNs aren't in the *Chemical Identifier Resolver*, the data base used by the code you sent me earlier.

```
# This file is a list of all 24 compounds mentioned across the 4 projects.
data <- read.csv('CompoundsList_Env316.csv', header = TRUE, fileEncoding="UTF-8-BOM")

CASRN <- as.vector(data$CAS)

cids <- get_cid(CASRN) # get's the pubchem Compound IDs (cids)
```

```
# Get properties listed on pubchem
x <- pc_prop(cids$cid,
             properties = c("MolecularFormula",
                           "MolecularWeight",
                           "CanonicalSMILES",
                           "IUPACName",
                           "XLogP",
                           "MonoisotopicMass",
                           "TPSA", # total polar surface area
                           "Charge",
                           "Volume3D"))

# Combining everything into one file for subsequent lectures
dataPubChem <- cbind(data,x)
write.csv(dataPubChem,
          file = "CmpdListPubChem_Env316.csv",
          row.names = FALSE)
```

Table 1: The first 6 columns of our compound chemical descriptors data frame.

Project	Compound.Name	Abbreviation	CAS	CID	MolecularFormula
1	Perfluorodecanoate	PFDA	73829-36-4	21895380	C10F19O2-
1	Perfluoroundecanoate	PFUnDA	196859-54-8	23533165	C11F21O2-
1	Perfluorododecanoate	PFDoDA	171978-95-3	22174013	C12F23O2-
1	Perfluorotridecanoate	PFTriDA	862374-87-6	23084971	C13F25O2-
1	Perfluorooctane sulfonate	PFOS	1763-23-1	74483	C8HF17O3S

The calculated values scrapped from PubChem are: CAS, CID, MolecularFormula, MolecularWeight, CanonicalSMILES, IUPACName, XLogP, MonoisotopicMass, TPSA, Charge, Volume3D.

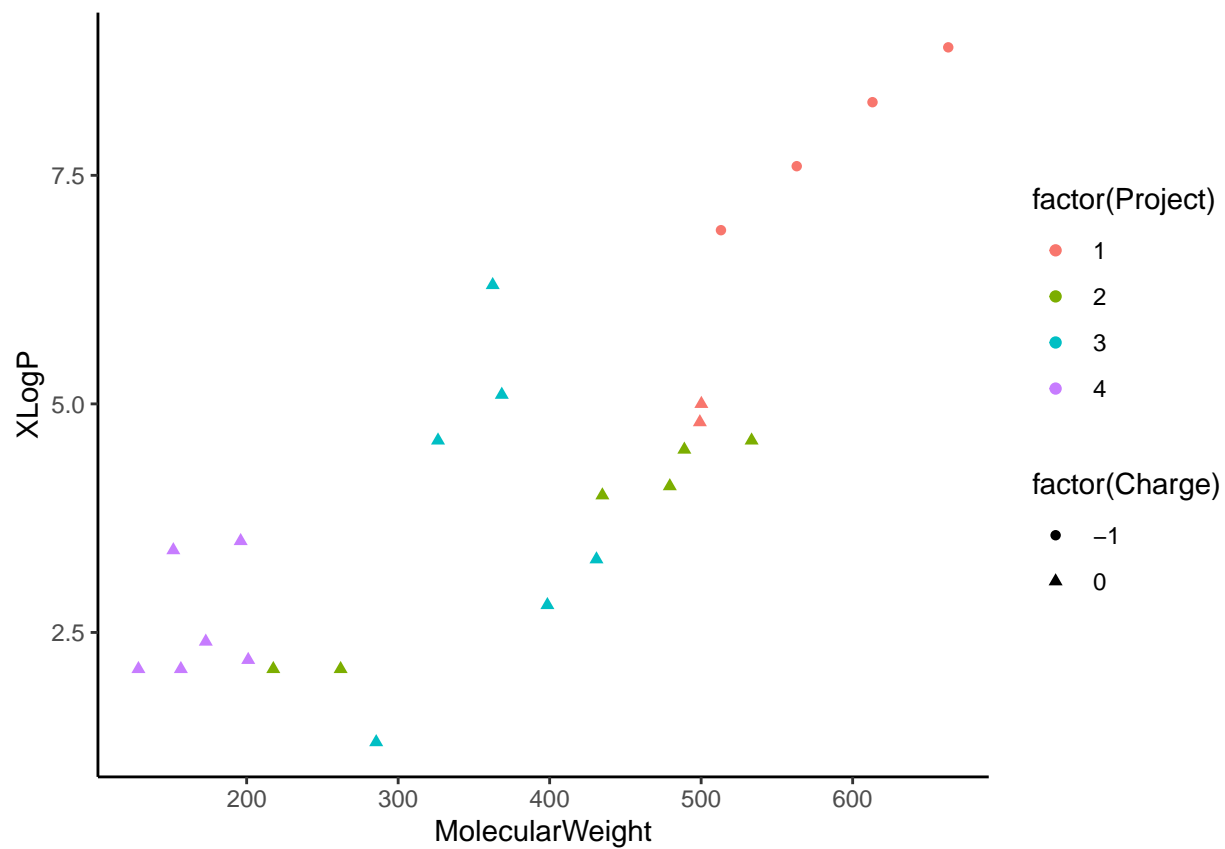
## Quick data exploration of our compounds

Need to expand, but essentially quick hypothesis to test is if polarity increases with molecular weight. The theory being that larger molecules are more complex, and hence have a lower chance of cancelling out dipole moments. We'll use XLogP as a stand-in for polarity and molecular weight as a standing for complexity. The higher the XLogP the more polar, the higher the molecular weight the greater the chemical complexity. **I know this needs a bit more justification, but i'm just trying to show them how to make plots.**

```
dataPubChem <- read.csv("CmpdListPubChem_Env316.csv", header = TRUE)

ggplot(dataPubChem, aes(x = MolecularWeight,
                        y = XLogP,
                        colour = factor(Project),
                        shape = factor(Charge))) +

geom_point() +
theme_classic()
```



From the plot in Figure 1, it appears there is a positive correlation between MW and LogP. You can see the impact of formal charges on LogP in the top right.