

Bioinformatics with python
practise proj 4/28/24

* Regression Model ~~use jupyter notebook~~

* Chembl database → can be used
for original boston projects

Pip install chembl resource client
(in python library)
→ do own due diligence on chembl
to pick what protein etc you want to
work on.

Import pandas, Import chembl client

→ Import Chembl (exact code →
From chembl_webresource_client.new_client

Import new-client

Search for target protein in
Search bar (like you search in search bar)

In this example we are searching for
'coronavirus' → (so like we're searching in
search bar coronavirus)

DATA
Collecting → TARGET SEARCH

~~(new client \Rightarrow target)~~

target search [protein/enzyme]
→ assigned as variable

then
targets is variable for

Search in dictionary

then put in data frame
print (refer to code in bioinfo
folder)

(Data frame \rightarrow with Pandas)

* You won't be able to view
the data with a py regular output
So make sure to do proj on the
jupyter notebook.

~~IMPORTANT~~

so you can load the data frame
properly, you just need to write it as
targets and NOT print

You are presented with a dataset
now you choose your protein

→ Select entry, with ID

→ returns ID

Now we retrieve bioactivity
data from the thing which
is classified by CP

So selected target → our targets ID
Activity → things which is being filtered they're

res → data [targeting things with our ID,
searching for things
considered IC50]

* It will take
a while to load so
but wait out
you can limit what appears of data frame
df. head(x)

We can check for unique types
In data frame, .unique()

"We want to define a particular standard type to make our dataset more uniform"
Standard value is potency of the drug So we want standard value to be as low as possible

IC₅₀, the standard value indicates how much of the drug you need to reach 50% inhibition effect

now we write out data frame into CSV file (commonly used to store tabular data)
we set the index to False in order to keep it from being displayed as data in CSV file. → leads to cluttering to be readable
There's just no need for it to be included | it saves it to the folder

Now we need to be able to access data from Google Drive

downloaded Google API Python client

*Important ⇒ We need to import file (CSV file) from system onto code
I couldn't find a way to import from Google Drive. but you use a pd.read (file path)

method and you define it as variable
so now we brought the data file which we downloaded onto the code, now we have filtered all standard columns which have no standard value

Data Preprocess

SORTING DATA

We Sort Data into 3 types

Active

These are all

Inactive

IC50's, so

Intermediate

We sort based on

Ex: K21000nm
is active

value (based on
nM amount)

We iterate through all instances of the
standard values in our second data frame

If --- then we append the

el.t--- label to it's class

else ---

Data set is comprised of many
compounds/molecules which are
drugs

These compounds exert some effect on
the target proteins

(ex: you take medication it exerts
effects on you)

Medicine works on that protein

So each compound has a unique
chembl ID, so some compounds

might have similar IDs, in that case we want to remove duplicates to prevent redundancy.

df2.molecule_chembl_ID →

displays chembl IDs
so we iterate through and add
to a list

do same for canonical Smiles

SMILES (standard molecular input line
Entry System)

way of
representing
molecules

overall molecules
ex: the SMILES struc-

for H₂O is "O" and "H H"
Save for standard value

We are making a new
data frame with only
the

SMILES

STANDARD IDs

and Chembl ID

Pd.Catcat → combines
multiple data frames into
a single one

Pd.Series → creates one dimensional
arrays that are tables

So we essentially made a new
data frame which sorts
the compounds to their appropriate
class based on nM, also has
Chemical IP, canonical smiles,
and standard values.
Now we create CSV file from
this