

دانشکده هوش مصنوعی، گروه کامپیوتر

عنوان

راهنمای پروژه تحلیل داده های ویروس کرونا و ارائه مناسب ترین راه مقابله دارویی با آن توسط مدل سازی داده ها و یادگیری ماشین

دانشجو

محمد حسین شوندی

استاد

دکتر امینه امینی

شماره دانشجویی

**39911541054042**

بهمن 1402

## فهرست

- بخش 1:

- مقدمه ..... 3

- بخش 2: طبقه بندی و نکته برداری

- نکته برداری های دستی ( به زبان انگلیسی) ..... 4

- توضیحات مختصر مراحل ..... 7

- هدف نهایی این پروژه ..... 8

- بخش 3: آموزش اجراء بخش 6

- توضیحات ..... 11

## مقدمه

با راهنمایی استاد محترم و علاقه‌ای که به حوزه بیوانفورماتیک دارم، خوشبختانه توانستم با این پروژه قدمی جدی برای ورود به این رشته جذاب و پر رمز و راز بردارم. این پروژه برایم فرصتی بود تا با موضوعاتی مثل یادگیری ماشین، مدل‌سازی، ساختارهای شیمیایی مولکول‌ها، و بیوتکنولوژی‌های مختلف بیشتر آشنا شوم و دانش خود را در این زمینه گسترش دهم.

این پروژه، در واقع بازسازی یکی از نمونه‌های اولیه و واقعی در حوزه بیوانفورماتیک است که توانستم با طبقه‌بندی و نکته‌برداری‌های خودم آن را به اجرا درآورم. برای درک بهتر و پیشبرد این پروژه، از منابع مختلفی همچون `chatgpt`، `GitHub`، `YouTube` و `Udemy` (لینک: [آموزش بیوانفورماتیک از پایه](#)) بهره گرفتم.

این تازه آغاز مسیر من در بیوانفورماتیک است و با توجه به دانش به‌دست‌آمده و حمایتی که از جانب استاد محترم دارم، مطمئن هستم که می‌توانم در آینده جایگاه شایسته‌ای در این رشته کسب کنم.

## با تشکر از شما.

### کتابخانه های استفاده شده:

Streamlit	pandas	base64	subprocess	os
pickle	padelpy	Numpy	seaborn	sklearn
lazypredict	scipy	chembl_webresource_client		

## نکته برداری های دستی ( به زبان انگلیسی )

در ادامه طبقه بندی ها و نکته برداری های خودم را در سه صفحه آینده به صورت عکس قرار می دهم.

~~Part 1: Retrieve & Download~~ ~~to construct ML models~~ ~~& quantitative structure-activity relationship~~  
**Part 1: Retrieve & Download to construct ML models & quantitative structure-activity relationship OR QSAR**

QSAR Allows for origin of biological activity and the interpretation of the model help to design better drug discovery

1. Import Chembl
2. search for target query
3. Get SARS Coronavirus & like Proteinase
4. filter for only IC50: the lower the number of standard value = the higher the potency of the drug. the number is like the amount of concentration needed for it to work the lower it is the better. (IC50)
5. saving ~~it~~ to CSV. with out index
6. managing missing data in standard value or Canonical smiles (Canonical smiles is the molecule structure)
7. Preprocessing Data IC50 > 10000 inactive, IC50 < 10000 Active, else intermediate
8. Molecule Chembl Id  $\Rightarrow$  comprised of many compounds or (molecules)  $\Rightarrow$  chemical structure that produces a modulatory activity or it exerts some effects on the target protein
9. final CSV with molecule Chembl Id, canonical smiles, standard value, class

**Part 2: ~~Retrieve & Download~~ Exploratory Data analysis**

1. install rdkit & Conda rdkit Allows to compute the molecular descriptors for the compounds in the dataset
2. Load CSV file
3. Lipinski rule: ① molecular weight < 500 Dalton ② Octanol-water partition Coefficient (LogP) < 5 ③ hydrogen bond donors < 5 ④ Hydrogen bond acceptors < 10

it shows molecular likeness & Absorption, distribution, metabolism, excretion or ADME  
Pharmacokinetic Profile

- ④ IC<sub>50</sub> to pIC<sub>50</sub> because the original IC<sub>50</sub> value has uneven distribution of the data points. we should Apply Negative Logarithmic transformation
- ⑤ Cap IC<sub>50</sub> value to 100,000,000 cause other wise -Log will be Negative value
- ⑥ Apply pIC<sub>50</sub> function and change pIC<sub>50</sub> to IC<sub>50</sub> value
- ⑦ Remove All intermediate bioactivity class
- ⑧ Data Analysis for chemical space, Active molecule is Like constellation, we have Constellation Plot for chemical space analysis where the more Active the molecule the Larger size it has
- 8-1 frequency plot for bioactivity class using Count Plot
- 8-2 Scatter plot of MW (molecular weight) vs LogP using Scatter plot
- 8-3 using boxplot for pIC<sub>50</sub> 6 > active 5 < inactive
- 8-4 using Mann-whitney U test for statistical significance of the difference
- 8-5 H<sub>0</sub> (Null Hypothesis) says there is no difference And difference is Accidental but our p-value is close to zero (less than Alpha=0.05) so it's Not valid.
- ⑨ using boxplot & Mann-whitney for all 4 of Lipinski values

### Part 3g

#### Descriptor calculation

- ① PaDEL-Descriptor Download & tool for calculating molecular descriptors And fingerprint from chemical structure it is to Analyze & represent the properties for molecules in quantitative manner its java-based APP that shows 1 molecular Descriptors 2 molecular fingerprint 3 Compatibility 4 customizable Descriptor types 5 Batch Processing
- pubchem's Binary-based value that encodes the presence or Absence of specific structural features. used for molecular similarity searching, clustering & classification
- ② Difference between Lipinski-Descriptors & pubchem's Lipinski is global features (the value of 1 or 0) but Pubchem is Local features & each molecule Like containing several Lego Blocks can have variety of structures And the way Lego blocks are connected will create a Unique set of properties And that's the essence of Drug Discovery.
- we have to find a way that the molecule provides the most potency and safe.





## در ادامه توضیح مختصری راجب پروژه به شما ارائه می دهیم:

این پروژه نحوه استفاده از پایتون و یادگیری ماشین در کشف دارو را ارائه می دهد و شامل جمع آوری داده، پیش پردازش، تحلیل اکتشافی، ساخت مدل و استقرار آن به عنوان یک وب اپلیکیشن است.:

همچنین به مقایسه مدل ها و انتخاب مناسب ترین مدل پرداخته و یک اینفوگرافیک خلاصه برای مرور آسان فراهم می شود.

(برای اجرای بخش 1 تا 5 با استفاده از Jupyter Notebook فقط کافی است فایل ها اجرا شود. کتابخانه های لازمه به صورت خودکار نصب می شوند).

معرفی پایگاه داده ChEMBL، که شامل بیش از 2 میلیون ترکیب شیمیایی از 76,000 سند است، و نصب ChEMBL web resource client برای دسترسی به داده های فعالیت زیستی.

کاوش در پایگاه داده ChEMBL، شامل ورژن و تعداد ترکیبات.

نصب ChEMBL client با استفاده از pip برای دانلود داده های فعالیت زیستی.

نحوه دریافت داده های بیواکتیویته برای پروتئین هایی مانند پروتئیناز ویروس کرونا SARS، با تمرکز بر قدرت تأثیر و یکپارچگی داده ها.

نمایش مراحل پیش پردازش داده های فعالیت زیستی، شامل دسته بندی ترکیبات بر اساس مقادیر IC50 به فعال، غیرفعال یا متوسط.

ایجاد پوشه ها، بررسی محتوای داده ها، و دسته بندی ترکیبات بر اساس IC50 برای استفاده در مدل های یادگیری ماشین در تحقیق دارویی.

محاسبه توصیف گرهای مولکولی با استفاده از RDKit و معرفی توصیف گرهای لیپینسکی، با توضیح اصول ارزیابی داروپذیری براساس ویژگی های فارماکوکینتیک.

توضیح مقیاس بندی مقادیر IC50 به pIC50 برای تحلیل در پروژه بیوانفورماتیک.

نمایش تکنیک های دست کاری و تحلیل داده ها در بیوانفورماتیک، شامل افزودن ستون های جدید، انجام تحلیل اکتشافی، و انجام آزمون های آماری روی کلاس های فعالیت زیستی.

بررسی داده‌های اکتشافی و تحلیل فضای شیمیایی با نمودارها و ذخیره نمودارها برای گزارش.

استفاده از ابزارهایی مانند Paddle برای محاسبه سریع‌تر توصیف‌گرهای مولکولی و بهبود دقت پروژه‌های بیوانفورماتیک.

دانلود، پردازش و آپلود فایل‌های محاسبه‌شده و آماده‌سازی ماتریس‌های داده برای مدل رگرسیون با الگوریتم جنگل تصادفی، با استفاده از اثر انگشت‌های PubChem برای بازدارنده‌های آنزیم استیل کولین استراز.

ساخت و مقایسه مدل‌های یادگیری ماشین با کتابخانه lazy predict، با تمرکز بر ارزیابی عملکرد مدل و ایجاد وب‌اپلیکیشن برای پیش‌بینی‌های کاربران.

نحوه بارگذاری فایل ورودی حاوی داده ساختار شیمیایی، تولید اثر انگشت‌های مولکولی، پیش‌بینی مقادیر فعالیت زیستی، و بهینه‌سازی مدل با کاهش ویژگی‌های تکراری برای تسریع در ساخت مدل.

ساخت مدل یادگیری ماشین با استفاده از رگرسیون جنگل تصادفی و ذخیره مدل با استفاده از Pickle برای استفاده‌های آینده.

ساخت یک ابزار بیوانفورماتیک با پایتون، شامل بارگذاری داده‌ها، محاسبه توصیف‌گرهای مولکولی، ساخت مدل و ارائه پیش‌بینی‌ها، با تأکید بر یادگیری علم داده از طریق کاربرد عملی.

### هدف پروژه:

هدف نهایی این پروژه توسعه یک مدل یادگیری ماشین است که بتواند فعالیت زیستی ترکیبات را در برابر یک پروتئین هدف خاص پیش‌بینی کند. این کار با استفاده از مقایسه IC50 و PIC50 که نشان دهنده قدرت اثر این ترکیب در مقابل پروتئین هدف است.

IC50 یا "غلظت بازدارنده نیمه‌حداکثری" یک معیار رایج در بیوشیمی و داروسازی است که نشان می‌دهد چه غلظتی از یک ترکیب (مثل دارو) لازم است تا 50٪ از فعالیت یک آنزیم یا گیرنده خاص را مهار کند. این معیار به‌عنوان شاخصی از قدرت و کارایی یک مهارکننده استفاده می‌شود؛ هرچه مقدار IC50 کمتر باشد، دارو یا ترکیب مؤثرتر است و می‌تواند با غلظت کمتری فعالیت زیستی موردنظر را مهار کند.



## اجرای بخش 6

اول از همه در فایل part 6 این فایل را اجرا کرده : bioactivity\_prediction\_app.ipynb تا نتایجی مانند پنج بخش قبلی را در csv برای اجرای webapp آماده کنیم.

سپس از ترمینال anaconda prompt که در بخش دو نصب کردیم دستورات زیر را وارد کرده:

```
Anaconda Prompt - streamlit run app.py
(base) C:\Users\pc>cd C:\Users\pc\Desktop\PROJECT\final\part 6\
(base) C:\Users\pc\Desktop\PROJECT\final\part 6>conda create -n my_env
Channels:
 - defaults
Platform: win-64
Collecting package metadata (repodata.json): done
Solving environment: done

## Package Plan ##

  environment location: C:\Users\pc\miniconda3\envs\my_env

Proceed ([y]/n)? y
Preparing transaction: done
Verifying transaction: done
Executing transaction: done
#
# To activate this environment, use
#
#   $ conda activate my_env
#
# To deactivate an active environment, use
#
#   $ conda deactivate

(base) C:\Users\pc\Desktop\PROJECT\final\part 6>conda activate my_env
(my_env) C:\Users\pc\Desktop\PROJECT\final\part 6>streamlit run app.py

You can now view your Streamlit app in your browser.

Local URL: http://localhost:8501
Network URL: http://192.168.12.66:8501
```

همانطور که در عکس میبیند با استفاده از streamlit برنامه app.py را برای اجرای بخش 6 اجرا می کنیم.

پس از اجرا با صفحه زیر مواجه می شویم:

>

3	CHEMBL426082	1	1	0	0	0	
---	--------------	---	---	---	---	---	--

(882 , 4)

### زیر مجموعه ای از توصیف گر ها از مدل های ساخته شده قبلی

	PubchemFP2	PubchemFP12	PubchemFP14	PubchemFP15	PubchemFP16	PubchemFP20	Pubche
0	0	0	1	1	0	0	
1	0	0	1	0	0	0	
2	0	1	1	0	0	1	
3	0	1	1	0	0	0	

(254 , 4)

### خروجی پیش بینی

	molecule_name	pIC50
0	CHEMBL187579	4.8953
1	CHEMBL188487	4.9333
2	CHEMBL185698	4.9768
3	CHEMBL426082	5.2424

[دانلود پیش بینی ها](#)

#### 1. فایل داده CSV خود را بارگذاری کنید

فایل ورودی خود را بارگذاری کنید

Drag and drop file here

Limit 200MB per file • TXT

Browse files

× example\_coronavirus.txt  
205.0B

[فایل ورودی نمونه](#)

پیش بینی

با استفاده از فایل ورودی نمونه در سمت راست می توانید به **github** من و فایل **txt** داخلش دسترسی پیدا کنید که داخل آن چهار نمونه پروتئین آزمایشی هست.

با اضافه کردن فایل **txt** و زدن دکمه پیش بینی می توان بهترین گزینه ها برای استفاده در دارو سازی که کمترین **PIC50** را دارند را به دست آورد.

(نکته: فرمت فایل ورودی باید مانند فایل نمونه باشد که شامل نام مولکول و ترکیب **chembl** و ترکیب شیمیایی آن مولکول است.)

با تشکر از صبر و لطف شما.