

# Predictive Analysis in Identifying Frequent Opioid Prescriber

Erjia Meng, Jiayi Hao, Shiyu Liu, Xinyi Liu, Yulin Feng

# Opioid & Opioid Epidemic

- Prescription Opioids - a class of drugs used treat moderate to severe pain
- Opioids are highly addictive, and also have serious risks and side effects
- Opioids were involved in 49,860 overdose deaths in 2019  
(70.6% of all drug overdose deaths).



# Opioid Over-Prescribers & “Pill Whales”

Pill Whale - Doctors who are over-prescribing Opioid drugs, and thus gain profit from medical representative.

First Step to find over-prescriber:

Identify Frequent Opioids Prescriber



Main Goal of our project





# Identify Frequent Opioids Prescriber

Predictive Model:

- Systematic method with Machine Learning techniques
- Dataset: 25,000 unique licensed medical professionals in the United States
- Independent Variable - Basic information of the medical professionals
- Target Variable - whether the professional is a 'Frequent Opioid Prescriber'

Frequent Opioid Prescriber  $\neq$  Over-Prescriber



Allow further investigation to identify over-prescribers

# Datasets





# Datasets

- Main dataset `prescriber-info.csv`: basic information of some medical professionals in 2014 in the United States and their prescription records for hundreds of common opioid and non-opioid drugs
- `overdoses.csv`: contains information on opioid related drug overdose fatalities

Source: <https://www.kaggle.com/datasets/apryor6/us-opiate-prescriptions>



# Datasets

## prescriber-info.csv

- **NPI** - unique National Provider Identifier number
- **Gender** - M/F
- **State** - U.S. State by abbreviation
- **Credentials** - set of initials indicative of medical degree
- **Specialty** - description of type of medicinal practice
- **A long list of drugs** - with numeric values indicating the total number of prescriptions written for the year by that individual
- **Opioid.Prescriber** - a boolean label indicating whether or not that individual prescribed opiate drugs more than 10 times in the year



# Datasets

## overdose.csv

- **State** - full name of U.S. State
- **Population** - population every state
- **Death** - number of death caused by overdose
- **Abbrev** - U.S. State by abbreviation



# Clean data





# prescriber-info

- **NPI** - unique National Provider Identifier number ← drop
- **Gender** - M/F ← To digit
- **State** - U.S. State by abbreviation ←
- **Credentials** - set of initials indicative of medical degree
- **Specialty** - description of type of medicinal practice
- **A long list of drugs** with numeric values indicating the total number of prescriptions written for the year by that individual
- **Opioid.Prescriber** - a boolean label indicating whether or not that individual prescribed opiate drugs more than 10 times in the year



# Overdose.csv

- **State** - full name of U.S. State
- **Population** - population every state
- **Death** - number of death caused by overdose
- **Abbrev** - U.S. State by abbreviation



**Death per capita**



# Dealing with “State” feature

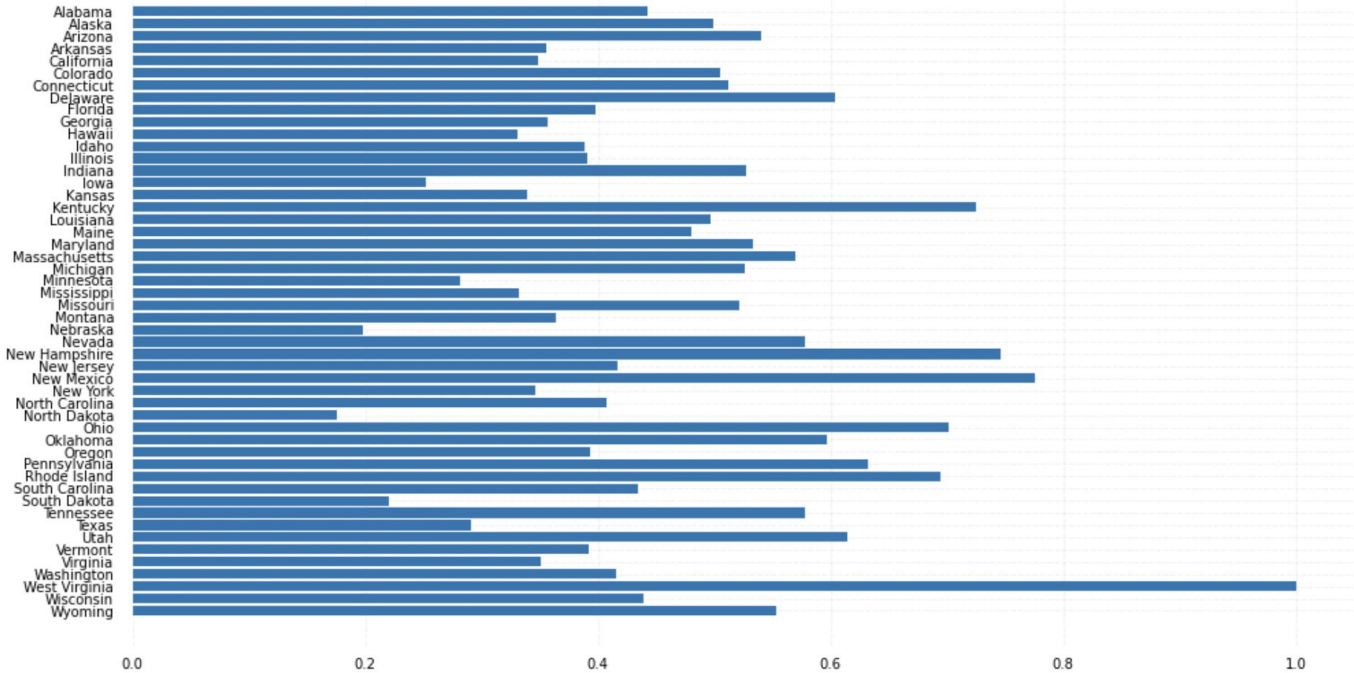
```
np.sort(arr1) State feature from prescriber_info.csv Washington D.C. 79 rows
```

```
array(['AA', 'AE', 'AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE',  
      'FL', 'GA', 'GU', 'HI', 'IA', 'ID', 'IL', 'IN', 'KS', 'KY', 'LA',  
      'MA', 'MD', 'ME', 'MI', 'MN', 'MO', 'MS', 'MT', 'NC', 'ND', 'NE',  
      'NH', 'NJ', 'NM', 'NV', 'NY', 'OH', 'OK', 'OR', 'PA', 'PR', 'RI',  
      'SC', 'SD', 'TN', 'TX', 'UT', 'VA', 'VI', 'VT', 'WA', 'WI', 'WV',  
      'WY', 'ZZ'], dtype=object) Puerto Rico 231 rows
```

```
np.sort(arr2) State feature from overdose.csv
```

```
array(['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DE', 'FL', 'GA', 'HI',  
      'IA', 'ID', 'IL', 'IN', 'KS', 'KY', 'LA', 'MA', 'MD', 'ME', 'MI',  
      'MN', 'MO', 'MS', 'MT', 'NC', 'ND', 'NE', 'NH', 'NJ', 'NM', 'NV',  
      'NY', 'OH', 'OK', 'OR', 'PA', 'RI', 'SC', 'SD', 'TN', 'TX', 'UT',  
      'VA', 'VT', 'WA', 'WI', 'WV', 'WY'], dtype=object)
```

Death per capita caused by opioid overdose per state (Normalized)





# Features

Gender (0/1)

State (str)

Credentials (str)

Specialty (str)

List of Drugs (0/1)



**Add a death\_Rate\_norm column**

Death\_Rate\_norm for “DC” and “PR” is filled by mean value



# Features

Gender (0/1)

State (str) ← **Categorical feature**

Credentials (str)

Specialty (str)

List of Drugs (0/1)

Death\_Rate\_norm (float between 0 and 1)

# One-hot encoding on State

```
features_df.head()
```

	Gender_is__M	State_is__AL	State_is__AR	State_is__AZ	State_is__CA	State_is__CO	State_is__CT
index							
0	1	0	0	0	0	0	0
1	0	1	0	0	0	0	0
2	0	0	0	0	0	0	0
3	1	0	0	1	0	0	0
4	1	0	0	0	0	0	0

5 rows × 346 columns





# Features

Gender (0/1)

List of State\_is\_ (0/1)

Credentials (str)


Specialty (str)

List of Drugs (0/1)

Death\_Rate\_norm (float between 0 and 1)



# prescriber-info

- **NPI** - unique National Provider Identifier number
- **Gender** - M/F
- **State** - U.S. State by abbreviation
- **Credentials** - set of initials indicative of medical degree 
- **Specialty** - description of type of medicinal practice
- **A long list of drugs** with numeric values indicating the total number of prescriptions written for the year by that individual
- **Opioid.Prescriber** - a boolean label indicating whether or not that individual prescribed opiate drugs more than 10 times in the year

# Credentials

```
df.Clean_Credentials.unique()
```

```
array(['DDS', 'DMD', 'MD', 'NoRecord', 'DDS MS', 'DDSPLLC', 'DMD MS',  
      'DDSMS', 'DDS PA', 'DDS MPH', 'DDSMAGD PC', 'DMDPC', 'DENTIST',  
      'DD S', 'DDS PC', 'DDS PHD', 'DR', 'DMDFACD', 'DDS MD',  
      'DDS M', 'ISTDDS',  
      'DDS M', 'MAGD',  
      'HERSC', 'DMD MMSC',  
      'BDS I', 'DDS MS PHD',  
      'DDS RN', 'APC DDS', 'BS DDS', 'DDS FACD', 'DMD MSD', 'EIN',  
      'DDSPRACTICELIMITE', 'DDS GREGSAWYER', 'DMDMS', 'DMD PC',  
      'DDS PHARMD', 'BDS', 'DDSPA', 'DDS PD', 'DO', 'DDSFAGD', 'DDSPS',  
      'DMD MDS', 'DMDDRMEEDDENT', 'DDS MS MSPH', 'MMD', 'DDS PS',  
      'DDS PLLC', 'DDA', 'DMD MPH', 'DMDBA', 'DMDPA', 'MSDMD', 'OD',  
      'ODPC', 'MD OD', 'DOCTOROFOPTOMETRY', 'OPTOMETRIST', 'DO DOS',  
      'OO', 'FAAO OD', 'MED OD', 'FAAO', 'MS OD', 'FAAO FOAA OD',  
      'MPH OD', 'OD PA', 'ODPS', 'MDPA', 'MDPC', 'MD MPH', 'MD',
```

```
len(df.Clean
```

```
ue())
```

Drop It

```
618
```

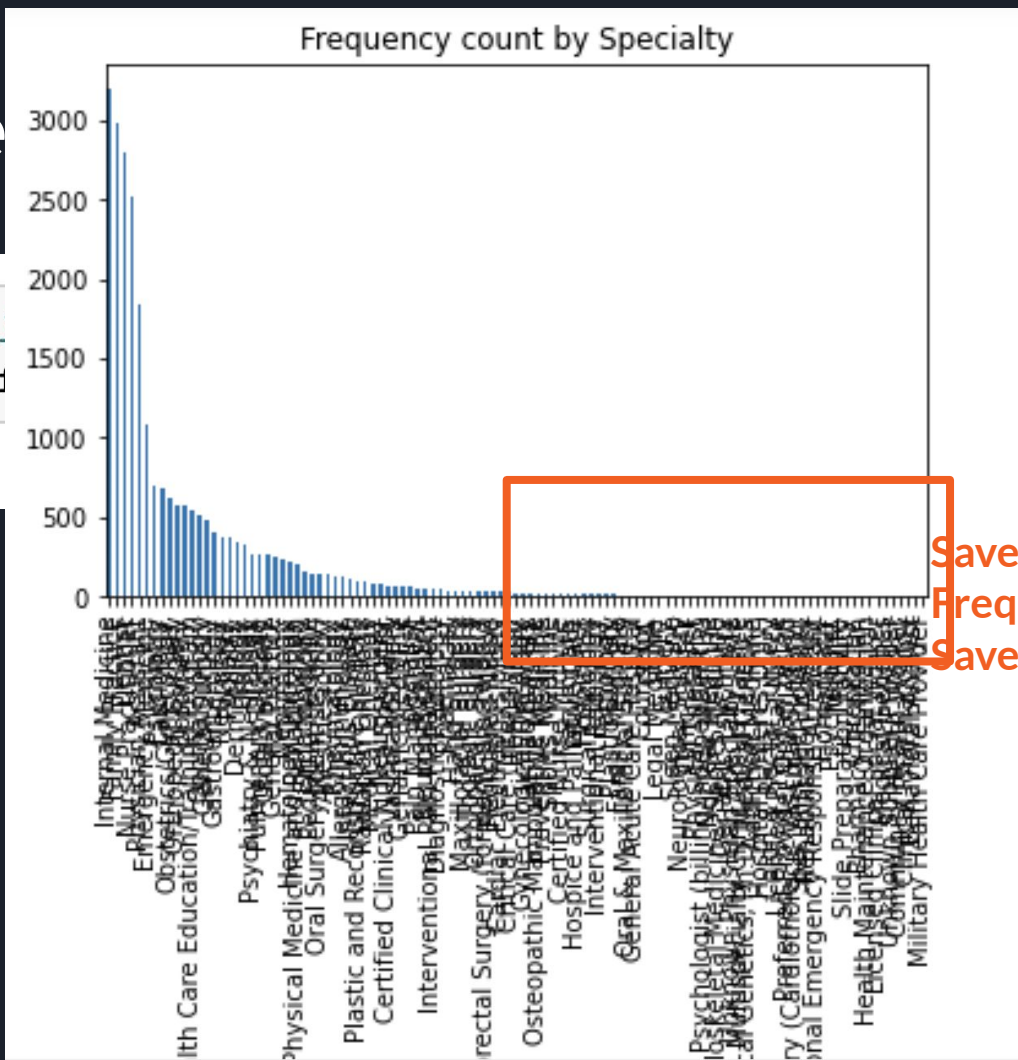


# prescriber-info

- **NPI** - unique National Provider Identifier number
- **Gender** - M/F
- **State** - U.S. State by abbreviation
- **Credentials** - set of initials indicative of medical degree
- **Specialty** - description of type of medicinal practice ←
- **A long list of drugs** with numeric values indicating the total number of prescriptions written for the year by that individual
- **Opioid.Prescriber** - a boolean label indicating whether or not that individual prescribed opiate drugs more than 10 times in the year

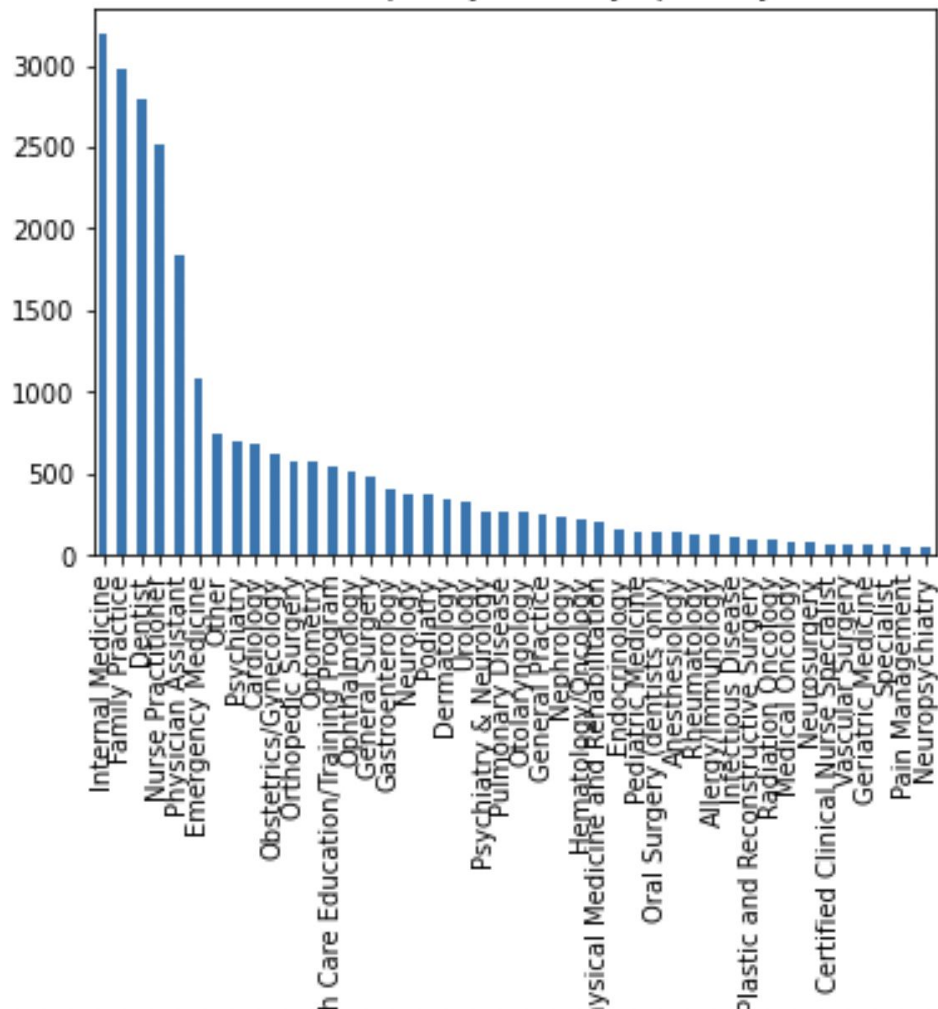
```
#now len(d)
```

109



SP

# New Frequency count by Specialty





# One-hot encoding on Specialty

```
features_list
```

```
'Specialty_is__Anesthesiology',  
'Specialty_is__Cardiology',  
'Specialty_is__Certified Clinical Nurse Specialist',  
'Specialty_is__Dentist',  
'Specialty_is__Dermatology',  
'Specialty_is__Emergency Medicine',  
'Specialty_is__Endocrinology',  
'Specialty_is__Family Practice',  
'Specialty_is__Gastroenterology',  
'Specialty_is__General Practice',  
'Specialty_is__General Surgery',  
'Specialty_is__Geriatric Medicine'
```



# Features

Gender (0/1)

List of State\_is\_ (0/1)

List of Specialty\_is\_ (0/1)

List of Drugs (0/1)

Death\_Rate\_norm (float between 0 and 1)



Now There is a total of

```
len(features_df.iloc[0,:])
```

346

Features

And is ready for the next step

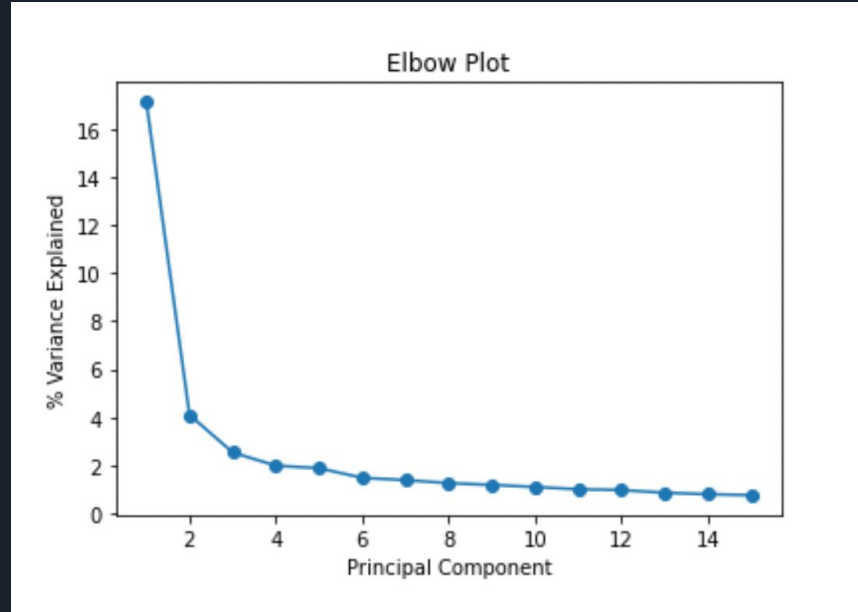
# Dimension Deduction



# Dimension Deduction

## Elbow Plot

- 2-3 variables

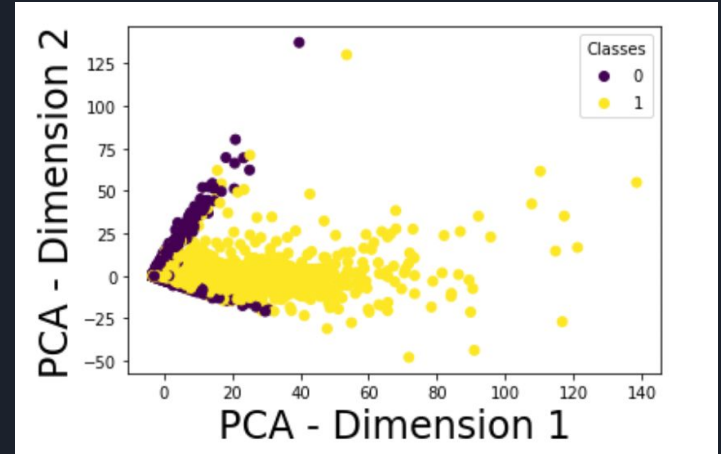


# Dimension Deduction

- Visualization

## 2D Map of Opioid Prescriber

- Obvious trend in two classes



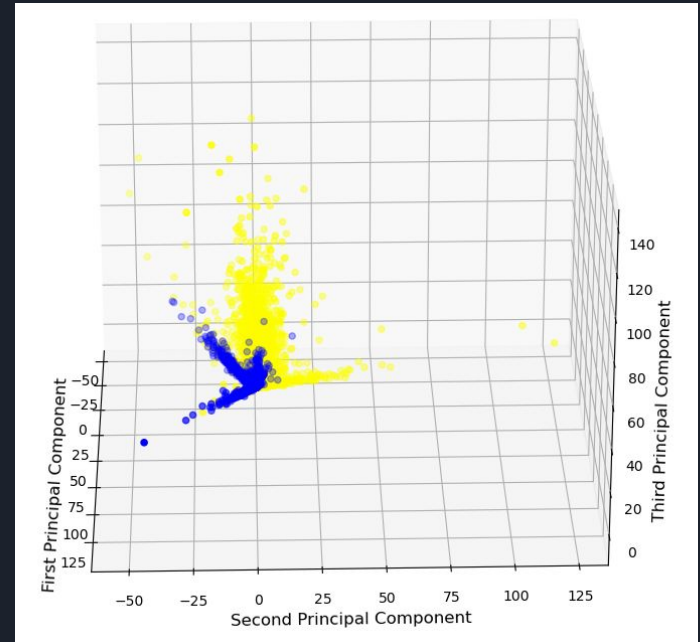
# Dimension Deduction

- Visualization

## 3D Map of Opioid Prescriber

- More obvious distribution

Yellow: 1    Blue: 0





# Dimension Deduction

- Eliminate Dimension
  - Reduce number of attributes to 211 with 90% of the variance in the target variable explained
  - Original attributes: 346
- Good Visualization
- Not very significant deduction with a loss of 10%

# Methods & Models





# Methods & Models

- Logistic Regression
- Decision Tree
- Random Forest (max\_depth = 100)
- KNN Classifier (n\_neighbors = 17)
- SVM Classifier (kernel = 'rbf')
- Gaussian Naive Bayes
- Classification\_report of confusion matrix
- ❖ Focus on the prediction of 1 of Opioid Prescriber





# Measures

- Target\_variable: Opioid.Prescriber (frequent (1) or not (0)).
- Overall precision: Percentage of correct prediction of both 1 and 0 .
- Precision of 1: Percentage of correct prediction of frequent opioid prescribers.
- Recall: Percentage of predicted frequent opioid prescribers out of all frequent opioid prescribers in real situation.
- F1 score: harmonic mean of precision and recall



# Precision of predictions

Model's name	Overall precision	Precision of 1
Logistic Regression	0.9213685474189676	<u>0.98</u>
Decision Tree	0.9039615846338536	<u>0.93</u>
Random Forest	<u>0.9227691076430572</u>	0.96
KNN Classifier	<u>0.8723489395758304</u>	0.95
SVM Classifier	0.8809523809523809	<u>0.98</u>
Gaussian Naive Bayes	0.9099639855942377	0.95



# Recall of predicting 1s

Model's name	Recall
Logistic Regression	0.88
Decision Tree	<u>0.91</u>
Random Forest	<u>0.91</u>
KNN Classifier	<u>0.82</u>
SVM Classifier	<u>0.82</u>
Gaussian Naive Bayes	0.89



## F-1 score of predicting 1s

Model's name	F-1 score
Logistic Regression	<u>0.93</u>
Decision Tree	0.92
Random Forest	<u>0.93</u>
KNN Classifier	<u>0.88</u>
SVM Classifier	0.89
Gaussian Naive Bayes	0.92



## Key features

HYDROCODONE.ACETAMINOPHEN	0.240765
OXYCODONE.ACETAMINOPHEN	0.061513
TRAMADOL.HCL	0.060548
OXYCODONE.HCL	0.028386
ACETAMINOPHEN.CODEINE	0.020295
PREDNISON	0.020028
GABAPENTIN	0.016958
Specialty_is_Emergency Medicine	0.014685
AMOXICILLIN	0.014468
Death_Rate_norm	0.013344



# Run Time

Model's name	Run Time (second)
Logistic Regression	0.6
Decision Tree	0.8
Random Forest	3.9
KNN Classifier	5.4
SVM Classifier	<u>48.6</u>
Gaussian Naive Bayes	<u>0.2</u>



# Conclusion

- Our goal: identify frequent opioid prescribers
- Data cleaning -> PCA -> Prediction
- Which is the optimum model? Why?
- Implications : systematic, efficient, narrow down the pool of potential overdoses
- Potential improvements



# Reference

1. CDC Opioids [Online] Available from: <https://www.cdc.gov/opioids/index.html>
2. "U.S. Opiate Prescriptions/Overdoses". Accessed on: Apr. 15, 2022. [Online]. Available: <https://www.kaggle.com/datasets/apryor6/us-opiate-prescriptions?select=prescriber-info.csv>
3. "Drug Overdose Deaths in the U.S. Top 100,000 Annually," Nov. 17, 2021. Accessed on: Apr. 15, 2022. [Online]. Available: [https://www.cdc.gov/nchs/pressroom/nchs\\_press\\_releases/2021/20211117.htm](https://www.cdc.gov/nchs/pressroom/nchs_press_releases/2021/20211117.htm)
4. "Detecting Frequent Opioid Prescription" [Online]. Available: <https://www.kaggle.com/code/apryor6/detecting-frequent-opioid-prescription>



**Thank you for  
listening!**

