

# **Abstract / Introduction**

With numerous initiatives to collect and make publicly available data sets in healthcare, a framework based on open source data mining technologies, primarily extending Python's SciPy ecosystem, is being designed to process them in search of novel correlations. [1] Using almost 900,000 practitioners connections and 3,000 hospitals, there is much to be found in the 50 fields present in the Center for Medicare and Medicaid's Hospital Value-Based Purchasing data. [2]

We use our framework to test a hypothesis about a correlation between age and performance. Visualizations were built to illustrate the findings. Linear regression analysis and Spearman's rank correlation analysis were used to find results for all hospitals and in individual specialties which suggest that there is no significant correlation between age and reported total performance score.

## Background

With a number of public datasets and public data analytic toolkits in Python, we've been able to collect, join, sanitize, and analyze a large amount of healthcare data accounting for approximately 7% of the US Healthcare Workforce. As we process these datasets we develop a high level interface to the open source components with a focus on automating common forms of processing in the healthcare field.

## **Datasets and Views**

- Medicare Hospital Value-Based Purchasing [2]: 3,090
- Medicare Practitioners [2]: 868,384
- Medicare Hospital Practitioner Connections: 1,067,937
- NY SPARCS Outcomes [3]: 1,907
- NY Practitioner Licenses [4]: 156
- NY Outcome-License connections: 1,865

## Tools in use

- SciPy framework
- Pandas [5]
- Matplotlib [6]
- Scikit-Learn [7]
- K-Means



**Figure 1.** Pie chart showing top specialties in dataset.

# Applying Big Data Analytics to Open Health Care Data: Exploring Relationships **Between Seniority and Performance in Healthcare Providers** Daniel Clarke, Dr. Ravi Rao, Prof. Maryelena Vargas Fairleigh Dickinson University

- Coherent wrapper around existing Python toolkits
- Some level of automatic data sanitation, encoding, scaling, and sampling
- High level visualizations of common models
- Abstract analysis for identifying significant trends or correlations in data
- Application of framework to public heathcare data
- Try to understand the relationship between hospital performance score and seniority

Initially finding a weak correlation by linear regression we dug deeper to validate the finding. We split the data by specialty and by graduation year to help understand any skewness in the data and performed a number of tests..





- Test for Normality with Skewness and Kurtosis
- Performance Scores
- Per year
- Per specialty
- Number of practitioners
- Per hospital • Per graduation year
- K-Means trend categorization Number of
- practitioner • Per
- graduation year



- Linear Regression and Spearman Correlation Coefficient
- Graduation year vs Performance Score
- Per Specialty





Figure 6. Plot showing recent quantitative trends in time for Nurse Practitioners.

	Spearman correlation coefficient	Spearman P-value	R²	m	μ	σ	Ν
TRY	0.164	0.04507	0.049	0.164	38.638	8.254	150
iΥ	0.11	0	0.006	0.075	37.643	9.151	1742
	0.077	0.00781	0.004	0.055	39.558	9.811	1191
.OGY	0.063	0.00087	0.003	0.041	39.625	9.368	2820
)GY	0.058	0.00001	0.003	0.049	39.724	9.979	5864
ASE	0.051	0	0.003	0.05	39.434	10.317	21587

**Table 1.** Table showing top specialty correlation results.