



# Group Project Presentation

**DSBA 6160: Big Data Design, Storage, and Provenance**

Alpha Group 1: Tony Bejos, Josh Peterson, Michael Lewis and Kevin Ovendorf

# Part I

Introduction

# Accessing and Loading MIMIC III Database

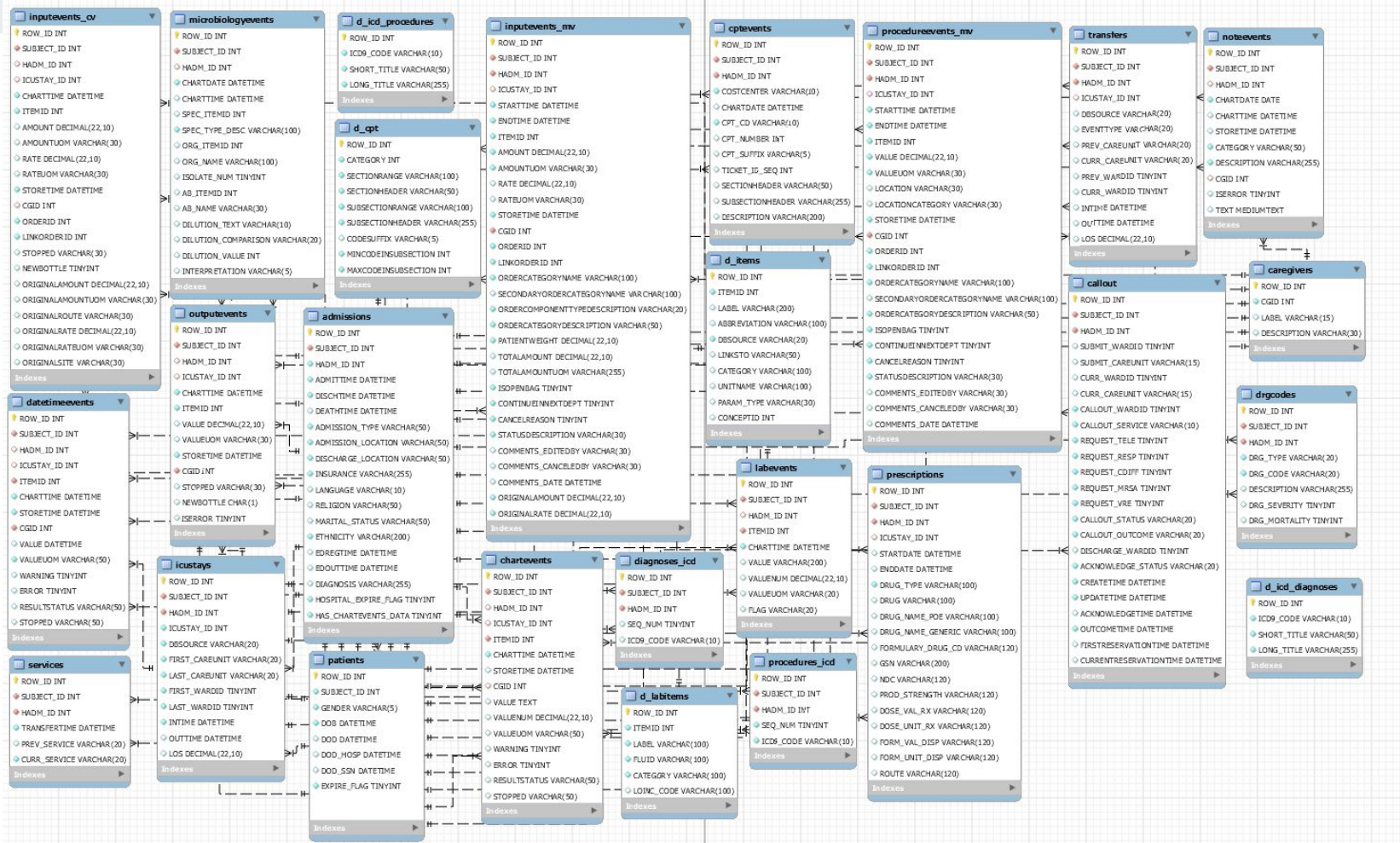


- Attained Training Certification through CITI online training program to gain access to MIMIC III database
  - 90% Score Required
  - Learned about human research practices and regulations
  - Patient Privacy remains essential
- Prior to loading the tables into MySQL, each table reviewed and discussed by the team
  - 27 Tables
  - Millions of Rows
  - Entity vs Association

# MIMIC III Database Build Overview

- Created a MYSQL db named *mimic\_iii\_project*
- Created table schema prior to load
- Extracted all zip files and utilized LOAD DATA LOCAL IN FILE code to load
- Created indexes on all tables for ease of use

```
LOAD DATA LOCAL INFILE 'ADMISSIONS.csv' INTO TABLE ADMISSIONS
FIELDS TERMINATED BY ',' ESCAPED BY '\\' OPTIONALLY ENCLOSED BY '"'
LINES TERMINATED BY '\n'
IGNORE 1 LINES
(@ROW_ID,@SUBJECT_ID,@HADM_ID,@ADMITTIME,@DISCHTIME,@DEATHTIME,@ADMISSION_TYPE,@ADMISSION_LOCATION,@DISCHARGE_LOCATION,@INSURANCE,@LANGUAGE,@RELIGION,@MARITAL_STATUS,@ETHNICITY,@EDREGTIME,@EDOUTTIME,@DIAGNOSIS,@HOSPITAL_EXPIRE_FLAG,@HAS_CHARTEVENTS_DATA)
SET
ROW_ID = @ROW_ID,
SUBJECT_ID = @SUBJECT_ID,
HADM_ID = @HADM_ID,
ADMITTIME = @ADMITTIME,
DISCHTIME = @DISCHTIME,
DEATHTIME = IF(@DEATHTIME='', NULL, @DEATHTIME),
ADMISSION_TYPE = @ADMISSION_TYPE,
ADMISSION_LOCATION = @ADMISSION_LOCATION,
DISCHARGE_LOCATION = @DISCHARGE_LOCATION,
INSURANCE = @INSURANCE,
LANGUAGE = IF(@LANGUAGE='', NULL, @LANGUAGE),
RELIGION = IF(@RELIGION='', NULL, @RELIGION),
MARITAL_STATUS = IF(@MARITAL_STATUS='', NULL, @MARITAL_STATUS),
ETHNICITY = @ETHNICITY,
EDREGTIME = IF(@EDREGTIME='', NULL, @EDREGTIME),
EDOUTTIME = IF(@EDOUTTIME='', NULL, @EDOUTTIME),
DIAGNOSIS = IF(@DIAGNOSIS='', NULL, @DIAGNOSIS),
HOSPITAL_EXPIRE_FLAG = @HOSPITAL_EXPIRE_FLAG,
HAS_CHARTEVENTS_DATA = @HAS_CHARTEVENTS_DATA;
```



# Part II

Notable Tables and Characteristics

# Table Descriptions and Characteristics



## DRGCODES

- Diagnosis Related Groups are codes used by the hospital for grouping diagnoses.
- There are two sources for these codes: MS and HCFA
- Can change over time, so must also use the description field when extracting these records

## ICUSTAYS

- Defines each ICU Stay, so patients can have multiple ICUSTay\_ID
- Provides the ICU database source relevant to that patient; they used CareVue from 2001-2008 and Metavision from 2008-2012

# Table Descriptions and Characteristics



## INPUTEVENTS CV & INPUTEVENTS MV

- These tables house the events that take place during a patient stay in either the CareVue or Metavision database.
- Each DB has similar, but different fields so must be familiar with each.

## LABEVENTS

- This table contains all lab measurements for a certain patients
- Lab measurements include the results from sampling a patient specimen (ie blood, urine, saliva)
- Time for these events is logged when specimen is acquired, not when data is available.



# Table Descriptions and Characteristics



## MICROBIOLOGYEVENTS

- This table is similar to the LABEVENTS table except it details situations where microbiology is required
- Provides details and results from various tests involving microbiology cultures

## NOTEEVENTS

- Details the patient notes taken by a doctor or nurse during the course of their ICU Stay
- Category is an important field as it can often tell us when a patient is discharged
- The TEXT field can be very long and should be extracted carefully as a result

# Table Descriptions and Characteristics



## ADMISSIONS

- Entity - All hospital admissions data
- Notably includes demographics of patients and admission/discharge times.
- Also include binary expiry flag that shows if patient is deceased.

## CALLOUT

- Entity - Provides information about ICU discharge planning
- Shows when a caregiver suggests callout, and if patient complied
- Also has “REQUEST” fields outlining specific diagnoses

## CAREGIVERS

- Entity table including caregiver information and CGID key field

# Table Descriptions and Characteristics



## CPTEVENTS

- Association table - outlines billing events
- One-to-many relationship with CAREGIVERS through CGID foreign key

## CHARTEVENTS

- Association table w/ 330 million rows - physical measurements and “chartable” events
- 5 different foreign keys (joins PATIENTS, ADMISSIONS, ICUSTAYS, D\_ITEMS, CAREGIVERS all with one-to-many relationships)

## DATETIMEEVENTS

- Association table w/ 4 million rows - datetime of all measurements of a patient
- 5 different foreign keys (joins PATIENTS, ADMISSIONS, ICUSTAYS, D\_ITEMS, CAREGIVERS all with one-to-many relationships)

# Table Descriptions and Characteristics



## PATIENTS

- Entity table - identifies the patient's sex, date of birth, date of death, and a mortality signifier
- One-to-many relationship with ADMISSIONS and ICUSTAYS

## OUTPUTEVENTS

- Entity table - identifies the chart times of the output events as well as details concerning the solutions and the caregivers responsible for the admission
- This table joins with the following tables:
  - PATIENTS
  - ADMISSIONS
  - ICUSTAYS
  - D\_ITEMS
  - CAREGIVERS

# Table Descriptions and Characteristics



## PRESCRIPTIONS

- Entity table - contains the identifiers and details about prescriptions written in the hospital
- Joins with the PATEINTS, ADMISSIONS, and ICUSTAYS tables

## PROCEDUREEVENTS\_MV

- Entity table - contains the identifiers and details regarding events surrounding procedures
- Joins with the following tables:
  - PATEINTS
  - ADMISSIONS
  - ICUSTAYS
  - D\_ITEMS

# Table Descriptions and Characteristics



## PROCEDURES ICD

- Association table - contains identifiers and the order in which procedures took place with the associated ICD-9 procedure code.
- Joins with the PATIENTS, ADMISSIONS, and D\_ICD\_PROCEDURES tables

## SERVICES

- Association table - Contains transfer times along with the previous service performed as well as the current service being performed
- Joins with the PATIENTS and ADMISSIONS tables

## TRANSFERS

- Entity table - contains the identifiers and details regarding the transfers of patients
- Joins with the PATIENTS, ADMISSIONS, and ICUSTAYS tables.

# Table Descriptions and Characteristics



## DICTIONARY TABLES ( 5 total 'D\_' tables)

- Association tables: provide a high-level dictionary of procedural terminology used in the medical field
  - D\_CPT
  - D\_ICD\_PROCEDURES
  - C\_ICD\_DIAGNOSES
  - D\_ITEMS
  - D\_LABITEMS
- ICD9\_CODE links DIAGNOSES\_ICD and PROCEDURES\_ICD tables separately with a one-to-one relationship
- ITEM\_ID links the following tables with a one-to-one relationship:
  - CHARTEVENTS, DATETIMEEVENTS, INPUTEVENTS\_MV, MICROBIOLOGYEVENTS, OUTPUTEVENTS, PROCEDUREEVENTS\_MV, LABEVENTS

# Table Descriptions and Characteristics



## DIAGNOSES ICD

- Association table - contains diagnoses which relate to a specific hospital's admissions through the use of the ICD9 Code system.
- Links to the PATIENTS table with a one-to-many relationship
- Links to the ADMISSIONS table with a one-to-one relationship

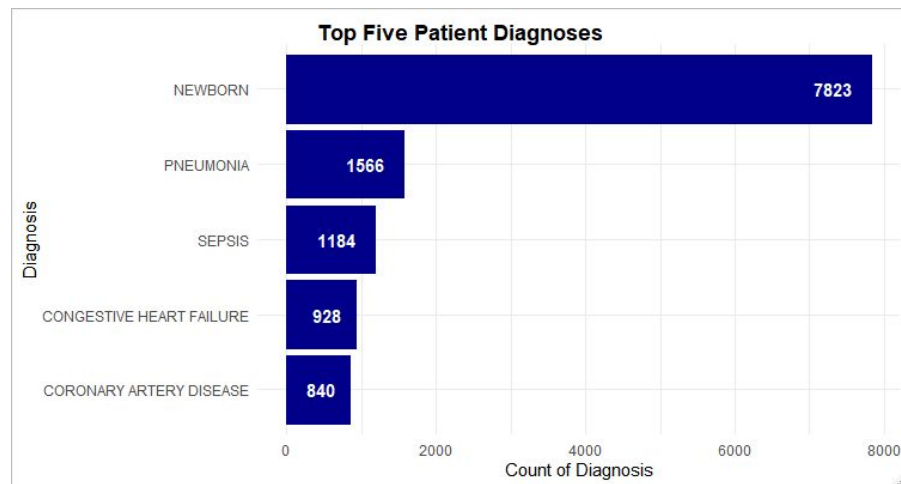


# PART III

Exploring Major Characteristics of the MIMIC III Database

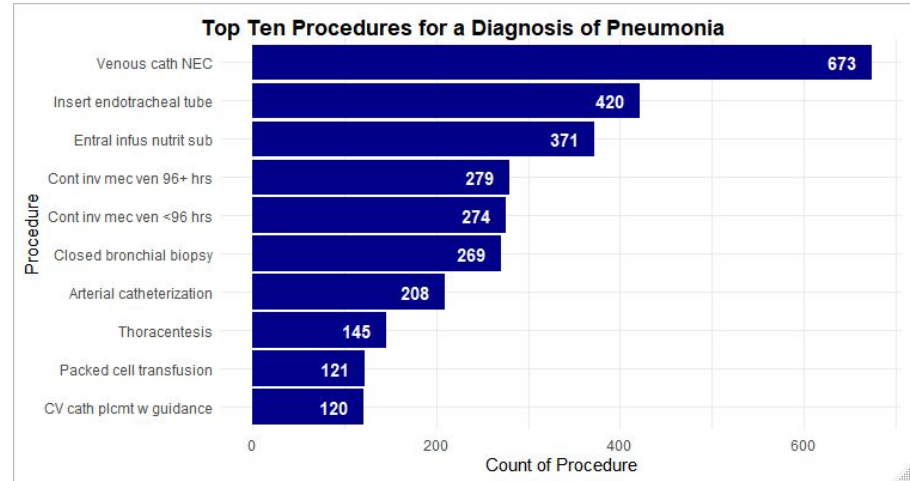
# Top Five Patient Diagnoses

- Displayed is the top five patient diagnostics from the database.
- The initial exploration was to find the most common diagnoses.
- This allows further drilling down because a baseline of what patients are dealing with has been presented



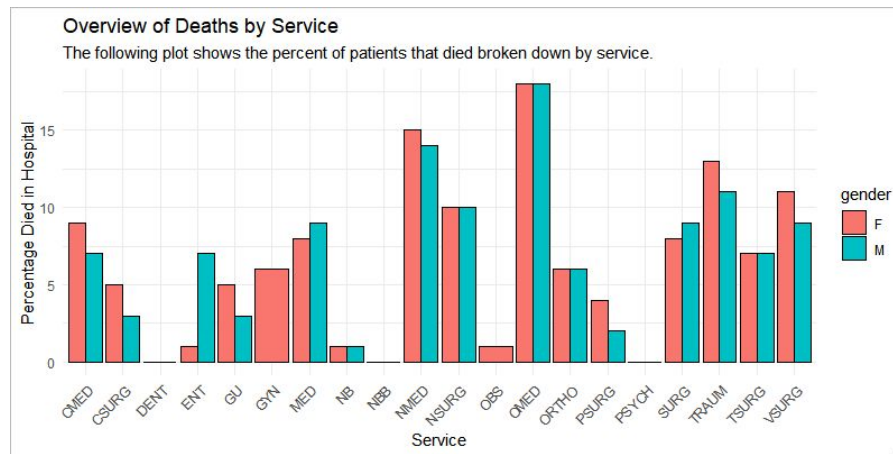
# Pneumonia Procedures

- Displays the number of each procedure performed on the pneumonia diagnostic ranked from high to low, Venous Cath NEC being the most commonly performed
- This plot was derived from the top five diagnostics, pneumonia being the second most common.



# Gender: Deaths By Service

- This displays the proportion of individuals, separated by gender, who died within the hospital during a particular visit.
- The largest difference between the genders was recorded at the ENT ( ear, nose, and throat) where there was a difference of > 5%.
- Beyond this plot, patients were bucketed into age groups for each service to further investigate the death rate of these services.



# Mortality Statistics



## Goals were to find:

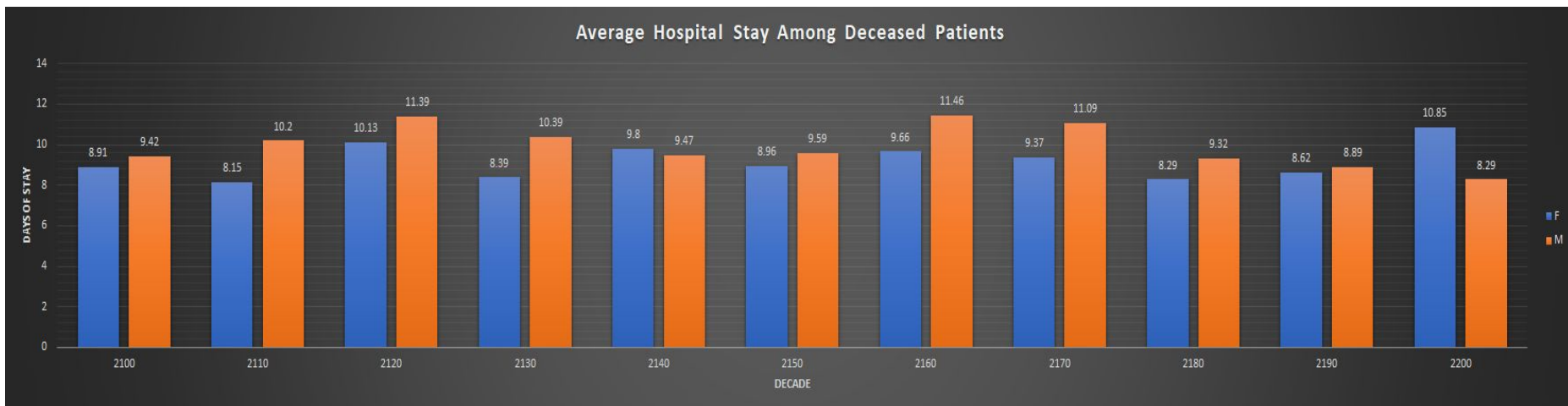
- Average hospital stay for deceased patients
- Most common diagnoses shared among deceased grouped by any time series and gender

## Achieved by creating a view joining:

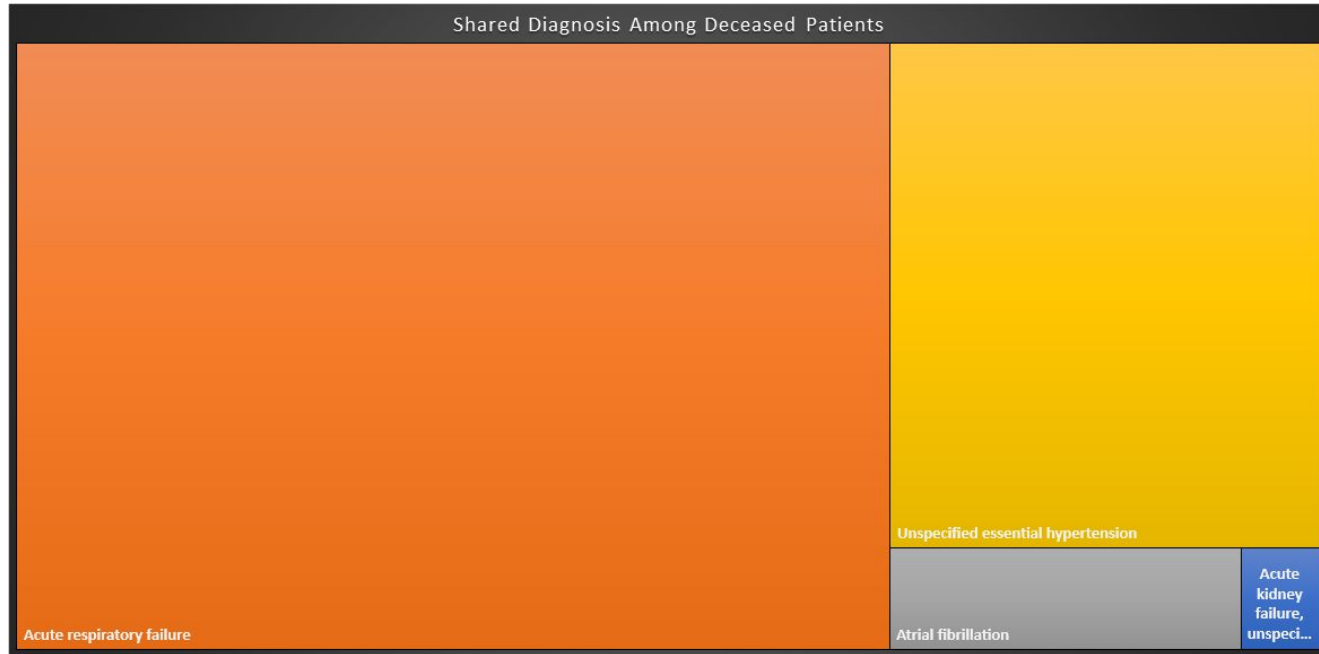
- ADMISSIONS
- PATIENTS
- DIAGNOSES\_ICD
- D\_ICD\_DIAGNOSES

**Condition - HOSPITAL\_EXPIRE\_FLAG = 1**

# Average Hospital Stay

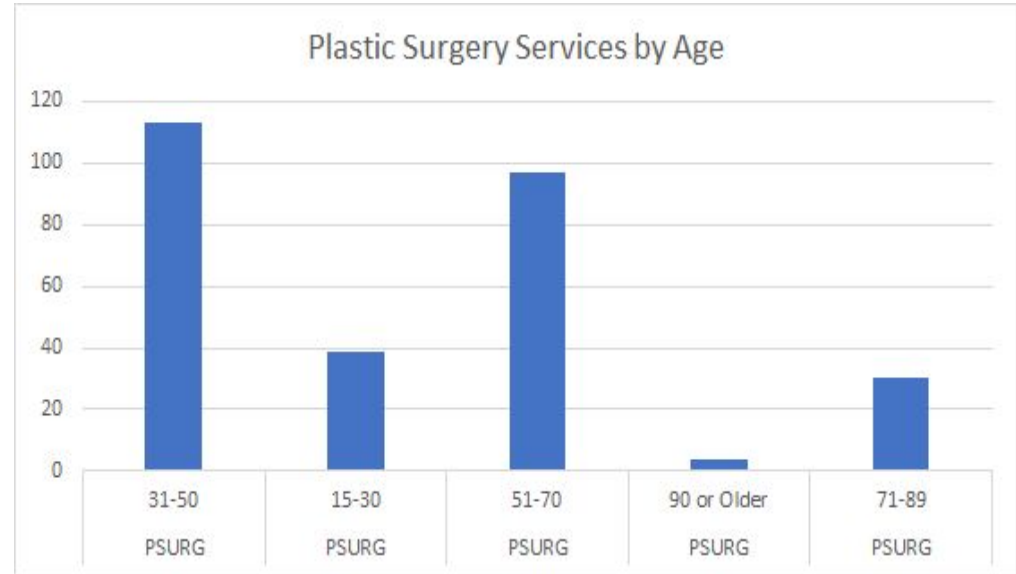


# Commonly Shared Diagnoses - Deceased Patients



# First-Visit Services Rendered by Age Group

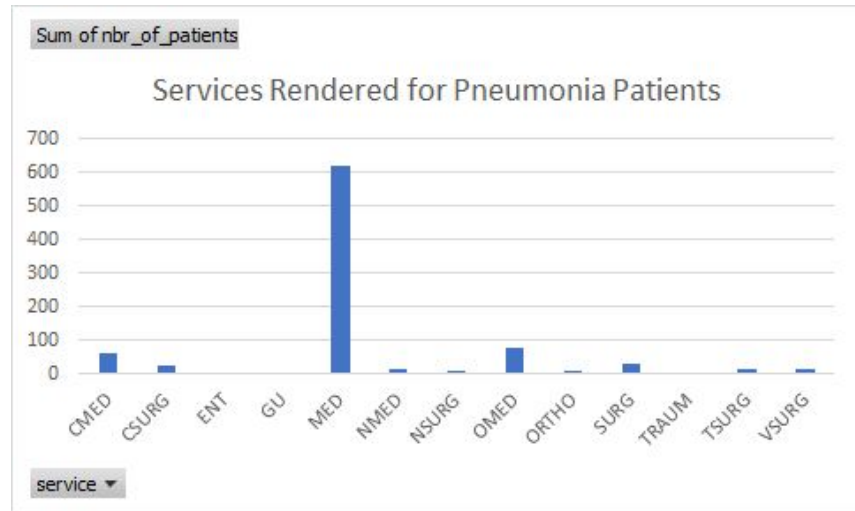
- Displays the services rendered to first-time ICU visits by age group
- Assists hospital with predicting the resources needed and helps prioritize patients
- This graph shows the age for Plastic Surgery, shows to be a more common service for younger people in relation to other services.





# First-Admit Services for Pneumonia Patients

- Search can filter by diagnosis to determine which services are commonly rendered in the case of a specific diagnosis.

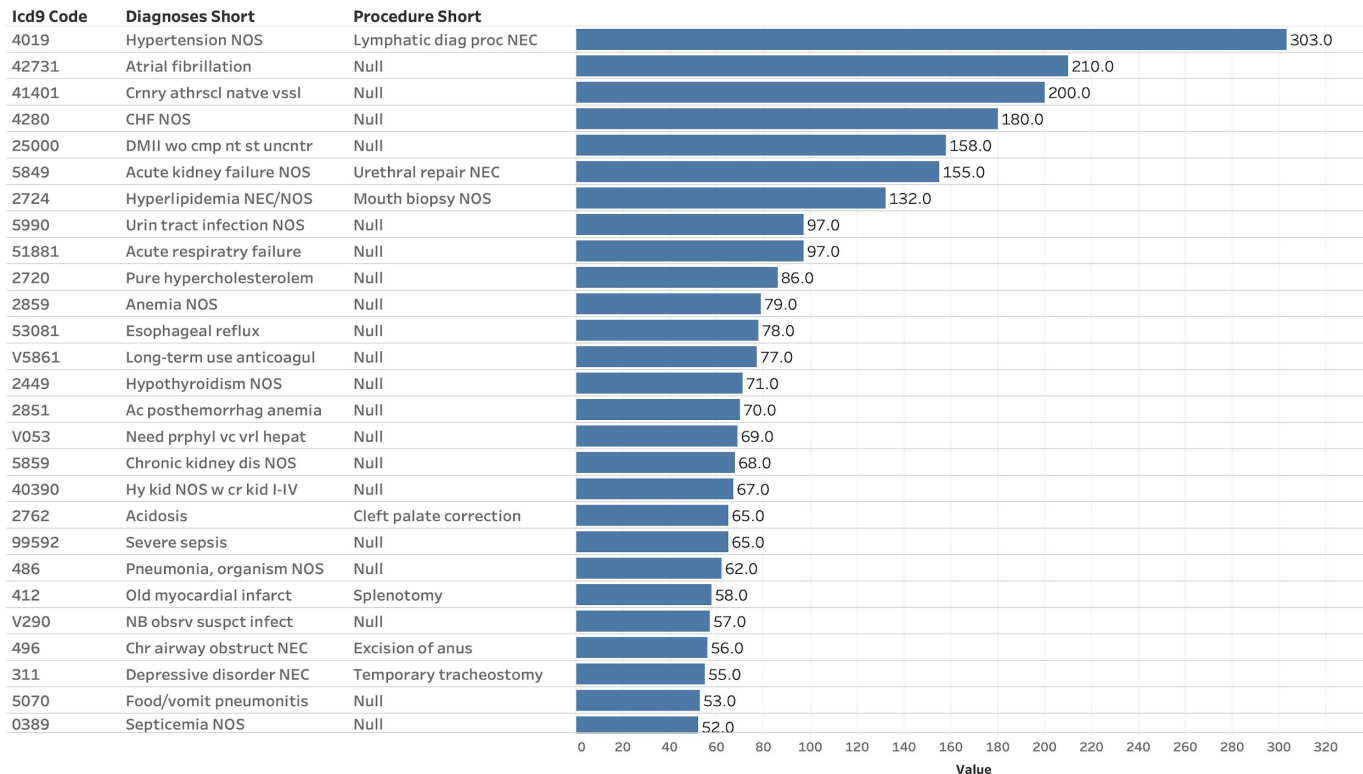


# ICD9 Codes and associated diagnoses and procedure

- ICD9 codes were joined with their respective short and long descriptions of its diagnosis and procedure.
- All ICD9 codes were extracted from the DIAGNOSES\_ICD and PROCEDURE\_ICD tables.
- This serves as a foundation to pursue future analytics that deal with features such datetime, gender, and mortality rate.

Icd9 Code	Diagnoses Short	Diagnoses Long	Procedure Short	Procedure Long
3572	Neuropathy in diabet...	Polyneuropathy in di...	Ventr septa def rep N...	Other and unspecifie...
41400	Cor ath unsp vsl ntv/...	Coronary atheroscler...	<i>null</i>	<i>null</i>
4019	Hypertension NOS	Unspecified essentia...	Lymphatic diag proc ...	Other diagnostic pro...
V4581	Aortocoronary bypass	Aortocoronary bypas...	<i>null</i>	<i>null</i>
41071	Subendo infarct, initial	Subendocardial infar...	<i>null</i>	<i>null</i>
5849	Acute kidney failure ...	Acute kidney failure, ...	Urethral repair NEC	Other repair of urethra
4280	CHF NOS	Congestive heart fail...	<i>null</i>	<i>null</i>
25060	DMII neuro nt st unc...	Diabetes with neurol...	<i>null</i>	<i>null</i>
41401	Crnry athrsc native v...	Coronary atheroscler...	<i>null</i>	<i>null</i>
4019	Hypertension NOS	Unspecified essentia...	Lymphatic diag proc ...	Other diagnostic pro...
2720	Pure hypercholester...	Pure hypercholester...	<i>null</i>	<i>null</i>
42731	Atrial fibrillation	Atrial fibrillation	<i>null</i>	<i>null</i>
V1582	History of tobacco use	Personal history of t...	<i>null</i>	<i>null</i>
2859	Anemia NOS	Anemia, unspecified	<i>null</i>	<i>null</i>
3572	Neuropathy in diabet...	Polyneuropathy in di...	Ventr septa def rep N...	Other and unspecifie...
78552	Septic shock	Septic shock	<i>null</i>	<i>null</i>

# ICD9 Codes and associated diagnoses and procedure



# Part IV

Analysis of Chapters 21 & 23

# Chapter 21: Modeling Mortality



- Extracting predictor variables with the intent of building various models to predict mortality in the ICU:
- Existing ICU performance metrics used to test the validity of the model:
  - APACHE Scores
  - The simplified acute physiology score
  - The mortality probability model
  - The sequential organ failure assessment score
- New performance metrics are extracted...

# Chapter 21: Extraction



- Only the first ICU admission of Adult patients is taken into account
- In order to successfully complete this task a number of adjustments were made to accommodate the requirements needed to work on the MIMIC-III database:
  - Identifiers used in the MIMIC-II database do not exist in the MIMIC-III, so a new query needed to be built to recreate the ICUSTAY\_DETAIL
    - This included a 'dense\_rank' function 'partitioned' over SUBJECT\_ID or HADM\_ID to get various ICU sequences
  - Further calculations were performed to determine the following features as well:
    - AGE
    - DAYS\_TILL\_DEATH\_AFTER\_DISCHARGE

# Chapter 21: Data Preparation



- The following steps were done in R after the data was extracted, this is to ensure the data is compatible with various mortality prediction models:
  - All patients older than 89 had their aged masked in MIMIC-III
    - These records were replaced with the median age of 91.4
  - All records with incomplete data were removed

# Chapter 23: Propensity Score Analysis



- Model that enables best selection for test and control groups
- Patient data must be extracted from MIMIC III using SQL to build model
  - The data to be extracted includes:
    - Demographics
    - Vital Signs
    - Basic Metabolic Panels
    - Past Medical Conditions
    - Disease Severity Scores
    - Types of Admission
    - Types of ICU
- **OUR GOAL** is to properly extract these pieces of information to feed the model for its propensity analysis and enable researchers to assign subjects to the correct group



# Chapter 23: Data Extraction Process



All pulled into single view to load into an analytics platform (i.e. *R*, *Python*, or *SAS*)

- Demographics
  - joined from ADMISSIONS and PATIENTS tables
- Vital Signs
  - created “pivot” subquery as *base\_events\_pivot* by pulling avg values for items in CHARTEVENTS table
- Basic Metabolic Panels
  - created “pivot” subquery as *labs\_pivot* by pulling max values for items in LABEVENTS table
- Past Medical Conditions
  - Separate table for reference by analyst
- Types of Admission
  - Joined admissions table
- Types of ICU
  - Joined ICUSTAYS table