

# FHIR-PIT

A tool to smooth the join of FHIR records with Environmental and Exposures data. It assumes a patient has lived in the same address for the study period. It joins FHIR records with exposure data through **patient's address**, and environmental data through **patient's address, study period** and **date-of-visit**.

## Technical Overview:

### Built using sbt

It expects a particular folder structure.

Build.sbt:

- Name, version, scalaVersion
- Library Dependencies
- AssemblyMergeStrategies
- Resolvers

All scala files should go under src/main/scala. Build.sbt acts similar to requirements.txt, lists all dependencies and compiles all scala code that matches the folder structure.

### Configured using the DHALL programming language

Example.dhall contains an example format configuration file for running FHIR-PIT. It calls pipeline.dhall.

The input parameters are:

**Report:**

**Progress:**

**Configdir:** Directory with extra configuration files (mainly icees\_features.yaml)

**Basedirinput:** Data directory. The following structure is expected:

**FHIR/:** Contains FHIR records

**EPR/** : EPR files

**ICEESPCD/**: Xwalk files

**other/spatial/**: Contains exposure data

**other/env/**: Contains environmental data.

**Basedir**: Directory to save intermediate results.

**Basediroutput**: Directory to save outputs.

**Fhirconfig**: JSON that configures the FHIR PIT run

**SkipList**: List configurations for the PreprocCSVTable transformation of each study period. (Poorly named)

**Pyexec**: Python exec.

The configuration file is YAML and defines a list of steps.

Each step has parameters: name, skip, arguments and function.

The function defines the object Step: It inherits the StepImpl trait and overrides a config type, a decoder and a step function. Afterwards, The main method parses the config.yaml file; creates a queue of steps; and executes each step in the queue. Results are saved in the directories specified by the arguments of each step.

### **Scala design-patterns/notions used:**

**Implicits**: Passing the “wrong” type. “Unknown” method call. It performs type conversion implicitly when the wrong type is passed. The programmer defines the conversion. It can also define methods that extend functionality with conversion

**Encoder/Decoder JSON**: Defines a case class object with attributes that match the fields in a JSON file. It then uses `io.circe.generic.semiauto.deriveDecoder` to parse the json into the respective class. The classes can be nested to match the nested nature of the JSON.

**Others:** Type classes, Case classes, companion objects, Nullary functions for exec. time, Mapper, HashMap with Multimaps, Match-case, foreach.

## Files used in demo:

Basedirinput = /FHIR-PIT/data/input

### [36M] FHIR data

{basedirinput}/FHIR/all

{basedirinput}/FHIR/all

### [4M] Environmental data:

{basedirinput}/other/env/merged\_cmaq\_2010.csv

### [1M] Environmental data:

{basedirinput}/other/env/cmaq2010

### [589M] Census\_data:

{basedirinput}/other/spatial/env/US\_Census\_Tracts\_LCC/US\_Census\_Tracts\_LCC.\*

### [1M] ACS\_data:

{basedirinput}/other/spatial/acs/ACS\_NC\_2016\_with\_column\_headers.csv

### [36M] ACSUR\_data:

{basedirinput}/other/spatial/acs/Appold\_trans\_geo\_cross\_02.10.10 - trans\_geo\_cross.csv

### [45M] geoid\_data:

{basedirinput}/other/spatial/acs/tl\_2016\_37\_bg\_lcc.shp

**[489M] NearestRoadTL:**

{basedirinput}/other/spatial/nearestRoadTL/tl\_2015\_allstates\_prisecro  
ads\_lcc.shp

**[2.6G] NearestRoadHPMS:**

{basedirinput}/other/spatial/nearestRoadHPMS/hpms2016\_major\_roa  
ds\_lcc.shp

**[8.3M] CAFO + Landfill:**

{basedirinput}/other/spatial/BDT\_PointDatasets/

**[2M] XWalkData:**

{basedirinput}/ICEESPCD/RegistryPtsXWalkForHao.csv

{basedirinput}/ICEESPCD/8000PtsXWalkForHao.csv

**[56K] EPR:**

{basedirinput}/EPR/UNC\_NIEHS\_XWalk\_for\_Hao\_shape\_h3.csv

{basedirinput}/EPR/TLR4\_AllData\_NewHash\_01292020 NO

PII\_no\_new\_line.csv

## Step description:

Each step (bold and underlined) describes the input files, output files and highlights the contribution/operational-role of each input file in the output schema. Each input file is colored in the description. Each output file is the result of transformations/joins of multiple input files. The output schema is listed under each filename in a box. Data-elements in the output file's schema are highlighted with a paler color to match its corresponding input-file.

## **PreprocFHIR**

**Takes in:** FHIR Data

### **Description:**

Process encounters and resources on a per-patient basis.

It reorganizes data into the following file structure:

FHIR\_processed/<Resource\_type>/<Patient\_num>/<Resource\_num>  
@<iter\_num>.json. Additionally it creates the frequently used geo.csv  
file.

### **Outputs:**

FHIR\_processed

“FHIR-PIT/data/output/FHIR\_processed/Patient/d0f9bf93-f99f-4544-b0  
95-1d3c5265b5bb”

geo.csv

“FHIR-PIT/data/output/FHIR\_processed/geo.csv”

Patient_num, lat, lon
-----------------------

## **EnvCoordinates**

**Takes in:** cmaq files, geo.csv, start date and end date

### **Description:**

For each patient in **geo.csv** and for all years between start date and  
end date, save the contents of **cmaq<year>/<Row><Col>Daily.csv**  
into other\_processed/env\_coordinates/<Patient\_num>.csv. Where  
Row and Col are estimated from the patient’s lat lon coordinates in  
**geo.csv**.

The result is all the daily exposures recorded throughout the study  
period in the Row,Col address of a patient.

### **Outputs:**

“FHIR-PIT/data/output/other\_processed/env\_coordinates/d0f9bf93-f99  
f-4544-b095-1d3c5265b5bb”

start_date,o3_avg,pm25_avg,o3_max,pm25_max,o3_min,pm25_mi
---

n,o3_stddev,pm25_stddev
-------------------------

**Notes:** It is assumed the patient lived in the same address (i.e. lat lon coordinates in geo.csv) during each year between start and end dates.

### **LatLongtoGEOID**

**Takes in:** US\_Census\_Tracts\_LCC.shp and geo.csv

**Description:**

For each patient in **geo.csv**, map their lat-lon coordinates to fips using **US\_Census\_Tracts\_LCC.shp**. Expand geo.csv with the FIPS column and save the result to: other\_processed/lat\_lon\_to\_geoid/geoids.csv.

To obtain GEOID, it uses the function getGeoidForLatLon. It creates an LCC point with lat, lon coordinates. Then it asks for the census block containing point. It does so by looking at the geometry and checking if it contains the point.

**Outputs:**

“FHIR-PIT/data/output/other\_processed/lat\_lon\_to\_geoid/geoids.csv”

patient_num,lat,lon,FIPS
--------------------------

### **PreprocEnvDataFIPS**

**Takes in:** merged\_cmaq\_<year>.csv, geoids.csv, start\_date, end\_date.

**Description:**

Union all **merged\_cmaq\_<year>.csv** files across all years between start\_date and end\_date. Then inner join the result on **FIPS** with **geoids.csv**. Expand the data frame with column: **start\_date** (i.e. yyyy/mm/dd) from Date. The result is all the daily exposures recorded throughout the study period in the FIPS address of a patient.

**Outputs:**

“FHIR-PIT/data/output/other\_processed/env\_FIPS/preagg”

FIPS,Date,Longitude,Latitude,CO_ppbv,NO_ppbv,NO2_ppbv,NOX_ppbv,SO2_ppbv,ALD2_ppbv,FORM_ppbv,pm25_daily_average,pm2
--

```
5_daily_average_stderr, ozone_daily_8hour_maximum, ozone_daily_8hour_maximum_stderr, BENZ_ppbv, start_date, patient_num
```

**Notes:** It is assumed the patient lived in the same address for every year between start and end dates. If columns in the env file are missing Mapper.envInputColumns2, the cols are aggregated as null (This causes an ERROR because writeDataframe does not accept null columns). start\_date column is created to match daily join later.

### **PreprocSplit**

**Takes in:** other\_processed/env\_FIPS/preagg.csv

#### **Description:**

Extract the rows from each patient in **preagg.csv** and save them to its own csv file.

#### **Outputs:**

“FHIR-PIT/data/output/other\_processed/env\_split\_FIPS/d0f9bf93-f99f-4544-b095-1d3c5265b5bb.csv”

```
FIPS, Date, Longitude, Latitude, CO_ppbv, NO_ppbv, NO2_ppbv, NOX_ppbv, SO2_ppbv, ALD2_ppbv, FORM_ppbv, pm25_daily_average, pm25_daily_average_stderr, ozone_daily_8hour_maximum, ozone_daily_8hour_maximum_stderr, BENZ_ppbv, start_date, patient_num
```

### **PreprocEnvDataAggregate**

**Takes in:** /env\_coordinates, indices, statistics, study\_periods, study\_period\_bounds

#### **Description:**

For each patient’s environmental exposures throughout the study period (i.e. /**env\_coordinates**), group exposures by the intervals [study\_period\_bounds(i), study\_period\_bounds(i+1)]. Then **aggregate/compute** statistics for a set of indices. Expand the original data with the computed set, and append the values from the **previous day**. Save the result to other\_processed/env\_agg\_coordinates/<patient\_num>.csv

## Outputs:

“FHIR-PIT/data/output/other\_processed/env\_agg\_coordinates/d0f9bf93-f99f-4544-b095-1d3c5265b5bb.csv”

```
start_date,pm25_max,pm25_avg,o3_max,o3_avg,pm25_max_avg,pm25_avg_avg,o3_max_avg,o3_avg_avg,pm25_max_max,pm25_avg_max,o3_max_max,o3_avg_max,pm25_max_prev_date,pm25_avg_prev_date,o3_max_prev_date,o3_avg_prev_date
```

## PreprocEnvDataAggregate

**Takes in:** /env\_split\_FIPS, indices, statistics, study\_periods, study\_period\_bounds

### Description:

Same as above but done on the /**env\_split\_FIPS** data.

For each patient's environmental exposures throughout the study period (i.e. /**env\_split\_FIPS**), group exposures by the intervals [study\_period\_bounds(i), study\_period\_bounds(i+1)]. Then **aggregate/compute** statistics for a set of indices. Expand the original data with the computed set, and append the values from the **previous day**.

## Outputs:

“FHIR-PIT/data/output/other\_processed/env\_agg\_FIPS/d0f9bf93-f99f-4544-b095-1d3c5265b5bb”

```
start_date,pm25_daily_average,ozone_daily_8hour_maximum,CO_ppbv,NO_ppbv,NO2_ppbv,NOX_ppbv,SO2_ppbv,ALD2_ppbv,FORM_ppbv,BENZ_ppbv,pm25_daily_average_avg,ozone_daily_8hour_maximum_avg,CO_ppbv_avg,NO_ppbv_avg,NO2_ppbv_avg,NOX_ppbv_avg,SO2_ppbv_avg,ALD2_ppbv_avg,FORM_ppbv_avg,BENZ_ppbv_avg,pm25_daily_average_max,ozone_daily_8hour_maximum_max,CO_ppbv_max,NO_ppbv_max,NO2_ppbv_max,NOX_ppbv_max,SO2_ppbv_max,ALD2_ppbv_max,FORM_ppbv_max,BENZ_ppbv_max,pm25_daily_average_prev_date,ozone_daily_8hour_maximum_prev_date,CO_ppbv_prev_date,NO_ppbv_prev_date,NO2_ppbv_prev_date,NOX_ppbv_prev_date,SO2_ppbv_prev_date,ALD2_ppbv
```



```
_prev_date,FORM_ppbv_prev_date,BENZ_ppbv_prev_date
```

### **PreprocPerPatSeriesACS**

**Takes**            **in:**            geo.csv,            tl\_2016\_37\_bg\_lcc.shp,  
ACS\_NC\_2016\_with\_column\_headers.csv, icees\_features.yaml/acs

#### **Description:**

Map the lat, lon columns of **geo.csv** to GEOID using tl\_2016\_37\_bg\_lcc.shp. Then inner join the result with **ACS\_NC\_2016\_with\_column\_headers.csv** on GEOID.

#### **Outputs:**

“FHIR-PIT/data/output/other\_processed/acs.csv”

```
patient_num,EstResidentialDensity,EstProbabilityHighSchoolMaxEd  
ducation,EstProbabilityNoHealthIns,EstProbabilityHouseholdNonHisp  
White,EstProbabilityESL,EstProbabilityNonHispWhite,EstHouseholdI  
ncome,EstResidentialDensity25Plus,EstProbabilityNoAuto
```

### **PreprocPerPatSeriesACSUR**

**Takes**            **in:**            geo.csv,            tl\_2016\_37\_bg\_lcc.shp,  
Appold\_trans\_geo\_cross\_02.10.10,            trans\_geo\_cross.csv,  
icees\_features.yaml/acsUR

#### **Description:**

Map the lat, lon columns of **geo.csv** to GEOID using tl\_2016\_37\_bg\_lcc.shp. Then inner join the result with **Appold\_trans\_geo\_cross\_02.10.10 - trans\_geo\_cross.csv** on GEOID.

#### **Outputs:**

“FHIR-PIT/data/output/other\_processed/acsUR.csv”

```
patient_num,ur
```

### **PreprocPerPatSeriesNearestRoad**

**Takes in:** geo.csv, tl\_2015\_allstates\_prisecroads\_lcc.shp, icees\_features.yaml/nearestRoadTL

**Description:**

For each patient in **geo.csv**, compute distance from patient's address to nearest road using **tl\_2015\_allstates\_prisecroads\_lcc.shp**.

**Outputs:**

"FHIR-PIT/data/output/other\_processed/nearestRoadTL.csv"

patient_num, MajorRoadwayHighwayExposure
--

**PreprocPerPatSeriesNearestRoadHPMS**

**Takes in:** geo.csv, hpms2016\_major\_roads\_lcc.shp, icees\_features.yaml/nearestRoadHPMS

**Description:**

For each patient in **geo.csv**, compute distance from patient's address to nearest road distance and features using **hpms2016\_major\_roads\_lcc.shp**.

**Outputs:**

"FHIR-PIT/data/output/other\_processed/nearestRoadHPMS.csv"

patient_num, RoadwayDistanceExposure, RoadwayType, RoadwayAADT, RoadwaySpeedLimit, RoadwayLanes
---

**PreprocPerPatSeriesNearestPointCafo**

**Takes in:** geo.csv, Permitted\_Animal\_Facilities-4-1-2020.shp, icees\_features.yaml/cafo

**Description:**

For each patient in **geo.csv**, find distance in meters between the patient's address and **Permitted\_Animal\_Facilities-4-1-2020.shp**.

**Outputs:**

"FHIR-PIT/data/output/other\_processed/cafo.csv"

patient_num, CAFO_Exposure
----------------------------

### **PreprocPerPatSeriesNearestPointLandfill**

**Takes in:** geo.csv, Active\_Permitted\_Landfills\_geo.shp,  
icees\_features.yaml/landfill

#### **Description:**

For each patient in **geo.csv**, find distance in meters between the patient's address and **Active\_Permitted\_Landfills\_geo.shp**.

#### **Outputs:**

"FHIR-PIT/data/output/other\_processed/landfill.csv"

patient_num, Landfill_Exposure
--------------------------------

### **PreprocPerPatSeriesToVector**

**Takes in:** FHIR\_processed/, start\_date, end\_date,  
icees\_features.yaml/

#### **Description:**

Load a patient json into a Patient object, which contains a list of Medication objects, Address objects, Encounter objects, Condition objects, Lab objects, Procedure objects, bmi objects.

An Encounter object also contains a list of: Condition, Lab, Medication, Procedure objects; in addition to an id, start\_date, end\_date. Encounters is a data structure to aid operations between all these pieces of information.

For each patient, group their encounters, medications (converted to Encounter), conditions (converted to Encounter), labs (converted to Encounter), procedures (converted to Encounter) by **day** using a HashMultiMap (i.e. A dictionary with day as key and set of encounters as value). Accordingly, encounters are grouped-by on a daily basis.

Then aggregate features from the encounters by: Count, First, Last.

Finally, collect all the vectors, union all the feature names collected each day, and populate a csv with missing values replaced by "".

**Outputs:**

"FHIR-PIT/data/output/FHIR\_vector/patient/d0f9bf93-f99f-4544-b095-1d3c5265b5bb.csv", which contains the demographic information of a patient.

patient_num,Ethnicity,Sex,birth_date,Race
---

"FHIR-PIT/data/output/FHIR\_vector/visit/d0f9bf93-f99f-4544-b095-1d3c5265b5bb.csv", which contains the patient's vectorized/aggregated encounters.

patient_num,start_date,Ethnicity,AgeVisit,ObesityBMIVisit,encounter_num,VisitType,Sex,birth_date,Race
---

## **PreprocPerPatSeriesToCSVTable**

**Takes in:** /FHIR\_vector/visit/<patient\_num>, landfill.csv, cafe.csv, nearestRoadHPMS.csv, nearestRoadTL.csv, acsur.csv, acs.csv, /other\_processed/env\_agg\_FIPS/<patient\_num>, /other\_processed/env\_agg\_coordinates/<patient\_num>, study\_periods, study\_period\_bounds.

### **Description:**

For each patient in /FHIR\_vector/visit/<patient\_num>.csv: Left join the **visits** **vector** with **environmental data** (i.e. /other\_processed/env\_agg\_coordinates/<patient\_num>) on start date. Thus attaching the aggregated environmental exposures that occurred the day of the encounter. Similarly, associate the **exposure** information through the patient\_num (Note exposures were associated to the patient\_num with the patient's address). Lastly, **expand** the data with **bucket**, **study\_period** and **year**. Extra columns that all carry the same information (i.e. Which study\_period\_bound is associated with a given row). Save the result to icees/<year>/<patient>.

### **Outputs:**

"FHIR-PIT/data/output/icees/2010/d0f9bf93-f99f-4544-b095-1d3c5265b5bb"

bucket,patient_num,start_date,Ethnicity,AgeVisit,ObesityBMIVisit,encounter_num,VisitType,Sex,birth_date,Race,pm25_daily_average,ozone_daily_8hour_maximum,CO_ppbv,NO_ppbv,NO2_ppbv,NOX_ppbv,SO2_ppbv,ALD2_ppbv,FORM_ppbv,BENZ_ppbv,pm25_daily_average_avg,ozone_daily_8hour_maximum_avg,CO_ppbv_avg,NO_ppbv_avg,NO2_ppbv_avg,NOX_ppbv_avg,SO2_ppbv_avg,ALD2_ppbv_avg,FORM_ppbv_avg,BENZ_ppbv_avg,pm25_daily_average_max,ozone_daily_8hour_maximum_max,CO_ppbv_max,NO_ppbv_max,NO2_ppbv_max,NOX_ppbv_max,SO2_ppbv_max,ALD2_ppbv_max,FORM_ppbv_max,BENZ_ppbv_max,pm25_daily_average_prev_date,ozone_daily_8hour_maximum_prev_date,CO_ppbv_prev_date,N
---

O\_ppbv\_prev\_date,NO2\_ppbv\_prev\_date,NOX\_ppbv\_prev\_date,SO2\_ppbv\_prev\_date,ALD2\_ppbv\_prev\_date,FORM\_ppbv\_prev\_date,BENZ\_ppbv\_prev\_date,EstResidentialDensity,EstProbabilityHighSchoolMaxEducation,EstProbabilityNoHealthIns,EstProbabilityHouseholdNonHispWhite,EstProbabilityESL,EstProbabilityNonHispWhite,EstHouseholdIncome,EstResidentialDensity25Plus,EstProbabilityNoAuto,ur,MajorRoadwayHighwayExposure,RoadwayDistanceExposure,RoadwayType,RoadwayAADT,RoadwaySpeedLimit,RoadwayLanes,CAFO\_Exposure,Landfill\_Exposure,study\_period,year

**Note:** Bucket, Study\_period, year information is redundant when study\_period\_bounds are yearly.

## PreprocCSVTable

### Takes

in:

“FHIR-PIT/data/output/icees/2010/d0f9bf93-f99f-4544-b095-1d3c5265b5bb.csv”,

“FHIR-PIT/data/output/FHIR\_Vector/patient/d0f9bf93-f99f-4544-b095-1d3c5265b5bb.csv”

### Description:

Per year, the visits of all patients are “union” into a single table, grouped by (patient\_num, study\_period) and aggregated with the following sql functions:

```
first(pm25_avg, false) AS `AvgDailyPM2.5Exposure`
first(pm25_avg_avg, false) AS `AvgDailyPM2.5Exposure_StudyAvg`
first(pm25_avg_max, false) AS `AvgDailyPM2.5Exposure_StudyMax`
first(pm25_max, false) AS `MaxDailyPM2.5Exposure`
first(pm25_max_avg, false) AS `MaxDailyPM2.5Exposure_StudyAvg`
first(pm25_max_max, false) AS `MaxDailyPM2.5Exposure_StudyMax`
first(o3_avg, false) AS `AvgDailyOzoneExposure`
first(o3_avg_avg, false) AS `AvgDailyOzoneExposure_StudyAvg`
first(o3_avg_max, false) AS `AvgDailyOzoneExposure_StudyMax`
first(o3_max, false) AS `MaxDailyOzoneExposure`
first(o3_max_avg, false) AS `MaxDailyOzoneExposure_StudyAvg`
first(o3_max_max, false) AS `MaxDailyOzoneExposure_StudyMax`
first(pm25_daily_average_avg, false) AS `AvgDailyPM2.5Exposure_2`
first(ozone_daily_8hour_maximum_avg, false) AS `MaxDailyOzoneExposure_2`
first(CO_ppbv_avg, false) AS `AvgDailyCOExposure_2`
first(NO_ppbv_avg, false) AS `AvgDailyNOExposure_2`
first(NO2_ppbv_avg, false) AS `AvgDailyNO2Exposure_2`
first(NOX_ppbv_avg, false) AS `AvgDailyNOxExposure_2`
first(SO2_ppbv_avg, false) AS `AvgDailySO2Exposure_2`
```

```
first(ALD2_ppbv_avg, false) AS `AvgDailyAcetaldehydeExposure_2`
first(FORM_ppbv_avg, false) AS `AvgDailyFormaldehydeExposure_2`
first(BENZ_ppbv_avg, false) AS `AvgDailyBenzeneExposure_2`
max(ObesityBMIVisit) AS `ObesityBMI`
totaltypevisits(VisitType, RespiratoryDx) AS `TotalEDVisits`
totaltypevisits(VisitType, RespiratoryDx) AS `TotalInpatientVisits`
(totaltypevisits(VisitType, RespiratoryDx) + totaltypevisits(VisitType, RespiratoryDx)) AS `TotalEDInpatientVisits`
first(Sex2, false) AS `Sex2`
first(birth_date, false) AS `birth_date`
first(Sex, false) AS `Sex`
first(Race, false) AS `Race`
first(Ethnicity, false) AS `Ethnicity`
first(MajorRoadwayHighwayExposure, false) AS `MajorRoadwayHighwayExposure`
first(RoadwayDistanceExposure, false) AS `RoadwayDistanceExposure`
first(RoadwayType, false) AS `RoadwayType`
first(RoadwayAADT, false) AS `RoadwayAADT`
first(RoadwaySpeedLimit, false) AS `RoadwaySpeedLimit`
first(RoadwayLanes, false) AS `RoadwayLanes`
```

```

first(CAFO_Exposure, false) AS `CAFO_Exposure`
first(Landfill_Exposure, false) AS `Landfill_Exposure`
first(EstResidentialDensity, false) AS `EstResidentialDensity`
first(EstProbabilityHighSchoolMaxEducation, false) AS
`EstProbabilityHighSchoolMaxEducation`
first(EstProbabilityNoHealthIns, false) AS `EstProbabilityNoHealthIns`
first(EstProbabilityHouseholdNonHispWhite, false) AS
`EstProbabilityHouseholdNonHispWhite`
first(EstProbabilityESL, false) AS `EstProbabilityESL`
first(EstProbabilityNonHispWhite, false) AS
`EstProbabilityNonHispWhite`
first(EstHouseholdIncome, false) AS `EstHouseholdIncome`
first(EstResidentialDensity25Plus, false) AS
`EstResidentialDensity25Plus`
first(EstProbabilityNoAuto, false) AS `EstProbabilityNoAuto`
first(ur, false) AS `ur`
CAST(sum(Propranolol) AS INT) AS `Propranolol`
CAST(sum(Cetirizine) AS INT) AS `Cetirizine`
CAST(sum(PregnancyDx) AS INT) AS `PregnancyDx`
CAST(sum(Fluoxetine) AS INT) AS `Fluoxetine`
CAST(sum(ObesityDx) AS INT) AS `ObesityDx`
CAST(sum(Fluticasone) AS INT) AS `Fluticasone`
CAST(sum(Mometasone) AS INT) AS `Mometasone`
CAST(sum(Leuprolide) AS INT) AS `Leuprolide`
CAST(sum(Albuterol) AS INT) AS `Albuterol`
CAST(sum(CroupDx) AS INT) AS `CroupDx`
CAST(sum(ReactiveAirwayDx) AS INT) AS `ReactiveAirwayDx`
CAST(sum(Ciclesonide) AS INT) AS `Ciclesonide`
CAST(sum(Ipratropium) AS INT) AS `Ipratropium`
CAST(sum(KidneyCancerDx) AS INT) AS `KidneyCancerDx`
CAST(sum(Diphenhydramine) AS INT) AS `Diphenhydramine`
CAST(sum(Escitalopram) AS INT) AS `Escitalopram`
CAST(sum(Fexofenadine) AS INT) AS `Fexofenadine`
CAST(sum(EndometriosisDx) AS INT) AS `EndometriosisDx`
CAST(sum(OvarianDysfunctionDx) AS INT) AS
`OvarianDysfunctionDx`
CAST(sum(Metaproterenol) AS INT) AS `Metaproterenol`
CAST(sum(TesticularCancerDx) AS INT) AS `TesticularCancerDx`
CAST(sum(Tamoxifen) AS INT) AS `Tamoxifen`
CAST(sum(Goserelin) AS INT) AS `Goserelin`
CAST(sum(DepressionDx) AS INT) AS `DepressionDx`
CAST(sum(AutismDx) AS INT) AS `AutismDx`
CAST(sum(Progesterone) AS INT) AS `Progesterone`
CAST(sum(CoughDx) AS INT) AS `CoughDx`
CAST(sum(Omalizumab) AS INT) AS `Omalizumab`
CAST(sum(AnxietyDx) AS INT) AS `AnxietyDx`
CAST(sum(Citalopram) AS INT) AS `Citalopram`
CAST(sum(Beclomethasone) AS INT) AS `Beclomethasone`
CAST(sum(Theophylline) AS INT) AS `Theophylline`
CAST(sum(FibromyalgiaDx) AS INT) AS `FibromyalgiaDx`
CAST(sum(Sertraline) AS INT) AS `Sertraline`
CAST(sum(Venlafaxine) AS INT) AS `Venlafaxine`
CAST(sum(DrugDependenceDx) AS INT) AS `DrugDependenceDx`
CAST(sum(Formoterol) AS INT) AS `Formoterol`
CAST(sum(AlcoholDependenceDx) AS INT) AS
`AlcoholDependenceDx`
CAST(sum(AlopeciaDx) AS INT) AS `AlopeciaDx`
CAST(sum(MenopauseDx) AS INT) AS `MenopauseDx`
CAST(sum(CervicalCancerDx) AS INT) AS `CervicalCancerDx`
CAST(sum(Mepolizumab) AS INT) AS `Mepolizumab`
CAST(sum(TesticularDysfunctionDx) AS INT) AS
`TesticularDysfunctionDx`
CAST(sum(Estropipate) AS INT) AS `Estropipate`
CAST(sum(Histrelin) AS INT) AS `Histrelin`
CAST(sum(Triptorelin) AS INT) AS `Triptorelin`
CAST(sum(Salmeterol) AS INT) AS `Salmeterol`
CAST(sum(Arformoterol) AS INT) AS `Arformoterol`
CAST(sum(Paroxetine) AS INT) AS `Paroxetine`
CAST(sum(Flunisolide) AS INT) AS `Flunisolide`
CAST(sum(Testosterone) AS INT) AS `Testosterone`
CAST(sum(Budesonide) AS INT) AS `Budesonide`
CAST(sum(DiabetesDx) AS INT) AS `DiabetesDx`
CAST(sum(Metformin) AS INT) AS `Metformin`
CAST(sum(Nandrolone) AS INT) AS `Nandrolone`
CAST(sum(Prasterone) AS INT) AS `Prasterone`
CAST(sum(AsthmaDx) AS INT) AS `AsthmaDx`
CAST(sum(Indacaterol) AS INT) AS `Indacaterol`
CAST(sum(Androstenedione) AS INT) AS `Androstenedione`
CAST(sum(Duloxetine) AS INT) AS `Duloxetine`
CAST(sum(Prednisone) AS INT) AS `Prednisone`
CAST(sum(PneumoniaDx) AS INT) AS `PneumoniaDx`
CAST(sum(UterineCancerDx) AS INT) AS `UterineCancerDx`
CAST(sum(Medroxyprogesterone) AS INT) AS
`Medroxyprogesterone`
CAST(sum(Hydroxyzine) AS INT) AS `Hydroxyzine`
CAST(sum(ProstateCancerDx) AS INT) AS `ProstateCancerDx`
CAST(sum(Estrogen) AS INT) AS `Estrogen`
CAST(sum(Trazodone) AS INT) AS `Trazodone`
CAST(sum(Estradiol) AS INT) AS `Estradiol`
CAST(sum(OvarianCancerDx) AS INT) AS `OvarianCancerDx`

```

The result of the aggregation is in `/icees2/2010/all_patient`. Intermediate files are created `/icees2/2010/visit_combined` and `/icees2/2010/all_visit` which contain the union of all patient's visit info. The only difference between the two are renames of columns and default missing values.

## Outputs:

“FHIR-PIT/data/output/icees2/2010/visit\_combined”

It contains all the icees/<year> information for all patients combined. The schema is the union of all columns from all files in icees/<year> with 0 as a default for missing values.

AgeVisit,ALD2\_ppbv,ALD2\_ppbv\_avg,ALD2\_ppbv\_max,ALD2\_ppbv\_prev\_date,BENZ\_ppbv,BENZ\_ppbv\_avg,BENZ\_ppbv\_max,BENZ\_ppbv\_prev\_date,birth\_date,bucket,CAFO\_Exposure,CO\_ppbv,CO\_ppbv\_avg,CO\_ppbv\_max,CO\_ppbv\_prev\_date,encounter\_num,EstHouseholdIncome,EstProbabilityESL,EstProbabilityHighSchoolMaxEducation,EstProbabilityHouseholdNonHispWhite,EstProbabilityNoAuto,EstProbabilityNoHealthIns,EstProbabilityNonHispWhite,EstResidentialDensity,EstResidentialDensity25Plus,Ethnicity,FORM\_ppbv,FORM\_ppbv\_avg,FORM\_ppbv\_max,FORM\_ppbv\_prev\_date,Landfill\_Exposure,MajorRoadwayHighwayExposure,NO2\_ppbv,NO2\_ppbv\_avg,NO2\_ppbv\_max,NO2\_ppbv\_prev\_date,NO\_ppbv,NO\_ppbv\_avg,NO\_ppbv\_max,NO\_ppbv\_prev\_date,NOX\_ppbv,NOX\_ppbv\_avg,NOX\_ppbv\_max,NOX\_ppbv\_prev\_date,ObesityBMIVisit,ozone\_daily\_8hour\_maximum,ozone\_daily\_8hour\_maximum\_avg,ozone\_daily\_8hour\_maximum\_max,ozone\_daily\_8hour\_maximum\_prev\_date,patient\_num,pm25\_daily\_average,pm25\_daily\_average\_avg,pm25\_daily\_average\_max,pm25\_daily\_average\_prev\_date,Race,RoadwayAADT,RoadwayDistanceExposure,RoadwayLanes,RoadwaySpeedLimit,RoadwayType,Sex,SO2\_ppbv,SO2\_ppbv\_avg,SO2\_ppbv\_max,SO2\_ppbv\_prev\_date,start\_date,study\_period,ur,VisitType,year

“FHIR-PIT/data/output/icees2/2010/all\_visit”

Renamed and less columns than /icees2/2010/visit\_combined

AgeVisit,Avg24hAcetaldehydeExposure\_2,Avg24hBenzeneExposure\_2,birth\_date,bucket,CAFO\_Exposure,Avg24hCOExposure\_2,encounter\_num,EstHouseholdIncome,EstProbabilityESL,EstProbabilityHighSchoolMaxEducation,EstProbabilityHouseholdNonHispWhite,EstProbabilityNoAuto,EstProbabilityNoHealthIns,EstProbabilityNonHispWhite,EstResidentialDensity,EstResidentialDensity25Plus,Ethnicity,Avg24hFormaldehydeExposure\_2,Landfill\_Exposure,MajorRoadwayHighwayExposure,Avg24hNO2Exposure\_2,Avg24hNOExposure\_2,Avg24hNOxExposure\_2,ObesityBMIVisit,Max24hOzoneExposure\_2,pa



tient\_num,Avg24hPM2.5Exposure\_2,Race,RoadwayAADT,RoadwayDistanceExposure,RoadwayLanes,RoadwaySpeedLimit,RoadwayType,Sex,Avg24hSO2Exposure\_2,start\_date,study\_period,ur,VisitType,year,BeclomethasoneVisit,EstradiolVisit,ProstateCancerDxVisit,AlopeciaDxVisit,UterineCancerDxVisit,AutismDxVisit,DepressionDxVisit,TheophyllineVisit,FormoterolVisit,OvarianDysfunctionDxVisit,AlcoholDependenceDxVisit,MetforminVisit,ObesityDxVisit,ParoxetineVisit,DiphenhydramineVisit,PregnancyDxVisit,MepolizumabVisit,MedroxyprogesteroneVisit,LeuprolideVisit,CiclesonideVisit,ReactiveAirwayDxVisit,EndometriosisDxVisit,AnxietyDxVisit,FibromyalgiaDxVisit,FluticasoneVisit,IndacaterolVisit,CetirizineVisit,HydroxyzineVisit,FexofenadineVisit,PrednisoneVisit,TamoxifenVisit,MometasoneVisit,TrazodoneVisit,PropranololVisit,CoughDxVisit,NandroloneVisit,DrugDependenceDxVisit,ArformoterolVisit,BudesonideVisit,OvarianCancerDxVisit,CroupDxVisit,VenlafaxineVisit,MetaproterenolVisit,HistrelinVisit,OmalizumabVisit,MenopauseDxVisit,PneumoniaDxVisit,DiabetesDxVisit,DuloxetineVisit,EstropipateVisit,GoserelinVisit,CitalopramVisit,CervicalCancerDxVisit,AndrostenedioneVisit,ProgesteroneVisit,TestosteroneVisit,SertralineVisit,EscitalopramVisit,EstrogenVisit,TriptorelinVisit,SalmeterolVisit,IpratropiumVisit,KidneyCancerDxVisit,TesticularCancerDxVisit,AsthmaDxVisit,FlunisolideVisit,PrasteroneVisit,TesticularDysfunctionDxVisit,FluoxetineVisit,AlbuterolVisit,Sex2,Avg24hPM2.5Exposure,Max24hPM2.5Exposure,Avg24hOzoneExposure,Max24hOzoneExposure

“FHIR-PIT/data/output/icees2/2010/all\_patient”

This file is the result of the aggregation.

**Note:** The last portion of the icees\_features.yaml file lists the Dx whose presence indicates a RespiratoryDx. That is RespiratoryDx is True if the visit had any of the Dx listed. Then, for each given (patient, study\_period),

**TotalEDVisits/TotalInpatientVisits/TotalEDInpatientVisits** count the total number of visits that indicate a RespiratoryDx and satisfy a given visitType. For instance, TotalInpatientVisits is the total count of “IMP”

type visits registered with “RespiratoryDx” flag for a given (patient, year); TotalEDVisits is the total count of “AMB”, “EMER” type visits registered with “RespiratoryDx” flag for a given (patient, year)

patient\_num, study\_period, AvgDailyPM2.5Exposure, AvgDailyPM2.5Exposure\_StudyAvg, AvgDailyPM2.5Exposure\_StudyMax, MaxDailyPM2.5Exposure, MaxDailyPM2.5Exposure\_StudyAvg, MaxDailyPM2.5Exposure\_StudyMax, AvgDailyOzoneExposure, AvgDailyOzoneExposure\_StudyAvg, AvgDailyOzoneExposure\_StudyMax, MaxDailyOzoneExposure, MaxDailyOzoneExposure\_StudyAvg, MaxDailyOzoneExposure\_StudyMax, AvgDailyPM2.5Exposure\_2, MaxDailyOzoneExposure\_2, AvgDailyCOExposure\_2, AvgDailyNOExposure\_2, AvgDailyNO2Exposure\_2, AvgDailyNOxExposure\_2, AvgDailySO2Exposure\_2, AvgDailyAcetaldehydeExposure\_2, AvgDailyFormaldehydeExposure\_2, AvgDailyBenzeneExposure\_2, ObesityBMI, TotalEDVisits, TotalInpatientVisits, TotalEDInpatientVisits, Sex2, birth\_date, Sex, Race, Ethnicity, MajorRoadwayHighwayExposure, RoadwayDistanceExposure, RoadwayType, RoadwayAADT, RoadwaySpeedLimit, RoadwayLanes, CAFO\_Exposure, Landfill\_Exposure, EstResidentialDensity, EstProbabilityHighSchoolMaxEducation, EstProbabilityNoHealthIns, EstProbabilityHouseholdNonHispWhite, EstProbabilityESL, EstProbabilityNonHispWhite, EstHouseholdIncome, EstResidentialDensity25Plus, EstProbabilityNoAuto, ur, Prednisone, OvarianDysfunctionDx, Hydroxyzine, Albuterol, Diphenhydramine, TesticularDysfunctionDx, Estradiol, PneumoniaDx, Formoterol, DiabetesDx, Metformin, Beclomethasone, Omalizumab, OvarianCancerDx, Venlafaxine, PregnancyDx, Fexofenadine, Arformoterol, Mometasone, Medroxyprogesterone, DepressionDx, Propranolol, Salmeterol, TesticularCancerDx, Budesonide, Prasterone, Trazodone, AlcoholDependenceDx, Androstenedione, Ipratropium, Tamoxifen, Nandrolone, Duloxetine, EndometriosisDx, CroupDx, Escitalopram, Metaproterenol, ObesityDx, CoughDx, Histrelin, AutismDx, Estrogen, Ciclesonide, UterineCancerDx, Sertraline, Testosterone, Estropipate, Indacaterol, AsthmaDx, Triptorelin, KidneyCancerDx, Progesterone, ReactiveAirwayDx, Ghrelin, ProstateCancerDx, AnxietyDx, Theophylline, DrugDependenceDx, Flunisolide, CervicalCancerDx, FibromyalgiaDx, Leuprolide, Fluoxetine, Citalopram, MenopauseDx, Cetirizine, Fluticasone, Paroxetine, Alopecia

ciaDx,Mepolizumab,AgeStudyStart,Active_In_Study_Period
--

## **BinICEES**

### **Takes in:**

Config\_file: "FHIR-PIT/spark/config/icees\_features.yaml",

Input\_dir: "FHIR-PIT/data/output/icees2",

Output\_dir: "FHIR-PIT/data/output/icees2\_bins",

Study\_periods: List of years

### **Description:**

Note: icees2/year has three files {all\_visit, all\_patient, visit\_combined}.

Add index information to dataframe.

The goal is to bin the features of all\_visit and all\_patient FHIR PIT outputs. Binning consists in categorizing the features based on bins. The bins are created to spread the data uniformly across each bin (i.e. qcut) or to bin the domain into evenly sized bins (i.e. cut). Note the difference between qcut and cut. Qcut may not produced evenly sized bins, and cut may not produced bins with the same amount of data in each.

The process for all\_patient and all\_visit is the same except for minor details. For each year get file {input\_dir}/{year}/{all\_patient, all\_visit} and:

1. PreProcAge bins age: Bins age into two groups, ['<5', '5-17', '18-44', '45-64', '65-89'] and ['0-2', '3-17', '18-34', '35-50', '51-69', '70-89'].
2. PreProcEnv bins environmental variables.
3. PreProcSocial bins exposure variables.
4. Cut\_col bins columns in the intersection between all\_patient/all\_visit cols and icees\_features.FHIR.keys()

The data is then deidentified by dropping: patient\_num and birthdate for all\_patient. And dropping a larger list (cols\_to\_drop) for all\_visit.

Save the result to: "/output/EPR\_binned/EPR\_binned"

### **Outputs:**

"FHIR-PIT/data/output/EPR\_binned/EPR\_binned"

### **BinEPR**

#### **Takes**

**in:**

"/home/jjgarcia/projects/fhir\_sample\_data/TLR4\_AllData\_NewHash\_0  
1292020 NO PII\_no\_new\_line.csv",  
"/home/jjgarcia/projects/fhir\_sample\_data/UNC\_NIEHS\_XWalk\_for\_H  
ao\_shape\_h3.csv",  
"FHIR-PIT/data/output/icees2\_bins/",  
"FHIR-PIT/data/output/FHIR\_processed/geo.csv",  
"FHIR-PIT/data/output/EPR\_binned/EPR\_binned",  
study\_periods

#### **Description:**

[Important] Cw has two columns [patient\_num, HASH\_VALUE]

Df has multiple columns, mostly on EPR\_binned.

Df\_pat\_geo is loaded from **geo.csv**

Preprocess df and write it to **EPR\_binned**.

JOIN LEFT df with cw on HASH\_VALUE; creates column In\_EPR;  
writes it to **EPR\_binned\_pat**. Save to df\_pat

JOIN LEFT df with cw on hash, then OUTER JOIN result with  
df\_pat\_geo[patient\_num]. Save to df\_pat\_ord.

The difference between df\_pat and df\_pat\_ord are the geo.csv  
columns.

For each year in study\_periods:

- Load Dfp from **icees2\_bins/{year}patient**.

- OUTER JOIN df\_pat with dfp on patient\_num. LEFT JOIN the result with df\_pat\_ord on “patient\_num” and “hash\_value”. This will keep all the null values from the initial outer join and attach the geo information to non-null matches.
- Save result to dfpe.
- Preprocess na values, sort, cast floats to ints from dfpe. Save to **EPR\_binned{year}patient**
- Drop [“patient\_num”, “hash\_value”] from dfpe. Save to **EPR\_binned{year}patient\_deidentified**.

Repeat the same process as above with two modifications: dfp is loaded from **icees2\_bins/{year}visit**; dfpe is the OUTER JOIN df\_pat with dfp on patient\_num. It essentially contains the visit information and the hash information from df\_pat (it contains patient\_num, hash\_value and all the original df information).

The results are saved to: **EPR\_binned{year}visit, EPR\_binned{year}visit\_deidentified**.

Lastly, drop “hash\_value” from df and save it to **EPR\_binned\_deidentified**.

### Outputs:

"FHIR-PIT/data/output/EPR\_binned/EPR\_binned"

**HASH\_VALUE**,IN\_FINAL\_SAMPLE,ORIGINAL\_ANALYTIC\_SAMPLE,TLR4\_AGE,GENDER,RACE,ETHNICITY,CURRENT\_AGE,RESPONDER\_STATUS,QXAGE,BMI\_CATEGORY,D28\_ASTHMA,D28A\_ASTHMA\_AD\_TEXT,D28B\_STILL\_HAVE\_ASTHMA,D28C\_ASTHMA\_EPISODE\_12M,D28D\_ASTHMA\_ER\_VISIT\_12M,D28E\_ASTHMA\_MED\_TAKE\_12M,D28F\_ASTHMA\_14D\_NUM\_NIGHTS\_TEXT,D28G\_ASTHMA\_14D\_LIMIT\_DAYS\_TEXT,D28H\_ASTHMA\_14D\_NUM\_WHEEZE\_TEXT,S177\_SMOKE\_100\_CIG\_LIFETIME,SMOKE\_CAT,SNP1,SNP2,SNP3,SNP4,ESTTOTALPOP,DISTANCE,AADT,RO

ADTYPE,SPEED,THROUGH\_LANES,O3\_N\_OBS,PM25\_N\_OBS,TLR4\_AGE2,CURRENT\_AGE2,QXAGE2,D28A\_ASTHMA\_AD\_TEXT2,TLR4\_DIST\_1X\_qcut,TLR4\_DIST\_1X\_cut,TLR4\_DIST\_2X\_qcut,TLR4\_DIST\_2X\_cut,TLR4\_DIST\_3X\_qcut,TLR4\_DIST\_3X\_cut,O3\_ANNUAL\_AVERAGE\_qcut,O3\_ANNUAL\_AVERAGE\_cut,PM25\_ANNUAL\_AVERAGE\_qcut,PM25\_ANNUAL\_AVERAGE\_cut,ESTTOTALPOP25PLUS\_qcut,ESTTOTALPOP25PLUS\_cut,ESTPROPPERSONSNONHISPWHITE\_qcut,ESTPROPPERSONSNONHISPWHITE\_cut,ESTPROPPERSONS25PLUSHSMAX\_qcut,ESTPROPPERSONS25PLUSHSMAX\_cut,ESTPROPHOUSEHOLDSNOAUTO\_qcut,ESTPROPHOUSEHOLDSNOAUTO\_cut,ESTPROPPERSONSNOEALTHINS\_qcut,ESTPROPPERSONSNOEALTHINS\_cut,ESTPROPPERSONS5PLUSNOENGLISH\_qcut,ESTPROPPERSONS5PLUSNOENGLISH\_cut,MEDIANHOUSEHOLDINCOME\_qcut,MEDIANHOUSEHOLDINCOME\_cut,DISTANCE2

"FHIR-PIT/data/output/EPR\_binned/EPR\_binned\_pat"

**HASH\_VALUE**,IN\_FINAL\_SAMPLE,ORIGINAL\_ANALYTIC\_SAMPLE,TLR4\_AGE,GENDER,RACE,ETHNICITY,CURRENT\_AGE,RESPONDER\_STATUS,QXAGE,BMI\_CATEGORY,D28\_ASTHMA,D28A\_ASTHMA\_AD\_TEXT,D28B\_STILL\_HAVE\_ASTHMA,D28C\_ASTHMA\_EPISODE\_12M,D28D\_ASTHMA\_ER\_VISIT\_12M,D28E\_ASTHMA\_MED\_TAKE\_12M,D28F\_ASTHMA\_14D\_NUM\_NIGHTS\_TEXT,D28G\_ASTHMA\_14D\_LIMIT\_DAYS\_TEXT,D28H\_ASTHMA\_14D\_NUM\_WHEEZE\_TEXT,S177\_SMOKE\_100\_CIG\_LIFETIME,SMOKE\_CAT,SNP1,SNP2,SNP3,SNP4,ESTTOTALPOP,DISTANCE,AADT,ROADTYPE,SPEED,THROUGH\_LANES,O3\_N\_OBS,PM25\_N\_OBS,TLR4\_AGE2,CURRENT\_AGE2,QXAGE2,D28A\_ASTHMA\_AD\_TEXT2,TLR4\_DIST\_1X\_qcut,TLR4\_DIST\_1X\_cut,TLR4\_DIST\_2X\_qcut,TLR4\_DIST\_2X\_cut,TLR4\_DIST\_3X\_qcut,TLR4\_DIST\_3X\_cut,O3\_ANNUAL\_AVERAGE\_qcut,O3\_ANNUAL\_AVERAGE\_cut,PM25\_ANNUAL\_AVERAGE\_qcut,PM25\_ANNUAL\_AVERAGE\_cut,ESTTOTALPOP25PLUS\_qcut,ESTTOTALPOP25PLUS\_cut,ESTPROPPERSONSNONHISPWHITE\_qcut,ESTPROPPERSONSNONHISPWHITE\_cut

t,ESTPROPPERSONS25PLUSHSMAX\_qcut,ESTPROPPERSONS25PLUSHSMAX\_cut,ESTPROPHOUSEHOLDSNOAUTO\_qcut,ESTPROPHOUSEHOLDSNOAUTO\_cut,ESTPROPPERSONSNOHEALTHINS\_qcut,ESTPROPPERSONSNOHEALTHINS\_cut,ESTPROPPERSONS5PLUSNOENGLISH\_qcut,ESTPROPPERSONS5PLUSNOENGLISH\_cut,MEDIANHOUSEHOLDINCOME\_qcut,MEDIANHOUSEHOLDINCOME\_cut,DISTANCE2,**patient\_num,IN\_EPR**

"FHIR-PIT/data/output/EPR\_binned/EPR\_binned2010patient"

HASH\_VALUE,patient\_num,index,IN\_FINAL\_SAMPLE,ORIGINAL\_ANALYTIC\_SAMPLE,TLR4\_AGE,GENDER,RACE,ETHNICITY,CURRENT\_AGE,RESPONDER\_STATUS,QXAGE,BMI\_CATEGORY,D28\_ASTHMA,D28A\_ASTHMA\_AD\_TEXT,D28B\_STILL\_HAVE\_ASTHMA,D28C\_ASTHMA\_EPISODE\_12M,D28D\_ASTHMA\_ER\_VISIT\_12M,D28E\_ASTHMA\_MED\_TAKE\_12M,D28F\_ASTHMA\_14D\_NUM\_NIGHTS\_TEXT,D28G\_ASTHMA\_14D\_LIMIT\_DAYS\_TEXT,D28H\_ASTHMA\_14D\_NUM\_WHEEZE\_TEXT,S177\_SMOKE\_100\_CIG\_LIFETIME,SMOKE\_CAT,SNP1,SNP2,SNP3,SNP4,ESTTOTALPOP,DISTANCE,AADT,ROADTYPE,SPEED,THROUGH\_LANES,O3\_N\_OBS,PM25\_N\_OBS,TLR4\_AGE2,CURRENT\_AGE2,QXAGE2,D28A\_ASTHMA\_AD\_TEXT2,TLR4\_DIST\_1X\_qcut,TLR4\_DIST\_1X\_cut,TLR4\_DIST\_2X\_qcut,TLR4\_DIST\_2X\_cut,TLR4\_DIST\_3X\_qcut,TLR4\_DIST\_3X\_cut,O3\_ANNUAL\_AVERAGE\_qcut,O3\_ANNUAL\_AVERAGE\_cut,PM25\_ANNUAL\_AVERAGE\_qcut,PM25\_ANNUAL\_AVERAGE\_cut,ESTTOTALPOP25PLUS\_qcut,ESTTOTALPOP25PLUS\_cut,ESTPROPPERSONSNONHISPWHITE\_qcut,ESTPROPPERSONSNONHISPWHITE\_cut,ESTPROPPERSONS25PLUSHSMAX\_qcut,ESTPROPPERSONS25PLUSHSMAX\_cut,ESTPROPHOUSEHOLDSNOAUTO\_qcut,ESTPROPHOUSEHOLDSNOAUTO\_cut,ESTPROPPERSONSNOHEALTHINS\_qcut,ESTPROPPERSONSNOHEALTHINS\_cut,ESTPROPPERSONS5PLUSNOENGLISH\_qcut,ESTPROPPERSONS5PLUSNOENGLISH\_cut,MEDIANHOUSEHOLDINCOME\_qcut,MEDIANHOUSEHOLDINCOME\_cut,DISTANCE2,IN\_EPR,study\_period,AvgDailyPM2.5Exposure,AvgDailyPM2.5Exposure\_StudyAvg,

AvgDailyPM2.5Exposure\_StudyMax,MaxDailyPM2.5Exposure,MaxDailyPM2.5Exposure\_StudyAvg,MaxDailyPM2.5Exposure\_StudyMax,AvgDailyOzoneExposure,AvgDailyOzoneExposure\_StudyAvg,AvgDailyOzoneExposure\_StudyMax,MaxDailyOzoneExposure,MaxDailyOzoneExposure\_StudyAvg,MaxDailyOzoneExposure\_StudyMax,AvgDailyPM2.5Exposure\_2,MaxDailyOzoneExposure\_2,AvgDailyCOExposure\_2,AvgDailyNOExposure\_2,AvgDailyNO2Exposure\_2,AvgDailyNOxExposure\_2,AvgDailySO2Exposure\_2,AvgDailyAcetaldehydeExposure\_2,AvgDailyFormaldehydeExposure\_2,AvgDailyBenzeneExposure\_2,ObesityBMI,TotalEDVisits,TotalInpatientVisits,TotalEDInpatientVisits,Sex2,Sex,Race,Ethnicity,MajorRoadwayHighwayExposure,RoadwayDistanceExposure,RoadwayType,RoadwayAADT,RoadwaySpeedLimit,RoadwayLanes,CAFO\_Exposure,Landfill\_Exposure,EstResidentialDensity,EstProbabilityHighSchoolMaxEducation,EstProbabilityNoHealthIns,EstProbabilityHouseholdNonHispWhite,EstProbabilityESL,EstProbabilityNonHispWhite,EstHouseholdIncome,EstResidentialDensity25Plus,EstProbabilityNoAuto,ur,PneumoniaDx,Ciclesonide,DrugDependenceDx,Leuprolide,CervicalCancerDx,AutismDx,TesticularCancerDx,Flunisolide,Escitalopram,AlopeciaDx,Propranolol,Ipratropium,Nandrolone,Testosterone,Prasterone,Trazodone,Venlafaxine,Estrogen,OvarianCancerDx,Sertraline,EndometriosisDx,CoughDx,Fluticasone,AnxietyDx,PregnancyDx,FibromyalgiaDx,Metformin,Citalopram,Paroxetine,Cetirizine,Androstenedione,Fluoxetine,Duloxetine,ObesityDx,AlcoholDependenceDx,Formoterol,Estropipate,OvarianDysfunctionDx,Budesonide,Omalizumab,Prednisone,Hydroxyzine,Albuterol,Metaproterenol,Theophylline,AsthmaDx,TesticularDysfunctionDx,Salmeterol,Goserelin,Tamoxifen,Indacaterol,KidneyCancerDx,CroupDx,Progesterone,MenopauseDx,Estradiol,Histrelin,Mometasone,Diphenhydramine,ReactiveAirwayDx,ProstateCancerDx,DiabetesDx,Medroxyprogesterone,Mepolizumab,DepressionDx,Fexofenadine,Triptorelin,Beclomethasone,UterineCancerDx,Arformoterol,AgeStudyStart,Active\_In\_Study\_Period,AgeStudyStart2,AvgDailyPM2.5Exposure\_StudyAvg\_qcut,AvgDailyPM2.5Exposure\_StudyMax\_qcut,AvgDailyPM2.5Exposure\_qcut,MaxDailyPM2.5Exposure\_StudyAvg\_qcut,MaxDailyPM2.5Exposure\_StudyMax\_qcut,MaxDailyPM2.5Exposure\_qcut,AvgDailyOzoneExposure\_StudyAvg\_qcut,AvgDailyOzoneExposure\_qcut



e\_StudyMax\_qcut,AvgDailyOzoneExposure\_qcut,MaxDailyOzoneExposure\_StudyAvg\_qcut,MaxDailyOzoneExposure\_StudyMax\_qcut,MaxDailyOzoneExposure\_qcut,AvgDailyPM2.5Exposure\_2\_qcut,MaxDailyOzoneExposure\_2\_qcut,AvgDailyCOExposure\_2\_qcut,AvgDailyNOExposure\_2\_qcut,AvgDailyNO2Exposure\_2\_qcut,AvgDailyNOxExposure\_2\_qcut,AvgDailySO2Exposure\_2\_qcut,AvgDailyAcetaldehydeExposure\_2\_qcut,AvgDailyFormaldehydeExposure\_2\_qcut,AvgDailyBenzeneExposure\_2\_qcut,MajorRoadwayHighwayExposure2,RoadwayDistanceExposure2,IN\_ICEES

"FHIR-PIT/data/output/EPR\_binned/EPR\_binned2010patient\_deidentified"

index,IN\_FINAL\_SAMPLE,ORIGINAL\_ANALYTIC\_SAMPLE,TLR4\_AGE,GENDER,RACE,ETHNICITY,CURRENT\_AGE,RESPONDER\_STATUS,QXAGE,BMI\_CATEGORY,D28\_ASTHMA,D28A\_ASTHMA\_AD\_TEXT,D28B\_STILL\_HAVE\_ASTHMA,D28C\_ASTHMA\_EPISODE\_12M,D28D\_ASTHMA\_ER\_VISIT\_12M,D28E\_ASTHMA\_MED\_TAKE\_12M,D28F\_ASTHMA\_14D\_NUM\_NIGHTS\_TEXT,D28G\_ASTHMA\_14D\_LIMIT\_DAYS\_TEXT,D28H\_ASTHMA\_14D\_NUM\_WHEEZE\_TEXT,S177\_SMOKE\_100\_CIG\_LIFETIME,SMOKE\_CAT,SNP1,SNP2,SNP3,SNP4,ESTTOTALPOP,DISTANCE,AADT,ROADTYPE,SPEED,THROUGH\_LANES,O3\_N\_OBS,PM25\_N\_OBS,TLR4\_AGE2,CURRENT\_AGE2,QXAGE2,D28A\_ASTHMA\_AD\_TEXT2,TLR4\_DIST\_1X\_qcut,TLR4\_DIST\_1X\_cut,TLR4\_DIST\_2X\_qcut,TLR4\_DIST\_2X\_cut,TLR4\_DIST\_3X\_qcut,TLR4\_DIST\_3X\_cut,O3\_ANNUAL\_AVERAGE\_qcut,O3\_ANNUAL\_AVERAGE\_cut,PM25\_ANNUAL\_AVERAGE\_qcut,PM25\_ANNUAL\_AVERAGE\_cut,ESTTOTALPOP25PLUS\_qcut,ESTTOTALPOP25PLUS\_cut,ESTPROPPERSONSNONHISPWHITE\_qcut,ESTPROPPERSONSNONHISPWHITE\_cut,ESTPROPPERSONS25PLUSHSMAX\_qcut,ESTPROPPERSONS25PLUSHSMAX\_cut,ESTPROPHOUSEHOLDSNOAUTO\_qcut,ESTPROPHOUSEHOLDSNOAUTO\_cut,ESTPROPPERSONSNOHEALTHINS\_qcut,ESTPROPPERSONSNOHEALTHINS\_cut,ESTPROPPERSONS5

PLUSNOENGLISH\_qcut,ESTPROPPERSONS5PLUSNOENGLISH\_qcut,MEDIANHOUSEHOLDINCOME\_qcut,MEDIANHOUSEHOLDINCOME\_cut,DISTANCE2,IN\_EPR,study\_period,AvgDailyPM2.5Exposure,AvgDailyPM2.5Exposure\_StudyAvg,AvgDailyPM2.5Exposure\_StudyMax,MaxDailyPM2.5Exposure,MaxDailyPM2.5Exposure\_StudyAvg,MaxDailyPM2.5Exposure\_StudyMax,AvgDailyOzoneExposure,AvgDailyOzoneExposure\_StudyAvg,AvgDailyOzoneExposure\_StudyMax,MaxDailyOzoneExposure,MaxDailyOzoneExposure\_StudyAvg,MaxDailyOzoneExposure\_StudyMax,AvgDailyPM2.5Exposure\_2,MaxDailyOzoneExposure\_2,AvgDailyCOExposure\_2,AvgDailyNOExposure\_2,AvgDailyNO2Exposure\_2,AvgDailyNOxExposure\_2,AvgDailySO2Exposure\_2,AvgDailyAcetaldehydeExposure\_2,AvgDailyFormaldehydeExposure\_2,AvgDailyBenzeneExposure\_2,ObesityBMI,TotalEDVisits,TotalInpatientVisits,TotalEDInpatientVisits,Sex2,Sex,Race,Ethnicity,MajorRoadwayHighwayExposure,RoadwayDistanceExposure,RoadwayType,RoadwayAADT,RoadwaySpeedLimit,RoadwayLanes,CAFO\_Exposure,Landfill\_Exposure,EstResidentialDensity,EstProbabilityHighSchoolMaxEducation,EstProbabilityNoHealthIns,EstProbabilityHouseholdNonHispWhite,EstProbabilityESL,EstProbabilityNonHispWhite,EstHouseholdIncome,EstResidentialDensity25Plus,EstProbabilityNoAuto,ur,PneumoniaDx,Ciclesonide,DrugDependenceDx,Leuprolide,CervicalCancerDx,AutismDx,TesticularCancerDx,Flunisolide,Escitalopram,AlopeciaDx,Propranolol,Ipratropium,Nandrolone,Testosterone,Prasterone,Trazodone,Venlafaxine,Estrogen,OvarianCancerDx,Sertraline,EndometriosisDx,CoughDx,Fluticasone,AnxietyDx,PregnancyDx,FibromyalgiaDx,Metformin,Citalopram,Paroxetine,Cetirizine,Androstenedione,Fluoxetine,Duloxetine,ObesityDx,AlcoholDependenceDx,Formoterol,Estropipate,OvarianDysfunctionDx,Budesonide,Omalizumab,Prednisone,Hydroxyzine,Albuterol,Metaproterenol,Theophylline,AsthmaDx,TesticularDysfunctionDx,Salmeterol,Goserelin,Tamoxifen,Indacaterol,KidneyCancerDx,CroupDx,Progesterone,MenopauseDx,Estradiol,Histrelin,Mometasone,Diphenhydramine,ReactiveAirwayDx,ProstateCancerDx,DiabetesDx,Medroxyprogesterone,Mepolizumab,DepressionDx,Fexofenadine,Triptorelin,Beclomethasone,UterineCancerDx,Arformoterol,AgeStudyStart,Active\_In\_Study\_Period,AgeStudyStart2,AvgDailyPM2.5Exposure\_StudyAvg\_qcut,AvgDai

lyPM2.5Exposure\_StudyMax\_qcut,AvgDailyPM2.5Exposure\_qcut,MaxDailyPM2.5Exposure\_StudyAvg\_qcut,MaxDailyPM2.5Exposure\_StudyMax\_qcut,MaxDailyPM2.5Exposure\_qcut,AvgDailyOzoneExposure\_StudyAvg\_qcut,AvgDailyOzoneExposure\_StudyMax\_qcut,AvgDailyOzoneExposure\_qcut,MaxDailyOzoneExposure\_StudyAvg\_qcut,MaxDailyOzoneExposure\_StudyMax\_qcut,MaxDailyOzoneExposure\_qcut,AvgDailyPM2.5Exposure\_2\_qcut,MaxDailyOzoneExposure\_2\_qcut,AvgDailyCOExposure\_2\_qcut,AvgDailyNOExposure\_2\_qcut,AvgDailyNO2Exposure\_2\_qcut,AvgDailyNOxExposure\_2\_qcut,AvgDailySO2Exposure\_2\_qcut,AvgDailyAcetaldehydeExposure\_2\_qcut,AvgDailyFormaldehydeExposure\_2\_qcut,AvgDailyBenzeneExposure\_2\_qcut,MajorRoadwayHighwayExposure2,RoadwayDistanceExposure2,IN\_ICEES

"FHIR-PIT/data/output/EPR\_binned/EPR\_binned2010visit"

HASH\_VALUE,IN\_FINAL\_SAMPLE,ORIGINAL\_ANALYTIC\_SAMPLE,TLR4\_AGE,GENDER,RACE,ETHNICITY,CURRENT\_AGE,RESPONDER\_STATUS,QXAGE,BMI\_CATEGORY,D28\_ASTHMA,D28A\_ASTHMA\_AD\_TEXT,D28B\_STILL\_HAVE\_ASTHMA,D28C\_ASTHMA\_EPISODE\_12M,D28D\_ASTHMA\_ER\_VISIT\_12M,D28E\_ASTHMA\_MED\_TAKE\_12M,D28F\_ASTHMA\_14D\_NUM\_NIGHTS\_TEXT,D28G\_ASTHMA\_14D\_LIMIT\_DAYS\_TEXT,D28H\_ASTHMA\_14D\_NUM\_WHEEZE\_TEXT,S177\_SMOKE\_100\_CIG\_LIFETIME,SMOKE\_CAT,SNP1,SNP2,SNP3,SNP4,ESTTOTALPOP,DISTANCE,AADT,ROADTYPE,SPEED,THROUGH\_LANES,O3\_N\_OBS,PM25\_N\_OBS,TLR4\_AGE2,CURRENT\_AGE2,QXAGE2,D28A\_ASTHMA\_AD\_TEXT2,TLR4\_DIST\_1X\_qcut,TLR4\_DIST\_1X\_cut,TLR4\_DIST\_2X\_qcut,TLR4\_DIST\_2X\_cut,TLR4\_DIST\_3X\_qcut,TLR4\_DIST\_3X\_cut,O3\_ANNUAL\_AVERAGE\_qcut,O3\_ANNUAL\_AVERAGE\_cut,PM25\_ANNUAL\_AVERAGE\_qcut,PM25\_ANNUAL\_AVERAGE\_cut,ESTTOTALPOP25PLUS\_qcut,ESTTOTALPOP25PLUS\_cut,ESTPROPPERSONSNONHISPWHITE\_qcut,ESTPROPPERSONSNONHISPWHITE\_cut,ESTPROPPERSONS25PLUSHSMAX\_qcut,ESTPROPPERSONS25PLUSHSMAX\_cut,ESTPROPHOUSEHOLDSNOAUTO\_qcut,ESTP

ROPHOUSEHOLDSNOAUTO\_cut,ESTPROPPERSONSNOHEALTHINS\_qcut,ESTPROPPERSONSNOHEALTHINS\_cut,ESTPROPPERSONS5PLUSNOENGLISH\_qcut,ESTPROPPERSONS5PLUSNOENGLISH\_cut,MEDIANHOUSEHOLDINCOME\_qcut,MEDIANHOUSEHOLDINCOME\_cut,DISTANCE2,patient\_num,IN\_EPR,AgeVisit,bucket,CAFO\_Exposure,EstHouseholdIncome,EstProbabilityESL,EstProbabilityHighSchoolMaxEducation,EstProbabilityHouseholdNonHispanicWhite,EstProbabilityNoAuto,EstProbabilityNoHealthIns,EstProbabilityNonHispanicWhite,EstResidentialDensity,EstResidentialDensity25Plus,Ethnicity,Landfill\_Exposure,MajorRoadwayHighwayExposure,ObesityBMIVisit,Race,RoadwayAADT,RoadwayDistanceExposure,RoadwayLanes,RoadwaySpeedLimit,RoadwayType,Sex,study\_period,ur,VisitType,year,BeclomethasoneVisit,EstradiolVisit,ProstateCancerDxVisit,AlopeciaDxVisit,UterineCancerDxVisit,AutismDxVisit,DepressionDxVisit,TheophyllineVisit,FormoterolVisit,OvarianDysfunctionDxVisit,AlcoholDependenceDxVisit,MetforminVisit,ObesityDxVisit,ParoxetineVisit,DiphenhydramineVisit,PregnancyDxVisit,MepolizumabVisit,MedroxyprogesteroneVisit,LeuprolideVisit,CiclesonideVisit,ReactiveAirwayDxVisit,EndometriosisDxVisit,AnxietyDxVisit,FibromyalgiaDxVisit,FluticasoneVisit,IndacaterolVisit,CetirizineVisit,HydroxyzineVisit,FexofenadineVisit,PrednisoneVisit,TamoxifenVisit,MometasoneVisit,TrazodoneVisit,PropranololVisit,CoughDxVisit,NandroloneVisit,DrugDependenceDxVisit,ArformoterolVisit,BudesonideVisit,OvarianCancerDxVisit,CroupDxVisit,VenlafaxineVisit,MetaproterenolVisit,HistrelinVisit,OmalizumabVisit,MenopauseDxVisit,PneumoniaDxVisit,DiabetesDxVisit,DuloxetineVisit,EstropipateVisit,GoserelinVisit,CitalopramVisit,CervicalCancerDxVisit,AndrostenedioneVisit,ProgesteroneVisit,TestosteroneVisit,SertralineVisit,EscitalopramVisit,EstrogenVisit,TriptorelinVisit,SalmeterolVisit,IpratropiumVisit,KidneyCancerDxVisit,TesticularCancerDxVisit,AsthmaDxVisit,FlunisolideVisit,PrasteroneVisit,TesticularDysfunctionDxVisit,FluoxetineVisit,AlbuterolVisit,Sex2,AgeVisit2,MajorRoadwayHighwayExposure2,RoadwayDistanceExposure2,IN\_ICEES,index

"FHIR-PIT/data/output/EPR\_binned/EPR\_binned2010visit\_deidentified"

N\_FINAL\_SAMPLE,ORIGINAL\_ANALYTIC\_SAMPLE,TLR4\_AGE,GENDER,RACE,ETHNICITY,CURRENT\_AGE,RESPONDER\_STATU  
S,QXAGE,BMI\_CATEGORY,D28\_ASTHMA,D28A\_ASTHMA\_AD\_T  
EXT,D28B\_STILL\_HAVE\_ASTHMA,D28C\_ASTHMA\_EPISODE\_12  
M,D28D\_ASTHMA\_ER\_VISIT\_12M,D28E\_ASTHMA\_MED\_TAKE\_1  
2M,D28F\_ASTHMA\_14D\_NUM\_NIGHTS\_TEXT,D28G\_ASTHMA\_1  
4D\_LIMIT\_DAYS\_TEXT,D28H\_ASTHMA\_14D\_NUM\_WHEEZE\_TE  
XT,S177\_SMOKE\_100\_CIG\_LIFETIME,SMOKE\_CAT,SNP1,SNP2,  
SNP3,SNP4,ESTTOTALPOP,DISTANCE,AADT,ROADTYPE,SPEED  
,THROUGH\_LANES,O3\_N\_OBS,PM25\_N\_OBS,TLR4\_AGE2,CUR  
RENT\_AGE2,QXAGE2,D28A\_ASTHMA\_AD\_TEXT2,TLR4\_DIST\_1  
X\_qcut,TLR4\_DIST\_1X\_cut,TLR4\_DIST\_2X\_qcut,TLR4\_DIST\_2X\_  
cut,TLR4\_DIST\_3X\_qcut,TLR4\_DIST\_3X\_cut,O3\_ANNUAL\_AVERA  
GE\_qcut,O3\_ANNUAL\_AVERAGE\_cut,PM25\_ANNUAL\_AVERAGE  
\_qcut,PM25\_ANNUAL\_AVERAGE\_cut,ESTTOTALPOP25PLUS\_qc  
ut,ESTTOTALPOP25PLUS\_cut,ESTPROPPERSONSNONHISPWHI  
TE\_qcut,ESTPROPPERSONSNONHISPWHITE\_cut,ESTPROPPER  
SONS25PLUSHSMAX\_qcut,ESTPROPPERSONS25PLUSHSMAX\_  
cut,ESTPROPHOUSEHOLDSNOAUTO\_qcut,ESTPROPHOUSEHÖ  
LDSNOAUTO\_cut,ESTPROPPERSONSNOHEALTHINS\_qcut,ESTP  
ROPPERSONSNOHEALTHINS\_cut,ESTPROPPERSONS5PLUSN  
OENGLISH\_qcut,ESTPROPPERSONS5PLUSNOENGLISH\_cut,ME  
DIANHOUSEHOLDINCOME\_qcut,MEDIANHOUSEHOLDINCOME\_  
cut,DISTANCE2,IN\_EPR,AgeVisit,bucket,CAFO\_Exposure,EstHous  
eholdIncome,EstProbabilityESL,EstProbabilityHighSchoolMaxEduca  
tion,EstProbabilityHouseholdNonHispWhite,EstProbabilityNoAuto,Es  
tProbabilityNoHealthIns,EstProbabilityNonHispWhite,EstResidential  
Density,EstResidentialDensity25Plus,Ethnicity,Landfill\_Exposure,Maj  
orRoadwayHighwayExposure,ObesityBMIVisit,Race,RoadwayAADT,  
RoadwayDistanceExposure,RoadwayLanes,RoadwaySpeedLimit,Ro  
adwayType,Sex,study\_period,ur,VisitType,year,BeclomethasoneVisit  
,EstradiolVisit,ProstateCancerDxVisit,AlopeciaDxVisit,UterineCancer  
DxVisit,AutismDxVisit,DepressionDxVisit,TheophyllineVisit,Formoter  
olVisit,OvarianDysfunctionDxVisit,AlcoholDependenceDxVisit,Metfor  
minVisit,ObesityDxVisit,ParoxetineVisit,DiphenhydramineVisit,Pregn

ancyDxVisit,MepolizumabVisit,MedroxyprogesteroneVisit,LeuprolideVisit,CiclesonideVisit,ReactiveAirwayDxVisit,EndometriosisDxVisit,AnxietyDxVisit,FibromyalgiaDxVisit,FluticasoneVisit,IndacaterolVisit,CetirizineVisit,HydroxyzineVisit,FexofenadineVisit,PrednisoneVisit,TamoxifenVisit,MometasoneVisit,TrazodoneVisit,PropranololVisit,CoughDxVisit,NandroloneVisit,DrugDependenceDxVisit,ArformoterolVisit,BudesonideVisit,OvarianCancerDxVisit,CroupDxVisit,VenlafaxineVisit,MetaproterenolVisit,HistrelinVisit,OmalizumabVisit,MenopauseDxVisit,PneumoniaDxVisit,DiabetesDxVisit,DuloxetineVisit,EstropipateVisit,GoserelinVisit,CitalopramVisit,CervicalCancerDxVisit,AndrostenedioneVisit,ProgesteroneVisit,TestosteroneVisit,SertralineVisit,EscitalopramVisit,EstrogenVisit,TriptorelinVisit,SalmeterolVisit,IpratropiumVisit,KidneyCancerDxVisit,TesticularCancerDxVisit,AsthmaDxVisit,FlunisolideVisit,PrasteroneVisit,TesticularDysfunctionDxVisit,FluoxetineVisit,AlbuterolVisit,Sex2,AgeVisit2,MajorRoadwayHighwayExposure2,RoadwayDistanceExposure2,IN\_ICEES,index

"FHIR-PIT/data/output/EPR\_binned/EPR\_binned\_deidentified"

IN\_FINAL\_SAMPLE,ORIGINAL\_ANALYTIC\_SAMPLE,TLR4\_AGE,GENDER,RACE,ETHNICITY,CURRENT\_AGE,RESPONDER\_STATUSS,QXAGE,BMI\_CATEGORY,D28\_ASTHMA,D28A\_ASTHMA\_AD\_TEXT,D28B\_STILL\_HAVE\_ASTHMA,D28C\_ASTHMA\_EPISODE\_12M,D28D\_ASTHMA\_ER\_VISIT\_12M,D28E\_ASTHMA\_MED\_TAKE\_12M,D28F\_ASTHMA\_14D\_NUM\_NIGHTS\_TEXT,D28G\_ASTHMA\_14D\_LIMIT\_DAYS\_TEXT,D28H\_ASTHMA\_14D\_NUM\_WHEEZE\_TEXT,S177\_SMOKE\_100\_CIG\_LIFETIME,SMOKE\_CAT,SNP1,SNP2,SNP3,SNP4,ESTTOTALPOP,DISTANCE,AADT,ROADTYPE,SPEED,THROUGH\_LANES,O3\_N\_OBS,PM25\_N\_OBS,TLR4\_AGE2,CURRENT\_AGE2,QXAGE2,D28A\_ASTHMA\_AD\_TEXT2,TLR4\_DIST\_1X\_qcut,TLR4\_DIST\_1X\_cut,TLR4\_DIST\_2X\_qcut,TLR4\_DIST\_2X\_cut,TLR4\_DIST\_3X\_qcut,TLR4\_DIST\_3X\_cut,O3\_ANNUAL\_AVERAGE\_qcut,O3\_ANNUAL\_AVERAGE\_cut,PM25\_ANNUAL\_AVERAGE\_qcut,PM25\_ANNUAL\_AVERAGE\_cut,ESTTOTALPOP25PLUS\_qcut,ESTTOTALPOP25PLUS\_cut,ESTPROPPERSONSNONHISPWHI

```
TE_qcut,ESTPROPPERSONSNONHISPWHITE_cut,ESTPROPPERSONS25PLUSHSMAX_qcut,ESTPROPPERSONS25PLUSHSMAX_cut,ESTPROPHOUSEHOLDSNOAUTO_qcut,ESTPROPHOUSEHOLDSNOAUTO_cut,ESTPROPPERSONSNOHEALTHINS_qcut,ESTPROPPERSONSNOHEALTHINS_cut,ESTPROPPERSONS5PLUSNOENGLISH_qcut,ESTPROPPERSONS5PLUSNOENGLISH_cut,MEDIANHOUSEHOLDINCOME_qcut,MEDIANHOUSEHOLDINCOME_cut,DISTANCE2
```

## YAML configuration file: icees\_features.yaml

The Mapper object consumes this file for various purposes:

1. To label multiple FHIR concept codes with the same feature name. This is important for users to know code to feature correspondence of integrated tables.
2. To relabel data source features for the integrated feature tables: ACS, UR, HPMS, TL, Cafo, Landfill. This is important for users to have easier to understand feature names.
3. To aid a use case. For the asthma use case, it list the diagnosis of interest to count as “Respiratory” visits at ED and inpatient locations. This is important for asthma researchers.

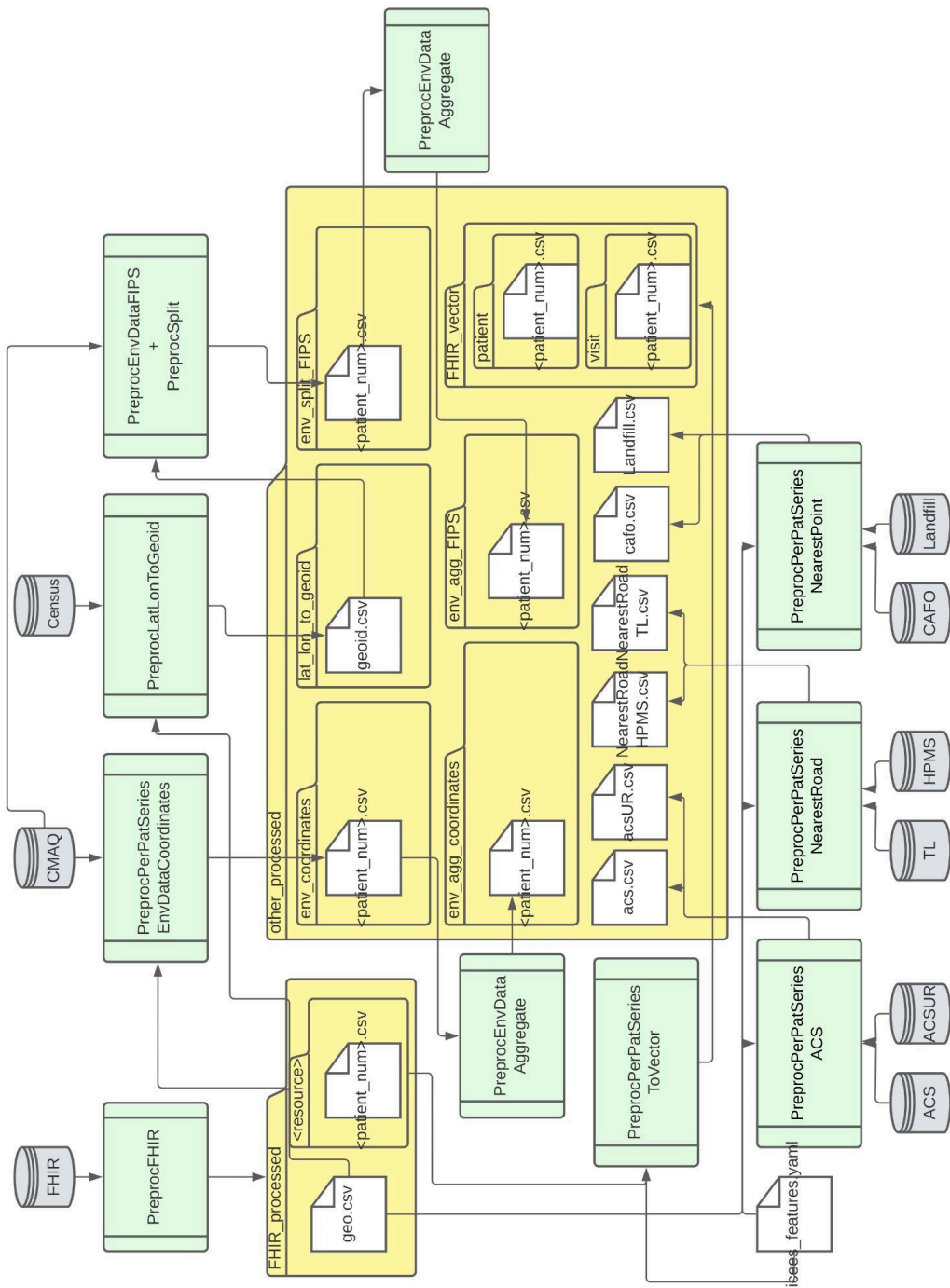


Figure 1: FHIR-PIT integrated feature tables pipeline. Arrows indicate dependency. Data sources are in gray, transformations in green, directories in yellow, and integrated feature tables in white. Please see the ‘Step description’ section for further details on the transformations.



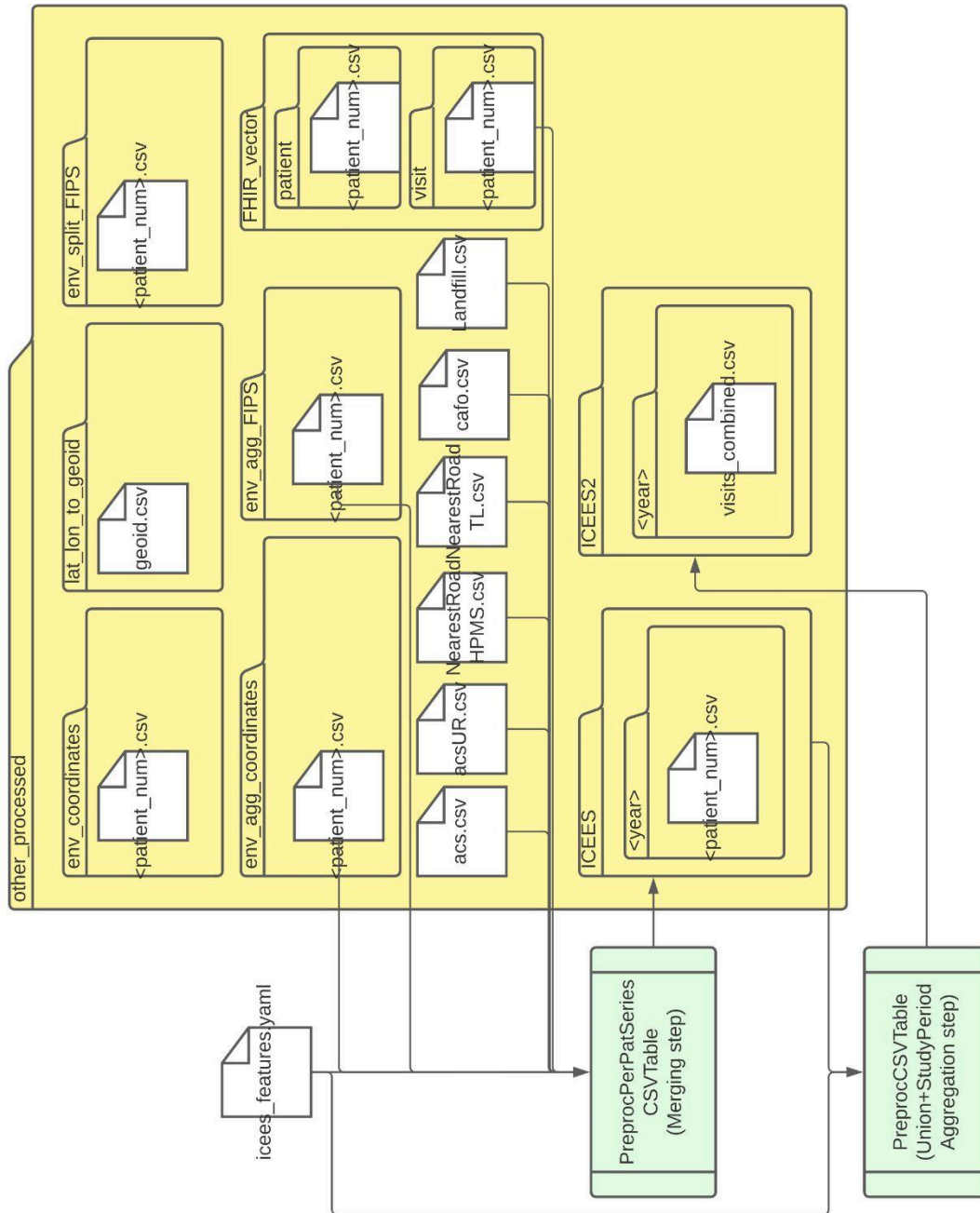


Figure 2: FHIR-PIT integrated feature tables merging and aggregation pipelines. Arrows indicate dependency. These transformations produce the final output to bin and deidentify. Data sources are in gray, transformations in green, directories in yellow, and integrated feature tables in white. Please see the ‘Step description’ section for further details on the transformations.