

ICEES – Environmental Exposures – Feature Variables

Table 1. ICEES feature variables – environmental exposures data: name, description, and binning strategy.*

Index Variables	Description and binning strategy
Index / PatientID	<i>Patient pseudo identifier, not available as feature variable, but allows for linkages across years when creating cohorts</i>
year	<i>Study period (calendar year, with a few exceptions)</i>
Active_In_Year	<i>Allows for selection of only those patients who are 'active' in a given study period or year, meaning that they had one or more visits with a healthcare provider</i>
Sex2	<i>Binned as "Male", "Female"</i>
Race_UNC Health* *Note that there are other race variables for cohorts that are not derived from UNC Health	<i>Binned as "Caucasian", "African American", "Asian", "Native Hawaiian/Pacific Islander", "American/Alaskan Native", "Other", "Unknown", "None"</i>
TotalEDVisits, TotalInpatientVisits, TotalEDInpatientVisits	<i>Total number of visits to ED, Inpatient Clinic, or Both over study period; binned as 0 ... >9</i>
Clinical Feature Variables	
Nomenclature format for medications: XXRX	<i>Medication XX prescribed or administered over study period, 0=no, 1=yes (one more more prescriptions/administrations over study period)</i>
Nomenclature format for diagnoses: XXDx	<i>Diagnosis XX made over study period, 0=no, 1=yes (one more more diagnoses over study period)</i>
Nomenclature format for procedures XX	<i>Procedure XX performed over study period, 0=no, 1=yes (one more more diagnoses over study period)</i>

Nomenclature format for laboratory measures	
XX	<i>Laboratory measurement (no units)</i>
XX _first_flag	<i>Flag (e.g., normal, above normal, below normal) for first laboratory measurement over study period</i>
XX _last_flag	<i>Flag (e.g., normal, above normal, below normal) for last laboratory measurement over study period</i>
US Census Bureau American Community Survey Data (years 2007-2011 & 2012-2016 survey samples)	
EstResidentialDensity	<i>Estimated total population [block group], binned according to US Census Bureau definitions (1=rural [0,2500), 2=urban cluster [2500,50000), 3=urbanized area [50000,inf))</i>
ur	<i>Estimated urban or rural residence, binned according to US Census Bureau definitions</i> <i>Note that u=rural, r=urban [known issue]</i>
EstResidentialDensity25Plus *2012-2016 sample only	<i>Estimated total population aged 25 years or older [block group], binned as quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
EstPropNonHispanicWhite	<i>Estimated proportion of persons who are non-Hispanic white [block group], binned as quartiles (pandas qcut) (1, 2, 3, 4)</i>
EstPropHouseholdNonHispanicWhite	<i>Estimated proportion of households that are non-Hispanic white [block group], binned as quartiles (pandas qcut) (1, 2, 3, 4)</i>
EstPropHighSchoolMaxEducation	<i>Estimated proportion persons aged 25 y or older with a HS diploma or less at their highest level of schooling [block group], binned as quartiles (pandas qcut) (1, 2, 3, 4)</i>
EstPropHighSchoolDropout	<i>Estimated proportion persons aged 16-19 y who are neither attending school nor HS graduates [block group], binned as quartiles (pandas qcut) (1, 2, 3, 4)</i>
EstPropHighSchoolDropoutNoWork	<i>Estimated proportion persons aged 16-19 y who are neither attending school nor HS graduates and are</i>

	<i>without work [block group], binned as quartiles (pandas qcut) (1, 2, 3, 4)</i>
EstPropHouseholdsNoAuto	<i>Estimated proportion of households without an automobile [block group], binned as quartiles (pandas qcut) (1, 2, 3, 4)</i>
EstPropHouseholdsNoHealthIns *2012-2016 sample only	<i>Estimated proportion of persons without health insurance [block group], binned as quartiles (pandas qcut) (1, 2, 3, 4)</i>
EstProp5PlusESL	<i>Estimated proportion of persons 5 years or older who sometimes speak a language other than English at home [block group], binned as quartiles (pandas qcut) (1, 2, 3, 4)</i>
EstMedianHouseholdIncome	<i>Estimated median household income [block group], binned as quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
EstPropMaleLittleWork	<i>Estimated proportion of males aged 16-64 years who worked less than 26 weeks in previous year [block group], binned as quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
EstPropHouseholdSSI	<i>Estimated proportion of households receiving Supplemental Security Income block group], binned as quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
EstPropHouseholdPA	<i>Estimated proportion of households receiving Public Assistance [block group], binned as quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
EstPropFemaleHouseholdNoSpouse	<i>Estimated proportion of family households headed by a female (no male partner present) [block group], binned as quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
EstPropFemaleHouseholdFamilyChild	<i>Estimated proportion of total households headed by a female with family children aged 18 y or less (no male partner present) [block group], binned as quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
EstPropFemaleHouseholdAnyChild	<i>Estimated proportion of total households headed by a female with any children aged 18 y or less (no male partner present) [block group], binned as quintiles (1, 2, 3, 4, 5)</i>

DeGauss Social Deprivation Index (2018)	
SocialDeprivationIndex	<p><i>DeGauss composite metric derived from six 2018 ACS variables: fraction_assisted_income, fraction_high_school_edu, median_income, fraction_no_health_ins, fraction_poverty, fraction_vacant_housing</i></p> <p><i>Not yet integrated into ICEES</i></p>
US Census Topologically Integrated Geographic Encoding and Referencing System (TIGERline) Roadway Data (years 2016, 2017)	
MajorRoadwayHighwayExposure	<i>Distance in meters from household to nearest major road/highway (1 = 0-49, 2 = 50-99, 3 = 100-199, 4 = 200-299, 5 = 300-499, 6 = >=500 meters)</i>
MajorRoadwayHighwayExposure2	<i>Distance in meters from household to nearest major road/highway (1 = 0-49, 2 = 50-99, 3 = 100-149, 4 = 150-199, 5 = 200-249, 6 = >=250 meters)</i>
US Department of Transportation, Federal Highway Administration, Highway Patrol Monitoring System Roadway Data (year 2016, 2017)	
RoadwayDistanceExposure	<i>Distance in meters from household to nearest roadway (1 = 0-49, 2 = 50-99, 3 = 100-199, 4 = 200-299, 5 = 300-499, 6 = >=500 meters)</i>
RoadwayDistanceExposure2	<i>Distance in meters from household to nearest roadway (1 = 0-49, 2 = 50-99, 3 = 100-149, 4 = 150-199, 5 = 200-249, 6 = >=250 meters)</i>
RoadwayType	<i>UNC DOT roadway classification (e.g., major highway)</i>
RoadwayAADT	<i>US DOT Annual average daily traffic estimate</i>
RoadwaySpeedLimit	<i>US DOT Roadway speed limit</i>
RoadwayLanes	<i>UNC DOT Roadway number of lanes</i>
NC Department of Environmental Quality Data (Years vary, see notes)	
CAFO_Distance	<i>NC Department of Environmental Quality distance in meters from household to nearest concentrated</i>

<p><i>*Years vary, based on latest permit, most valid through 2024: 06/1998-03/2020</i></p>	<p><i>animal farming operation (1 = <500, 2 = 500-1000, 3 = 1000-2000, 4 = 2000-4000, 5 = >4000)</i></p> <p><i>Includes data on:</i></p> <ul style="list-style-type: none"> • Swines • Cattle • Poultry
<p>LandfillDistance</p> <p><i>*Pre-regulatory landfills (1983 and before) and active permitted landfills (1983 and after)</i></p> <p><i>**No dates, data pulled on 4/1/2020</i></p>	<p><i>NC Department of Environmental Quality distance in meters from household to nearest landfill (1 = <500, 2 = 500-1000, 3 = 1000-2000, 4 = 2000-4000, 5 = >4000)</i></p> <p><i>Includes data on:</i></p> <ul style="list-style-type: none"> • Unregulated landfills (pre-1983) • Regulated landfills • Superfund sites
<p>National Center for Education Statistics Public School Data (2018)</p>	
<p>PublicSchoolDistance</p>	<p><i>Distance from nearest public school</i></p> <p><i>Not yet integrated into ICEES</i></p>
<p>US Environmental Protection Agency Community Multiscale Air Quality (CMAQ) model exposure estimates: PM2.5, ozone</p>	
<p><i>CMAQ exposure estimates: PM2.5, ozone (years 2010 & 2011) [§]</i></p>	

Variable nomenclature format:	
AvgDailyXXExposure_StudyAvg	<i>UNC IE average of estimated average daily PM2.5 exposure over 'study' period, binned by data values (pandas cut) (1, 2, 3, 4, 5)</i>
MaxDailyXXExposure_StudyAvg	<i>UNC IE average of estimated maximum daily exposure over 'study' period, binned by data values (pandas cut) (1, 2, 3, 4, 5)</i>
AvgDailyXXExposure_StudyMax	<i>UNC IE maximum of estimated average daily PM2.5 exposure over 'study' period, binned by data values (pandas cut) (1, 2, 3, 4, 5)</i>
MaxDailyXXExposure_StudyMax	<i>UNC IE maximum of estimated maximum daily exposure over 'study' period, binned by data values (pandas cut) (1, 2, 3, 4, 5)</i>
AvgDailyXXExposure_StudyAvg_qcut	<i>UNC IE average of estimated average daily PM2.5 exposure over 'study' period, binned by quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
MaxDailyXXExposure_StudyAvg_qcut	<i>UNC IE average of estimated maximum daily exposure over 'study' period, binned by quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
AvgDailyXXExposure_StudyMax_qcut	<i>UNC IE maximum of estimated average daily PM2.5 exposure over 'study' period, binned by quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
MaxDailyXXExposure_StudyMax_qcut	<i>UNC IE maximum of estimated maximum daily exposure over 'study' period, binned by quintiles (pandas qcut) (1, 2, 3, 4, 5)</i>
<i>US Environmental Protection Agency continental USA (conUS) CMAQ model exposure estimates: PM2.5, ozone (years 2002-2019), CO, NO, NO2, NOx, SO2, acetaldehyde, formaldehyde, benzene (2002-2019, estimates are not consistently available for all years)**</i>	

Variable nomenclature format:	
AvgDailyXXExposure_2	<i>US EPA conUS CMAQ daily exposure estimates for PM2.5, CO, NO, NO2, NOx, SO2, acetaldehyde, formaldehyde, or benzene, averaged over 'study' period, binned by data values (pandas cut) (1,2,3,4,5)</i>
MaxDailyXXExposure_2	<i>US EPA conUS CMAQ 8-hour ozone maximum exposure estimate, averaged over 'study' period, binned by data values (pandas cut) (1, 2, 3, 4, 5)</i>
AvgDailyXXExposure_2_qcut	<i>US EPA conUS CMAQ daily exposure estimates for PM2.5, CO, NO, NO2, NOx, SO2, acetaldehyde, formaldehyde, or benzene, averaged over 'study' period, binned by quintiles (pandas.qcut) (1,2,3,4,5)</i>
MaxDailyXXExposure_2_qcut	<i>US EPA conUS CMAQ 8-hour PM2.5 exposure estimate, averaged over 'study' period, binned by quintiles (pandas.qcut) (1,2,3,4,5)</i>
<i>US EPA National Air Quality System Pollutant Data (2021-2022)</i>	<p><i>Data on exposure estimates for pollen, nitrogen dioxide, sulfur dioxide, ozone, and particulates</i></p> <p><i>Not yet integrated into ICEES</i></p>
<i>US EPA National Emissions Inventory (NEI) Data (2017)</i>	
<i>Variables not yet named</i>	<p><i>Estimates of air emissions of criteria pollutants, criteria precursors, and hazardous air pollutants from air emissions sources (nonpoint wagon wheel (HAP, chromium, PM); CAP and lead emissions for wildland (wild and prescribed) fires; commercial marine vehicle emission estimates; airport-related and aircraft emission estimates; locomotive and railyard emission estimates)</i></p> <p><i>Not yet integrated into ICEES</i></p>
<i>US EPA Unregulated Contaminant Monitoring Rule (UCMR) Public Water Supply Data, UCMR3 (2015-2015) & ICMR4 (2018-2020)</i>	

<i>Variables not yet named</i>	<i>Data for contaminants suspected to be present in public water supplies, but that do not have regulatory standards set under the Safe Drinking Water Act (SDWA), e.g., PFAS contaminants</i> <i>Not yet integrated into ICEES</i>
--------------------------------	--

Abbreviations: PM_{2.5} = particulate matter $\leq 2.5 \mu\text{m}$ in diameter

*The feature variables listed in the table are those for the patient-level tables, which include data on each patient for each year of available data (i.e., data on individual patients are represented as rows in the table). Similar feature variables are available for the visit-level tables, although the variables are sometimes treated differently. For example, PM_{2.5} exposures in visit-level tables are expressed in relation to the 24-hour and two-week period before visits, not in relation to the one-year 'study' period, as was done for the patient-level tables. Additional feature variables (e.g., laboratory measures) are available for select years. Further information can be accessed at <https://github.com/NCATSTranslator/Translator-All/wiki/Exposures-Provider-ICEES>.

§From first batch of CMAQ output, derived from UNC Institute for the Environment, hourly estimates, 36-km (2010) or 12-km (2011) resolutions

**From second batch of CMAQ output, derived from US EPA continental USA (conUS); US Census tract resolution, 2002-2015 for PM_{2.5} and ozone, 2002 only for other chemicals