# ESTIMATING EMPLOYMENT WITH NIGHTTIME LIGHTS AND TRANSPORTATION DATA USING MACHINE LEARNING

MARK FOLDEN
SENIOR PREDICTIVE ANALYTICS SPECIALIST
NCTCOG

MAY 13, 2025

# I LOVE IT WHEN A PLAN COMES TOGETHER

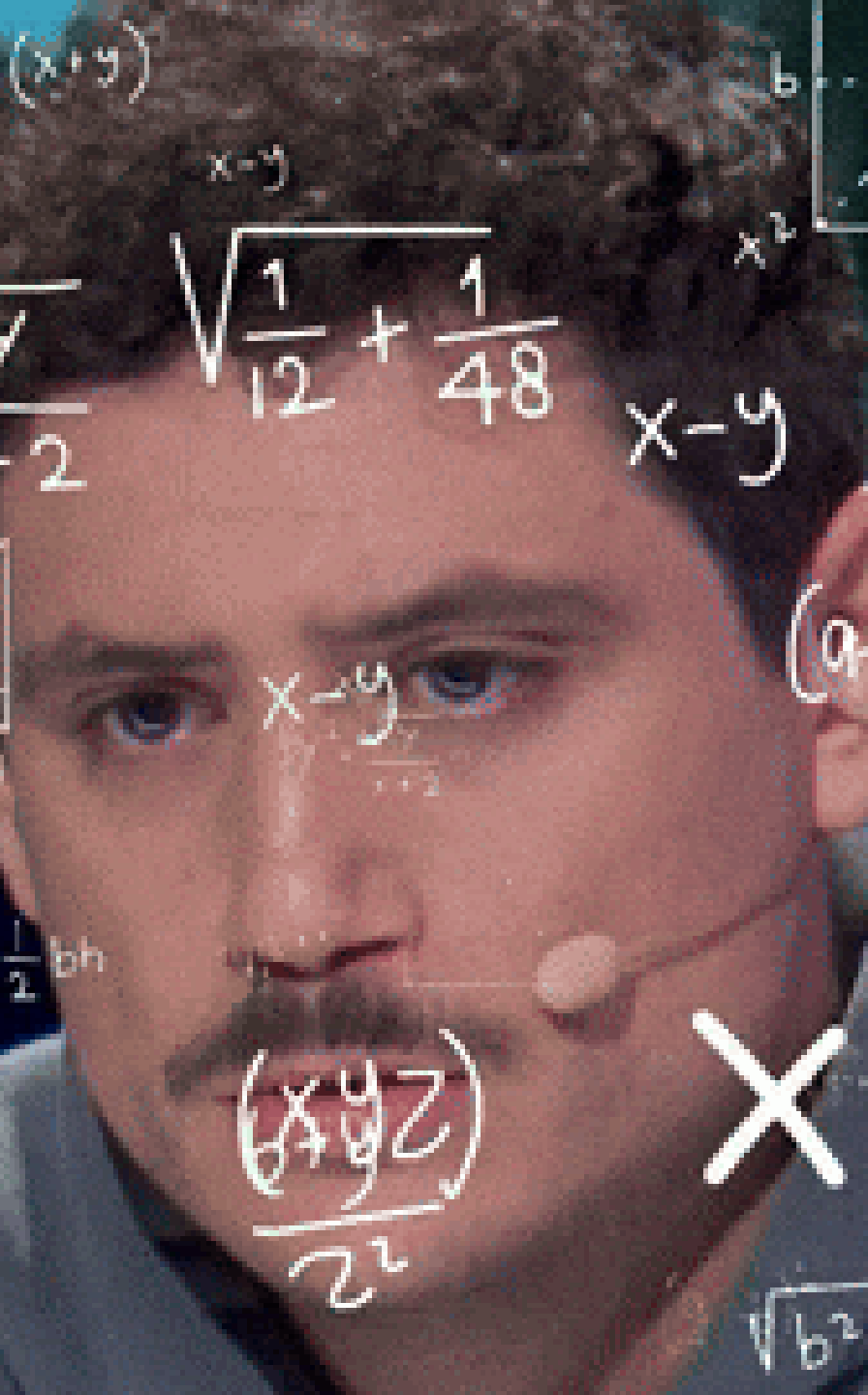PROBLEM TO SOLVE

DATA & GIS PROCESSES

DATA SCIENCE PROCESS

TAKEAWAYS

FUTURE RESEARCH

# PROBLEM TO SOLVE

- Small Area Estimates

- Sub-county estimates of Households, Population, Employment

- LODES is only published source of employment estimates; use as a starting point, but need to overcome several limitations

- Need a way to allocate to 30x30 meter grid; outputs should sum to the inputs, but be independent of zone structure

- Serve as starting point for Forecast

- Apply to past data and improve temporal consistency for Forecast model validation

- Apply in the future as new data becomes available, increased efficiency in generating updated data

# DATA & GIS PROCESSES

**Small Area Estimates**
- NCTCOG
- Already on 30x30m, but needs cleaning up

**Landuse**
- NCTCOG

**Nighttime Lights**
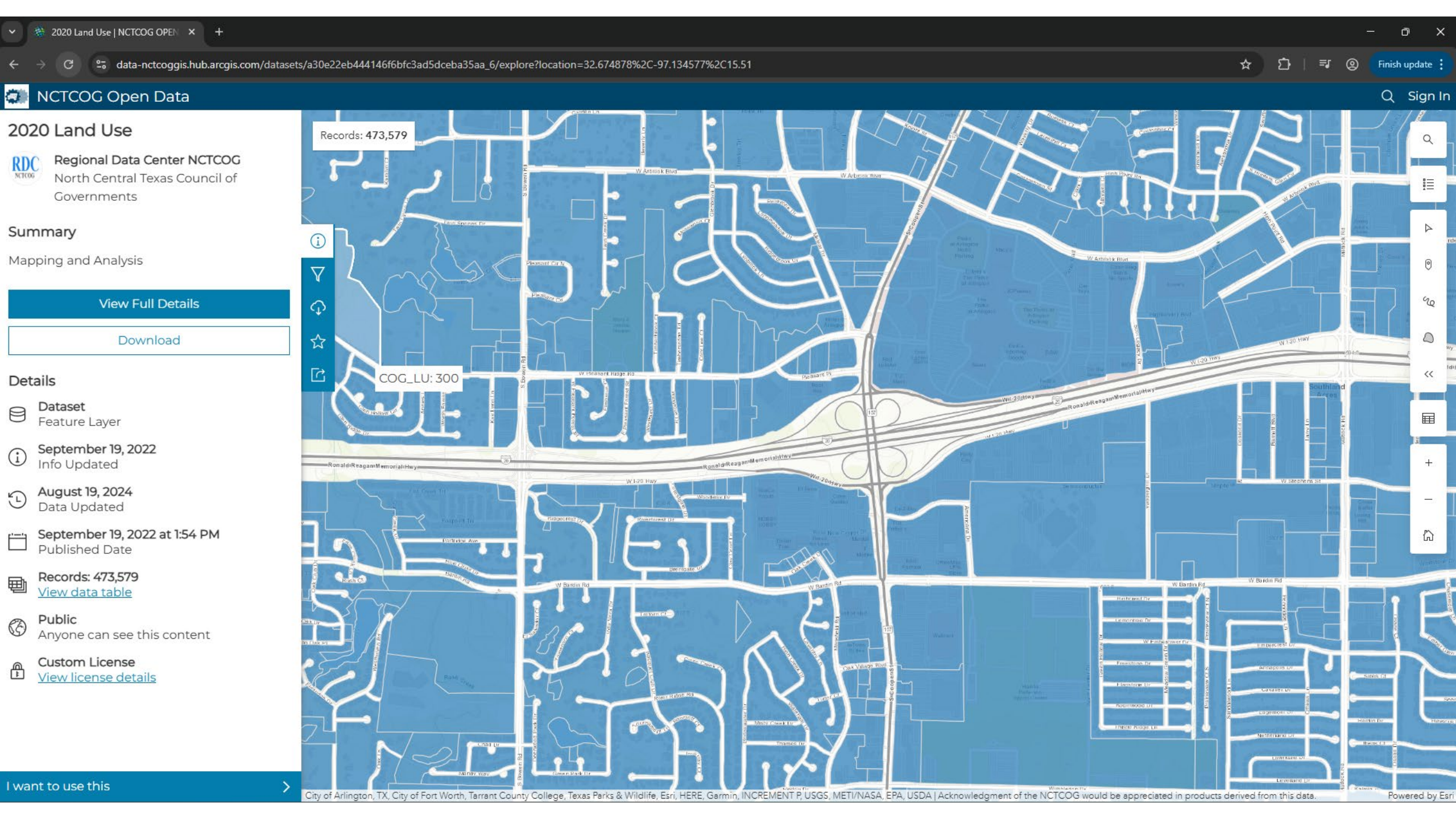- Earth Observation Center, Colorado School of Mines
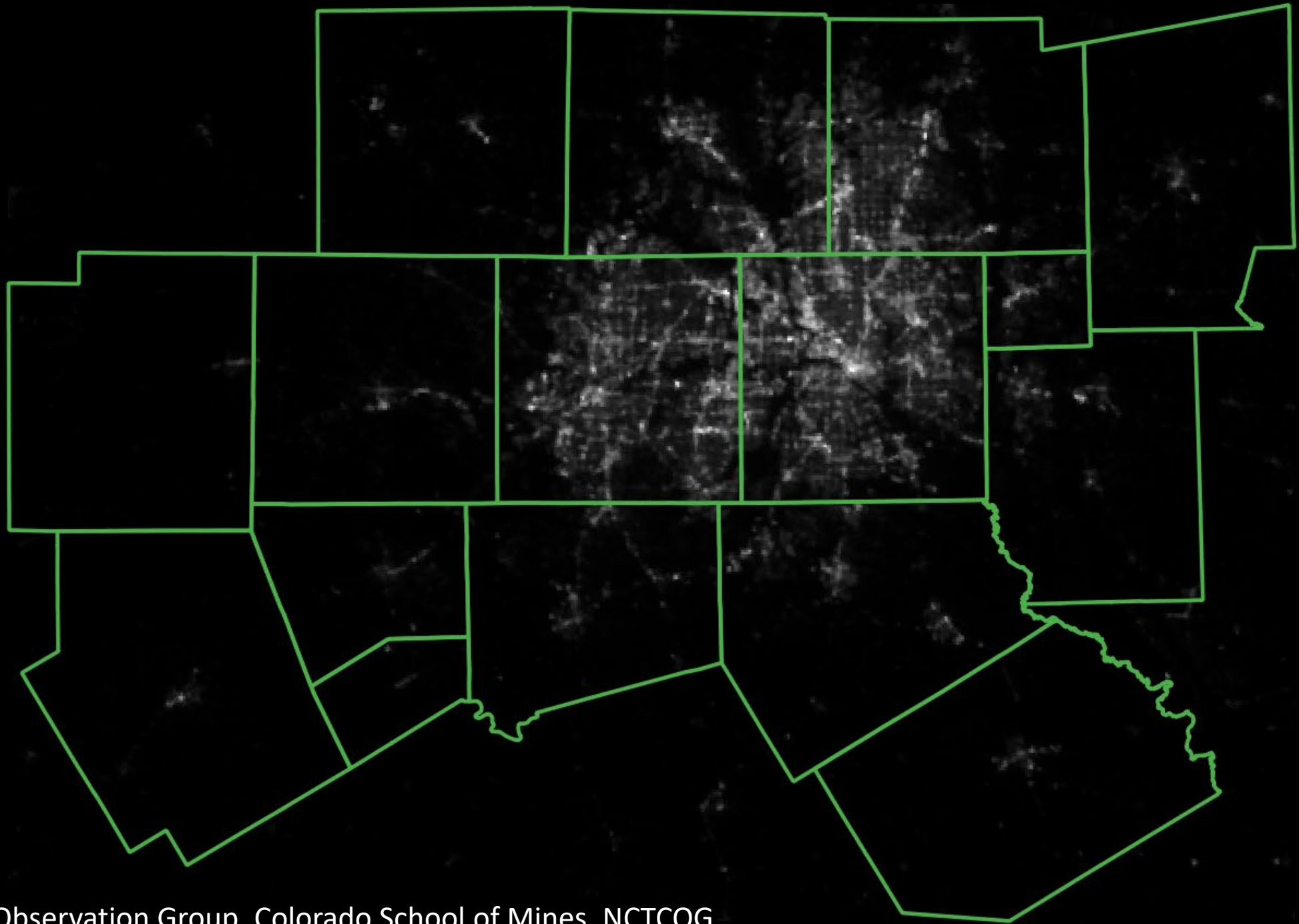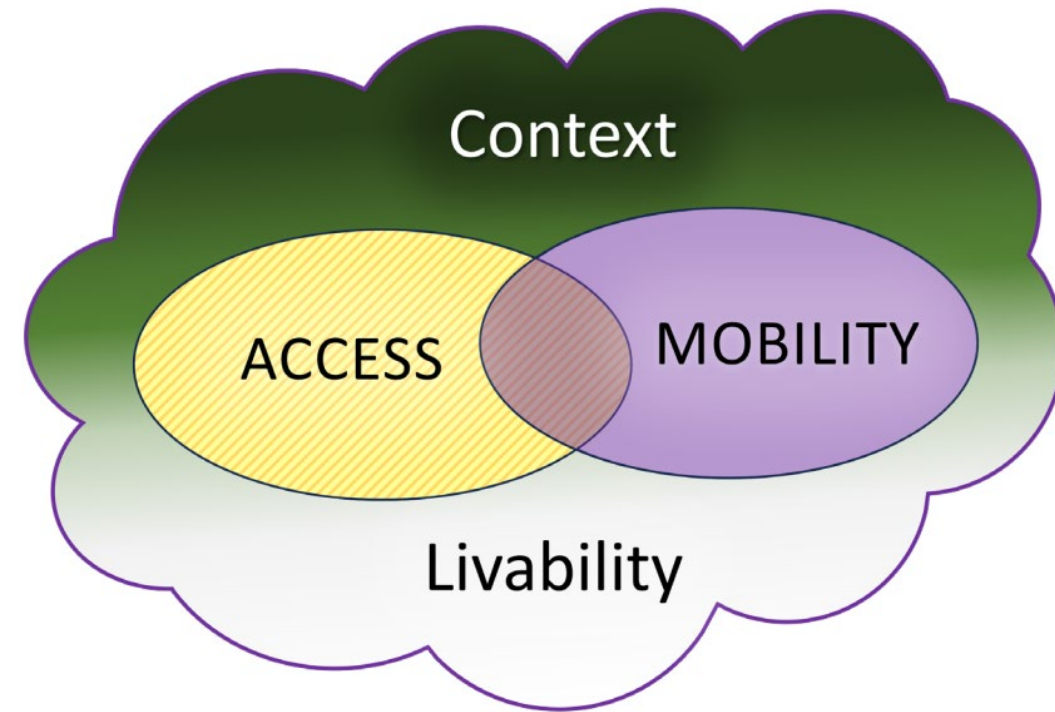
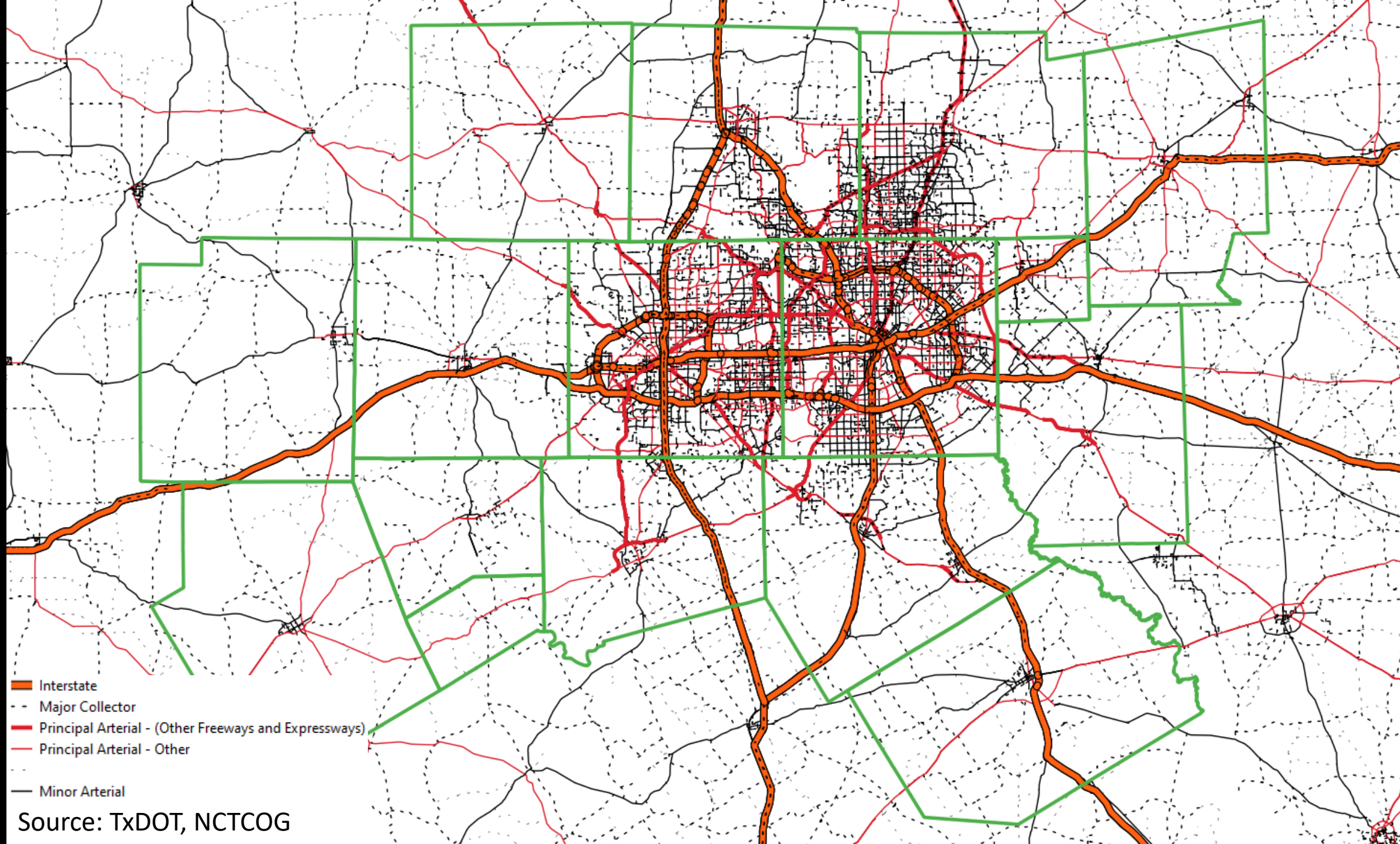**Roadway Network**
- TxDOT

**Traffic Counts**
- TxDOT

# ROADWAY NETWORK: FUNCTIONAL CLASSIFICATION

- Transportation facilities can be classified on a continuum between interrelated goals

  - Access vs. Mobility

- For our purposes:

  - Mobility

    - Facility cannot have a driveway connected to it

    - Interstates

    - "Major Arterials – Other Freeways & Expressways"

  - Access

    - Facility can have a driveway connected directly to it
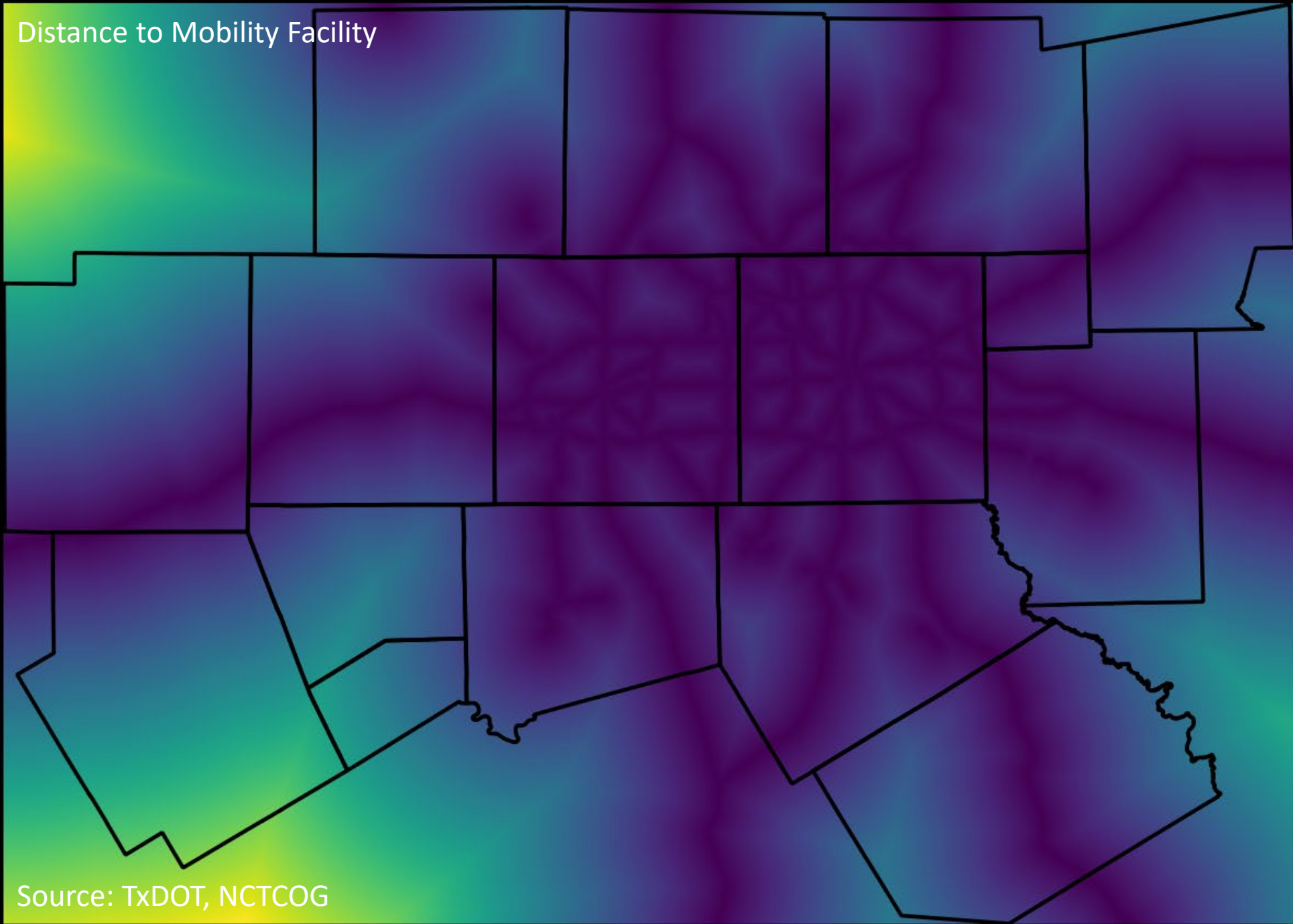


Source: FHWA

Interstate

Major Collector

Principal Arterial - (Other Freeways and Expressways)

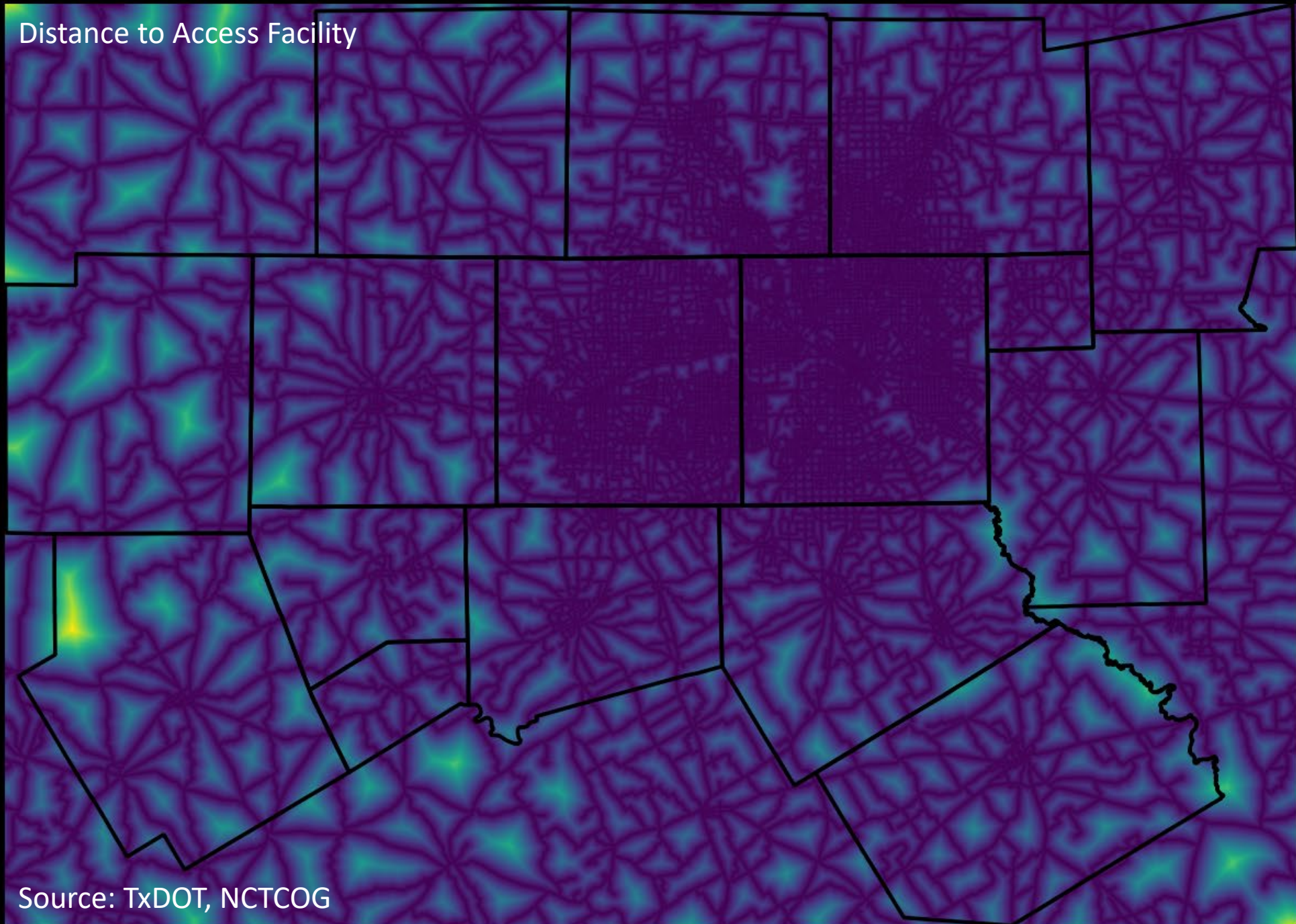Principal Arterial - Other

Minor Arterial

Source: TxDOT, NCTCOG

Distance to Mobility Facility

Source: TxDOT, NCTCOG
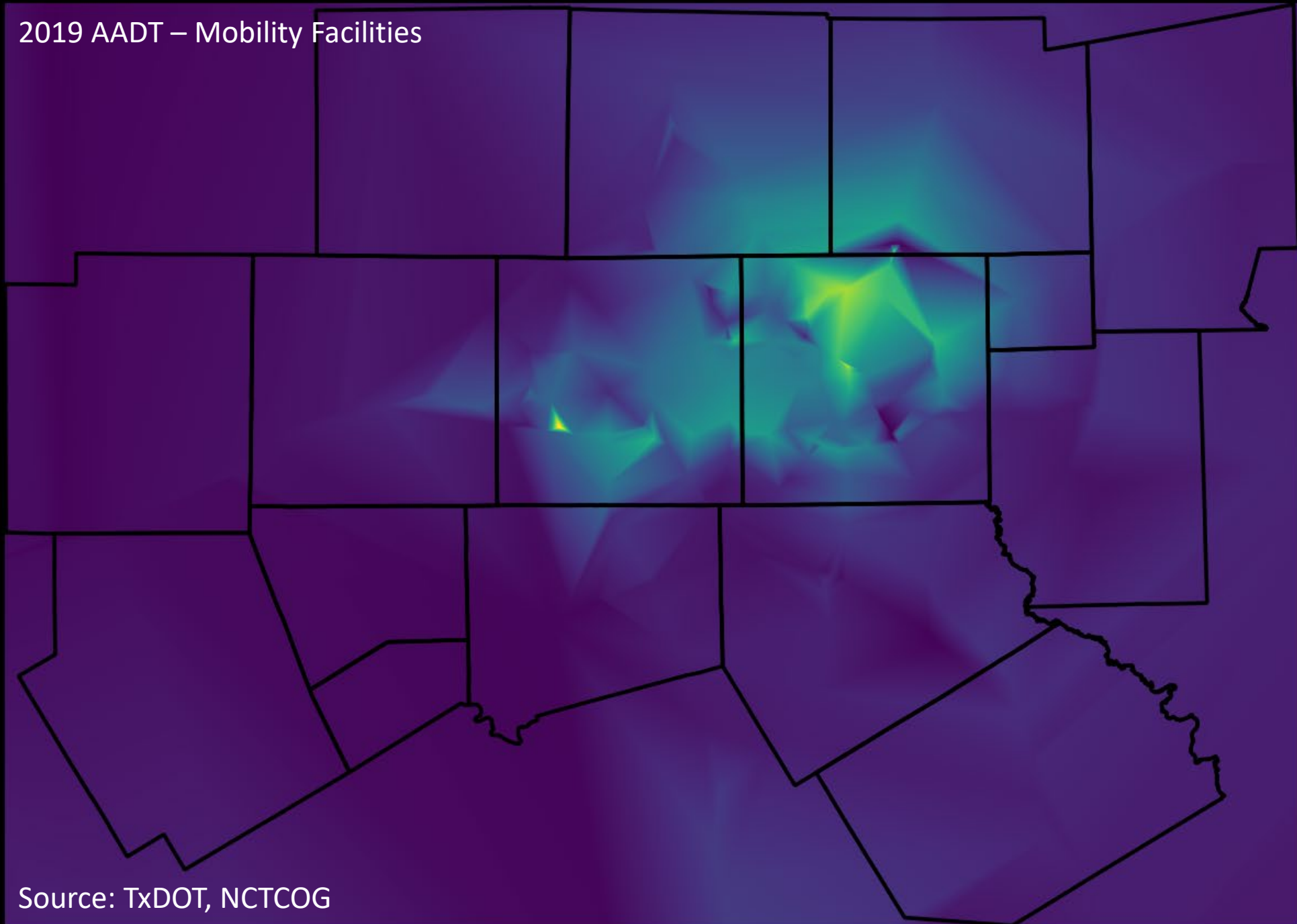
Distance to Access Facility

Source: TxDOT, NCTCOG

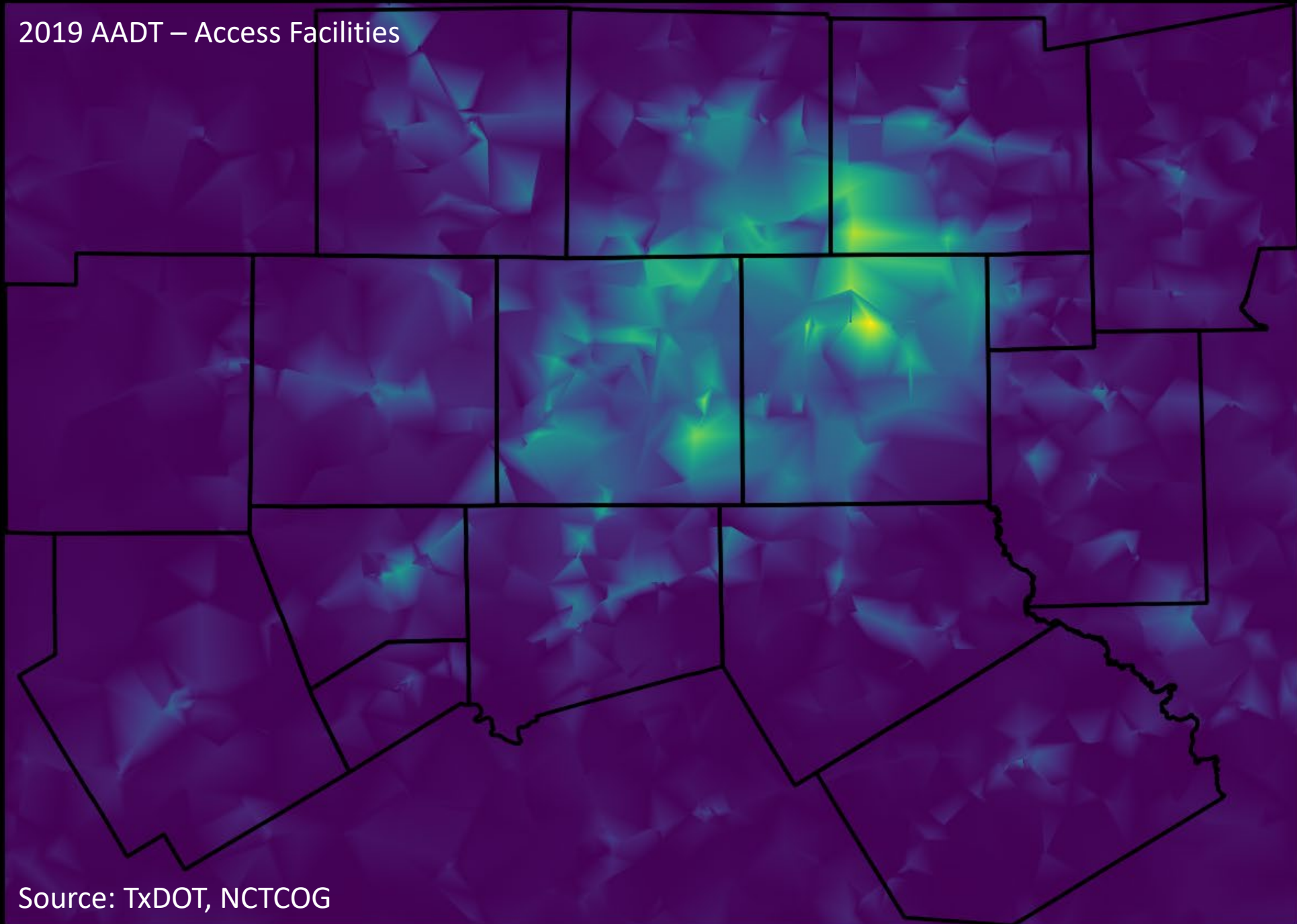Traffic Count Locations

Source: TxDOT, NCTCOG

2019 AADT – Mobility Facilities

Source: TxDOT, NCTCOG

2019 AADT – Access Facilities

Source: TxDOT, NCTCOG

# DATA SCIENCE PROCESS

CLEAN DATA

CHECK FOR NORMALITY

RESCALE (IF NECESSARY)

CHECK FOR CORRELATION

TRAIN ALGORITHM

VALIDATE & TEST

# CORRELATION



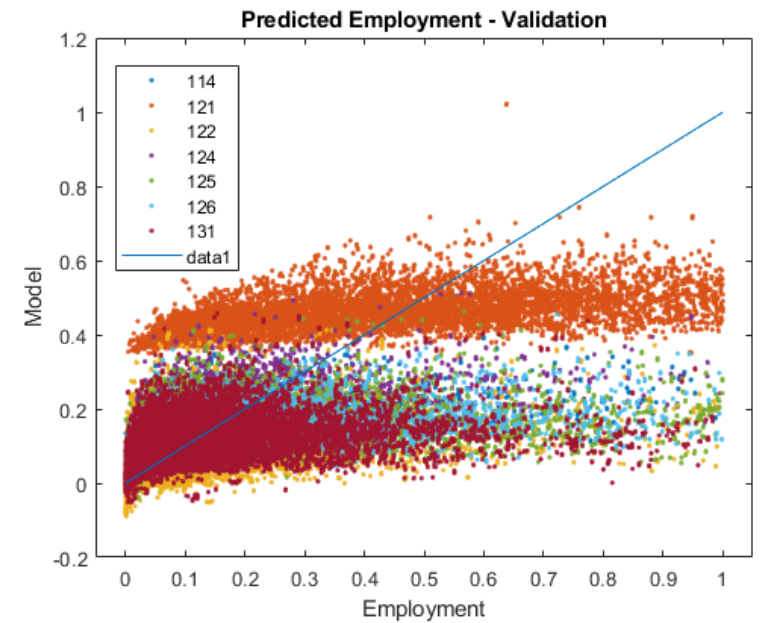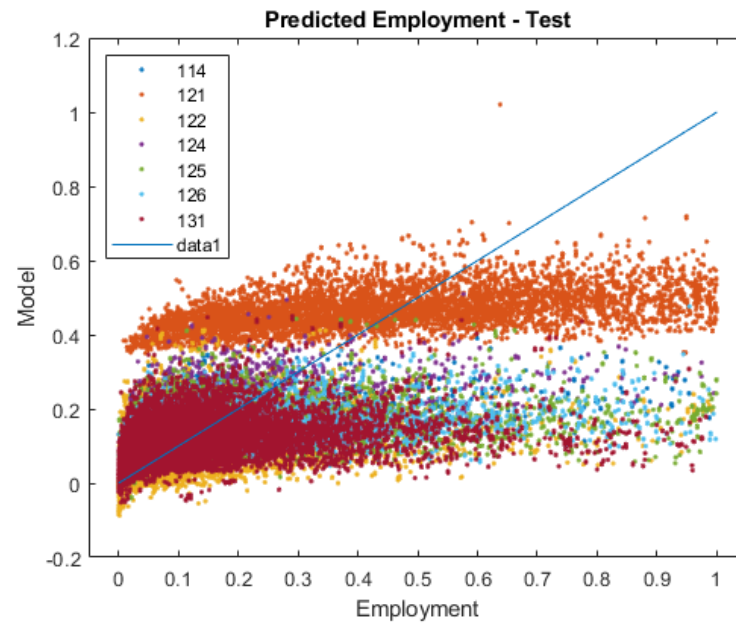|  | ntl | imp_sfc | dstaccess | cntaccess | dstmobility | cntmobility |
|---|---|---|---|---|---|---|
| ntl | 1 | 0.33645 | -0.17821 | 0.26478 | -0.33711 | 0.39738 |
| imp_sfc | 0.33645 | 1 | -0.09948 | 0.11172 | -0.16217 | 0.19147 |
| dstaccess | -0.17821 | -0.09948 | 1 | -0.12049 | 0.094757 | -0.15932 |
| cntaccess | 0.26478 | 0.11172 | -0.12049 | 1 | -0.066314 | 0.25983 |
| dstmobility | -0.33711 | -0.16217 | 0.094757 | -0.066314 | 1 | -0.3409 |
| cntmobility | 0.39738 | 0.19147 | -0.15932 | 0.25983 | -0.3409 | 1 |

# ALGORITHMS

MATLAB 2022a with Statistics and Machine Learning Toolbox

- Multivariate Linear Regression (Plain Vanilla)

- Support Vector Machine

- LS Boosted Ensemble

- Treebagger

- Bagged Ensemble
  - https://www.ibm.com/think/topics/bagging
  - https://www.mathworks.com/help/stats/ensemble-algorithms.html

n = 1,143,918

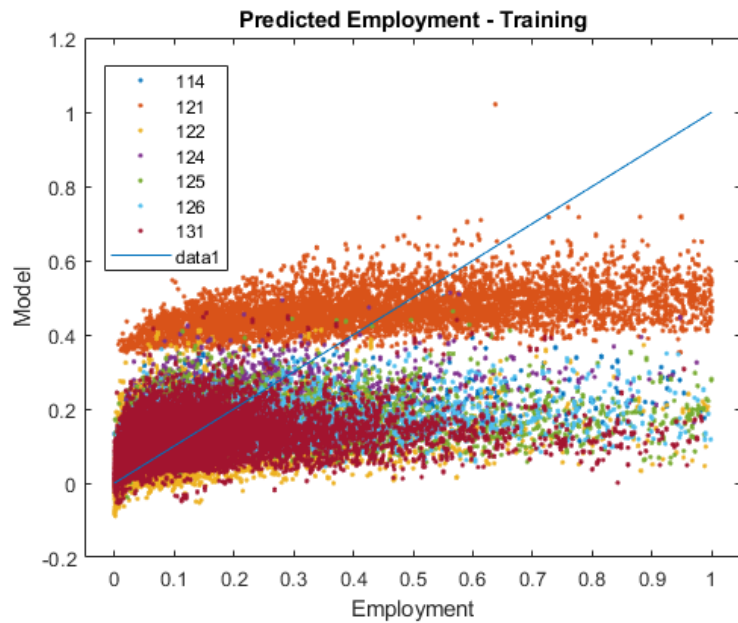- Training
  - 30% : 342,967

- Test
  - 20%: 228,935

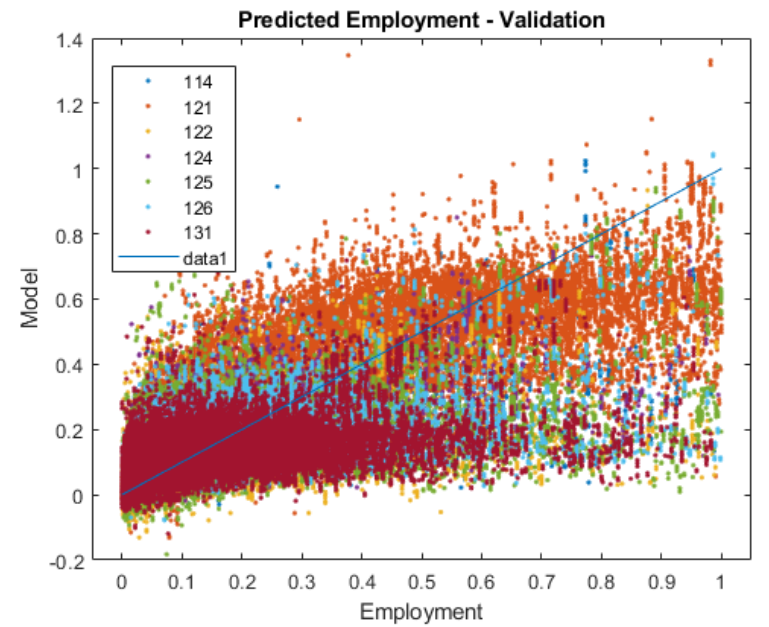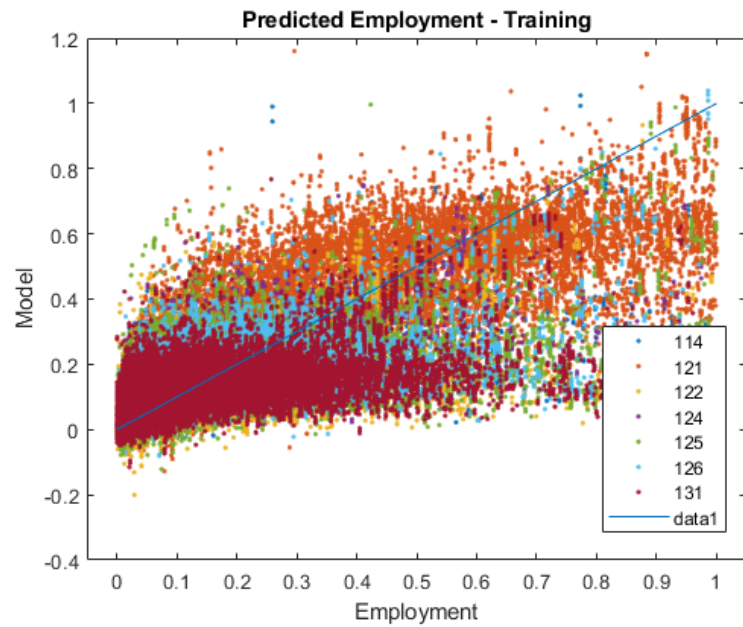- Validation
  - 50%: 572,016

# MULTIVARIATE LINEAR REGRESSION

# SUPPORT VECTOR MACHINE

# LS BOOSTED ENSEMBLE
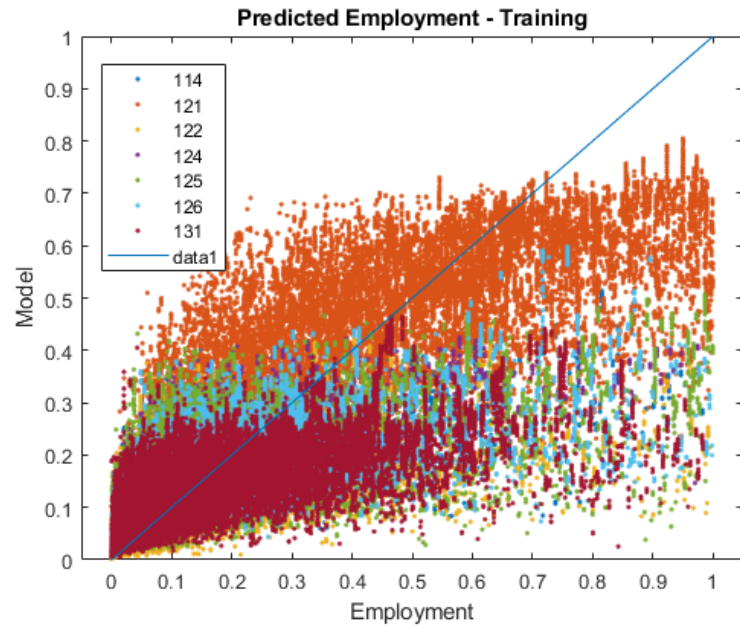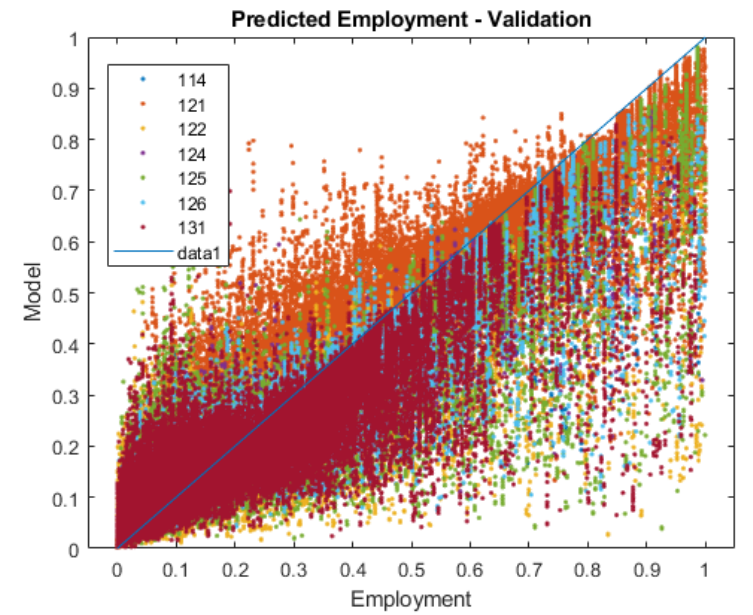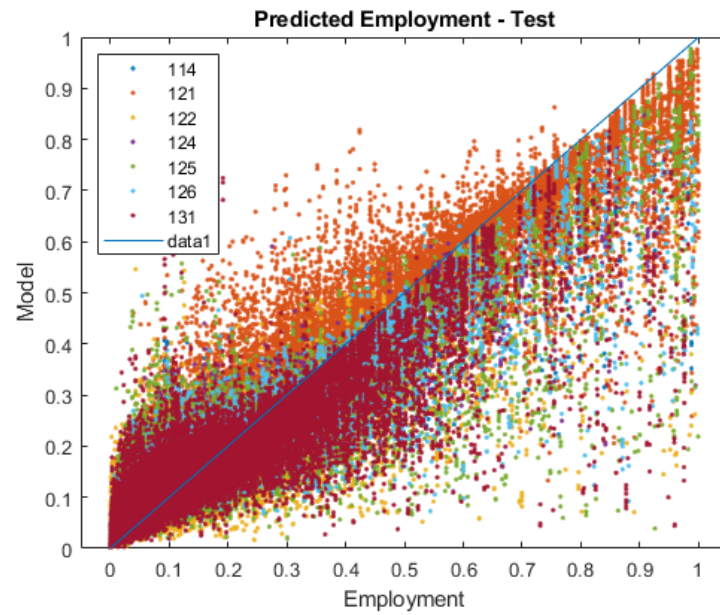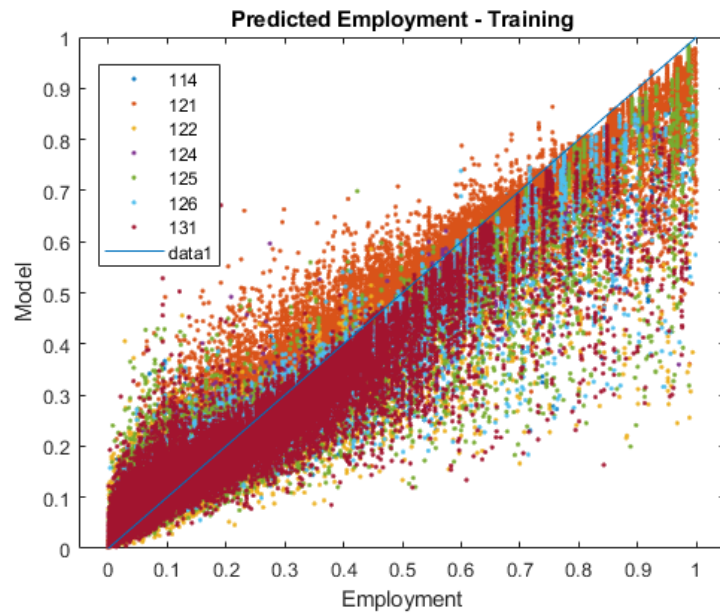
# TREEBAGGER

# BAGGED ENSEMBLE

# SUMMARY
## R-SQUARED

- $R^2$ closer to 1.0 the better
- Don't want a large drop from Training to Test and Validation

| Algorithm | Training | Test | Validation |
|---|---|---|---|
| Multivariate Linear Regression | 0.65027 | 0.64930 | 0.65293 |
| Support Vector Machine | 0.64824 | 0.64711 | 0.65109 |
| LS Boosted Ensemble | 0.75222 | 0.74138 | 0.74352 |
| Treebagger | 0.80157 | 0.78811 | 0.78892 |
| Bagged Ensemble | 0.97062 | 0.93480 | 0.93508 |

# TAKEAWAYS

Machine Learning techniques compared to traditional Multinomial Least Squares regression is a tradeoff between model performance and interpretability
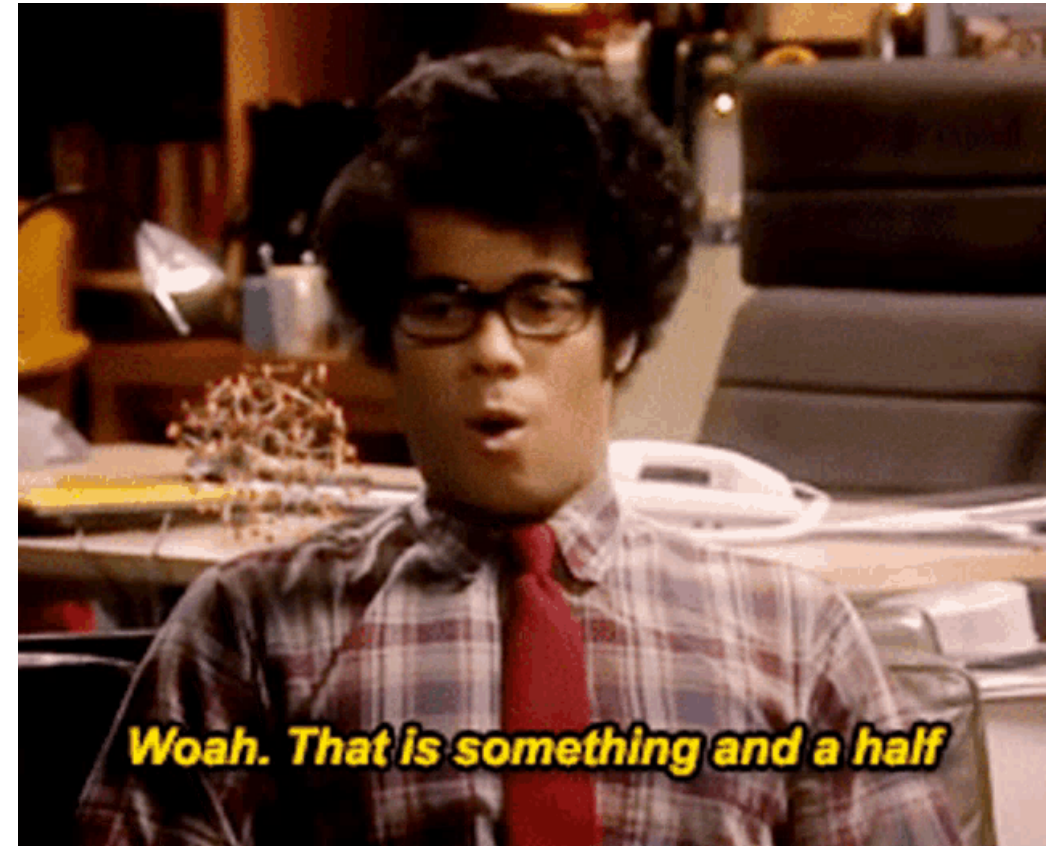
You can offset some of that tradeoff with additional analysis and trying multiple methods

You might not know <u>exactly</u> how an algorithm works, but go ahead and try it and see if you get an improvement

Even two algorithms that work similarly can get you different results
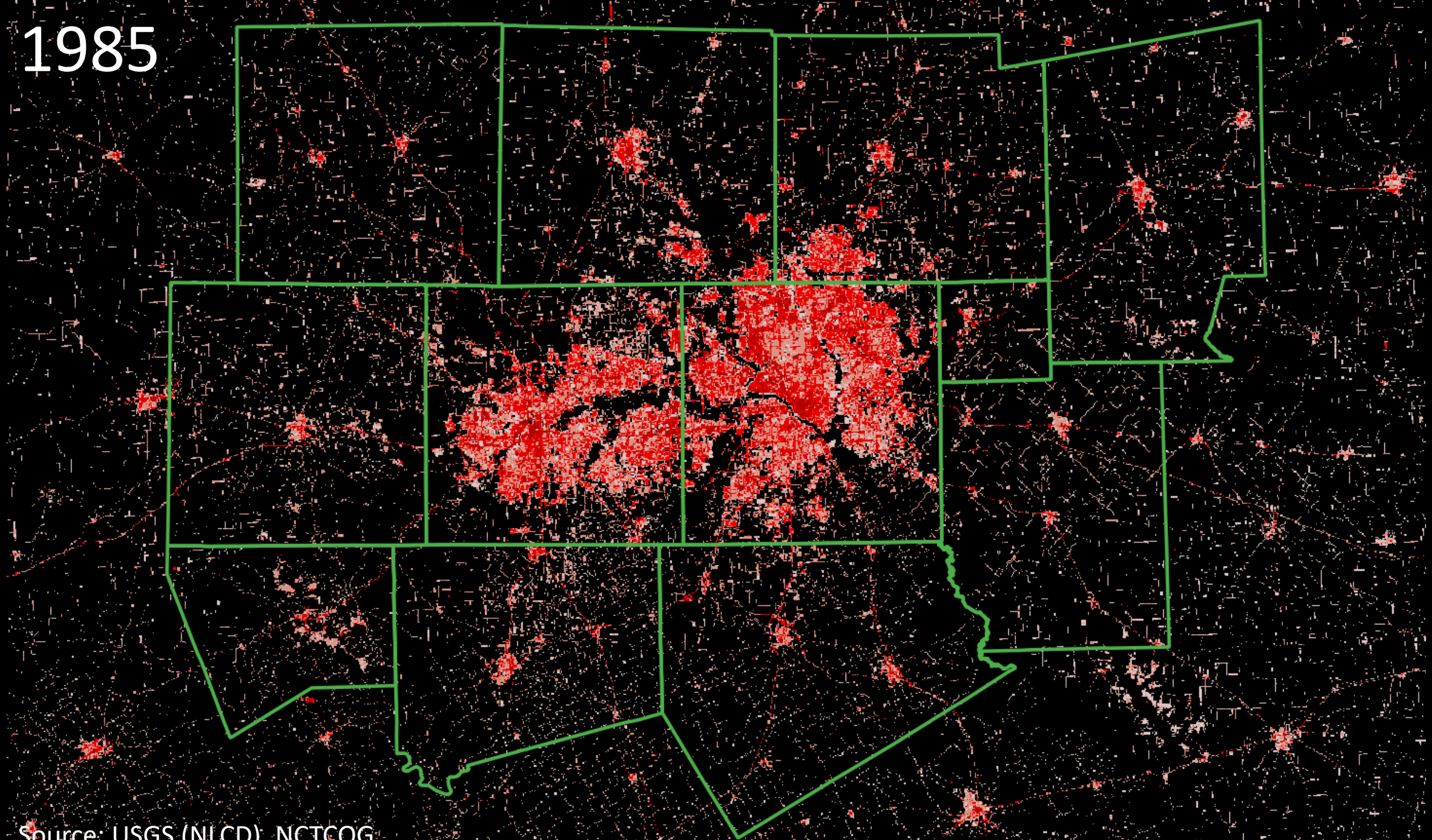
Categorical variables can be VERY powerful

More data != better
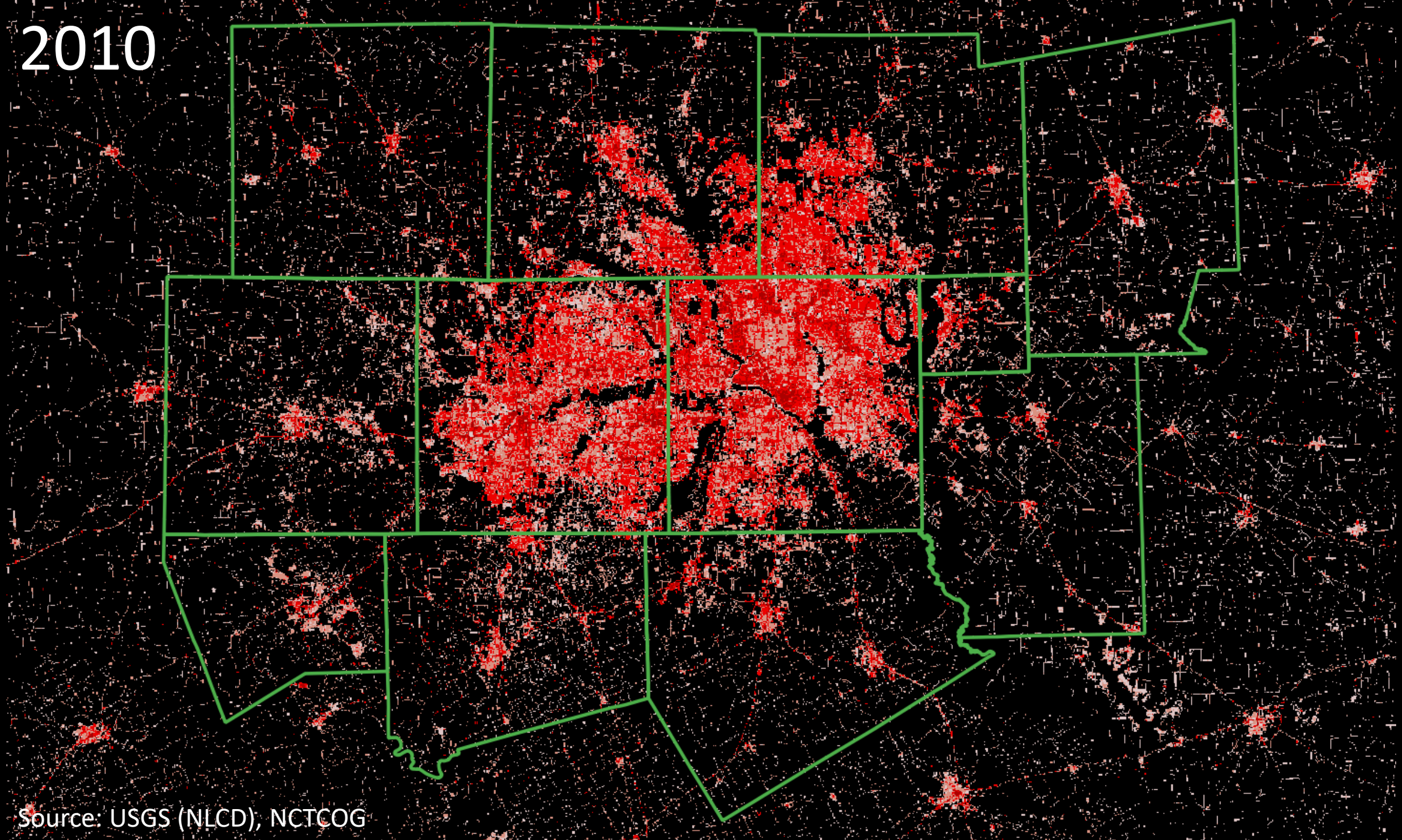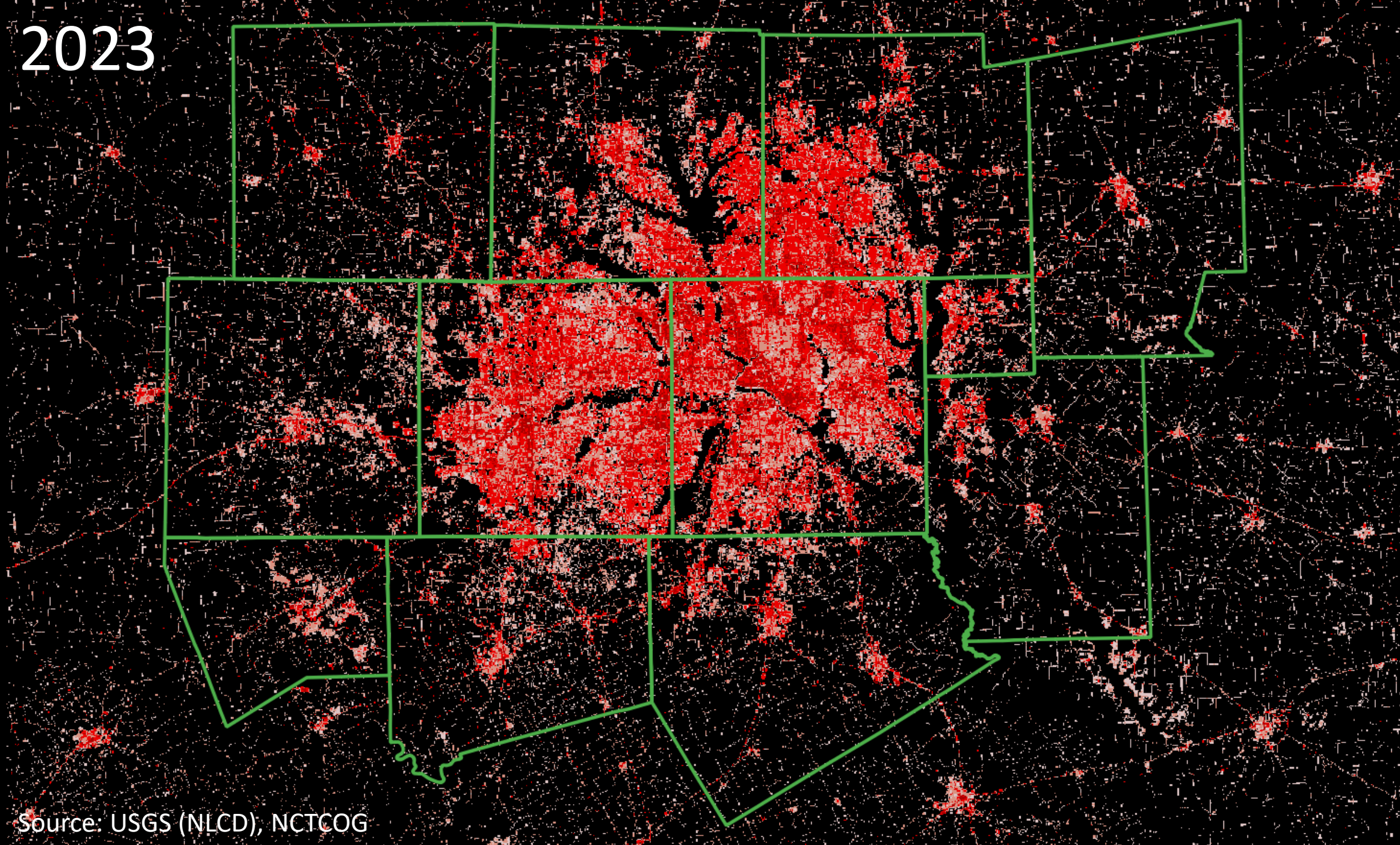


Woah. That is something and a half
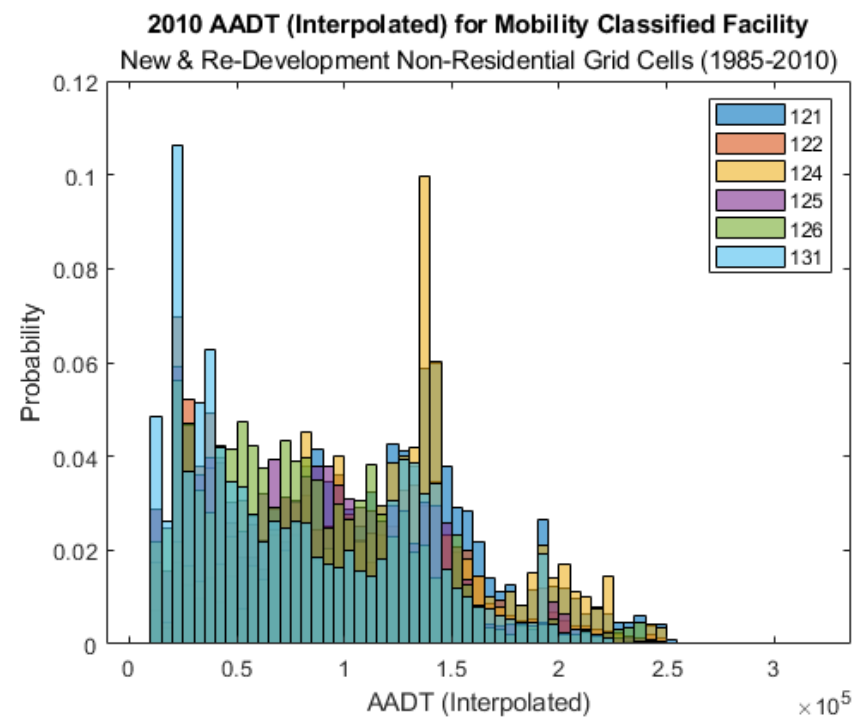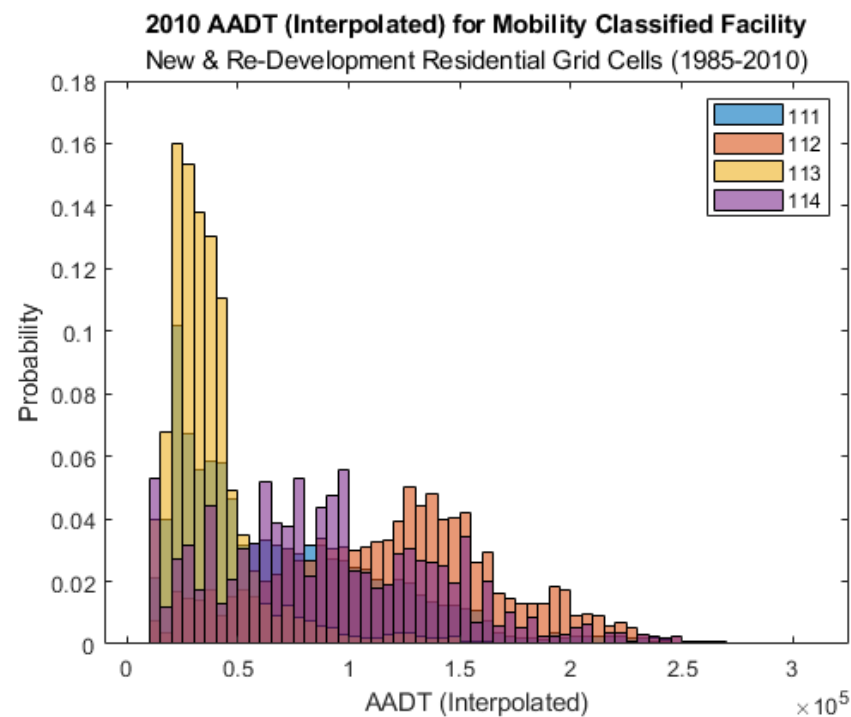
1985

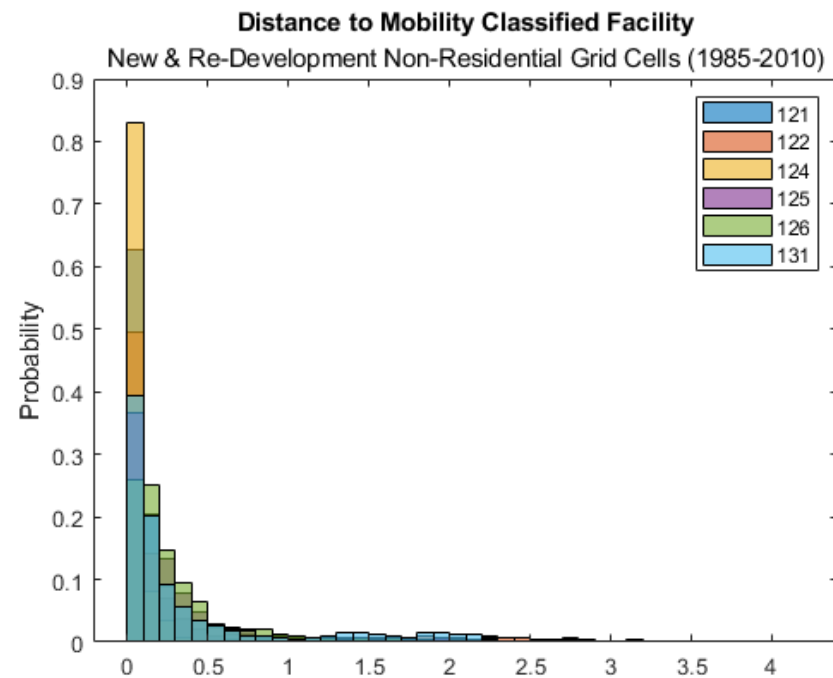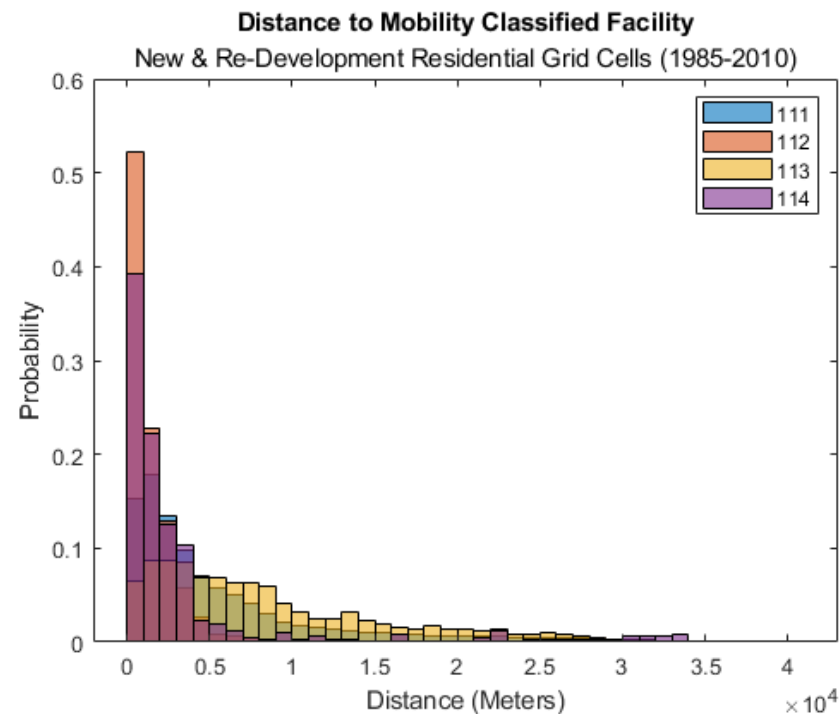Source: USGS (NLCD), NCTCOG

2010

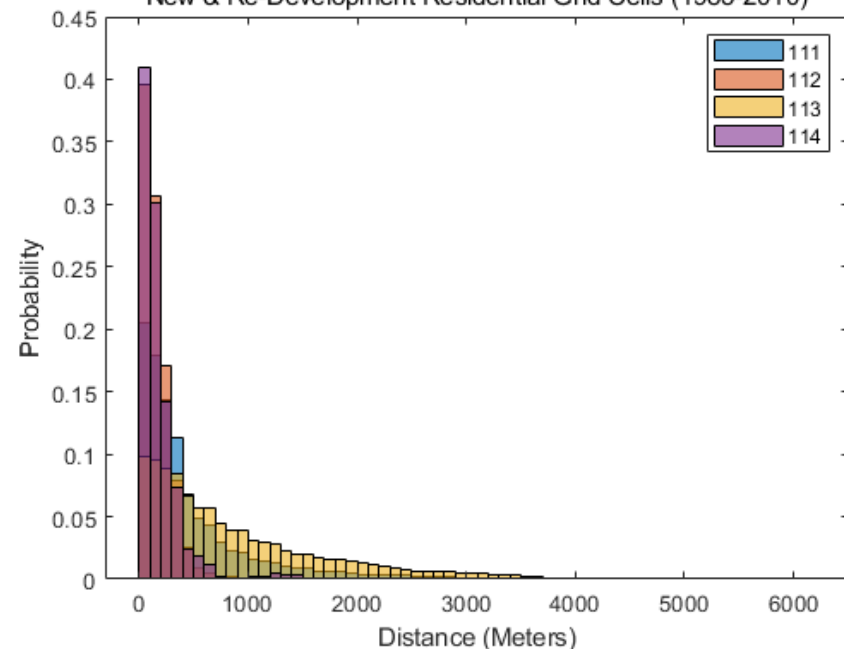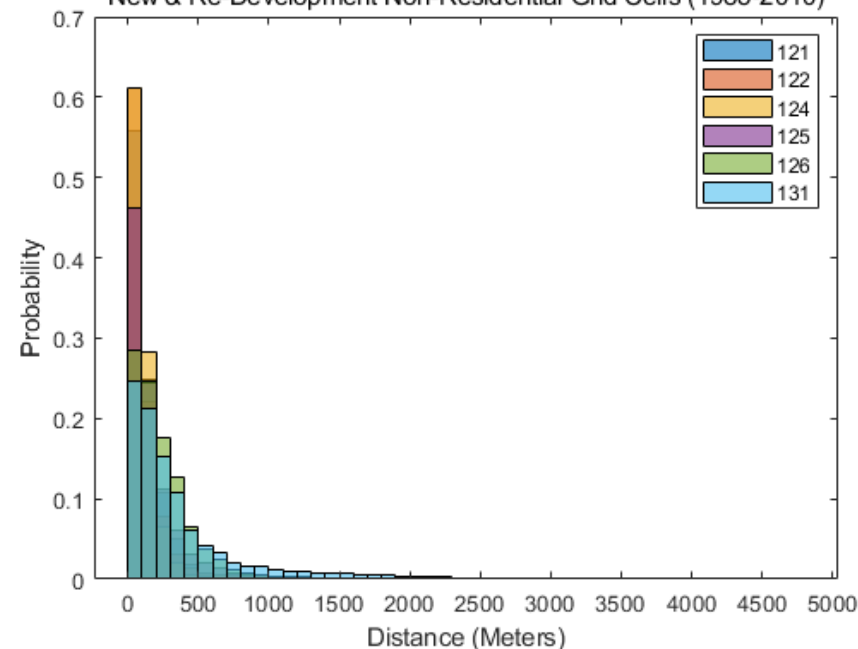Source: USGS (NLCD), NCTCOG

2023

Source: USGS (NLCD), NCTCOG

**Distance to Access Classified Facility**
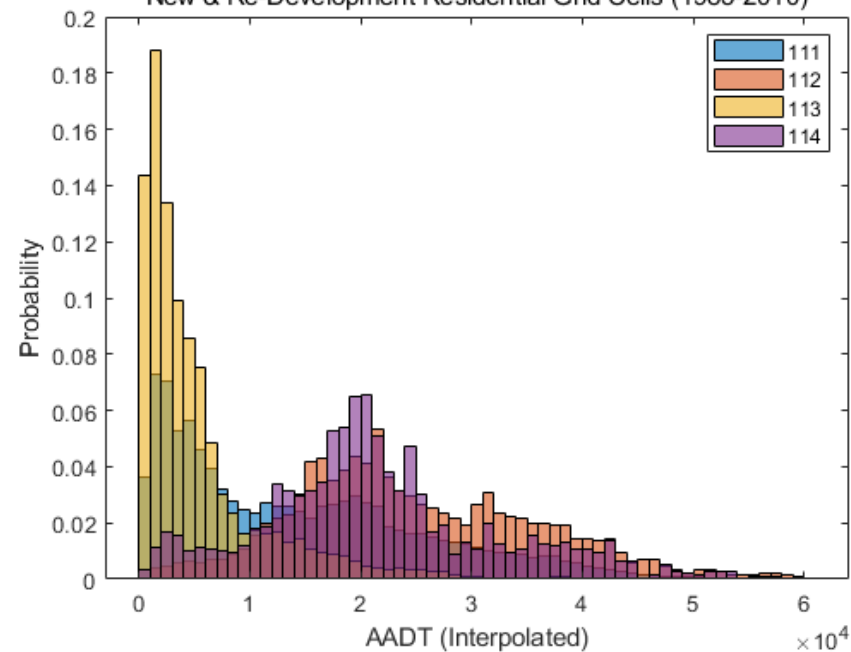New & Re-Development Residential Grid Cells (1985-2010)

**Distance to Access Classified Facility**
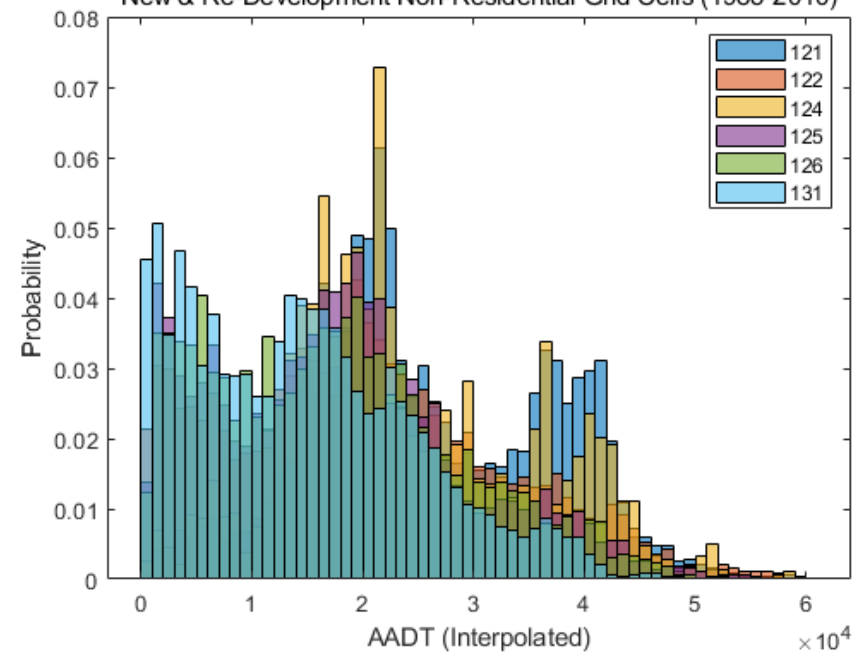New & Re-Development Non-Residential Grid Cells (1985-2010)

**2010 AADT (Interpolated) for Access Classified Facility**
New & Re-Development Residential Grid Cells (1985-2010)

**2010 AADT (Interpolated) for Access Classified Facility**
New & Re-Development Non-Residential Grid Cells (1985-2010)

QUESTIONS?