# Predicting Building Energy Usage

w207 section 3 - C. Ilin Donahoe, Meier, Rosenberg https://github.com/bennnyys/w207-sec3-final-proj-DMR

# Motivation

# What question are we working on? Why is it interesting?

- Can we predict energy usage from building characteristics?
  - o Energy use intensity site
- Current modeling requires many inputs, lots of expensive time, specific people, and computational power
  - EnergyPlus, DOE2, eQUEST, etc.
- Reduced order modeling can more easily inform building system planning for new construction as well as retrofit decision making
- If we change one thing in a building, will it save energy?
- Target within 20% error

#### Estimated "Best Features"

- year of construction
- building type
- window-wall ratio
- wall r-value
- floor r-value
- ceiling r-value
- window surface area

- climate zone
- heating type
- cooling type
- number of stories
- number of occupants
- estimated air leakage



### What data are we using?

- U.S. Department of Energy Building Performance Database
  - Energy consumption and building characteristics of actual buildings from 16 states/municipalities in the U.S. (New York, San Francisco, Chicago, D.C., Seattle, Philadelphia, Austin, etc.)
- Up to 10 years of panel data per building
- Number of rows: 296,065
- Number of features: 33
- Main features
  - Site EUI (energy use intensity: kBtu per square feet, 91k NaN)
    - Also: year, electric EUI (95k NaN), fuel EUI (135k NaN), GHG emissions (152k NaN)
  - Climate region (combination of temperature and humidity, 823 NaN)
  - Residential or Commercial (0 NaN)
  - Facility type (office, education, warehouse, industrial, lodging, etc., 156k NaN)
  - Floor Area (square feet, 0 NaN)
  - Year built (4,779 NaN)
  - Energy star rating (206k NaN)

## State/Municipality Dataset Analysis

- Austin 8,390 (~1,600 rows with roof/ceiling and window glass data)
- Berkeley 414
- Boston 6,835 (~4,500 rows with energy star ratings)
- CA Building Energy Benchmarking Program 3,881
- CA Prop 39 K-12 Program 1,501
- Cambridge 4,122
- Chicago 11,937
- Fannie Mae 857
- Gainesville 154,528 (includes cooling, roof\_ceiling type, but only 75k with site\_eui)
- NYC Ordinance 54,298
- NY Residential 439
- Philadelphia 6,457
- San Francisco 9,647
- Seattle <u>17,002</u>
- Syracuse 139
- D.C. 15,618

#### Fairness

- Location skew data collected at state/municipal level mostly large coastal cities
- Income skew buildings which receive energy audits may be skewed to affluent communities
  - Health and safety usually prioritized over energy efficiency in less affluent communities
- Consistency Commercial vs Residential
  - Commercial buildings have significant variation in construction compared to residential buildings

## Descriptive Statistics: numeric

	year	floor_area	year_built	energy_star_rating	electric_eui	fuel_eui	site_eui	source_eui	ghg_emissions_int
count	296065.000000	2.960650e+05	291286.000000	89543.000000	200966.000000	160800.000000	204210.000000	196936.000000	143854.000000
mean	2013.978505	7.469098e+04	1971.887276	61.809511	29.602798	23.291331	63.231037	126.836547	4.937401
std	2.609692	1.796745e+05	28.742165	28.910037	28.690731	32.060086	54.413696	110.338840	4.202703
min	2010.000000	5.000000e+02	1649.000000	0.000000	0.000000	0.000000	1.001169	1.075772	0.000000
25%	2012.000000	1.635000e+03	1961.000000	40.000000	16.253904	8.182419	32.009412	72.190118	2.881654
50%	2014.000000	3.260000e+03	1978.000000	68.000000	23.384545	15.138294	50.051977	102.748173	4.026421
75%	2016.000000	8.036400e+04	1993.000000	86.000000	33.968446	26.041885	79.382514	146.732770	5.665176
max	2020.000000	6.385382e+06	2020.000000	100.000000	987.466930	936.379589	997.866120	3133.315574	109.708218

### Descriptive Statistics: categorical

```
CLIMATE - number of distinct vals: 12 -----
2A Hot - Humid (Houston-TX)
                                      162918
4A Mixed - Humid (Baltimore-MD)
                                       76484
5A Cool - Humid (Chicago-IL)
                                       23269
4C Mixed - Marine (Salem-OR)
                                       17039
3C Warm - Marine (San Francisco-CA)
                                       11171
3B Warm - Dry (El Paso-TX)
                                        3957
NaN
                                         823
6A Cold - Humid (Burlington-VT)
                                         183
4B Mixed - Dry (Albuquerque-NM)
                                          95
5B Cool - Dry (Boise-ID)
                                          67
Name: climate, dtype: int64
BUILDING CLASS - number of distinct vals: 2 -----
Residential
             223601
Commercial 72464
Name: building class, dtype: int64
```

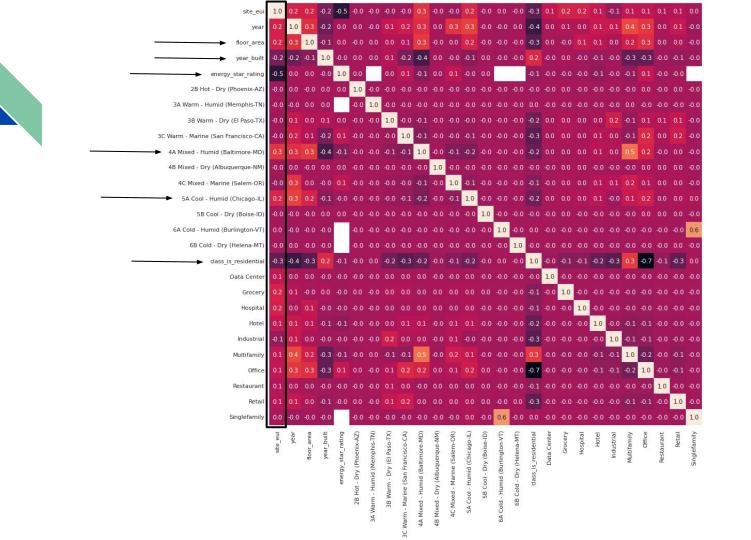
```
FACILITY TYPE - number of distinct vals: 82 -----
NaN
                                     156941
Multifamily - Uncategorized
                                      56268
Office - Uncategorized
                                      22030
Education - Other classroom
                                       8350
2-4 Unit Building
                                       5018
5+ Unit Building
                                       4508
Lodging - Hotel
                                       4327
Retail - Uncategorized
                                       3362
Education - College or university
                                       3118
Commercial - Other
                                       2586
Name: facility type, dtype: int64
```

### Descriptive Statistics: correlation matrix

- Most promising features
  - Features with low NaN counts
  - Turn categorical features into multi-hot
  - Calculate correlation between features

site_eui			100000																			
year floor area	Water Steel																					0.1 -0.0 0.0 -0.0
year built																						0.1 -0.0
energy_star_rating	-0.5 0				0.0					0.1 -0.					0.0 -0							
2B Hot - Dry (Phoenix-AZ)	-0.0 0	0.0 -0.0	0.0	0.0	1.0 -0	0.0 -0.0	0.0-	-0.0 -	-0.0 -0	0.0 -0.	0 -0.0	0 -0.0	-0.0	-0.0 -	0.0 0.		0 -0.0	-0.0	-0.0	0.0	-0.0	0.0 -0.0
3A Warm - Humid (Memphis-TN)	-0.0 -0	0.0 0.0	0.0	4	0.0 1.	.0 -0.0	-0.0	-0.0 -	-0.0 -0	0.0 -0.	0 -0.0	0 -0.0	-0.0	0.0	0.0 -0		0 -0.0	-0.0	0.0	-0.0	-0.0 -	0.0 -0.0
3B Warm - Dry (El Paso-TX)	-0.0 0	0.1 0.0	0.1	0.0	0.0 -0	0.0 1.0	-0.0	-0.1	-0.0 -0	0.0 -0.	0 -0.0	0 -0.0	-0.0	-0.2	0.0 0.	.0 0.0	0.0	0.2	-0.1	0.1	0.1	0.1 -0.0
3C Warm - Marine (San Francisco-CA)	-0.0 0	0.2 0.1	-0.2	0.1 -	0.0 -0	0.0 -0.0	1.0	-0.1	-0.0 -0	0.0 -0.	1 -0.0	0 -0.0	-0.0	-0.3	0.0 0.	0 0.0	0 0.1	0.0	-0.1	0.2	0.0	0.2 -0.0
4A Mixed - Humid (Baltimore-MD)	0.3 0	0.3	-0.4	-0.1 -0	0.0 -0	0.0 -0.1	-0.1	1.0	-0.0 -0	0.1 -0.	2 -0.0	0 -0.0	-0.0	-0.2	0.0 0.	.0 0.0	0 0.1	0.0	0.5	0.2	-0.0	0.0 -0.0
4B Mixed - Dry (Albuquerque-NM)	-0.0 0	0.0 -0.0	0.0	0.0	0.0 -0	0.0 -0.0	-0.0	-0.0	1.0	0.0 -0.	0 -0.0	0 -0.0	-0.0	-0.0 -	0.0 0.	.0 0.0	0.0	-0.0	-0.0	0.0	0.0	0.0 -0.0
4C Mixed - Marine (Salem-OR)	-0.0 0	0.0	0.0	0.1	0.0 -0	0.0 -0.0	-0.0	-0.1	-0.0 1	1.0 -0.	1 -0.0	0 -0.0	-0.0	-0.1	0.0 0.	.0 0.0	0.1	0.1	0.2	0.1	0.0	0.0 -0.0
5A Cool - Humid (Chicago-IL)	0.2 0	0.3 0.2	-0.1	-0.0 -	0.0 -0	0.0 -0.0	-0.1	-0.2	-0.0 -0	0.1 1.0	-0.0	0 -0.0	-0.0	-0.2	0.0 0.	0.0	0 0.1	-0.0	0.1	0.2	0.0	0.0 0.0
5B Cool - Dry (Boise-ID)	-0.0 -0	0.0 -0.0	0.0	0.0 -	0.0 -0	0.0 -0.0	-0.0	-0.0	-0.0 -0	0.0 -0.	0 1.0	0.0-	-0.0	-0.0 -	0.0 0.		0 -0.0	-0.0	-0.0	0.0	0.0	0.0 -0.0
6A Cold - Humid (Burlington-VT)	0.0 -0	0.0 -0.0	0.0	-(	0.0 -0	0.0 -0.0	-0.0	-0.0 +	-0.0 -0	0.0 -0.	0 -0.0	0 1.0	-0.0	0.0 -	0.0 -0		0 -0.0	-0.0	-0.0	-0.0	-0.0 -	0.0 0.6
6B Cold - Dry (Helena-MT)	-0.0 0	0.0 -0.0	0.0-	-	0.0 -0	0.0 -0.0	-0.0	-0.0 -	-0.0 -0	0.0 -0.	0 -0.0	0 -0.0	1.0	-0.0 -	0.0 -0		0 -0.0	-0.0	-0.0	-0.0	-0.0	0.0 -0.0
dass_is_residential	-0.3 -0	0.4 -0.3	0.2	-0.1 -0	0.0 0.	.0 -0.2	-0.3	-0.2	-0.0 -0	0.1 -0.	2 -0.0	0.0	-0.0	1.0	0.0 -0		1 -0.2	-0.3	0.3	-0.7	-0.1 -	0.3 0.0
Data Center	0.1 0	0.0	0.0	-0.0 -0	0.0 -0	0.0 0.0	0.0	0.0	-0.0 -0	0.0 0.0	0.0-	0 -0.0	-0.0	-0.0	1.0 -0	.0 -0.	0 -0.0	-0.0	-0.0	-0.0	-0.0 -	0.0 -0.0
Grocery	0.2 0	0.1 -0.0	0.0	-0.0	0.0 -0	0.0 0.0	0.0	0.0	0.0 0	0.0 0.0	0.0	0.0-	-0.0	-0.1	0.0 1.	.0 -0.	0 -0.0	-0.0	-0.0	-0.0	-0.0 -	0.0 -0.0
Hospital	0.2 0	0.0 0.1	-0.0	-0.0 -0	0.0 -0	0.0 0.0	0.0	0.0	0.0 0	0.0 0.0	0.0-	0 -0.0	-0.0	-0.1 -	0.0 -0	.0 1.0	0.0	-0.0	-0.0	-0.0	-0.0 -	0.0 -0.0
Hotel	0.1 0	0.1 0.1	-0.1	-0.1	0.0 -0	0.0 0.0	0.1	0.1	-0.0 0	0.1 0.	1 -0.0	0.0-	-0.0	-0.2	0.0 -0		0 1.0	-0.0	-0.1	-0.1	-0.0 -	0.0 -0.0
Industrial	-0.1 0	0.1 0.0	0.0-	-0.0 -	0.0 -0	0.0 0.2	0.0	0.0	-0.0 0	0.1 -0.	0 -0.0	0 -0.0	-0.0	-0.3	0.0 -0		0 -0.0	1.0	-0.1	-0.1	-0.0 -	0.0 -0.0
Multifamily																				0.000		0.1 -0.0
Office	0.1 0	0.3 0.3	-0.3	0.1	0.0 -0	0.0 0.1	0.2	0.2	0.0 0	0.1 0.	2 0.0	0.0-	-0.0	-0.7	0.0 -0		0 -0.1	-0.1	-0.2	1.0	-0.0 -	0.1 -0.0
Restaurant	0.1 0	0.0	0.0-	-0.0 -	0.0 -0	0.0 0.1	0.0	-0.0	0.0 0	0.0 0.0	0.0	0.0-	-0.0	-0.1	0.0 -0		0 -0.0	-0.0	-0.0	-0.0	1.0	0.0 -0.0
Retail	0.1 0	0.1 0.0	-0.1	-0.0	0.0 -0	0.0 0.1															-	1.0 -0.0
Singlefamily	0.0 -0	0.0 -0.0	0.0-	-(	0.0 -0	0.0 -0.0	0.0-	-0.0	-0.0 -0	0.0 0.0	0.0-	0.6	-0.0	0.0 -	0.0 -0	.0 -0.	0 -0.0	-0.0	-0.0	-0.0	-0.0 -	0.0 1.0
	e eui	year	ar_built	ating	oenix-AZ)	o-TX)	0-CA)	-WD)	-NM)	1-OR)	(e-ID)	n-VT)	3-MT)	ential	enter	pital	Hotel	dustrial	Itifamily	Office	ırant	Retail
	site	floor	year	energy_star_rating	hoeni	El Pas	incisco	timore	rerque	(Salem Chicag	(Bois	-	(Helena-MT)	is_residential	Data Cente	· 호	-	Indus	Multifa	J	Restau	nglefa
				ergy	ot - Dry (Ph Humid (Me	Dry (	an Fran	d (Balti	bng	arine imid (		(Bur	Dry (F	class is							_	i/s
					Hot - I	/arm -	ine (Sa		> .	ed - Mar ool - Hur	SB Cool		6B Cold -	G								
				Ì	ZB H	38 W	- Marin	- pax	- pa	M O	2	Cold -	99									
					3.4 V		Warm		48 Mi	D 75		6A 0										
							3C )	0														

site_eu	1.0	0.2	0.2 -	0.2 -0	.5 -0.0	0.0- 0	0.0-	-0.0	0.3	-0.0	-0.0	0.2	-0.0	0.0 -0	.0 -0.3	0.1	0.2	0.2	0.1	0.1	0.1 0.1	1 0.1	0.1	0.0	
yea	0.2	1.0	0.3 -	0.2 0.	0 0.0	0.0-	0.1	0.2	0.3	0.0	0.3		-0.0	0.0 0	.0 -0.4	0.0	0.1	0.0	0.1	0.1		0.0	0.1	-0.0	
floor_are	0.2	0.3	1.0 -	0.1 0.	0 -0.	0.0	0.0	0.1	0.3	-0.0	0.0	0.2	-0.0	-0.0 -0	.0 -0.3	0.0	-0.0	0.1	0.1	0.0	0.2 0.3	0.0	0.0	-0.0	
year_buil	-0.2	-0.2	-0.1	1.0 -0	.0 0.0	0.0	0.1	-0.2	-0.4	0.0	-0.0	-0.1	0.0	-0.0 -0	.0 0.2	-0.0	0.0	-0.0	-0.1 -	0.0	0.3 -0.	3 -0.0	-0.1	-0.0	
energy_star_ratin	-0.5	0.0	0.0 -	0.0 1.	0.0		0.0	0.1	-0.1	0.0	0.1	-0.0	0.0		-0.1	-0.0	-0.0	-0.0	-0.1 -	0.0 -	0.1 0.1	1 -0.0	-0.0		
2B Hot - Dry (Phoenix-AZ	-0.0	0.0	-0.0	0.0 0.	0 1.0	-0.0	0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0 -0	.0 -0.0	-0.0	0.0	-0.0	-0.0	0.0 -	0.0 0.0	0.0-	0.0	-0.0	
3A Warm - Humid (Memphis-TM	-0.0	-0.0	0.0	0.0	-0.0	0 1.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0 -0	.0 0.0	-0.0	-0.0	-0.0	-0.0 -	0.0	0.0 -0.	0 -0.0	-0.0	-0.0	
3B Warm - Dry (El Paso-T)	-0.0	0.1	0.0	0.1 0.	0 -0.0	0.0-	1.0	-0.0	-0.1	-0.0	-0.0	-0.0	-0.0	-0.0 -0	.0 -0.2	0.0	0.0	0.0	0.0	0.2	0.1 0.1	1 0.1	0.1	-0.0	
3C Warm - Marine (San Francisco-CA	-0.0	0.2	0.1	0.2 0.	1 -0.0	0.0-	-0.0	1.0	-0.1	-0.0	-0.0	-0.1	-0.0	-0.0 -0	.0 -0.3	0.0	0.0	0.0	0.1	0.0 -	0.1 0.2	2 0.0	0.2	-0.0	
4A Mixed - Humid (Baltimore-MD	0.3	0.3	0.3 -	0.4 -0	.1 -0.0	0.0- 0	0 -0.1	-0.1	1.0	-0.0	-0.1	-0.2	-0.0	-0.0 -0	.0 -0.2	0.0	0.0	0.0	0.1	0.0	0.5 0.2	2 -0.0	0.0	-0.0	
4B Mixed - Dry (Albuquerque-NM	-0.0	0.0	-0.0	0.0 0.	0 -0.	0.0-	0.0	-0.0	-0.0	1.0	-0.0	-0.0	-0.0	-0.0 -0	.0 -0.0	-0.0	0.0	0.0	-0.0 -	0.0	0.0 0.0	0.0	0.0	-0.0	
4C Mixed - Marine (Salem-OF	-0.0	0.3	0.0 -	0.0 0.	1 -0.0	0.0-	0.0	-0.0	-0.1	-0.0	1.0	-0.1	-0.0	-0.0 -0	.0 -0.1	-0.0	0.0	0.0	0.1	0.1 (	0.2 0.1	1 0.0	0.0	-0.0	
5A Cool - Humid (Chicago-Il	0.2	0.3	0.2 -	0.1 -0	.0 -0.	0.0-	0.0-	-0.1	-0.2	-0.0	-0.1	1.0	-0.0	-0.0 -0	.0 -0.2	0.0	0.0	0.0	0.1 -	0.0	0.1 0.2	2 0.0		0.0	
5B Cool - Dry (Boise-ID	-0.0	-0.0	-0.0	0.0 0.	0 -0.0	0.0-	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	1.0	-0.0 -0	.0 -0.0	-0.0	0.0	-0.0	-0.0 -	0.0 -	0.0 0.0	0.0	0.0	-0.0	
6A Cold - Humid (Burlington-V1	0.0	-0.0	-0.0 -	0.0	-0.0	0.0-	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	1.0 -0	.0 0.0	-0.0	-0.0	-0.0	-0.0	0.0 -	0.0 -0.	0 -0.0	-0.0	0.6	
6B Cold - Dry (Helena-MT	-0.0	0.0	-0.0 -	0.0	-0.0	0.0-	0.0	-0.0	-0.0	-0.0	-0.0	-0.0	-0.0	0.0 1	.0 -0.0	-0.0	-0.0	-0.0	-0.0 -	0.0 -	0.0 -0.	0 -0.0	0.0	-0.0	
dass_is_residentia	-0.3	-0.4	-0.3	0.2 -0	.1 -0.0	0.0	-0.2	-0.3	-0.2	-0.0	-0.1	-0.2	-0.0	0.0 -0	.0 1.0	-0.0	-0.1	-0.1	-0.2 -	0.3	0.3 -0.	7 -0.1	-0.3	0.0	
Data Cente	0.1	0.0	0.0	0.0 -0	.0 -0.	0.0-	0.0	0.0	0.0	-0.0	-0.0	0.0	-0.0	-0.0 -0	.0 -0.0	1.0	-0.0	-0.0	-0.0	0.0	0.0 -0.	0.0-	-0.0	-0.0	
Grocer	0.2	0.1	-0.0	0.0 -0	.0 0.0	-0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	-0.0 -0	.0 -0.1	-0.0	1.0	-0.0	-0.0	0.0	0.0 -0.	0.0-	-0.0	-0.0	
Hospita	0.2	0.0	0.1	0.0 -0	.0 -0.	0.0-	0.0	0.0	0.0	0.0	0.0	0.0	-0.0	0.0 -0	.0 -0.1	-0.0	-0.0	1.0	-0.0	0.0 -	0.0 -0.	0 -0.0	-0.0	-0.0	
Hote	0.1	0.1	0.1 -	0.1 -0	.1 -0.0	0.0-	0.0	0.1	0.1	-0.0	0.1	0.1	-0.0	-0.0 -0	.0 -0.2	-0.0	-0.0	-0.0	1.0	0.0	0.1 -0.	1 -0.0	-0.0	-0.0	
Industria	-0.1	0.1	0.0 -	0.0 -0	.0 -0.	0.0- 0	0.2	0.0	0.0	-0.0	0.1	-0.0	-0.0	-0.0 -0	.0 -0.3	-0.0	-0.0	-0.0	-0.0	1.0	0.1 -0.	1 -0.0	-0.0	-0.0	
Multifamil	0.1	0.4	0.2	0.3 -0	.1 -0.0	0.0	-0.1	-0.1	0.5	-0.0	0.2	0.1	-0.0	0.0 -0	.0 0.3	-0.0	-0.0	-0.0	-0.1 -	0.1	1.0 -0.	2 -0.0	-0.1	-0.0	
Offic	0.1	0.3	0.3 -	0.3 0.	1 0.0	-0.0	0.1	0.2	0.2	0.0	0.1	0.2	0.0	0.0 -0	.0 -0.7	-0.0	-0.0	-0.0	-0.1 -	0.1	0.2 1.0	0.0	-0.1	-0.0	
Restauran	0.1	0.0	0.0 -	0.0 -0	.0 -0.	0.0-	0.1	0.0	-0.0	0.0	0.0	0.0	0.0	-0.0 -0	.0 -0.1	-0.0	-0.0	-0.0	-0.0 -	0.0	0.0 -0.	0 1.0	-0.0	-0.0	
Reta	0.1	0.1	0.0 -	0.1 -0	.0 0.0	-0.0	0.1	0.2	0.0	0.0	0.0	0.0	0.0	-0.0 0	.0 -0.3	-0.0	-0.0	-0.0	-0.0	0.0 -	0.1 -0.	1 -0.0	1.0	-0.0	
Singlefamil	0.0	-0.0	-0.0 -	0.0	-0.0	0.0-	0.0-	-0.0	-0.0	-0.0	-0.0	0.0	-0.0	0.6	.0 0.0	-0.0	-0.0	-0.0	-0.0	0.0 -	0.0 -0.	0.0-	-0.0	1.0	
	.in	year	rea	# S	AZ)	Ñ.	ίχ	(S	4D)	(M)	OR)	<u>-</u>	(Q)	5 5	tial	iter	ery	Ital	Hotel	lal	family	ant	Retail	ylic	
	site_eui	\$	floor_area	year_built	- Dry (Phoenix-AZ)	-siydi	Paso-TX)		nore-N	dne-N	lem-(	icago	Dry (Boise-ID)	- Humid (Burlington-VT)	dass_is_residential	Data Center	Grocery	Hospital	운	Industrial	Multifamily Office	stauran	-Re	Singlefamily	
			ij	* +	, (Pho	(Mem	Dry (El	(San Francisco	Baltin	ndner	e (Sa	d (Ch	Dry (E	Burlin	is re	Dat				-	₩.	S.		Sing	
				9	- 1	pimn	ū.	(San	) pim	/ (Alb	Marin	Hum	Cool -	mid (	dass										
					28 Hot - Dry (Phoenix-AZ	Warm - Humid (Memphis-TN)	3 Warm -	Marine	d - Hu	Mixed - Dry (Albuquerque-NM	Mixed - Marine (Salem-OR)	Cool - Humid (Chicago-IL)	580	Id - Humid (Burlington-VT)	5										
					124		38		4A Mixed - Humid (Baltimore-MD)	Mixe	4C Mi	5A (		6A Cold	,										
						3A		3C Warm	4A	48				9											
								ñ																	



# Experiments

### Prediction Algorithm

#### **Linear Regression**

• The feature we are trying to predict is a continuous, numerical variable which makes it a good candidate for predicting with a linear model.

#### Random Forest

• A random forest regression may be best able to handle the number of categorical and missing variables.

#### Feed Forward Neural Network

 Tried to improve on our initial linear model by building a feed forward neural network to try and achieve any possible marginal gains.

#### How will we evaluate our results?

- We chose to use Mean Absolute Error (MAE) as our evaluation metric.
- It is easy to interpret since it tells us how many energy intensity units our predictions were off by.
- MAE is less sensitive to outliers than Mean Squared Error and our data have observations with site EUI values multiple standard deviations above the mean.

$$MAE = \frac{1}{N} \sum_{i=1}^{N} \left| y_i - \widehat{y}_i \right|$$

#### Baseline Model

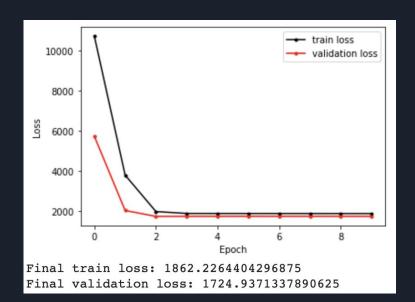
MSE and MAE from the mean

MSE Baseline to Train: 1563.4135692561551 MSE Baseline to Test: 1349.2001965901652

MAE Baseline to Train: 29.96627005205943 MAE Baseline to Test: 28.480106853244518

# Experiments: Linear Regression

#### Baseline Linear Model - Year Built

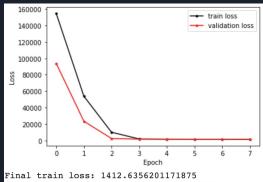


train mae validation mae

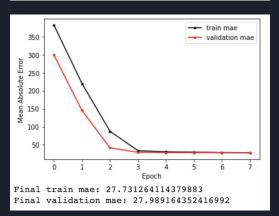
Train mae validation mae

Train mae

### Linear Model - Full Feature



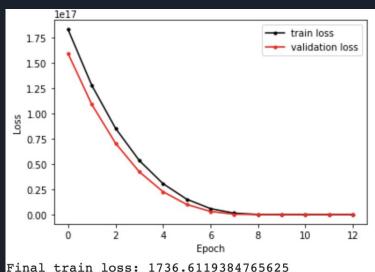
Final validation loss: 1533.863037109375



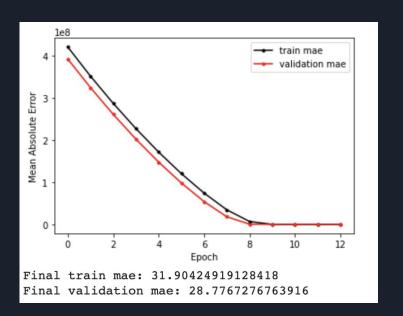
#### Features include:

- Floor area
- Year of construction
- Year of data collection
- Climate zone one hot
- Building type one hot
- Window glass layers one hot
- Roof type one hot
- Cooling type one hot
- Heating type one hot
- Wall type one hot

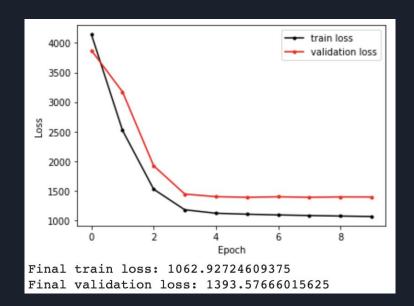
#### Normalized Baseline Linear Model - Year Built

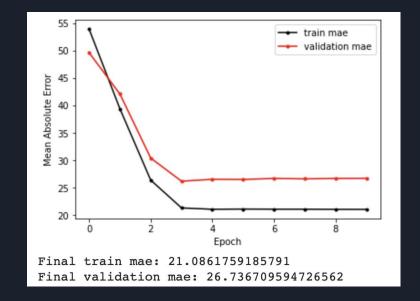


Final train loss: 1736.6119384765625
Final validation loss: 1636.8787841796875



#### Normalized Linear Model - Full Feature

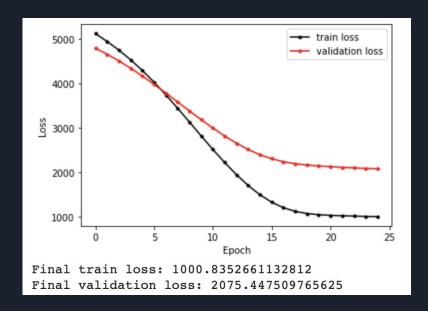


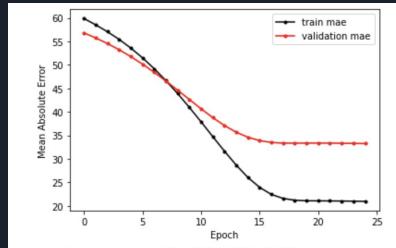


### Deep Neural Network Regression

- Normalization layer
- Dense non-linear layer with ReLU activation
- Dense linear single-output layer

Model: "sequential_1"		
Layer (type)	Output Shape	Param #
normalization (Normalizatio n)	(None, 55)	111
dense_2 (Dense)	(None, 64)	3584
dense_3 (Dense)	(None, 1)	65
Total params: 3,760 Trainable params: 3,649 Non-trainable params: 111		





Final train mae: 20.94822120666504 Final validation mae: 33.24668502807617

#### Best Linear Model (Normalized Full-Feature) Weights

Final bias: 6.33 year 222.7343 floor area 127.4397 year built 175.1638 clim 2A Hot - Humid 0.3617 clim 3C Warm - Marine 0.0001 clim 5A Cool - Humid 0.0 clim 3B Warm - Dry 0.0184 clim 2B Hot - Dry 0.0415 clim 4C Mixed - Marine 0.361 clim 6B Cold - Dry 0.0002 clim 4B Mixed - Dry 0.0893 clim 5B Cool - Dry 0.1233 clim 4A Mixed - Humid 0.0 clim 6A Cold - Humid 0.0011 clim 3A Warm - Humid 0.0 BC Residential 0.314 BC Commercial 0.686 FT Multifamily 0.0003 FT Industrial 0.0073 FT Office 0.0056 FT Retail 0.0442 FT Other 0.036 FT Hotel 0.3233 FT Restaurant 0.1804 FT Grocery 0.0113 FT Hospital 0.0024 FT Data Center 0.037 FT Singlefamily 0.0027

```
WGL Single-pane 0.0012
WGL Double-pane 0.0058
WGT Low-e 0.0002
roof Shingles 0.3233
roof Built-up 0.0168
roof Slate or tile shingles 0.0
roof Metal surfacing 0.0023
roof Other Or Combination 0.0031
roof Asphalt/fiberglass/other shingles 0.0002
roof Wood shingles/shakes/other wood 0.0071
roof Plastic/rubber/synthetic sheeting 0.0006
roof Green Roof 0.0
cool Central AC 0.3305
cool No cooling 0.0
cool Other 0.0102
cool Cooling Heat Pump 0.0
cool Split AC 0.0
heat Boiler 0.0002
heat Resistance Heating 0.0
heat Heating Heat pump 0.0
heat Other 0.0
heat Heating Furnace 0.0
wall Other 0.0013
wall wall Wood 0.0004
wall wall Metal 0.0
wall wall Brick stone 0.0003
wall wall Concrete 0.0007
```

#### Linear Model Results - Test Set

Full featured normalized model has best performance and had most ideal loss plot

	Year_Built	Full_Feature	Year_Built_Norm	Full_Feature_Norm	DNN_Model
MSE	1617.91	1175.92	1477.95	803.35	965.60
MAE	31.59	25.49	29.90	19.12	20.75

# Experiments: Random Forest

#### Random Forest Model

#### Features:

- Facility type, floor area, year\_built, roof/ceiling, window glass layers, window glass type, energy star label and rating, cooling ,heating, residential, dryness, temperature
- Data cleaning:
  - Drop rows without site EUI
  - Drop duplicate rows for same building (keep most recent)
  - o Train: 34k, Validation: 8k, Test: 19k



#### Random Forest

#### Model:

```
# Specify the model.
model = tfdf.keras.RandomForestModel(
    task = tfdf.keras.Task.REGRESSION,
    temp directory= path + 'output ben/',
    verbose = 2,
    allow na conditions=True,
    max depth=25,
    min examples=3,
    num trees=300,
    random seed=3000,
    name='Random Forest Regressor',
```

#### Train:

mse: 540.2357 mae: 10.4744

mape: 32.1941

RMSE: 23.24297004985123

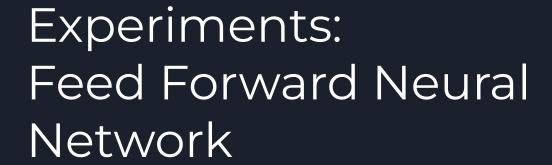
Test:

mse: 1218.7836

mae: 16.1808

mape: 47.1652

RMSE: 34.9110808961272

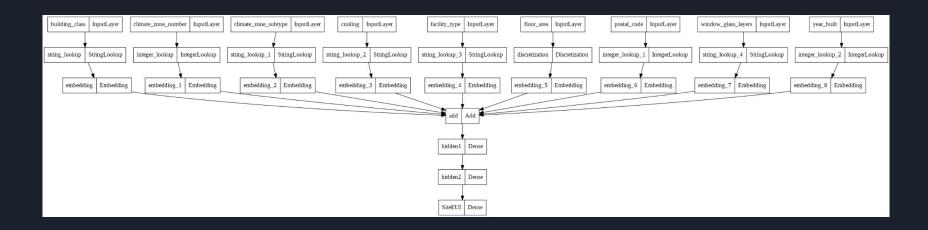


#### **FFNN**

- Extract climate zone number and subtype as individual features
- Convert all features to embeddings
- Use Tensorflow Functional API with a simple model architecture of 2 fully connected layers with 32 nodes each.
- Exponential decay learning rate schedule

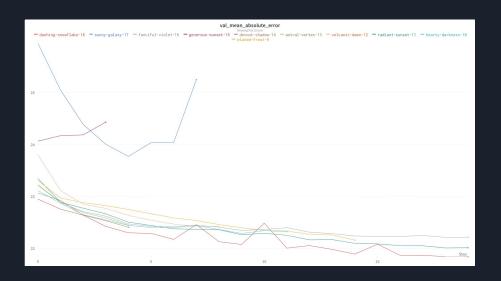
```
# Tensorflow input lavers
building_class = tf.keras.layers.Input(shape=(1,), dtype=tf.string, name='building_class')
climate_zone_number = tf.keras.layers.Input(shape=(1,), dtype=tf.int64, name='climate_zone_number')
climate_zone_subtype = tf.keras.layers.Input(shape=(1,), dtype=tf.string, name='climate_zone_subtype')
cooling = tf.keras.layers.Input(shape=(1,), dtype=tf.string, name='cooling')
facility type = tf.keras.layers.Input(shape=(1,), dtype=tf.string, name='facility type')
floor_area = tf.keras.layers.Input(shape=(1,), dtype=tf.int64, name='floor_area')
postal_code = tf.keras.layers.Input(shape=(1,), dtype=tf.int64, name='postal_code')
window_glass_layers = tf.keras.layers.Input(shape=(1,), dtype=tf.string, name='window_glass_layers')
year_built = tf.keras.layers.Input(shape=(1,), dtype=tf.int64, name='year_built')
# Tensorflow categorical feature layers
building_class_binned = tf.keras.layers.StringLookup(vocabulary=['Commercial', 'Residential'])(building_class)
climate_zone_number_binned = tf.keras.layers.IntegerLookup(vocabulary=[2, 3, 4, 5, 6])(climate_zone_number)
climate zone subtype binned = tf.keras.layers.StringLookup(vocabulary=['A', 'B', 'C'])(climate zone subtype)
cooling vocab = get feature vocab(df train.cooling)
cooling binned = tf.keras.layers.StringLookup(vocabulary=cooling vocab)(cooling)
facility type vocab = get feature vocab(df train.facility type)
facility_type_binned = tf.keras.layers.StringLookup(vocabulary=facility_type_vocab)(facility_type)
floor_area_bins = list(np.percentile(df_train_int["floor_area"],[10,20,30,40,50,50,70,80,90]))
floor area binned = tf.keras.layers.Discretization(bin boundaries=floor area bins)(floor area)
postal code vocab = get feature vocab(df train.postal code)
postal code binned = tf.keras.layers.IntegerLookup(vocabulary=postal code vocab)(postal code)
window_glass_layers_binned = tf.keras.layers.StringLookup(vocabulary=['Single-pane', 'Double-pane'])(window_glass_layers)
year_built_vocab = get_feature_vocab(df_train.year_built)
year_built_binned = tf.keras.layers.IntegerLookup(vocabulary=year_built_vocab)(year_built)
# Tensorflow embedding layers
building_class_embed = tf.keras.layers.Embedding(
    input dim = 3, output dim = embed dim, input length = 1)(building class binned)
climate zone number embed = tf.keras.layers.Embedding(
    input dim = 6. output dim = embed dim. input length = 1)(climate zone number binned)
climate zone subtype embed = tf.keras.layers.Embedding(
    input_dim = 4, output_dim = embed_dim, input_length = 1)(climate_zone_subtype_binned)
cooling_embed = tf.keras.layers.Embedding(
    input_dim = len(cooling_vocab) + 1, output_dim = embed_dim, input_length = 1)(cooling_binned)
facility_type_embed = tf.keras.layers.Embedding(
    input dim = len(facility type vocab) + 1 , output dim = embed dim, input length = 1)(facility type binned)
floor area embed = tf.keras.layers.Embedding(
    input_dim = len(floor_area_bins)+1, output_dim = embed_dim, input_length = 1)(floor_area_binned)
postal_code_embed = tf.keras.layers.Embedding(
    input_dim = len(postal_code_vocab)+1, output_dim = embed_dim, input_length = 1)(postal_code_binned)
window_glass_layers_embed = tf.keras.layers.Embedding(
    input_dim = 3, output_dim = embed_dim, input_length = 1)(window_glass_layers_binned)
year_built_embed = tf.keras.layers.Embedding(
    input_dim = len(year_built_vocab) + 1, output_dim = embed_dim, input_length = 1)(year_built_binned)
```

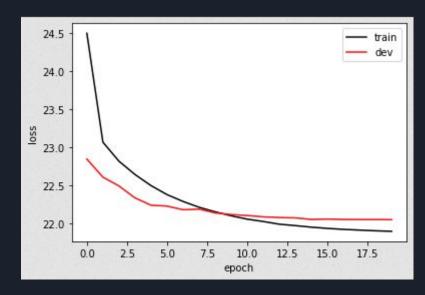
#### FFNN



#### FFNN

It proved to be very difficult to beat the performance of the full-feature, normalized linear regression model Extensive experimentation and hyperparameter tuning could not return a MAE below 21.





# Conclusions

#### Conclusions

- More and better data is needed
- Need more representative buildings from the communities this model will be used in
- Residential and commercial buildings should potentially be separately modeled to reduce complexity and more effectively select features

# Code and Contributions

### Code and Contributions

#### Code

• GitHub Repo: <a href="https://github.com/bennnyys/w207-sec3-final-proj-DMR">https://github.com/bennnyys/w207-sec3-final-proj-DMR</a>

	Ben	Sam	CJ
Theoretical Research		X	
Data cleaning	X	X	X
Data splitting	X		
Hyperparameter tuning	X	X	X
Augmentations	(random forest)	(linear regression)	(FFNN)
Presentation Slides	X	X	X