

2016 CampusLab

Detecting Unusual Trends in Electricity Usage Data

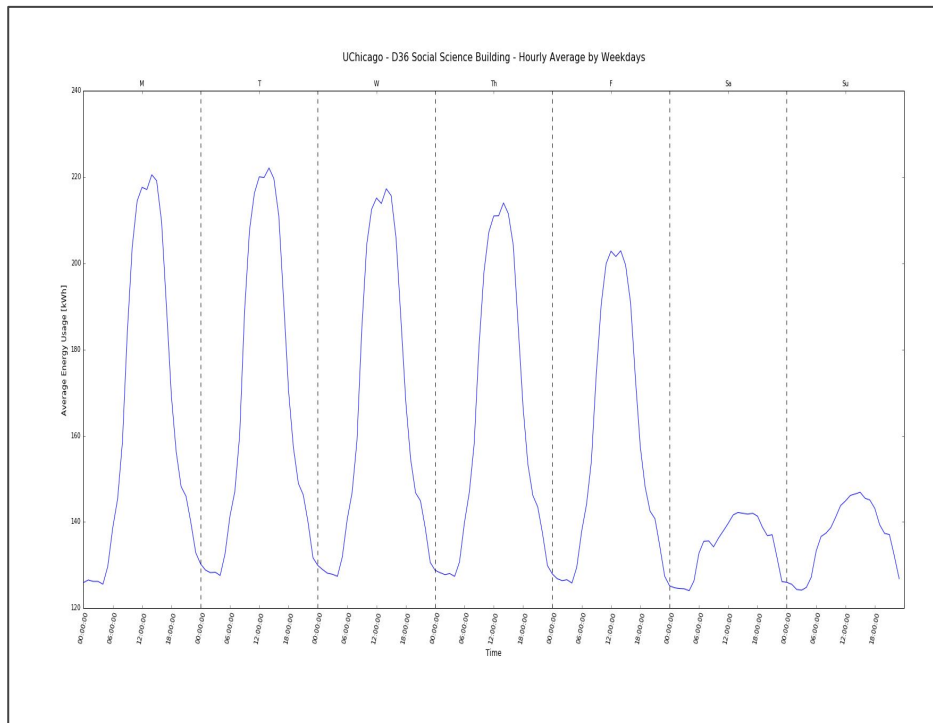
CS with Applications III (CS 12300 / CAPP 30123)
The University of Chicago, Spring 2016
Christine Chung and Leith McIndewar



University of Chicago Energy Data

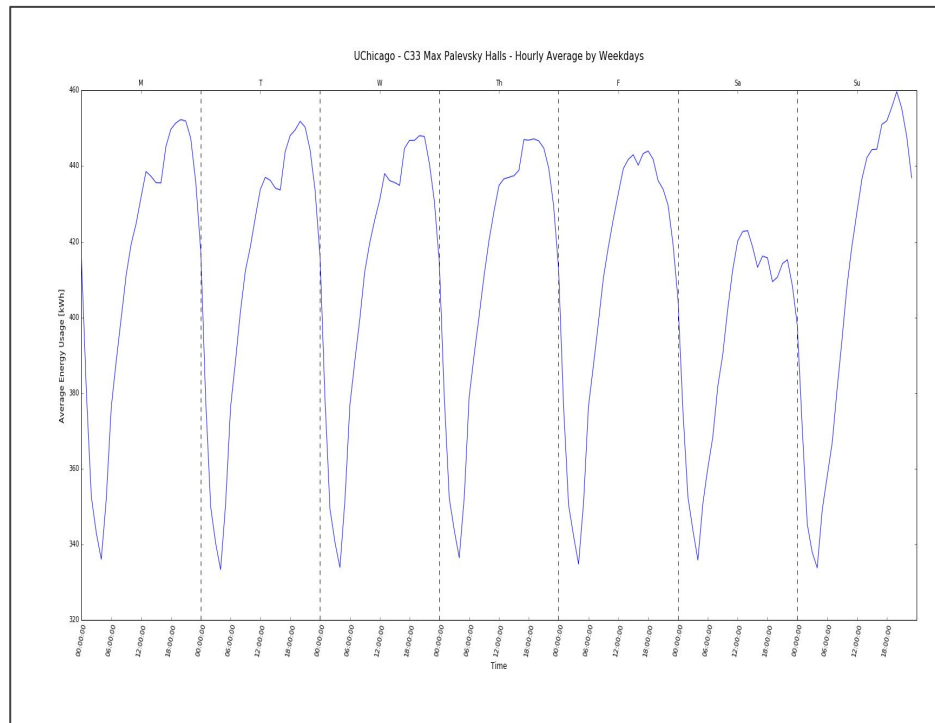
Data

- UChicago is one of the largest energy consumers in Chicago (with O'Hare)
- ~120 buildings
- Over 300 meters
- Meters and buildings do not always have a clear relationships



University of Chicago Energy Data

- Right now data is available at half hour intervals, but moving to more granular readings
- Purchase power hourly
- How would we detect if usage in a building is abnormal?
- **Goals:** Pattern & Anomaly Detection



Tools and Implementation

Tools

- PostgreSQL
- Amazon EC2, EMR, and S3
- ElasticSearch & Kibana
- MRJob
- Python Threading and Queue Class

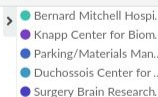
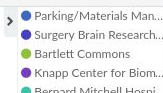
Algorithms

For anomaly detection:

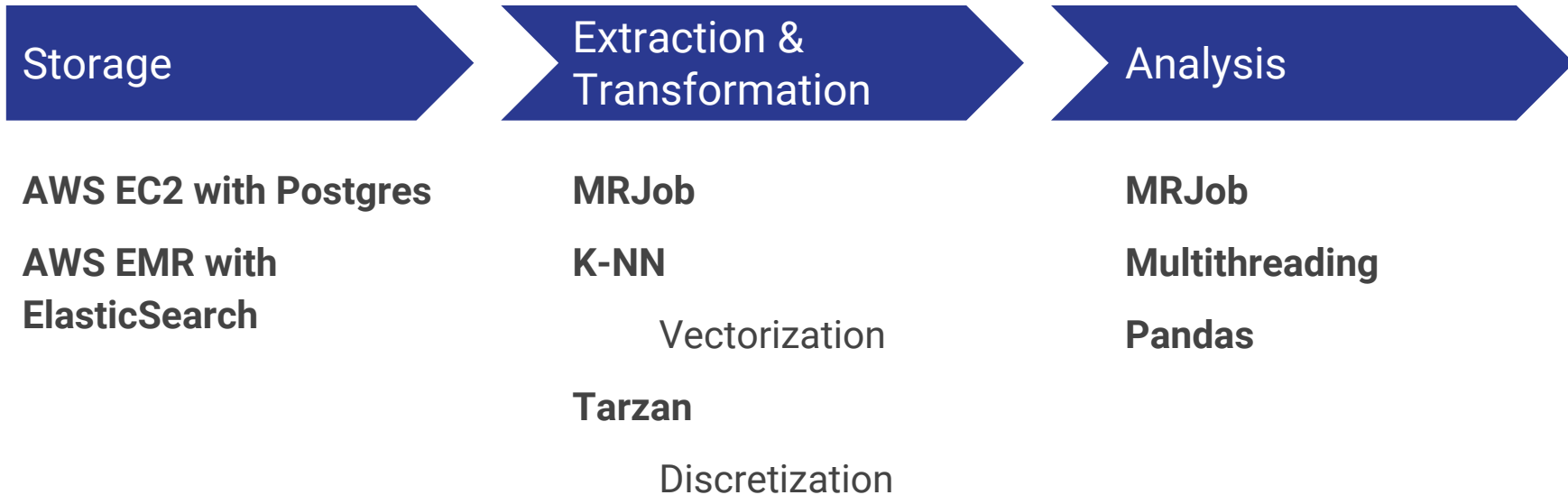
- Tarzan (Suffix Trees, Markov Models)
- K-Nearest Neighbors

Count

Average usage

Data Flow



Algorithms: Tarzan

```
void TARZAN (time_series  $R$ , time_series  $X$ ,  
             int  $l_1$ , int  $a$ , int  $l_2$ , real  $c$ )  
    let  $x = \text{DISCRETIZE\_TIME\_SERIES}(X, l_1, a)$   
    let  $r = \text{DISCRETIZE\_TIME\_SERIES}(R, l_1, a)$   
    let  $T_x = \text{PREPROCESS}(r, x)$   
    for  $i = 1, |x| - l_2 + 1$   
        let  $w = x_{[i, i+l_2-1]}$   
        retrieve  $z(w)$  from  $T_x$   
        if  $|z(w)| > c$  then print  $i, z(w)$ 
```

Table 4: Outline of the Tarzan algorithm: l_1 is the feature window length, a is the alphabet size for the discretization, l_2 is the scanning window length and c is the threshold

Tarzan Algorithm

Finding Surprising Patterns in a Time Series Database in Linear Time and Space – Keogh et al. *SIGKDD 2002*

Step 1:

Discretize time series into strings

Step 2:

Preprocess strings with Suffix Trees

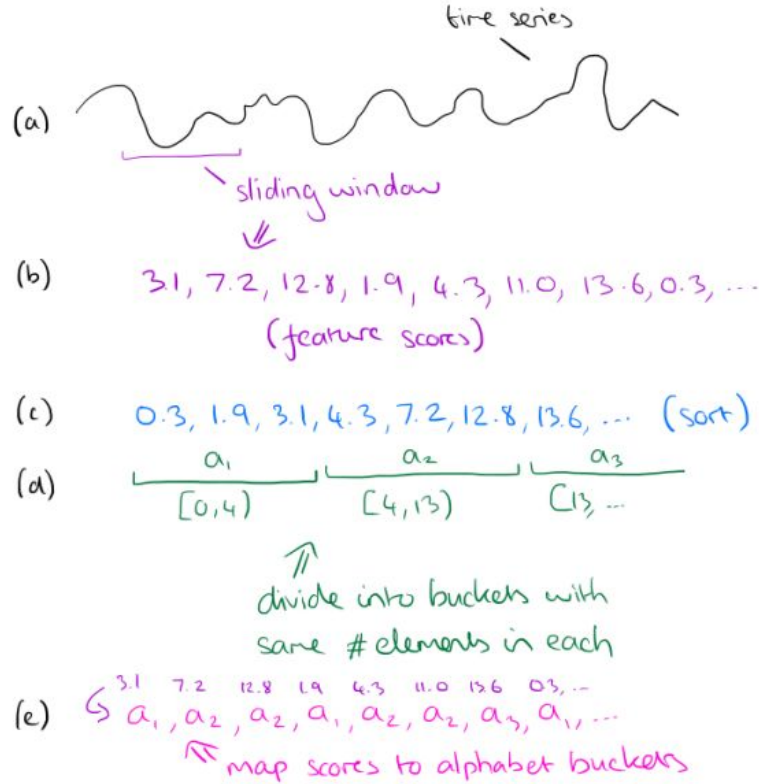
Step 3:

Score strings with Markov Models

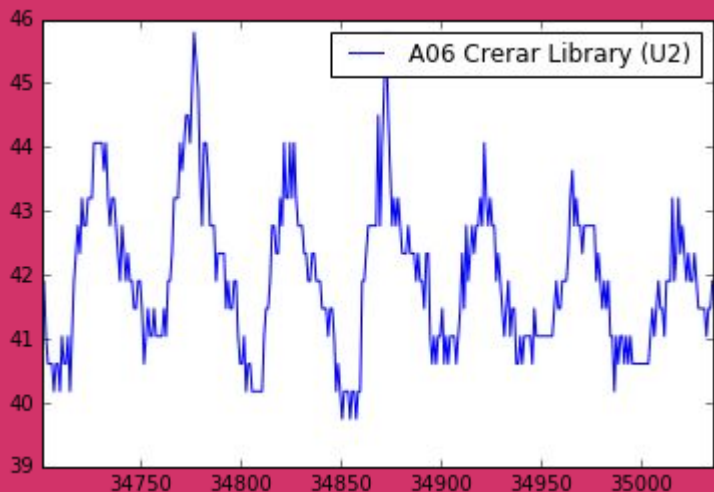


Discretization

1. Slide window over feature scores.
2. Generate one real number (slope)
3. Bucket slopes by quantiles
4. Assign symbol to each bucket
5. Replace numbers in original time series with label from bucket



Step 1: time series -----> string



=

edaafmhkkkjgklmonkkippooohhafdeelib
bgkgccebahphdejneileahplnnmmoonn
oojklkgfggkkdeccdgphecdhpfedjggfhba
hhhommpmmlalnjjhjkgefaaabaipjabaab
aelidedccmpbfmjmflkloophhkjjkjggfdb
bbapleabfcebknngffcnmoggdmpafnmlli
kjoomplljgebedddcbcnjggfcfmnmphakk
kahhhhhhpoofhpppmkgdebmfphhhcd
dccbnddeggnmhfcnceahhhhhhpononof
bknnljhljdccceilcbbhaelp

Algorithms: K-Nearest Neighbors

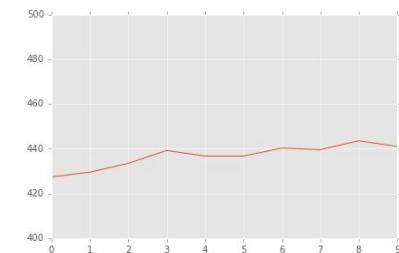
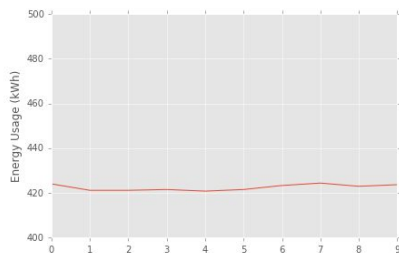
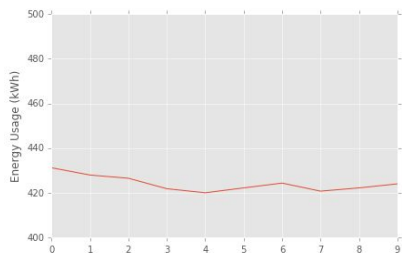
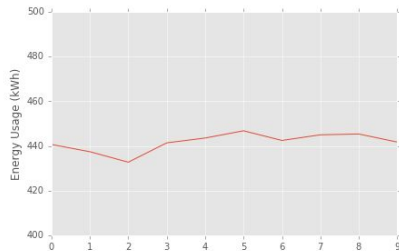
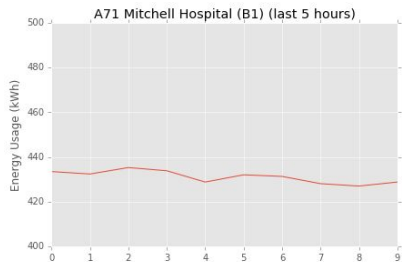
K Nearest Neighbors

- Define a window for comparison and k parameter
- For each building, compare the readings from all other buildings to find the k nearest neighbors
- Find what that building is doing during its 'nearest neighbors'
- Calculate distance to between those readings and the current readings
- Repeat for each building

RY	All Other Buildings						
6.5	2.2	4.5	7.6	3.4	9.1
3.2	2.4	5.7	8.4	2.8	1.0
.							
.							
.							
.							

6.5	2.2	4.5	7.6	3.4	9.1
3.2	2.4	5.7	8.4	2.8	1.0
.							
.							
.							
.							

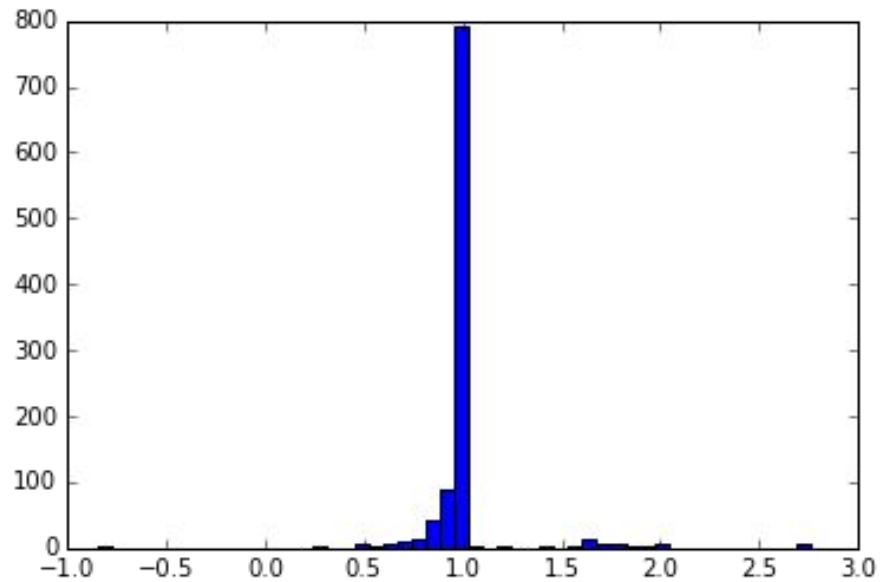
K Nearest Neighbors

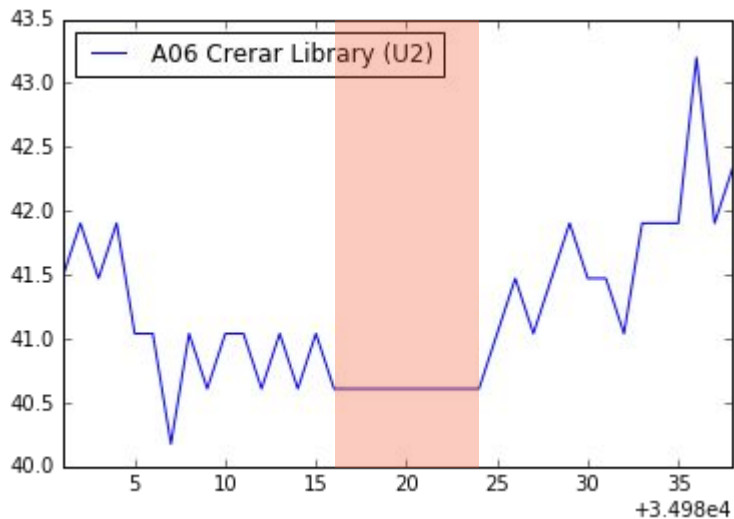


- Nearest Neighbors for Mitchell Hospital meter
- Distance between current readings (upper left) and average of 5 nearest neighbors:
 - 14.146

Results

Distribution of “Surprise Scores”





	datetime	A06 Crerar Library (U2)
34995	2016-04-10 03:00:00	41.040
34996	2016-04-10 03:30:00	40.608
34997	2016-04-10 04:00:00	40.608
34998	2016-04-10 04:30:00	40.608
34999	2016-04-10 05:00:00	40.608
35000	2016-04-10 05:30:00	40.608
35001	2016-04-10 06:00:00	40.608
35002	2016-04-10 06:30:00	40.608
35003	2016-04-10 07:00:00	40.608
35004	2016-04-10 07:30:00	40.608

“Surprise Result”

April 10, 2016 | Crerar Library

4 hours of unchanging usage from
3:30am - 7:30am

Run Times

K-NN	~ 31 min
K-NN with MR Job	~ 26 min
Tarzan	~ 17 min
Tarzan w/ Multithreading	~ ∞ min

Next Steps

Sharding with EMR and ElasticSearch

Effectively use MR Job and Multithreading with both Algorithms

Reduce Run Times



Questions?