# K-Means and Gaussian Mixtures

Dr. Chelsea Parlett-Pelleriti

# Unsupervised Machine Learning

# Clustering

# K-Means
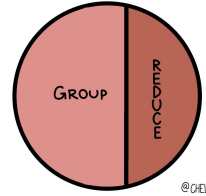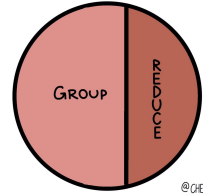
1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the centers
4. Repeat 2 and 3 until either:
   a. Cluster membership does not change
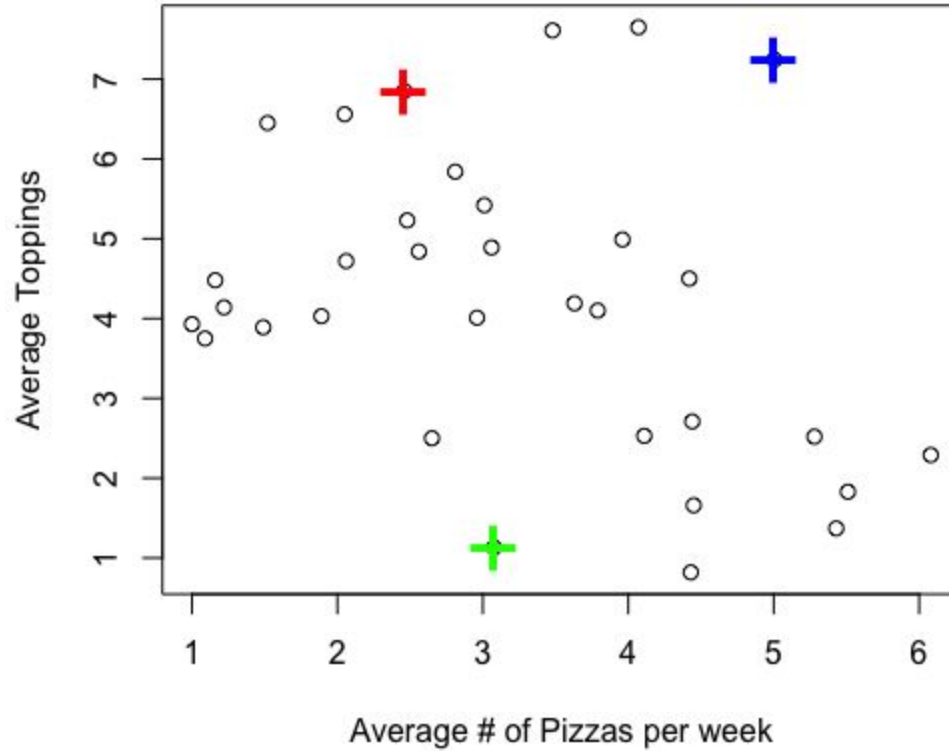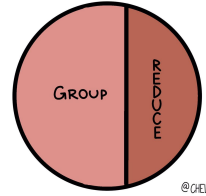   b. Centers change only a tiny amount

# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the centers



SIMPLIFY

GROUP | REDUCE

@CHELSEAPARLETT

**1**

# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
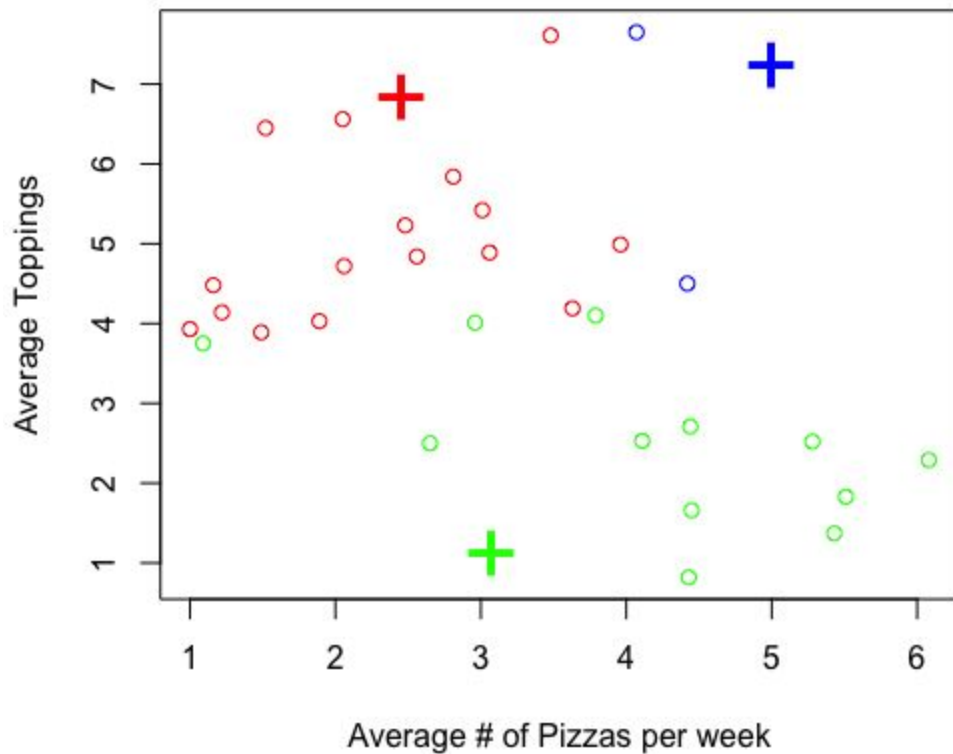3. Using these assignments, recalculate the centers

SIMPLIFY

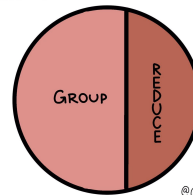GROUP REDUCE

@CHELSEAPARLETT

**2**

# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the centers

**3**

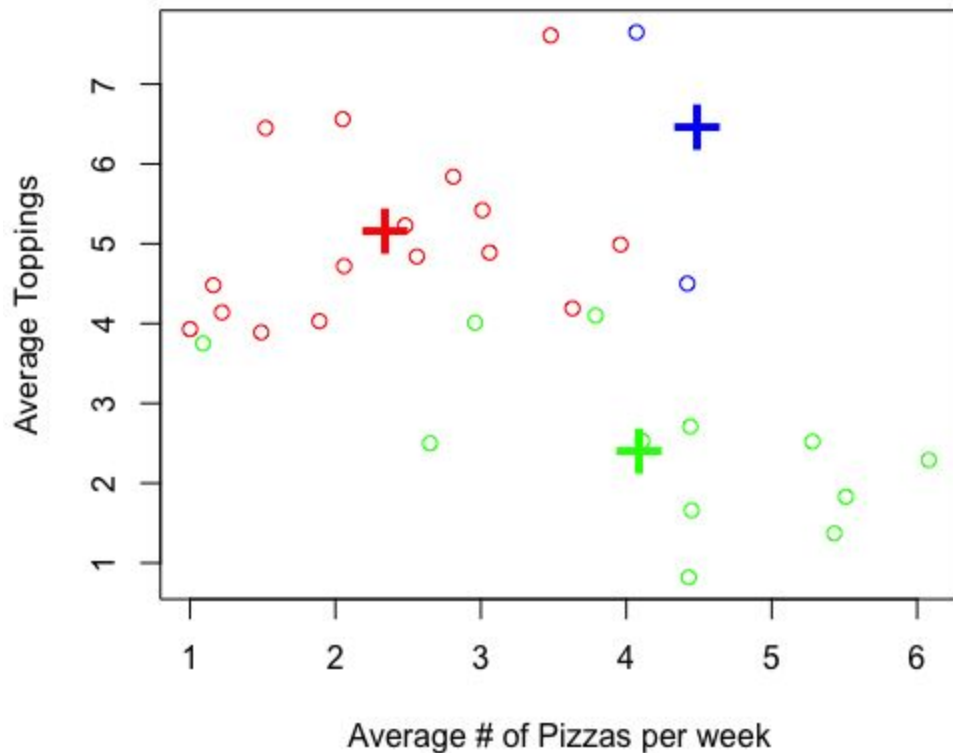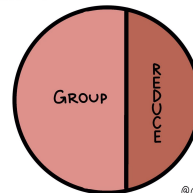# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the centers



**2**
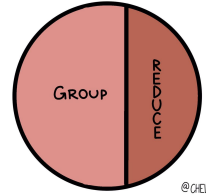
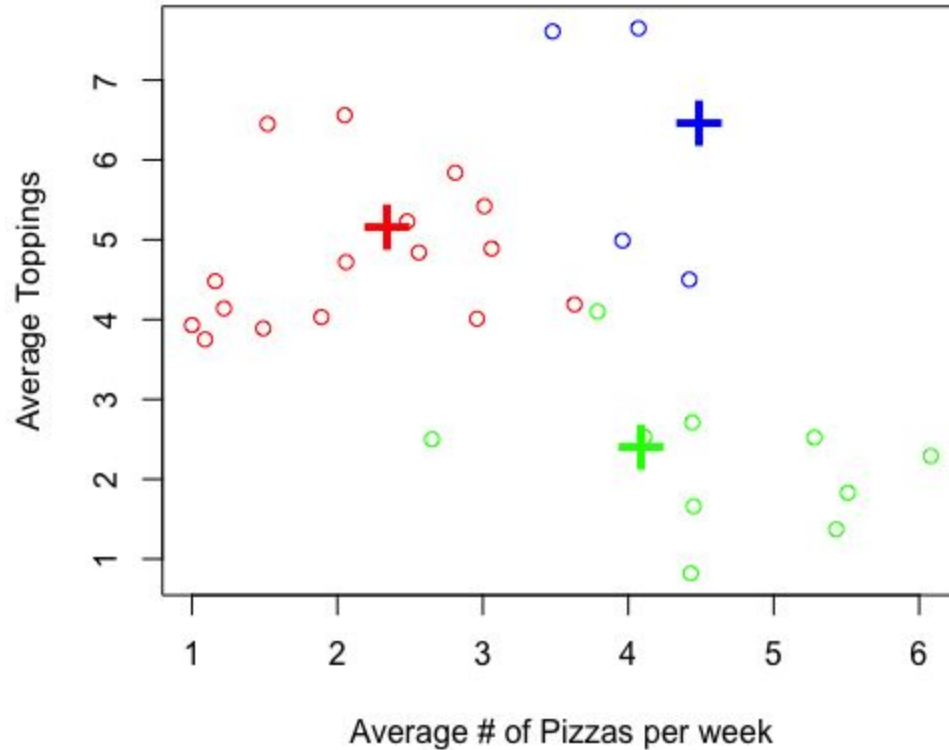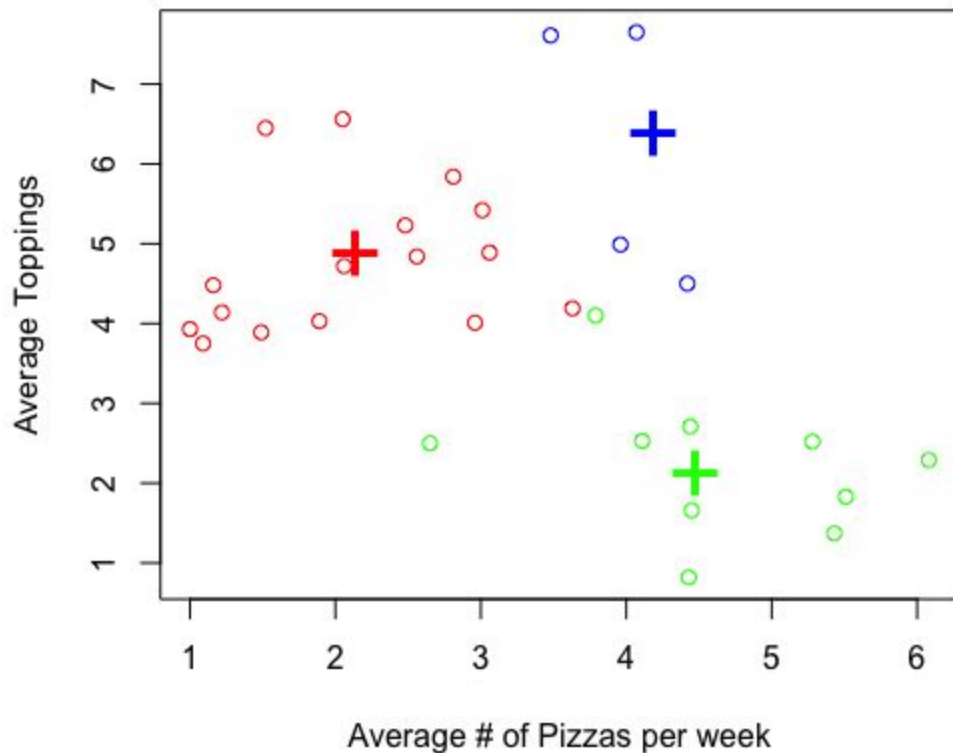SIMPLIFY

GROUP | REDUCE

@CHELSEAPARLETT

# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the centers
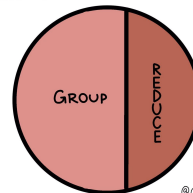
**3**

# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the centers

**2**

# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
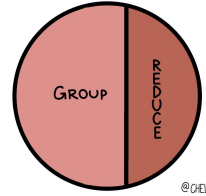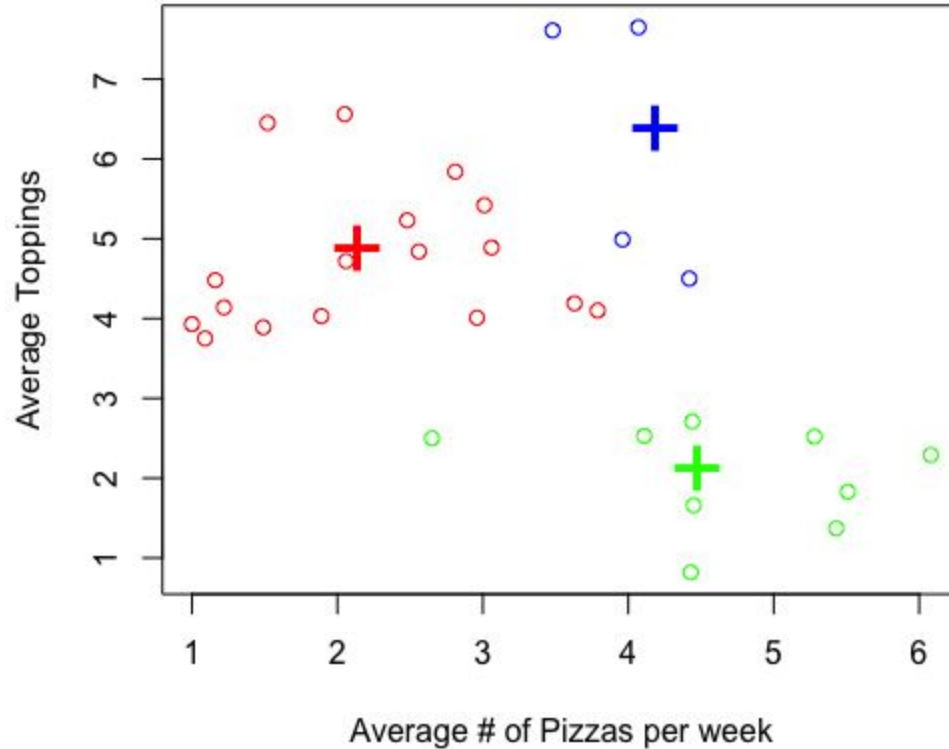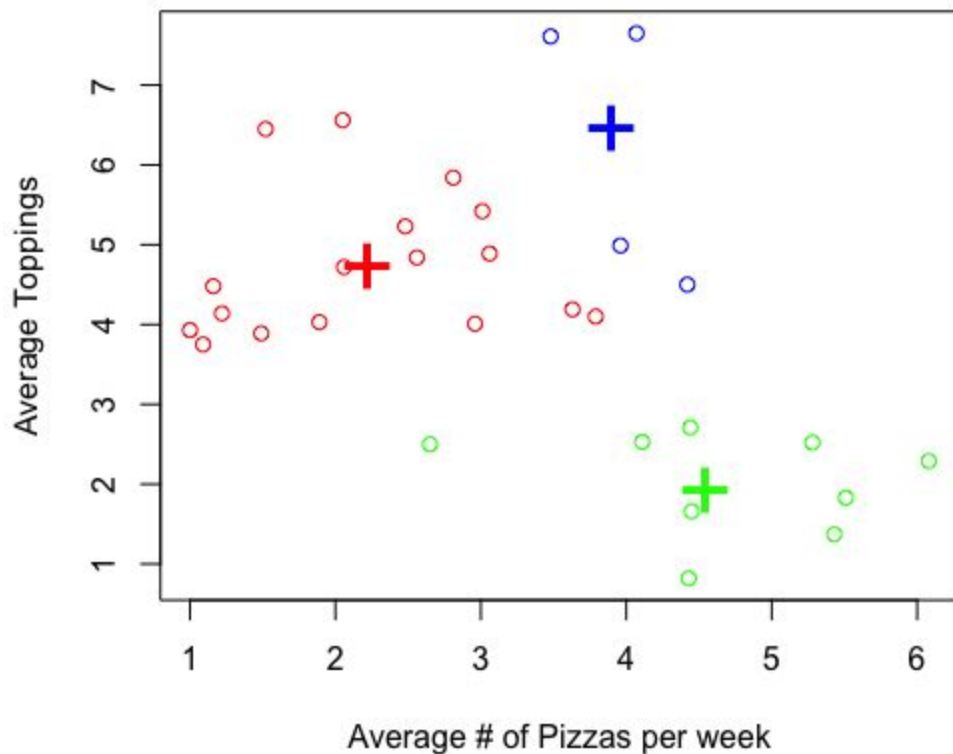3. Using these assignments, recalculate the centers

**SIMPLIFY**
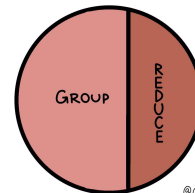
GROUP | REDUCE

@CHELSEAPARLETT

**3**

# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
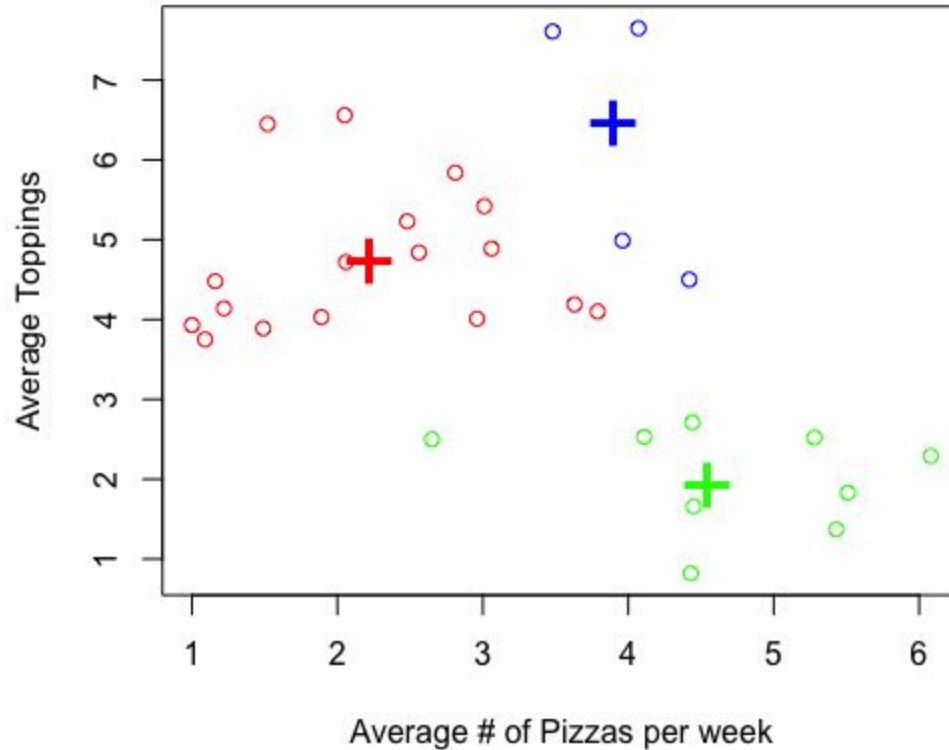3. Using these assignments, recalculate the centers

# K-Means

1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
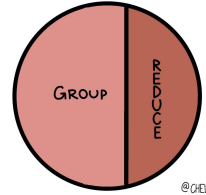3. Using these assignments, recalculate the centers

# K-Means
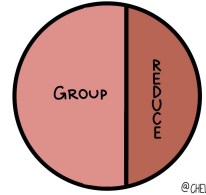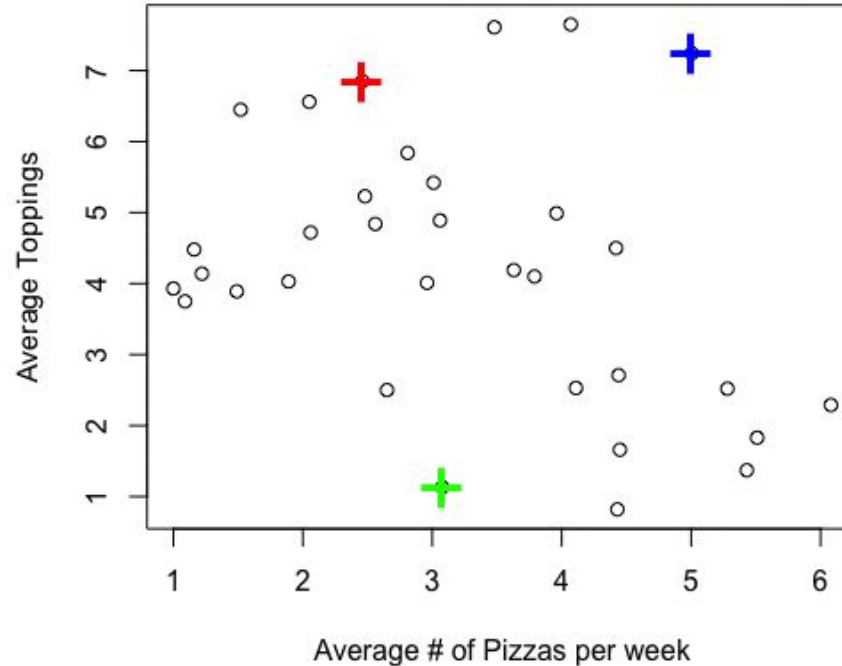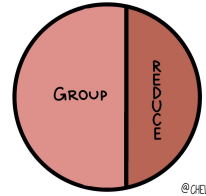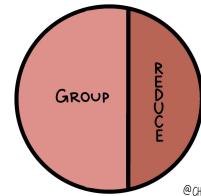
Assumptions

Spherical Clusters

Roughly the same # in each cluster

# Evaluating Unsupervised Models

Cohesion:

Separation:

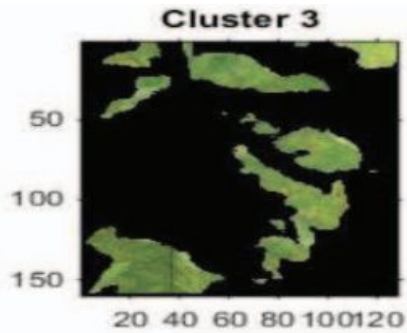$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

# Applications
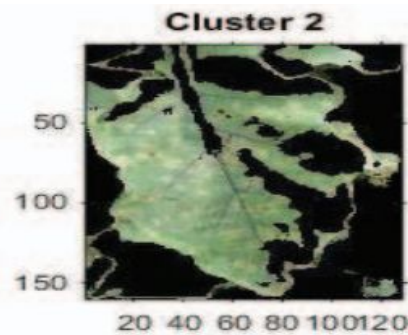

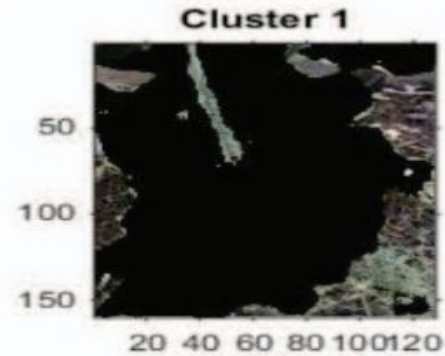
Fig. 8: Only Leaf



Fig. 7: Both Brinjal and Leaf



Fig. 6: Only Brinjal

# Applications
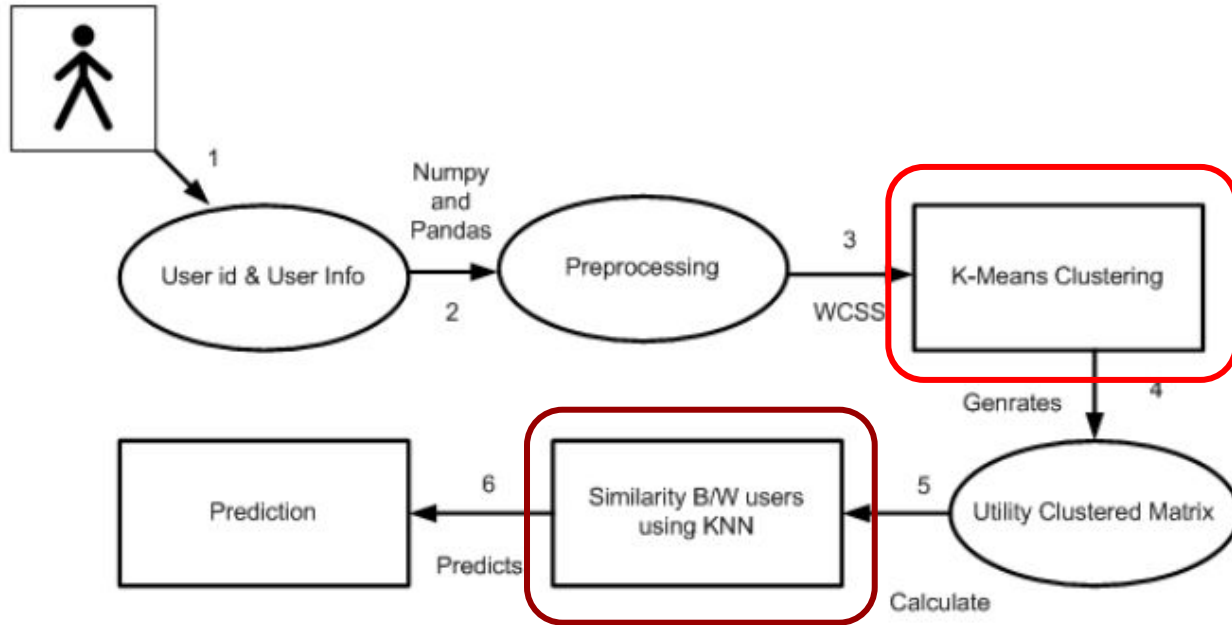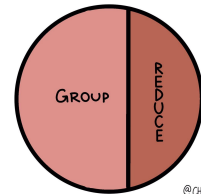


Fig. 2: Process Flow Diagram

# K-Means and Gaussian Mixtures

Dr. Chelsea Parlett-Pelleriti

# Normal (Gaussian) Distribution

The Normal Distribution

f(x)

μ

x

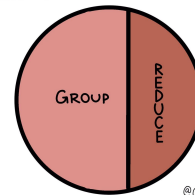$$y = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\mu =$ Mean
$\sigma =$ Standard Deviation
$\pi \approx 3.14159\cdots$
$e \approx 2.71828\cdots$

# GMM



$ Spent on Fast Food per Week

# GMM

$ Spent on Fast Food per Week

# Multivariate Normal Distributions



(a) 3 components Gaussian mixture density

(b) Data from 3 components Gaussian mixture density

# GMM

### K means

- Hard Assignment
- All Variances the Same

### GMM

- Soft (probabilistic) Assignment
- Variances can be different

GMM

# K-Means Review

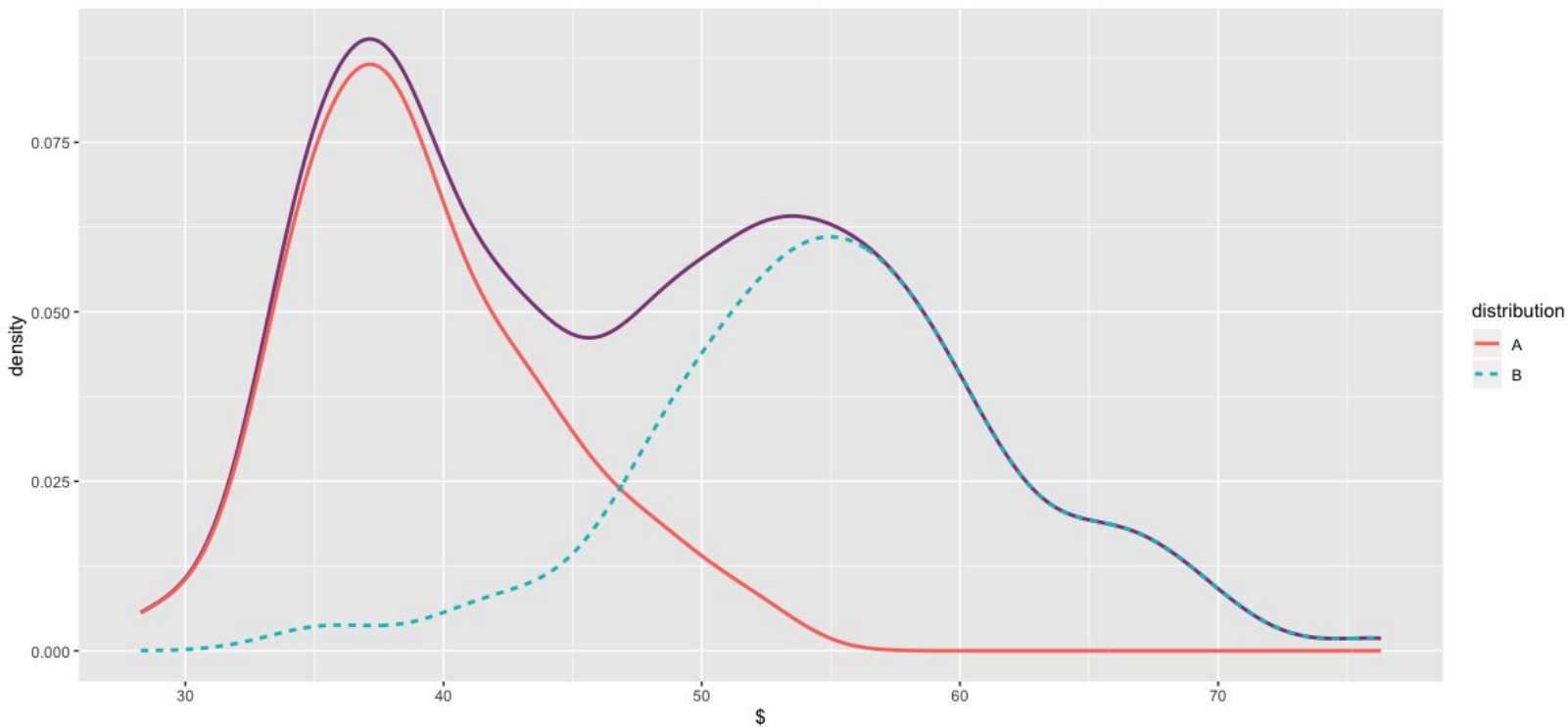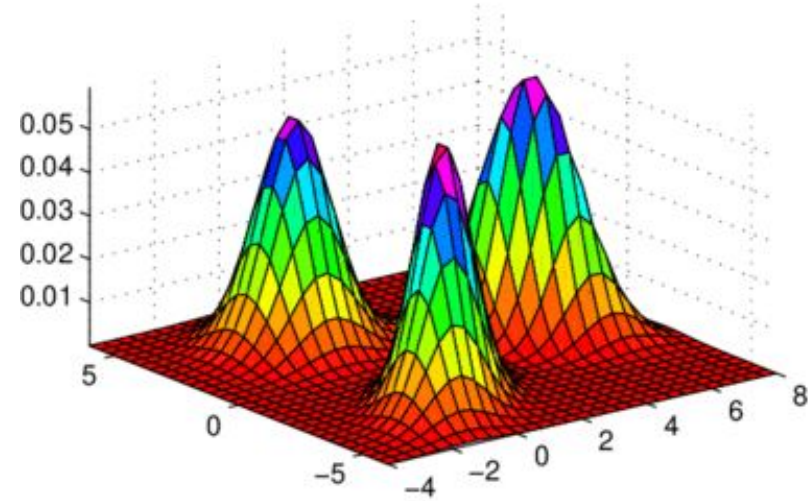1. Choose **k** random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the **centers**
4. Repeat 2 and 3 until either:
   a. Cluster membership does not change
   b. Centers change only a tiny amount

# GMM

1. Choose **k** random points to be cluster centers (**or estimate using k-means...etc)**
2. For each data point, calculate the **probability** of belonging to each cluster
3. Using these probability weights, recalculate the **means + variances**
4. Repeat 2 and 3 until **distributions converge.**

# Formulas (E-Step)

$$p_k(\underline{x}|\theta_k) = \frac{1}{(2\pi)^{d/2}|\Sigma_k|^{1/2}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu}_k)^t \Sigma_k^{-1}(\underline{x}-\underline{\mu}_k)}$$

# Formulas (M-Step)

# Applications

# Applications



**Normalized vectors**
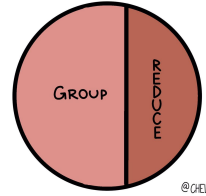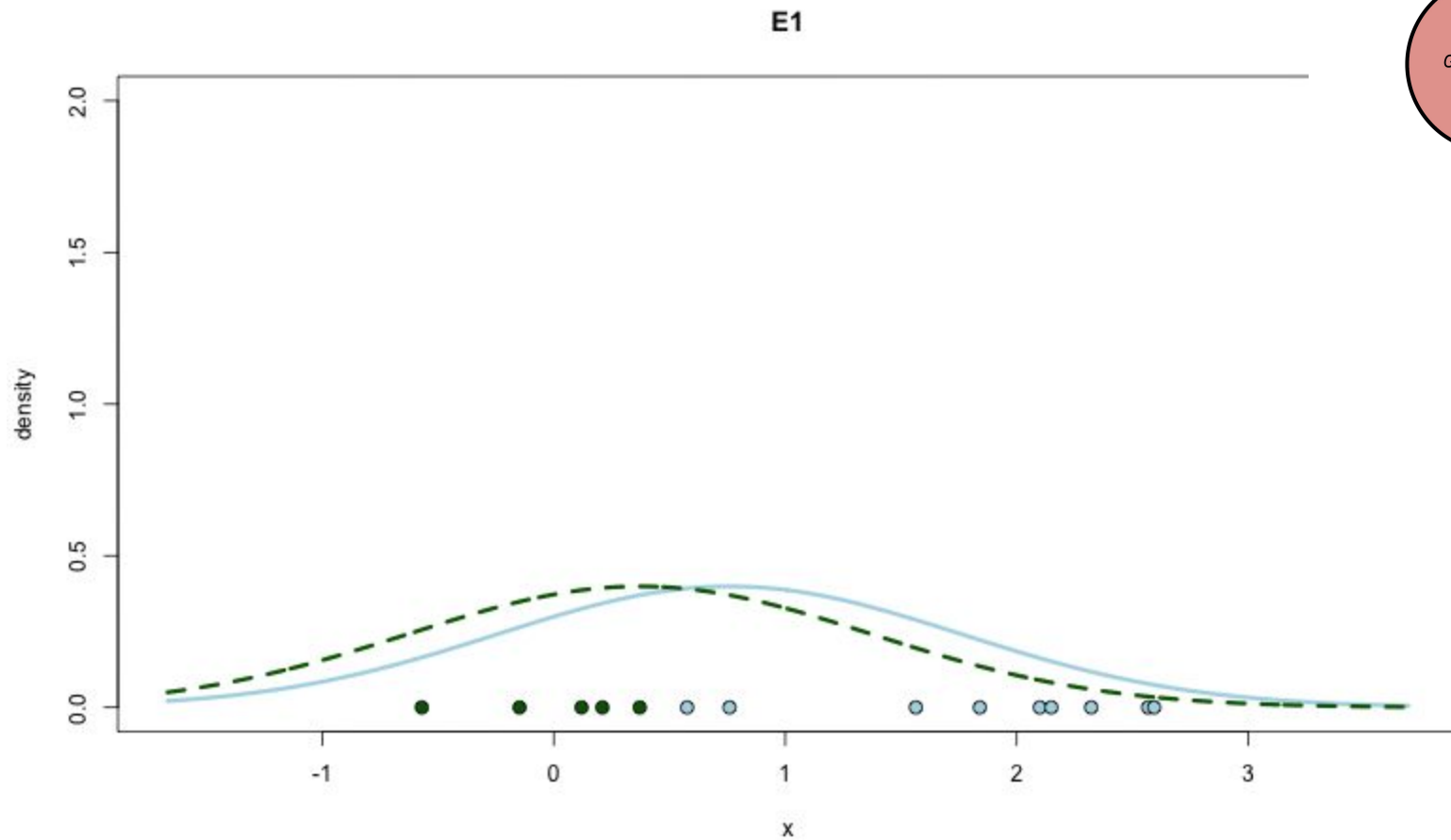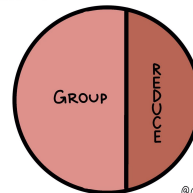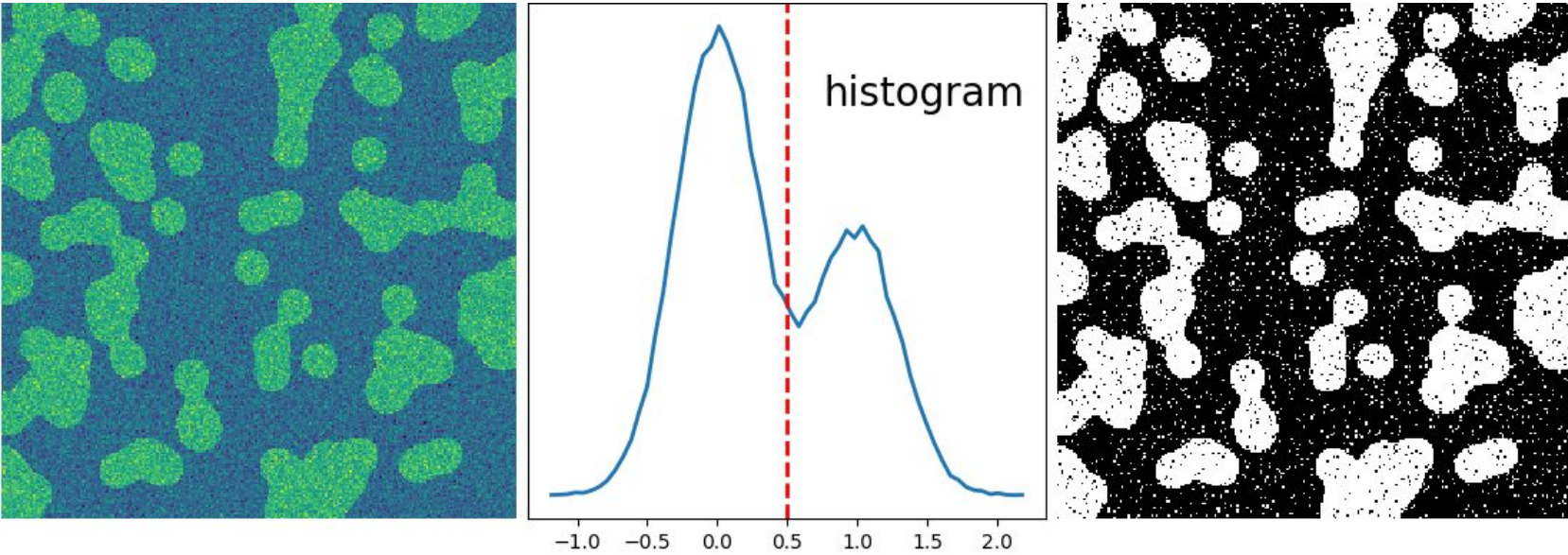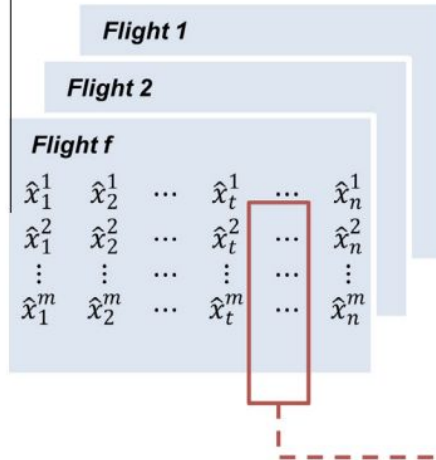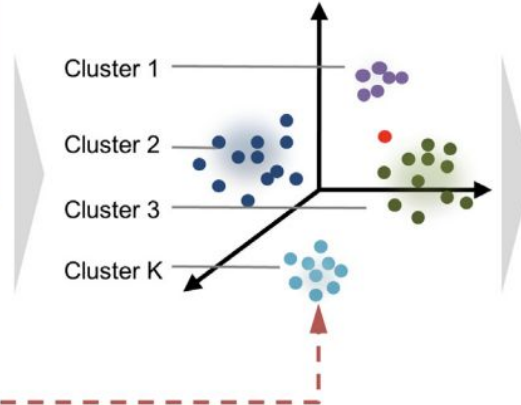Every flight parameter is normalized to have "zero mean and unit variance"

**Clusters**
GMM clustering is performed on normalized vectors; each cluster represent a typical operation of aircraft

**Temporal distribution of clusters**
The temporal distribution of clusters is summarized by observation frequency of each cluster along the temporal reference

*Flight 1*

*Flight 2*

*Flight f*

$$\begin{matrix} \hat{x}_1^1 & \hat{x}_2^1 & \cdots & \hat{x}_t^1 & \cdots & \hat{x}_n^1 \\ \hat{x}_1^2 & \hat{x}_2^2 & \cdots & \hat{x}_t^2 & \cdots & \hat{x}_n^2 \\ \vdots & \vdots & \cdots & \vdots & \cdots & \vdots \\ \hat{x}_1^m & \hat{x}_2^m & \cdots & \hat{x}_t^m & \cdots & \hat{x}_n^m \end{matrix}$$

Cluster 1
Cluster 2
Cluster 3
Cluster K

Cluster index

Distance before touchdown (nm)

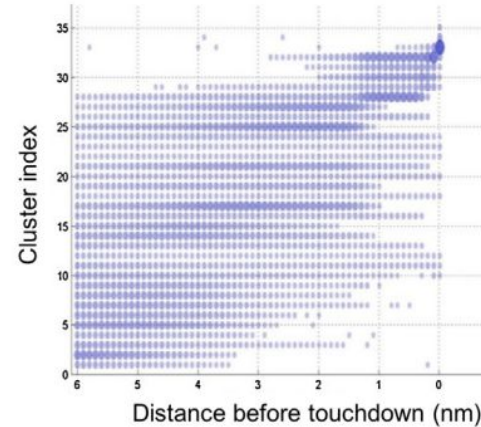Larger circle size and darker color indicates a higher observation frequency

**Fig. 3.** Cluster analysis: identify typical operations and temporal distribution.
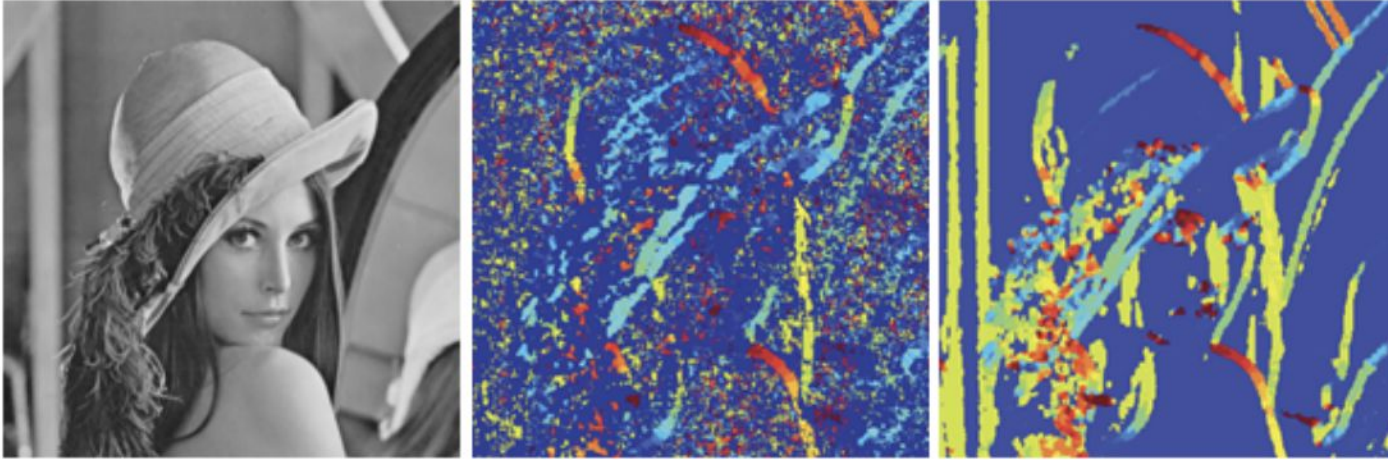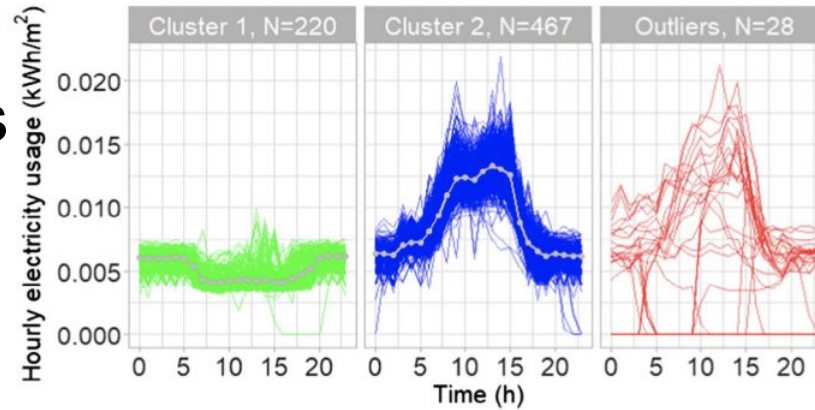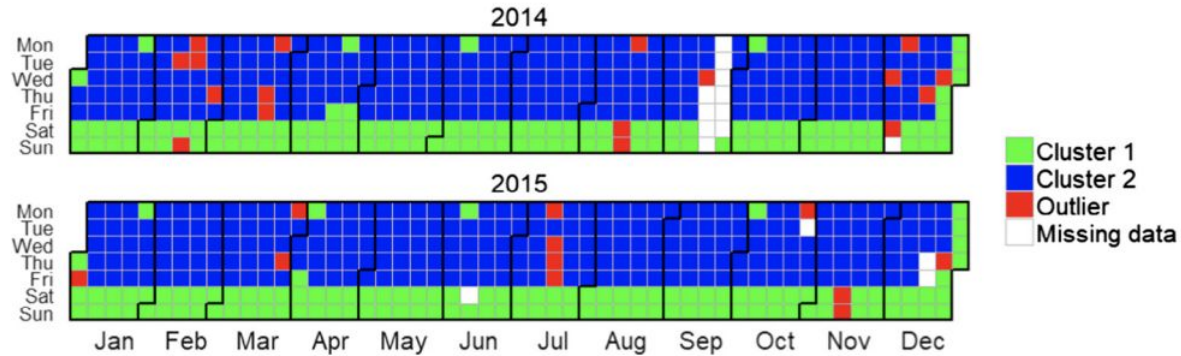
# Applications



Fig. 1. Illustration of clustering of patches in the PLE method for the Lena image. LEFT: Original image; RIGHT: Clustered image; The pixels in the same color indicate that $8 \times 8$ patches around them are in the same cluster. It can be seen that patches from different parts of image are grouped into one cluster [17].

# Applications



a) TDEU profiles and outliers



b) Distribution of the TDEU profiles

**Fig. 7.** Visualisation of the intra-building clustering result of Building #16.