

# OLAC

-

## Label Buying Decision Engine

### Problem setting:

Compute the optimal allocation of labels acquisitions that maximize, under resource constraints, long run utility.

### Definitions:

- Let  $\mathbf{X} \in \mathbb{R}^{nm}$  be the observed data.
- Let  $\mathbf{f} \in \mathbb{R}^m$  be the vector of features.
- Let  $\mathbf{H} \in \{0, 1\}^{nk_t}$  be the cluster mask.
- Let  $\mathbf{y} \in \{0, 1\}^n$  be the vector of true labels.
- Let  $\theta \in \{0, 1\}^n$  be the vector of acquired labels.
- Let  $\mathcal{K}_t$  be the set of clusters of size  $k_t$  at time  $t$ .
- Let  $\mathbf{k} \in \{1, 2, \dots, n\}^t$  be the vector of optimal number of clusters.
- Let  $\mathbf{u} \in \mathbb{R}^n \sim LN(\mu, \sigma)$  the vector of utility earned by detecting fraud drawn from a log-normal distribution.

### Information K-means

For each block of time  $t$ ,  $t \in \mathbb{Z}$ , the newly observed data is compared to the historical data and evaluate whether the number of clusters needs to be adjusted.

Let  $k_t = g(\mathbf{f})$ ; where  $g$  is a function that approximates the optimal number of clusters. Let  $\mathcal{K}_t = f(k_t, \mathbf{X})$  where  $f$  is the K-means algorithm that approximates the optimal centroid location and assignment of data points.

*cluster assignment:*

$$K_\eta^{(q)} = \{x_i : \|x_i - m_\eta\|^2 \leq \|x_i - m_j\|^2 \forall j, 1 \leq j \leq k\} \quad (1)$$

*centroid update:*

$$m_\eta^{(q+1)} = \frac{1}{|K_\eta^{(q)}|} \sum_{x_i \in K_\eta^{(q)}} x_i \quad (2)$$

## Utility function

Utility consists of the cost of buying a label,  $c \in \mathbb{R}_+$ , the revenue gained by detecting fraud,  $u \in \mathbb{R}_+$  with the assumption that the utility gained from determining fraud is greater than the cost,  $c < u$ . Additional elements can be added such as cost of true negatives, missing fraud, or additional costs for investigating false positives.

$$U(\mathbf{u}, \mathbf{y}, \theta, c) = \mathbf{u} \times (\mathbf{y} \times \theta) - c\theta = \sum_{i=0}^n u_i y_i \theta_i - c\theta_i \quad (3)$$

## Revenue

In order to determine the monetary gain of buying a label in a particular cluster we calculate the expected revenue of allocating all labels to a particular cluster. The matrix  $\mathbf{H} \in \{0, 1\}^{n \times k}$  contains the masks for each cluster,  $K_\eta \subset \mathcal{K}$  where, i.e.:

$$h_{i,\eta} = \begin{cases} 1, & \text{if } \mathbf{x}_i \in K_\eta \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

where  $i \in \{0, \dots, n-1\}$ ,  $\eta \in \{1, \dots, k\}$ .

The expected revenue for cluster  $\eta$  is defined as:

$$\pi_\eta = \mathbf{h}_\eta \cdot \mathbb{E}[U(\mathbf{u}, \mathbf{y}, \theta, c)] = \mathbf{h}_\eta \cdot (\mathbb{E}[\mathbf{u}] \times (\mathbf{y} \times \theta) - c\theta) \quad (5)$$

## ADW-MAB

Rough idea:

1. Determine the number of clusters (Information K-means, ...)
2. Determine the windows size dependent of measure of change (Adaptive windowing - ADWIN, ...)
3. Determine optimal allocation label (Multi-armed Bandit, ...)