

Supplemental material

K-means clustering

K-means clustering can be stated as finding a set of K cluster centers $\{\mu_i\}_{i=1}^K$ that for a dataset of N examples \mathbf{x}_j satisfies,

$$\arg \min_{\{\mu_i\}_{i=1}^K} \sum_{j=1}^N \min_{i \in \{1, \dots, K\}} \|\mu_i - \mathbf{x}_j\|^2.$$

Subgroup discovery

Algorithm 1 Subgroup discovery procedure

- 1: N - number of examples.
 - 2: $\alpha \in (0, 1)$ - size as percentage of total, for subgroup to have significant size.
 - 3: $\beta > 1$ - size of effect for interesting subgroup.
 - 4: N_x - number of clusters for modality X .
 - 5: N_y - number of clusters for modality Y .
 - 6: Perform K-Means on modality X with clusters C_x .
 - 7: Perform K-Means on modality Y with clusters C_y .
 - 8: **for** $c_x \in C_x$ **do**
 - 9: **for** $c_y \in C_y$ **do**
 - 10: **if** reject hypothesis c_x is independent to c_y **then**
 - 11: Classify subgroup of $c_x \wedge c_y$ as interesting.
 - 12: **if** Number of examples in cluster c_x and $c_y > \alpha N$ **then**
 - 13: Classify subgroup as significant in size.
 - 14: **end if**
 - 15: **if** $N|c_x \wedge c_y|/|c_x||c_y| > \beta$ **then**
 - 16: Classify subgroup effect as significant.
 - 17: **end if**
 - 18: **end if**
 - 19: **end for**
 - 20: **end for**
-

Gini impurity

Given data containing N examples each labeled with 1 of J categorical labels, where N_i of the examples are labeled with category i , the Gini impurity is defined as,

$$I_G = \sum_{i=1}^J \frac{N_i}{N} \left(1 - \frac{N_i}{N}\right). \quad (2)$$

Gaussian process regression

A Gaussian process is a collection of random variables where the joint distribution of any finite subset of the collection is a multivariate Gaussian distribution (Rasmussen and Williams 2006). A Gaussian process is defined by a mean function, $m(\mathbf{x})$, and covariance function, $k(\mathbf{x}, \mathbf{x}')$ and can be thought of as a prior distribution over functions. For simplicity it is often assumed that $m(\mathbf{x}) = 0$. Given this prior knowledge, and observations X with labels \mathbf{y} , the rules of probability can be applied to provide an *a posteriori* prediction of a label y_* for a new observation \mathbf{x}_* (a "posterior"). Assuming Gaussian noise on the observations and a Gaussian process prior on the function to be learned, the posterior is also given by a Gaussian process. Assuming Gaussian noise, the

posterior is given by,

$$P(y_* | \mathbf{x}_*, X, \mathbf{y}) = \mathcal{N}(\mu_*, \Sigma_*), \quad (3)$$

where $\mathcal{N}()$ is a multivariate Gaussian distribution with mean μ_* , and covariance Σ_* , where,

$$\mu_* = K(\mathbf{x}_*, X) [K(X, X) + \sigma^2 I]^{-1} \mathbf{y}, \quad (4)$$

$$\Sigma_* = K(\mathbf{x}_*, \mathbf{x}_*) \quad (5)$$

$$- K(\mathbf{x}_*, X) [K(X, X) + \sigma^2 I]^{-1} K(X, \mathbf{x}_*). \quad (6)$$

Here $K(X, X')$ is the covariance matrix, where the entry $K_{i,j}$ corresponds to the covariance function applied to the i example in X and the j example in X' , and σ^2 is the measurement noise in the observations. This noise is assumed to be constant. The values of the diagonal of the covariance (the variance) of the posterior gives an indicator of how uncertain the prediction for that input is. The lower the variance the more confident the prediction.