# EXHIBIT 4
# PART

unnatural assumption if we are exploring an essentially unknown set of data. Thus, a constantly recurring problem in cluster analysis is that of deciding just how many clusters are present.

When clustering is done by extremizing a criterion function, a common approach is to repeat the clustering procedure for $c = 1$, $c = 2$, $c = 3$, etc., and to see how the criterion function changes with $c$. For example, it is clear that the sum-of-squared-error criterion $J_e$ must decrease monotonically with $c$, since the squared error can be reduced each time $c$ is increased merely by transferring a single sample to the new cluster. If the $n$ samples are really grouped into $\hat{c}$ compact, well separated clusters, one would expect to see $J_e$ decrease rapidly until $c = \hat{c}$, decreasing much more slowly thereafter until it reaches zero at $c = n$. Similar arguments have been advanced for hierarchical clustering procedures, the usual assumption being that large disparities in the levels at which clusters merge indicate the presence of natural groupings.

A more formal approach to this problem is to devise some measure of goodness of fit that expresses how well a given $c$-cluster description matches the data. The chi-square and Kolmogorov-Smirnov statistics are the traditional measures of goodness of fit, but the curse of dimensionality usually demands the use of simpler measures, such as a criterion function $J(c)$. Since we expect a description in terms of $c + 1$ clusters to give a better fit than a description in terms of $c$ clusters, we would like to know what constitutes a statistically significant improvement in $J(c)$.

A formal way to proceed is to advance the *null hypothesis* that there are exactly $c$ clusters present, and to compute the sampling distribution for $J(c + 1)$ under this hypothesis. This distribution tells us what kind of apparent improvement to expect when a $c$-cluster description is actually correct. The decision procedure would be to accept the null hypothesis if the observed value of $J(c + 1)$ falls within limits corresponding to an acceptable probability of false rejection.

Unfortunately, it is usually very difficult to do anything more than crudely estimate the sampling distribution of $J(c + 1)$. The resulting solutions are not above suspicion, and the statistical problem of testing cluster validity is still essentially unsolved. However, under the assumption that a suspicious test is better than none, we include the following approximate analysis for the simple sum-of-squared-error criterion.

Suppose that we have a set $\mathcal{X}$ of $n$ samples and we want to decide whether or not there is any justification for assuming that they form more than one cluster. Let us advance the null hypothesis that all $n$ samples come from a normal population with mean $\mu$ and covariance matrix $\sigma^2 I$. If this hypothesis were true, any clusters found would have to have been formed by chance, and any observed decrease in the sum of squared error obtained by clustering would have no significance.

242      UNSUPERVISED  LEARNING  AND  CLUSTERING

The sum of squared error $J_e(1)$ is a random variable, since it depends on the particular set of samples:

$$J_e(1) = \sum_{x \in \mathcal{X}} \|x - m\|^2,$$

where $m$ is the mean of the $n$ samples. Under the null hypothesis, the distribution for $J_e(1)$ is approximately normal with mean $nd\sigma^2$ and variance $2nd\sigma^4$.

Suppose now that we partition the set of samples into two subsets $\mathcal{X}_1$ and $\mathcal{X}_2$ so as to minimize $J_e(2)$, where

$$J_e(2) = \sum_{i=1}^{2} \sum_{x \in \mathcal{X}_i} \|x - m_i\|^2,$$

$m_i$ being the mean of the samples in $\mathcal{X}_i$. Under the null hypothesis, this partitioning is spurious, but it nevertheless results in a value for $J_e(2)$ that is smaller than $J_e(1)$. If we knew the sampling distribution for $J_e(2)$, we could determine how small $J_e(2)$ would have to be before we were forced to abandon a one-cluster null hypothesis. Lacking an analytical solution for the optimal partitioning, we cannot derive an exact solution for the sampling distribution. However, we can obtain a rough estimate by considering the suboptimal partition provided by a hyperplane through the sample mean. For large $n$, it can be shown that the sum of squared error for this partition is approximately normal with mean $n(d - 2/\pi)\sigma^2$ and variance $2n(d - 8/\pi^2)\sigma^4$.

This result agrees with our statement that $J_e(2)$ is smaller than $J_e(1)$, since the mean of $J_e(2)$ for the suboptimal partition—$n(d - 2/\pi)\sigma^2$—is less than the mean for $J_e(1)$—$nd\sigma^2$. To be considered significant, the reduction in the sum of squared error must certainly be greater than this. We can obtain an approximate critical value for $J_e(2)$ by assuming that the suboptimal partition is nearly optimal, by using the normal approximation for the sampling distribution, and by estimating $\sigma^2$ by

$$\hat{\sigma}^2 = \frac{1}{nd} \sum_{x \in \mathcal{X}} \|x - m\|^2 = \frac{1}{nd} J_e(1).$$

The final result can be stated as follows: Reject the null hypothesis at the $p$-percent significance level if

$$\frac{J_e(2)}{J_e(1)} < 1 - \frac{2}{\pi d} - \alpha \sqrt{\frac{2(1 - 8/\pi^2 d)}{nd}}, \tag{44}$$

where $\alpha$ is determined by

$$p = 100 \int_{\alpha}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-1/2 u^2} \, du.$$

Thus, this provides us with a test for deciding whether or not the splitting of a cluster is justified. Clearly, the $c$-cluster problem can be treated by applying the same test to all clusters found.

## 6.13 LOW-DIMENSIONAL REPRESENTATIONS AND MULTIDIMENSIONAL SCALING

Part of the problem of deciding whether or not a given clustering means anything stems from our inability to visualize the structure of multidimensional data. This problem is further aggravated when similarity or dissimilarity measures are used that lack the familiar properties of distance. One way to attack this problem is to try to represent the data points as points in some lower-dimensional space in such a way that the distances between points in the lower-dimensional space correspond to the dissimilarities between points in the original space. If acceptably accurate representations can be found in two or perhaps three dimensions, this can be an extremely valuable way to gain insight into the structure of the data. The general process of finding a configuration of points whose interpoint distances correspond to dissimilarities is often called *multidimensional scaling*.

Let us begin with the simpler case where it is meaningful to talk about the distances between the $n$ samples $x_1, \ldots, x_n$. Let $y_i$ be the lower-dimensional *image* of $x_i$, $\delta_{ij}$ be the distance between $x_i$ and $x_j$, and $d_{ij}$ be the distance between $y_i$ and $y_j$. Then we are looking for a *configuration* of image points $y_1, \ldots, y_n$ for which the $n(n-1)/2$ distances $d_{ij}$ between image points are as close as possible to the corresponding original distances $\delta_{ij}$. Since it will usually not be possible to find a configuration for which $d_{ij} = \delta_{ij}$ for all $i$ and $j$, we need some criterion for deciding whether or not one configuration is better than another. The following sum-of-squared-error functions are all reasonable candidates:

$$J_{ee} = \frac{1}{\sum_{i<j} \delta_{ij}^2} \sum_{i<j} (d_{ij} - \delta_{ij})^2 \tag{45}$$

$$J_{ff} = \sum_{i<j} \left( \frac{d_{ij} - \delta_{ij}}{\delta_{ij}} \right)^2 \tag{46}$$

$$J_{ef} = \frac{1}{\sum_{i<j} \delta_{ij}} \sum_{i<j} \frac{(d_{ij} - \delta_{ij})^2}{\delta_{ij}}. \tag{47}$$

Since these criterion functions involve only the distances between points, they are invariant to rigid-body motions of the configurations. Moreover,

244    UNSUPERVISED LEARNING AND CLUSTERING

they have all been normalized so that their minimum values are invariant to dilations of the sample points. $J_{ee}$ emphasizes the largest errors, regardless whether the distances $\delta_{ij}$ are large or small. $J_{ff}$ emphasizes the largest fractional errors, regardless whether the errors $|d_{ij} - \delta_{ij}|$ are large or small. $J_{ef}$ is a useful compromise, emphasizing the largest product of error and fractional error.

Once a criterion function has been selected, an optimal configuration $y_1, \ldots, y_n$ is defined as one that minimizes that criterion function. An optimal configuration can be sought by a standard gradient-descent procedure, starting with some initial configuration and changing the $y_i$'s in the direction of greatest rate of decrease in the criterion function. Since

$$d_{ij} = \|\mathbf{y}_i - \mathbf{y}_j\|,$$

the gradient of $d_{ij}$ with respect to $y_i$ is merely a unit vector in the direction of $y_i - y_j$. Thus, the gradients of the criterion functions are easy to compute:*

$$\nabla_{\mathbf{y}_k} J_{ee} = \frac{2}{\sum\limits_{i<j} \delta_{ij}^2} \sum\limits_{j \neq k} (d_{kj} - \delta_{kj}) \frac{\mathbf{y}_k - \mathbf{y}_j}{d_{kj}}$$

$$\nabla_{\mathbf{y}_k} J_{ff} = 2 \sum\limits_{j \neq k} \frac{d_{kj} - \delta_{kj}}{\delta_{kj}^2} \frac{\mathbf{y}_k - \mathbf{y}_j}{d_{kj}}$$

$$\nabla_{\mathbf{y}_k} J_{ef} = \frac{2}{\sum\limits_{i<j} \delta_{ij}} \sum\limits_{j \neq k} \frac{d_{kj} - \delta_{kj}}{\delta_{kj}} \frac{\mathbf{y}_k - \mathbf{y}_j}{d_{kj}}.$$

The starting configuration can be chosen randomly, or in any convenient way that spreads the image points about. If the image points lie in a $\hat{d}$-dimensional space, then a simple and effective starting configuration can be found by selecting those $\hat{d}$ coordinates of the samples that have the largest variance.

The following example illustrates the kind of results than can be obtained by these techniques.† The data consist of thirty points spaced at unit intervals along a three-dimensional helix:

$$x_1(k) = \cos x_3$$
$$x_2(k) = \sin x_3$$
$$x_3(k) = k/\sqrt{2}, \qquad k = 0, 1, \ldots, 29.$$

---

* Second partial derivatives can also be computed easily, so that Newton's algorithm can be used. Note that if $\mathbf{y}_i = \mathbf{y}_j$, the unit vector from $\mathbf{y}_i$ to $\mathbf{y}_j$ is undefined. Should that situation arise, $(\mathbf{y}_i - \mathbf{y}_j)/d_{ij}$ can be replaced by an arbitrary unit vector.

† This example was taken from J. W. Sammon, Jr., "A nonlinear mapping for data structure analysis," *IEEE Trans. Comp.*, **C-18**, 401–409 (May 1969).

# DOCKET ALARM

# Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

## Real-Time Litigation Alerts

Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

## Advanced Docket Research

With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

## Analytics At Your Fingertips

Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

## API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

### LAW FIRMS

Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

### FINANCIAL INSTITUTIONS

Litigation and bankruptcy checks for companies and debtors.

### E-DISCOVERY AND LEGAL VENDORS

Sync your system to PACER to automate legal marketing.