

Okay, here is a study note based on the PDF content, including key concepts in both Chinese and English, along with explanations in Chinese:

## Machine Learning: Unsupervised Learning - Density Estimation (无监督学习：密度估计)

### 1. Introduction (简介)

- This lecture covers density estimation, a type of unsupervised machine learning.
  - The lecture is part of the CBU5201 Machine Learning course at Queen Mary University of London.
- **Bell's bit of scruff (贝尔的“一点杂乱”)**
  - The observation of a pulsar (CP 1919) signal led to a famous quote, "I logged it with a question mark and moved on," highlighting the importance of re-examining data. (重新检查你的数据！)

### 2. Machine Learning Taxonomy (机器学习分类法)

- **Machine Learning (机器学习)** is broadly divided into:
  - **Supervised Learning (监督学习)**: Uses labeled data to train models.
    - **Classification (分类)**: Predicts categorical labels (e.g., spam/not spam).
    - **Regression (回归)**: Predicts continuous values (e.g., house prices).
  - **Unsupervised Learning (无监督学习)**: Uses unlabeled data to find patterns and structures.
    - **Density Estimation (密度估计)**: Models the probability distribution of data.
    - **Structure Analysis (结构分析)**: Discovers relationships and patterns in data.

### 3. Probability Densities (概率密度)

- **Data Distribution (数据分布)**: Datasets are collections of samples distributed in the attribute space (属性空间), but they don't occupy the entire space. 数据集是在属性空间中分布的样本集合，但它们并不占据整个空间。
- **Key Questions (关键问题)**:
  - Where are most of my samples? 我的样本主要分布在哪里？
  - Should I expect a sample in this region? 我是否应该期望在这个区域找到样本？
  - What is the probability of finding a sample here? 在这个区域找到样本的概率是多少？
  - Given a probability, where will I find my next sample? 给定一个概率，我会在哪里找到下一个样本？
- **Divide and Count (分割与计数)**: Partitioning the space and counting samples in each region is a simple way to describe the observed distribution (观察分布). 将空间划分为多个区域，并计算每个区域中的样本数量，是描述观察分布的一种简单方法。
- **Probability Densities (概率密度)**: Models that describe the underlying true distribution (真实分布) of data. 概率密度是描述数据潜在真实分布的模型。
  - They quantify the probability of a sample being in a region. 它们量化了样本在某个区域内的概率。
  - "The probability of finding a sample anywhere in the attribute space is 1. This probability is not uniform: it is denser in some regions than others." (在属性空间中任意位置找到样本的概率为1。这个概率不是均匀的：它在某些区域比其他区域更密集。)
- **Density Estimation (密度估计)**: The task of building probability densities from data. 从数据中构建概率密度的任务。
  - "In machine learning we use data to build probability densities. We call this task density estimation." (在机器学习中，我们使用数据来构建概率密度。我们将此任务称为密度估计。)
  - **Marginal Probability Densities (边缘概率密度)**: Densities considering a subset of attributes. 考虑属性子集的密度。
  - Notation:  $p(x_1, x_2)$  or  $p(x)$  for joint density,  $p(x_1)$  and  $p(x_2)$  for marginal densities. 联合密度的符号表示为  $p(x_1, x_2)$  或  $p(x)$ ，边缘密度的符号表示为  $p(x_1)$  和  $p(x_2)$ 。

### 4. Non-parametric Methods (非参数方法)

- "Non-parametric methods for density estimation do not assume any specific shape for the probability density." (用于密度估计的非参数方法不假定概率密度的任何特定形状。) 非参数方法不预先假设概率密度的具体形状。
- **Histogram (直方图)**:
  - "The histogram is the simplest and best known non-parametric method for density estimation." (直方图是用于密度估计的最简单且最著名的非参数方法。)
  - Divides the feature space into equal-sized bins (等宽的箱子). 将特征空间划分为大小相等的箱子。
  - Density is approximated by the fraction of samples in each bin. 密度由每个箱子中样本的比例来近似。
  - **Bin Size (箱子大小)**: Small bins lead to spiky (尖峰) estimates, large bins lead to flat (平坦) estimates and loss of structure. 小的箱子会导致尖峰估计，大的箱子会导致平坦估计和结构损失。
- **Kernel Methods (核方法)**:
  - "Kernel methods proceed by building an individual density around each sample first and then combining all the densities together." (核方法首先围绕每个样本构建一个单独的密度，然后将所有密度组合在一起。)

- Build individual densities (e.g., Gaussian) around each sample and combine them. 围绕每个样本构建单独的密度（例如高斯密度）并将它们组合起来。
- Individual densities have the same shape (the kernel). 每个单独的密度都具有相同的形状（核）。
- **Curse of Dimensionality (维度灾难):** Histograms can have one sample per bin; kernel densities can become isolated in high dimensions. 直方图可能每个箱子只有一个样本；核密度在高维度中可能会变得孤立。

## 5. Parametric Density Estimation (参数密度估计)

- "Parametric approaches specify the shape of the probability density. The problem of density estimation consists of estimating its parameters." (参数方法指定概率密度的形状。密度估计的问题在于估计其参数。) 参数方法预先指定概率密度的形状，密度估计问题就变成了估计这些形状的参数。
- **Common Models (常见模型):**
  - **Gaussian Distribution (高斯分布/正态分布):**  $N(\mu, \Sigma)$  (denoted by  $N(\mu, \Sigma)$ , 通常用  $N(\mu, \Sigma)$  表示)
  - Log-normal Distribution (对数正态分布)
  - Uniform Distribution (均匀分布)
  - Gamma Distribution (伽马分布)
- **Gaussian Distribution (高斯分布):**
  - Defined by mean (均值)  $\mu$  and standard deviation (标准差)  $\sigma$  (or variance (方差)  $\sigma^2$ ) in 1D. 一维情况下，由均值  $\mu$  和标准差  $\sigma$  (或方差  $\sigma^2$ ) 定义。
  - " $\mu$  is known as the mean and  $\sigma$  is the standard deviation. The value  $\sigma^2$  is known as the variance." ( $\mu$  被称为均值， $\sigma$  是标准差。值  $\sigma^2$  被称为方差。)
  - Formula (1D):  $p(x_1) = (1 / (\sigma\sqrt{2\pi})) * e^{(-(x_1 - \mu)^2 / (2\sigma^2))}$
  - **Multivariate Gaussian (多变量高斯分布):** Extends to higher dimensions with mean vector  $\mu$  and covariance matrix ( $\Sigma$ ). 使用均值向量  $\mu$  和协方差矩阵  $\Sigma$  扩展到更高维度。
  - Formula:  $p(x) = (1 / ((2\pi)^{k/2} * |\Sigma|^{1/2})) * e^{(-(1/2)(x - \mu)^T \Sigma^{-1} (x - \mu))}$
  - **Independent Attributes (独立属性):** If attributes are independent, the covariance matrix is diagonal, and the joint density is the product of marginal densities. 如果属性是独立的，则协方差矩阵是对角的，联合密度是边缘密度的乘积。
- **Central Limit Theorem (CLT, 中心极限定理):** The sum of many independent random variables tends to be Gaussian. 许多独立随机变量的总和趋向于高斯分布。
  - "The CLT is perhaps the main reason why the Gaussian distribution is one of our favourite density models." (中心极限定理可能是高斯分布成为我们最喜欢的密度模型之一的主要原因。)
  - Explains why noise (噪声) is often Gaussian. 解释了为什么噪声通常是高斯分布的。
- **Estimation (估计):** Parameters are estimated using maximum likelihood (最大似然法). 使用最大似然法估计参数。
  - 1D:  $\mu = (1/N)\sum x_i$ ,  $\sigma^2 = (1/N)\sum (x_i - \mu)^2$
  - Multivariate:  $\mu = (1/N)\sum x_i$ ,  $\Sigma = (1/N)\sum (x_i - \mu)(x_i - \mu)^T$
- **High Dimensionality (高维度):** Can lead to overfitting (过拟合). 可能导致过拟合。
- **Mixture Models (混合模型):** Used when data has multiple modes (clumps, 多个模态/团). 当数据具有多个模态（团）时使用。
  - **Gaussian Mixture Models (GMM, 高斯混合模型):** Combines multiple Gaussian densities. 组合多个高斯密度。
  - "A GMM probability density is formulated as a combination of Gaussian densities  $g_m(x)$  with their own mean  $\mu_m$  and covariance matrix  $\Sigma_m$ ." (GMM 概率密度表示为具有各自均值  $\mu_m$  和协方差矩阵  $\Sigma_m$  的高斯密度  $g_m(x)$  的组合。)
  - Formula:  $p(x) = \sum \pi_m * g_m(x)$  (where  $\pi_m$  are mixing coefficients, 其中  $\pi_m$  是混合系数)
  - **Expectation-Maximization (EM) algorithm (期望最大化算法):** Iterative process to fit GMMs, similar to K-means. 类似于 K 均值的迭代过程，用于拟合 GMM。

## 6. Applications (应用)

- **Anomaly Detection (异常检测):**
  - "The main idea behind an anomaly detection algorithm is to quantify the probability of observing samples some distance away from the general pattern." (异常检测算法背后的主要思想是量化观察到远离一般模式的样本的概率。)
  - Identify outliers (异常值) based on low probability. 根据低概率识别异常值。
  - "If  $p(x_i) < T$ ,  $x_i$  is an anomaly." (如果  $p(x_i) < T$ ，则  $x_i$  是一个异常值。)
- **Classification (分类):**
  - "Classifiers that apply Bayes rule turn posterior probabilities into priors and class densities." (应用贝叶斯规则的分类器将后验概率转换为先验概率和类密度。)
  - Estimate class densities  $p(x|C)$  for each class. 估计每个类的类密度  $p(x|C)$ 。
  - **Naive Bayes (朴素贝叶斯):** Assumes independent predictors, simplifying calculations. 假设预测变量是独立的，从而简化计算。

- **Clustering (聚类):**

- GMMs can be used for clustering. GMM 可用于聚类。
- "GMM can be used as a clustering method that produces ellipsoidal clusters." (GMM 可以用作生成椭圆簇的聚类方法。)

## 7. Summary (总结)

- "Probability densities are models that allow us to calculate the probability of finding a sample in a region of the attribute space." (概率密度是允许我们计算在属性空间的某个区域中找到样本的概率的模型。)
- "Non-parametric methods do not assume any particular shape for the probability density, whereas parametric methods do." (非参数方法不假定概率密度有任何特定的形状，而参数方法则假定。)
- "The Central Limit Theorem and other mathematical properties, make the Gaussian distribution one of the most popular choices." (中心极限定理和其他数学性质使高斯分布成为最受欢迎的选择之一。)
- "Probability densities can be used in many machine learning problems, such as anomaly detection, classification and clustering." (概率密度可用于许多机器学习问题，例如异常检测、分类和聚类。)

I hope this comprehensive study note is helpful for your understanding of the PDF content. Let me know if you have any other questions.