

1 Normal distribution

Univariate: A random variable ξ is said to have a normal distribution with mean μ and variance σ^2 if its probability density function (pdf) is given by

$$f_{\xi}(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\} \quad (1)$$

where μ is the mean and σ^2 is the variance of the distribution. More compactly, it can be written as

$$\xi \sim \mathcal{N}(\mu, \sigma^2) \quad (2)$$

Uncorrelated multivariate: A random vector $\boldsymbol{\xi} = \begin{pmatrix} \xi_1 \\ \vdots \\ \xi_k \end{pmatrix}$ is said to have an uncorrelated multivariate normal distribution with mean $\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_k \end{pmatrix}$ and variances $\sigma_1^2, \dots, \sigma_k^2$ if the pdf of every random component of $\boldsymbol{\xi}$ is given by

$$f_{\xi_j}(x) = \frac{1}{\sigma_j\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu_j}{\sigma_j}\right)^2\right\} \quad (3)$$

where μ_j is the mean and σ_j^2 is the variance of the j -th component.

All components of $\boldsymbol{\xi}$ are assumed to be independent, so the joint pdf of $\boldsymbol{\xi}$ is the product of the pdfs of its components:

$$\begin{aligned} f_{\boldsymbol{\xi}}(x_1, \dots, x_k) &= \prod_{i=1}^k f_{\xi_i}(x_i) \\ &= \prod_{i=1}^k \frac{1}{\sigma_i\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x_i-\mu_i}{\sigma_i}\right)^2\right\} \end{aligned} \quad (4)$$

Covariance matrix: All variance parameters $\sigma_1^2, \dots, \sigma_k^2$ can be combined into a covariance matrix Σ . The covariance matrix is a symmetric positive definite matrix that describes the covariance between the components of $\boldsymbol{\xi}$.

$$\Sigma = \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_k^2 \end{pmatrix} \quad (5)$$

Here, the covariance matrix is diagonal (all off-diagonal elements are zero), because we assumed that the components of $\boldsymbol{\xi}$ are uncorrelated, *i.e.*, $\text{Cov}[\xi_i, \xi_j] = 0$ for all $i \neq j$.

The pdf of the multivariate normal distribution can be written in terms of the covariance matrix:

$$f_{\boldsymbol{\xi}}(x_1, \dots, x_k) = \frac{\exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}}{\sqrt{(2\pi)^k \det \Sigma}} \quad (6)$$

The covariance matrix Σ above is a diagonal matrix, but in general, it's a symmetric positive definite matrix that describes the covariance between the components of $\boldsymbol{\xi}$:

$$\Sigma := \begin{pmatrix} \text{Cov}[\xi_1, \xi_1] & \dots & \text{Cov}[\xi_1, \xi_k] \\ \vdots & \ddots & \vdots \\ \text{Cov}[\xi_k, \xi_1] & \dots & \text{Cov}[\xi_k, \xi_k] \end{pmatrix}. \quad (7)$$

If we substitute the non-diagonal covariance matrix Σ into the pdf, we get the general form of the multivariate normal distribution.

Technically, each component of Σ is the covariance between the corresponding components

$$\Sigma_{i,j} := \text{Cov}[\xi_i, \xi_j] = \mathbb{E}[(\xi_i - \mu_i)(\xi_j - \mu_j)]. \quad (8)$$

The term $\det \Sigma$ is the generalized variance.

The uncorrelated multivariate normal distribution is a special case of the general multivariate normal. When components are uncorrelated, the covariance matrix is diagonal, which simplifies many calculations.

For a sample $X = \{x_1, \dots, x_{\ell}\} \subset \mathbb{R}$, the variance is the average of the squared differences from the mean:

$$\mathbb{D}[X] := \frac{1}{\ell} \sum_{i=1}^{\ell} (x_i - \bar{x})^2.$$

Given another sample $Y = \{y_1, \dots, y_{\ell}\} \subset \mathbb{R}$, the *co*-variance between two samples is characterized by how much they vary together:

$$\text{Cov}[X, Y] := \frac{1}{\ell} \sum_{i=1}^{\ell} (x_i - \bar{x}) \cdot (y_i - \bar{y}).$$

Both per sample variance and two samples covariance can be combined into a covariance matrix.

$$\Sigma = \begin{pmatrix} \text{Cov}[x, x] & \text{Cov}[x, y] \\ \text{Cov}[y, x] & \text{Cov}[y, y] \end{pmatrix} = \begin{pmatrix} \mathbb{D}[x] & \text{Cov}[x, y] \\ \text{Cov}[y, x] & \mathbb{D}[y] \end{pmatrix}$$

Mahalanobis distance: The distance between a point \mathbf{x} and the distribution $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$ can be measured using the Mahalanobis distance.

The premise is that the covariance matrix Σ captures the correlations between the components of $\boldsymbol{\xi}$. The Mahalanobis distance is a measure of how many standard deviations away a point \mathbf{x} is from the mean $\boldsymbol{\mu}$, taking into account the correlations between the components of $\boldsymbol{\xi}$.

We can define a quadratic form

$$\begin{aligned} Q(\mathbf{x}) &:= (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\ &= \sum_{i,j} (x_i - \mu_i) (\Sigma^{-1})_{i,j} (x_j - \mu_j). \end{aligned} \quad (9)$$

The square root of this quadratic form $\sqrt{Q(\mathbf{x})}$ is the Mahalanobis distance between a point \mathbf{x} and the distribution $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$.

Quadratic form $Q(\mathbf{x})$ is a scalar function of a vector \mathbf{x} that can be expressed as a weighted sum of the squares of the components of \mathbf{x} :

$$Q(\mathbf{x}) = \sum_{i,j} w_{i,j} x_i x_j.$$

These weights can be gathered into a matrix W , and the quadratic form can be written as a matrix product:

$$Q(\mathbf{x}) = \mathbf{x}^\top W \mathbf{x}.$$

2 Multivariate Normal Distribution

General form: The probability density function of the multivariate normal distribution is given by:

$$f_{\boldsymbol{\xi}}(x_1, \dots, x_k) := \frac{\exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\}}{\sqrt{(2\pi)^k \det \Sigma}} \quad (10)$$

where:

- * $\boldsymbol{\xi} = \begin{pmatrix} \xi_1 \\ \vdots \\ \xi_k \end{pmatrix}$ is the random vector
- * $\boldsymbol{\mu} := \mathbb{E}\boldsymbol{\xi} = \begin{pmatrix} \mathbb{E}\xi_1 \\ \vdots \\ \mathbb{E}\xi_k \end{pmatrix}$ is the mean vector
- * $\Sigma_{i,j} := \text{Cov}[\xi_i, \xi_j] = \mathbb{E}[(\xi_i - \mu_i)(\xi_j - \mu_j)]$ is the covariance matrix (symmetric positive definite)
- * $\det \Sigma$ is the generalized variance
- * $Q(\mathbf{x}) = (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})$ is a quadratic form
- * $\sqrt{Q(\mathbf{x})}$ is the Mahalanobis distance between a point \mathbf{x} and the distribution $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$

In compact notation, this is written as: $\boldsymbol{\xi} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$

Uncorrelated components: When the components of the distribution are uncorrelated: $\text{Cor}[\xi_i, \xi_j] = 0 \quad \forall i \neq j$

This has geometric implications - the axes of the probability density ellipsoid are parallel to the coordinate axes.

The covariance matrix simplifies to a diagonal matrix: $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$

For uncorrelated components, the multivariate normal pdf can be expressed as a product of univariate normal pdfs:

$$\begin{aligned} f_{\boldsymbol{\xi}}(x_1, \dots, x_k) &:= \frac{\exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\}}{\sqrt{(2\pi)^k \det \Sigma}} \\ &= \prod_{i=1}^k \frac{\exp\left\{-\frac{1}{2}\left(\frac{x_i - \mu_i}{\sigma_i}\right)^2\right\}}{\sigma_i \sqrt{2\pi}} \end{aligned} \quad (11)$$

This factorization enables simpler parameter estimation methods.

Decorrelation transformation:

For correlated components, we need to find a transformation that makes the components uncorrelated. When components of a multivariate normal distribution are correlated, the covariance matrix is not diagonal.

Decorrelation is a critical preprocessing step in many machine learning applications. It simplifies the data by removing linear dependencies between features.

To decorrelate the components:

1. Apply spectral decomposition to the covariance matrix (a special case of SVD for symmetric matrices):

$$\Sigma = VSV^T, \quad S = \text{diag}(\lambda_1^2, \dots, \lambda_k^2)$$

2. The quadratic form can be rewritten as:

$$\begin{aligned} Q &= (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^T (VSV^T)^{-1}(\mathbf{x} - \boldsymbol{\mu}) \\ &= (\mathbf{x} - \boldsymbol{\mu})^T V^{-T} S^{-1} V^{-1}(\mathbf{x} - \boldsymbol{\mu}) \\ &= (\mathbf{x} - \boldsymbol{\mu})^T V S^{-1} V^T (\mathbf{x} - \boldsymbol{\mu}) \\ &= (V^T (\mathbf{x} - \boldsymbol{\mu}))^T S^{-1} (V^T (\mathbf{x} - \boldsymbol{\mu})) \end{aligned} \quad (12)$$

3. Define the decorrelation transformation:

$$\mathbf{x}' := V^T \mathbf{x}$$

4. The transformed parameters are:

$$\begin{aligned}\boldsymbol{\mu}' &= V^{\top} \boldsymbol{\mu} \\ \Sigma' &= S\end{aligned}\tag{13}$$

5. To transform parameters back to the original space:

$$\begin{aligned}\boldsymbol{\mu} &= V \boldsymbol{\mu}' \\ \Sigma &= V S V^{\top}\end{aligned}\tag{14}$$

For orthogonal (rotation) matrices: $V^{-1} = V^{\top}$, $V^{-T} = V$

3 Parameter Estimation for Normal Distribution

Maximum likelihood estimation: When a sample is generated from a Gaussian distribution, we can estimate its parameters using maximum likelihood estimation:

$$\begin{aligned}\frac{\partial}{\partial \mu} \ln L(\mu, \Sigma \mid X^\ell) &= 0 \Rightarrow \hat{\mu} = \frac{1}{\ell} \sum_{x \in X^\ell} x \\ \frac{\partial}{\partial \Sigma} \ln L(\mu, \Sigma \mid X^\ell) &= 0 \Rightarrow \hat{\Sigma} = \frac{1}{\ell} \sum_{x \in X^\ell} (x - \hat{\mu})(x - \hat{\mu})^\top\end{aligned}\tag{15}$$

Multicollinearity issues:

When estimating parameters from data, multicollinearity can cause problems:

- * The sample covariance matrix $\hat{\Sigma}$ is constructed from ℓ rank-1 matrices, so $\text{rank } \hat{\Sigma} \leq \ell$
- * If the number of features exceeds the sample size, $\hat{\Sigma}$ will be singular ($\det = 0$)
- * A singular covariance matrix cannot be inverted, making it impossible to evaluate the density function

Solutions to multicollinearity include:

1. Reducing the number of features through feature selection methods (PCA, SFS, *etc.*)
2. Increasing the sample size
3. Adding regularization to the covariance matrix:

$$\hat{\Sigma}' \leftarrow \hat{\Sigma} + \tau I$$

4. Assuming uncorrelated features (diagonal covariance matrix)

The MLE estimates are asymptotically unbiased and efficient, making them optimal for large samples. For smaller samples, however, the covariance estimate is biased.

Multicollinearity is a statistical phenomenon where two or more predictor variables in a model are highly correlated. This creates redundant information that doesn't contribute uniquely to the model's explanatory power.

4 Non-parametric Density Estimation

Parametric vs non-parametric approaches:

Parametric methods for density estimation assume a specific distribution shape (*e.g.*, normal), which simplifies calculations and requires less data. However, they can be inaccurate if the actual data distribution differs significantly from the assumed form and are limited in capturing complex structures.

Non-parametric methods offer greater flexibility and can provide more accurate estimates for complex and multimodal distributions without assuming a specific form. However, they typically require larger data sets, are more computationally intensive, and their results are harder to interpret due to the lack of explicit parameters.

The choice between parametric and non-parametric methods involves a bias-variance tradeoff. Parametric methods generally have higher bias but lower variance, while non-parametric methods have lower bias but higher variance.

Empirical density estimator: The simplest non-parametric estimator for a probability density function is:

$$\hat{f}(x) := \frac{1}{\ell} \sum_{x' \in X^\ell} \mathbb{I}[x = x'], \quad (16)$$

and for the cumulative distribution function:

$$\hat{F}(x) := \frac{1}{\ell} \sum_{x' \in X^\ell} \mathbb{I}[x'_1 \leq x_1] \cdots \mathbb{I}[x'_k \leq x_k] \quad (17)$$

where:

- * f is a pdf/pmf
- * F is a cdf
- * $x' = (x'_1, \dots, x'_k)$ is a point from the sample X^ℓ
- * $x = (x_1, \dots, x_k)$ is a new point

Histogram density estimator: A more practical approach divides the feature space into bins:

$$\hat{f}(x) := \frac{1}{\ell} \cdot \#(B(x) \cap X^\ell), \quad (18)$$

and for the CDF:

$$\hat{F}(x) := \frac{1}{\ell} \sum_B \mathbb{I}[B \leq B(x)] \cdot \#B \quad (19)$$

where:

- * $X^\ell \subseteq \bigsqcup_B B$ is the partition into bins
- * $\#B$ is the bin size
- * $B(x)$ is the specific bin containing x
- * $B \leq B'$ iff $\bigwedge_{j=1}^k \sup B_j \leq \sup B'_j$ allows ordering of bins
- * n_j is the number of bins for the j th feature
- * $h_j := \frac{f_j^{\max} - f_j^{\min}}{n}$ is the corresponding bin width

Window averaging:

Instead of bins, we can use a window function centered at each data point:

$$\hat{f}(x) := \frac{1}{\ell} \cdot \frac{1}{2h} \sum_{x' \in X^\ell} \mathbb{I}\left[\frac{\|x - x'\|}{h} < 1\right], \quad (20)$$

and for the CDF:

$$\hat{F}(x) := \frac{1}{\ell} \sum_{x' \in X^\ell} \mathbb{I}[x' \leq x \oplus h] \quad (21)$$

where:

The choice of window width h is critical in kernel density estimation and presents the classic bias-variance tradeoff.

- * h is the window width (radius)
- * $\mathbf{x}' \leq \mathbf{x}$ means all components are less or equal
- * \oplus means componentwise addition
- * $K(r) = \frac{1}{2} \mathbb{I}[|r| < 1]$ is a kernel function

5 Kernel Density Estimation

Parzen-Rosenblatt window method:

There are two main approaches for multivariate density estimation:

1. Product kernel approach (assumes local independence):

$$\hat{f}(\mathbf{x}) := \frac{1}{\ell} \sum_{\mathbf{x}' \in X^\ell} \prod_{j=1}^k \frac{1}{h_j} \cdot K\left(\frac{\mathbf{x}^j - \mathbf{x}'^j}{h_j}\right) \quad (22)$$

2. Multivariate kernel approach:

$$\hat{f}(\mathbf{x}) := \frac{1}{\ell} \cdot \frac{1}{V(h)} \sum_{\mathbf{x}' \in X^\ell} K\left(\frac{\rho(\mathbf{x}, \mathbf{x}')}{h}\right) \quad (23)$$

where:

- * $K(r)$ is a kernel function satisfying:
 - > $\int K(r) dr = 1$
 - > $K(r) > 0$
 - > For all $r > 0$: $K(r) \searrow$ (non-increasing)
- * $V(h) := \int_{-\infty}^{+\infty} K\left(\frac{\rho}{h}\right) d\rho$

Kernel functions: Several standard kernel functions are used in practice:

- * Epanechnikov kernel: $E(r) = \frac{3}{4}(1 - r^2)_+$
- * Quartic (biweight) kernel: $Q(r) = \frac{15}{16}(1 - r^2)^2 \cdot \mathbb{I}[|r| > 0]$
- * Triangle kernel: $T(r) = (1 - |r|)_+$
- * Gaussian kernel: $G(r) = \frac{1}{\sqrt{2\pi}} e^{-\frac{r^2}{2}}$
- * Uniform kernel: $\Pi(r) = \frac{1}{2} \cdot \mathbb{I}[|r| \leq 1]$

Bandwidth selection: The optimal bandwidth can be found by minimizing the cross-validation criterion:

$$Q(h) = - \sum_{\mathbf{x} \in X^\ell} \ln \hat{f}(\mathbf{x} \mid X^\ell \setminus \mathbf{x}, h) \rightarrow \min_h \quad (24)$$

where:

- * h is the window width (radius)
- * \hat{f} is the estimated density function
- * The notation indicates leave-one-out cross-validation

This approach:

1. Excludes each point from the training set
2. Estimates the density at that point using the remaining points
3. Minimizes the negative log-likelihood of these estimates

Relationship to other methods:

Kernel methods form a common framework that includes:

1. Density estimation: $a_1(\mathbf{x}) = \frac{1}{\ell V(h)} \sum_{\mathbf{x}' \in X^\ell} K\left(\frac{\rho(\mathbf{x}, \mathbf{x}')}{h}\right)$
2. Classification (Parzen window): $a_2(\mathbf{x}) = \arg \max_{y \in Y} \sum_{\mathbf{x}' \in X^\ell} \mathbb{I}[y(\mathbf{x}') = y] \cdot K\left(\frac{\rho(\mathbf{x}, \mathbf{x}')}{h}\right)$
3. Regression (Nadaraya-Watson): $a_3(\mathbf{x}) = \frac{\sum_{\mathbf{x}' \in X^\ell} y(\mathbf{x}') \cdot K\left(\frac{\rho(\mathbf{x}, \mathbf{x}')}{h}\right)}{\sum_{\mathbf{x}'} K\left(\frac{\rho(\mathbf{x}, \mathbf{x}')}{h}\right)}$

In these methods:

- * ρ is a distance function
- * $K(\rho)$ is a similarity function (larger distance means smaller similarity)

Kernel density estimation is often called the Parzen-Rosenblatt window method after its developers Emanuel Parzen and Murray Rosenblatt.

The Epanechnikov kernel is theoretically optimal in terms of mean integrated squared error, but in practice, the choice of kernel usually has less impact than the bandwidth selection.

Cross-validation provides a data-driven approach to bandwidth selection, avoiding both oversmoothing and undersmoothing.

The connection between these methods highlights the unified theoretical foundation underlying non-parametric approaches.

6 Mixture Models

Definition:

A mixture model combines multiple probability distributions:

$$f(\mathbf{x}) := \sum_{n=1}^N w_n \cdot f_n(\mathbf{x} \mid \boldsymbol{\theta}_n) \quad (25)$$

where:

- * $\mathbf{x} \sim \mathcal{D}_1, \dots, \mathcal{D}_N$ indicates data generated by N different distributions
- * $w_n := \mathbb{P}(\mathbf{x} \sim \mathcal{D}_n)$ is the probability of being generated by the n th distribution
- * $\sum_{n=1}^N w_n = 1, \quad w_n \geq 0$
- * $f_n(\mathbf{x} \mid \boldsymbol{\theta}_n)$ is the pdf/pmf of the n th distribution

Mixture models provide a flexible framework that can approximate any continuous distribution with arbitrary precision.

Parameter estimation: The log-likelihood function for a mixture model is:

$$l(\mathbf{w}, \boldsymbol{\theta}) = \sum_{\mathbf{x} \in X^\ell} \ln \sum_{n=1}^N w_n \cdot f_n(\mathbf{x} \mid \boldsymbol{\theta}_n) \rightarrow \max_{\mathbf{w}, \boldsymbol{\theta}} \quad (26)$$

where:

- * $\sum_n w_n = 1, \quad w_n \geq 0$ are constraints
- * $f_n(\mathbf{x} \mid \boldsymbol{\theta}_n)$ is the pdf/pmf of the n th distribution

Direct optimization is challenging because the logarithm of a sum doesn't simplify easily. The EM algorithm provides an iterative solution.

Fixed point method:

Before introducing EM, it's helpful to understand the fixed point iteration method:

$$x_{n+1} = f(x_n)$$

This method converges if: $|f'(x^*)| < 1$

where x^* is the fixed point such that $f(x^*) = x^*$.

The fixed point method underlies many iterative algorithms in machine learning, including the EM algorithm.

7 Expectation-Maximization Algorithm

EM algorithm:

The EM algorithm iteratively optimizes parameters w_n and θ_n for mixture components. Each iteration consists of two steps:

1. Expectation step (E-step): Calculate the posterior probability for each data point:

$$w'_n(x) := \mathbb{P}[x \sim f_n \mid x] = \frac{w_n \cdot f_n(x \mid \theta_n)}{\sum_{m=1}^N w_m f_m(x, \theta_m)} \quad (27)$$

2. Maximization step (M-step): Update parameters:

$$\theta_n \leftarrow \arg \max_{\theta} \sum_{x \in X^\ell} w'_n(x) \cdot \ln f(x \mid \theta_n) \quad (28)$$

$$w_n \leftarrow \frac{1}{\ell} \sum_{x \in X^\ell} w'_n(x) \quad (29)$$

Theoretical foundation:

The EM algorithm optimizes a Lagrangian function:

$$Q(w, \Theta) = \sum_{x \in X^\ell} \ln \left(\sum_{n=1}^N w_n \cdot f_n(x \mid \theta_n) \right) - \lambda \cdot \left(\sum_{n=1}^N w_n - 1 \right) \rightarrow \max_{w, \Theta} \quad (30)$$

where:

- * $\Theta = [\theta_1, \dots, \theta_N]$ is the parameters matrix
- * λ is the Lagrange multiplier
- * f_n is the pdf/pmf of the n th distribution
- * $w_n := \mathbb{P}[x \sim f_n]$ is the probability of coming from the n th distribution

Variants of EM:

Generalized EM (GEM) relaxes the maximization requirement:

- * Standard EM: $\theta^{(t+1)} \leftarrow \arg \max_{\theta} \ell(\theta)$
- * GEM: $\theta^{(t+1)} \leftarrow \theta^* : \ell(\theta^*) > \ell(\theta^{(t)})$

Stochastic EM (SEM) optimizes parameters using sampled subsets:

- * Generate samples from the estimated distributions
- * Optimize parameters independently for each component
- * This often accelerates convergence and can use standard maximum likelihood methods

Determining the number of components: Several approaches can determine the optimal number of distributions N in a mixture:

1. Greedy addition: Start with fewer components; add new components if the likelihood for some data points is below a threshold
2. Greedy deletion: Start with more components; remove components with small weights
3. AddDel: Combination of both approaches
4. Regularization: Use cross-entropy regularization to encourage sparsity in component weights

Hierarchical EM: Hierarchical EM extends the standard algorithm to restore hierarchical relationships in the data. It operates by greedily adding components and splitting components with low likelihood. This approach is useful for clustering and enhancing data understanding by revealing structure where single clusters may have multiple subclusters.

The EM algorithm was formalized by Dempster, Laird, and Rubin in 1977, though similar approaches had been used earlier.

The EM algorithm can be derived using the Lagrangian function and Karush-Kuhn-Tucker conditions.

Information criteria like AIC or BIC can also be used to select the optimal number of components, balancing model complexity with goodness-of-fit.

8 Comparing Density Estimation Methods

Framework comparison: Different approaches to density estimation include:

1. Parametric: Assumes a specific functional form

$$\hat{f}(\mathbf{x}) = f(\mathbf{x} \mid \boldsymbol{\theta})$$

2. Non-parametric kernel: Based on local estimations around each training point

$$\hat{f}(\mathbf{x}) = \frac{1}{\ell} \sum_{\mathbf{x}' \in X^\ell} \frac{1}{V(h)} K\left(\frac{\rho(\mathbf{x}, \mathbf{x}')}{h}\right)$$

3. Mixture models: Combines multiple distributions

$$\hat{f}(\mathbf{x}) = \sum_{n=1}^N w_n \cdot f_n(\mathbf{x} \mid \boldsymbol{\theta}_n)$$

These approaches represent a spectrum: mixture models generalize both parametric methods (when $N = 1$) and non-parametric approaches (when N equals the sample size).

Mixture models provide a unified framework - when $N = 1$, they reduce to parametric methods, and as N approaches the sample size, they approximate non-parametric methods.