# Question

Discuss the importance of discriminant functions in classification. How do they influence decision boundaries? (Short answer)

---

# Answer

In classification, **discriminant functions** $(g_k(x))$ are like **scoring systems** for each class $C_k$. For any input data $x$, you calculate a score $g_k(x)$ for every class.

**Importance:**

- They give each class a score based on the data $x$.
- The decision is simple: classify $x$ into the class with the **highest score**.

**Influence on Decision Boundaries:**

- The **decision boundary** between two classes (say $C_i$ and $C_j$) is where their scores are equal: $g_i(x) = g_j(x)$.
- The **type of mathematical function** used for $g_k(x)$ directly controls the **shape** of this boundary:
    - If $g_k(x)$ are simple **linear** functions, the boundary is a **straight line** (or flat plane).
    - If $g_k(x)$ are **more complex** (like quadratic or non-linear), the boundary can be **curvy or wiggly**.
- Training a classifier is essentially learning the numbers inside these functions $(g_k(x))$ so that the boundary is placed correctly to separate the data.

Basically, discriminant functions score how well data fits each class, and where scores are tied between two classes, that's your decision boundary. The math inside the score function determines the boundary's shape.

# Question

Consider a system for detecting spam emails. Formulate the classification problem using Bayesian decision theory. Specify the discriminant function and decision surface for classifying emails as spam or non-spam.

# Answer

Let's formulate the spam detection problem using the framework of Bayesian decision theory.

## 1. Problem Formulation

- **Classes:** We have two classes:
    - $C_1$: Spam email
    - $C_2$: Non-spam email (often called "Ham")
- **Feature Vector ($x$):** An email needs to be represented as a feature vector $x$. This vector could contain information extracted from the email, such as:
    - Word frequencies (e.g., count of the word "free", "money", etc.)
    - Presence/absence of specific phrases (e.g., "click here")
    - Sender's address characteristics
    - Email structure (e.g., use of excessive capitalization, exclamation marks)
    - Metadata (e.g., time sent)
- **Goal:** Given the feature vector $x$ of an incoming email, we want to classify it as either $C_1$ or $C_2$. Bayesian decision theory aims to make the decision that minimizes the expected risk (or minimizes the probability of error, assuming equal costs for misclassification).

## 2. Bayesian Decision Rule (Minimum Error Rate)

The decision rule for minimizing the probability of error is to choose the class $C_k$ that has the highest posterior probability given the observed features $x$:
Classify $x$ as $C_1$ if $P(C_1 \mid x) > P(C_2 \mid x)$
Classify $x$ as $C_2$ if $P(C_2 \mid x) > P(C_1 \mid x)$

Using Bayes' Theorem, $P(C_k \mid x) = \frac{P(x|C_k)P(C_k)}{P(x)}$, where:

- $P(C_k)$ is the **prior probability** of class $C_k$ (e.g., the overall proportion of spam/non-spam emails you receive).
- $P(x \mid C_k)$ is the **likelihood**, representing the probability of observing the feature vector $x$ given that the email belongs to class $C_k$.
- $P(x)$ is the **evidence** (probability of observing $x$), which is the same for all classes for a given $x$.

The decision rule $P(C_1 \mid x) > P(C_2 \mid x)$ is equivalent to comparing the unnormalized posteriors:
Classify $x$ as $C_1$ if $P(x \mid C_1)P(C_1) > P(x \mid C_2)P(C_2)$
Classify $x$ as $C_2$ if $P(x \mid C_1)P(C_1) < P(x \mid C_2)P(C_2)$

## 3. Discriminant Function

A common choice for the **discriminant function** $g_k(x)$ in Bayesian classification is proportional to the posterior probability. A widely used form, especially when likelihoods are multiplied (like in models assuming conditional independence), is the logarithm of the unnormalized posterior:

$$g_k(x) = \log(P(x \mid C_k)P(C_k)) = \log P(x \mid C_k) + \log P(C_k)$$

For our spam/non-spam problem ($C_1$ = Spam, $C_2$ = Non-spam), the discriminant functions are:

$$g_{\text{Spam}}(x) = \log P(x \mid \text{Spam}) + \log P(\text{Spam})$$

$$g_{\text{Non-spam}}(x) = \log P(x \mid \text{Non-spam}) + \log P(\text{Non-spam})$$

To build the classifier, we need to estimate the prior probabilities ($P(\text{Spam})$ and $P(\text{Non-spam})$ from data frequencies) and, more importantly, model the likelihoods ($P(x \mid \text{Spam})$ and $P(x \mid \text{Non-spam})$). The form of $P(x \mid C_k)$ depends on the chosen model (e.g., Naive Bayes, Gaussian models, etc.).

## 4. Decision Surface

The **decision surface** (or decision boundary) is the set of points $x$ in the feature space where the decision changes. For two classes, this occurs where the discriminant functions are equal:

$$g_{\text{Spam}}(x) = g_{\text{Non-spam}}(x)$$

$$\log P(x \mid \text{Spam}) + \log P(\text{Spam}) = \log P(x \mid \text{Non-spam}) + \log P(\text{Non-spam})$$

This equation can be rearranged to:

$$\log \left( \frac{P(x \mid \text{Spam})}{P(x \mid \text{Non-spam})} \right) = \log \left( \frac{P(\text{Non-spam})}{P(\text{Spam})} \right)$$

Or, in terms of the likelihood ratio:

$$\frac{P(x \mid \text{Spam})}{P(x \mid \text{Non-spam})} = \frac{P(\text{Non-spam})}{P(\text{Spam})}$$

This is the classic likelihood ratio test for the decision boundary in a two-class problem.

The **shape** of the decision surface defined by $g_{\text{Spam}}(x) = g_{\text{Non-spam}}(x)$ depends entirely on the models used for the likelihoods $P(x \mid \text{Spam})$ and $P(x \mid \text{Non-spam})$. For example:

- If $P(x \mid C_k)$ are assumed to be multivariate Gaussians with equal covariance matrices, the boundary is linear.
- If $P(x \mid C_k)$ are multivariate Gaussians with unequal covariance matrices, the boundary is quadratic.
- If a Naive Bayes assumption is made ($P(x \mid C_k) = \prod_i P(x_i \mid C_k)$) and features are binary or counts, the $\log P(x \mid C_k)$ term involves sums of individual feature log-probabilities, leading to potentially complex boundaries depending on the feature distributions.

In summary, Bayesian spam classification involves modeling the probability of features given the spam/non-spam classes ($P(x \mid C_k)$) and using these, along with prior probabilities ($P(C_k)$), in discriminant functions to make decisions by comparing weighted likelihoods. The boundary is where these weighted likelihoods are equal.

# Question

High-dimensional feature spaces often lead to the "curse of dimensionality." Discuss this problem and propose techniques to mitigate its effects in automatic pattern recognition systems, keeping the theory simple and concise.

---

# Answer

### The Problem: The "Curse of Dimensionality"

When you have data with a very high number of features (dimensions), the space that the data lives in becomes incredibly vast. The core problem is:

- **Data becomes incredibly sparse:** Your data points, no matter how many you have, get spread out extremely thin across this huge space. They become very far away from each other.
- **Models struggle:** Classifiers or other pattern recognition systems find it hard to identify reliable patterns because they don't see enough data points close together. They might find patterns in random noise by mistake.
- **Overfitting is easy:** Models can become too complex and fit the noise in the training data instead of the true underlying signal. They won't work well on new data.
- **Everything slows down:** Computations involving high-dimensional data take much longer and require more memory.

Essentially, in high dimensions, the data is sparse, distances are less meaningful, and models need exponentially more data to perform well, which is usually impossible to get.

### Mitigation Techniques (Solutions)

The main way to fight the curse of dimensionality is to **reduce the number of features** before or during model training.

  1. **Dimensionality Reduction:**
   - **Feature Selection:** Choose the most important original features and discard the rest. (Example: Pick features that are highly correlated with the class you want to predict). This keeps the original meaning of the features.

- **Feature Extraction:** Create a smaller set of entirely new features that are combinations of the original ones, capturing most of the important information in fewer dimensions. (Example: Principal Component Analysis (PCA) is a popular method that finds directions of maximum data variation).

2. **Regularization:** Use techniques during model training that penalize complexity. This helps prevent models from relying too much on any single feature or finding overly complex, noisy patterns. (Lasso, or L1 regularization, is a type that can also help by pushing the weights of less important features to zero).

3. **Choose Appropriate Models:** Some machine learning models are less affected by high dimensions than others, although none are completely immune.

In short, the curse means high dimensions make data sparse and models overfit. You mitigate it mainly by reducing the number of features using selection or extraction methods, and by using regularization during training to keep models simple.

# Question

Explain the key design concepts in an automatic pattern recognition system. How does feature selection impact system performance?

---

# Answer

An automatic pattern recognition system is designed to take raw data, process it, and make a decision or identify a pattern (like classifying an object in an image or detecting spam). Building such a system involves several key design concepts or stages:

**Key Design Concepts in Automatic Pattern Recognition Systems:**

1. **Data Acquisition:** This is the initial step where raw data is collected from sensors, databases, files, etc. The quality and quantity of this data are crucial. (e.g., getting images from a camera, audio from a microphone).

2. **Preprocessing:** Raw data is often noisy, incomplete, or in a format not suitable for direct analysis. Preprocessing cleans and prepares the data. This can include:
   - Noise reduction (e.g., smoothing images, filtering audio).
   - Handling missing values.
   - Normalization or scaling (e.g., ensuring all feature values are within a similar range).
   - Data transformation (e.g., converting audio to spectrograms).

3. **Feature Extraction / Representation:** This is a critical step where meaningful characteristics (features) are extracted from the preprocessed data. The goal is to transform the data into a

numerical vector or representation that captures the essential information for distinguishing between different patterns while discarding irrelevant variations. (e.g., extracting edges, textures from images; calculating spectral features from audio). This reduces the dimensionality from the raw data space.

4. **Pattern Recognition / Modeling (Classification or Regression):** This is the core intelligence of the system. A machine learning model or algorithm is chosen and trained using labeled data (in supervised learning) to learn the relationship between the feature vectors and the target patterns/classes. This model uses the features to classify new, unseen data points. (e.g., training a Support Vector Machine, a Neural Network, or a Decision Tree).

5. **Evaluation:** After training, the model's performance is rigorously evaluated on a separate dataset (test set) that it has not seen before. Common metrics include accuracy, precision, recall, F1-score, etc. This step is essential to ensure the model generalizes well to new data and to fine-tune parameters.

6. **Postprocessing (Optional):** The output of the model might be further processed. For instance, if the model outputs a probability score, postprocessing might involve applying a threshold. Or, combining outputs from multiple models.

## Impact of Feature Selection on System Performance

**Feature selection** is a technique typically applied during the **Preprocessing** or **Feature Extraction** phase. It involves choosing a *subset* of the most relevant, informative, and non-redundant features from the initial set (often the features extracted in step 3).

Feature selection significantly impacts system performance in several ways:

1. **Combats the "Curse of Dimensionality":** By reducing the number of dimensions, it helps mitigate issues like data sparsity and the increased risk of overfitting that occur in high-dimensional spaces.

2. **Improves Model Accuracy and Generalization:** Removing irrelevant or noisy features helps the learning algorithm focus on the truly discriminative information. This often leads to models that learn better patterns and generalize more effectively to unseen data.

3. **Reduces Training Time:** Training machine learning models is computationally less expensive and faster when dealing with fewer features.

4. **Decreases Memory Usage:** Storing and processing datasets with fewer features requires less memory.

5. **Enhances Model Interpretability:** A model based on a smaller, carefully selected set of features can be easier to understand and explain.

However, it's crucial to note that **poor feature selection** can degrade performance if important discriminative features are mistakenly removed. Therefore, feature selection techniques must be applied carefully, often guided by domain knowledge or systematic evaluation methods.

# Question

Differentiate between Minimum-Error-Rate classification and Maximum A Posteriori (MAP) classification. Provide examples where each is preferable.

---

# Answer

Both Minimum-Error-Rate classification and Maximum A Posteriori (MAP) classification are decision rules derived from Bayesian decision theory. However, they differ in what they prioritize regarding the costs of mistakes.

The key difference is how they handle the **costs associated with different types of misclassification**.

## Maximum A Posteriori (MAP) Classification

- **Goal:** To choose the class that is most probable given the observed data. It aims to maximize the **posterior probability** $P(C_k \mid x)$ for a given data point $x$.
- **Decision Rule:** Classify $x$ into class $C_k$ such that $P(C_k \mid x)$ is maximized.
- **Assumed Costs:** MAP classification is implicitly minimizing the error rate *assuming that the cost of any misclassification is the same*, regardless of which incorrect class is chosen, and that the cost of a correct classification is zero. If falsely classifying a cat as a dog costs the same as falsely classifying a dog as a cat, MAP is the appropriate rule.

## Minimum-Error-Rate Classification

- **Goal:** To minimize the overall **probability of making a classification error**, taking into account that different types of errors might have different consequences or costs. More generally, it minimizes the **expected risk**.
- **Decision Rule:** This involves calculating the expected risk $R(\alpha_i \mid x)$ for classifying $x$ into each possible class $C_i$, considering the potential true classes $C_j$ and the specific costs $\lambda(i \mid j)$ (cost of deciding $C_i$ when the true class is $C_j$). The rule is to choose the class $C_i$ that minimizes $R(\alpha_i \mid x)$.
- **Considers Costs:** This method explicitly incorporates a **loss function** or **cost matrix** $\lambda(i \mid j)$ that assigns potentially different penalties to different types of misclassification errors.

## Key Differentiation:

MAP is a **special case** of Minimum-Error-Rate classification that arises when all misclassification costs $\lambda(i \mid j)$ for $i \neq j$ are equal. If the costs are *unequal*, the Minimum-Error-

Rate rule (based on minimizing expected risk using the specific cost matrix) will generally lead to a different decision boundary than the simple MAP rule.

**Examples of When Each is Preferable:**

- **MAP Classification is preferable when:**
    - All misclassification errors have roughly **equal consequences**.
    - You simply want the most **statistically probable class** assignment based on the data and priors, without specific information about the costs of errors.
    - Examples: Standard benchmarking problems (like classifying digits in MNIST or species in Iris) where the "cost" of any wrong answer is treated equally. General object recognition in a photo collection where classifying a 'car' as 'truck' is no worse than classifying it as 'bus'.
- **Minimum-Error-Rate (with unequal costs) is preferable when:**
    - Different types of misclassification errors have **significantly different costs** or consequences.
    - Examples:
        - **Medical Diagnosis:** Falsely diagnosing a healthy person with a disease (False Positive) has a cost (stress, unnecessary tests), but falsely telling a sick person they are healthy (False Negative) can have a much higher cost (delayed treatment, health deterioration). Minimizing overall errors might not be best; minimizing the *cost* of False Negatives is often critical.
        - **Security Systems:** Falsely identifying an authorized person as a threat (False Positive) is costly (delays, inconvenience), but falsely identifying a threat as authorized (False Negative) can be vastly more expensive (security breach). Minimizing the cost of a False Negative is paramount.

In summary, use MAP when all errors cost the same; use the general Minimum-Error-Rate approach when different errors have different associated costs, allowing you to build a classifier that minimizes the total expected cost of its decisions.

# Question

Consider an application where self-driving cars use pattern recognition to detect stop signs. Describe the role of Bayesian decision theory in this scenario and formulate a discriminant function based on prior probabilities and likelihood estimates.

# Answer

In the context of a self-driving car detecting stop signs, the task is to classify whether a **stop sign is present** or **not present** based on sensor data, typically images from cameras. Bayesian decision theory provides a formal framework for making this decision optimally under uncertainty.

## Role of Bayesian Decision Theory in Stop Sign Detection

Self-driving cars operate in a complex, dynamic, and uncertain environment. Sensor data can be noisy, lighting conditions vary, objects can be partially occluded, and other signs might look similar to stop signs. Bayesian decision theory helps the system make robust decisions by explicitly combining:

1. **Evidence from Sensors (Likelihood):** What does the sensor data (e.g., features extracted from a camera image $x$) tell us about the presence of a stop sign? This is captured by the **likelihood** $P(x \mid \text{Stop Sign Present})$ and $P(x \mid \text{Stop Sign Not Present})$. These probabilities come from the pattern recognition model trained to recognize stop signs based on their visual features.
2. **Contextual Belief (Prior Probability):** How likely is it, independently of the current sensor data, that a stop sign should be present at the vehicle's current location? This is the **prior probability** $P(\text{Stop Sign Present})$. A car's internal state, GPS location, and map data provide this prior information. For instance, $P(\text{Stop Sign Present})$ would be very high if the car is approaching an intersection known to have a stop sign according to its map, and very low if it's in the middle of a highway. $P(\text{Stop Sign Not Present}) = 1 - P(\text{Stop Sign Present})$.

Bayesian decision theory allows the system to calculate the **posterior probability** $P(\text{Stop Sign Present} \mid x)$, which is the probability that a stop sign is actually present given the observed sensor data $x$. The core rule is to choose the class with the highest posterior probability (or minimize expected risk, considering different costs for hitting vs. unnecessarily stopping).

## Formulating a Discriminant Function

Let $C_1$ be the class "Stop Sign Present" and $C_2$ be the class "Stop Sign Not Present". Let $x$ be the feature vector derived from sensor data.

A common choice for the **discriminant function** $g_k(x)$ in Bayesian classification is the logarithm of the unnormalized posterior probability:

$$g_k(x) = \log P(x \mid C_k) + \log P(C_k)$$

For our two classes, the discriminant functions are:

1. Discriminant function for "Stop Sign Present" ($C_1$):

$$g_1(x) = \log P(x \mid \text{Stop Sign Present}) + \log P(\text{Stop Sign Present})$$

Here, $P(x \mid \text{Stop Sign Present})$ is the likelihood given the sign is present (provided by the visual detection model), and $P(\text{Stop Sign Present})$ is the prior probability (from maps, location, etc.).

2. Discriminant function for "Stop Sign Not Present" ($C_2$):

$$g_2(x) = \log P(x \mid \text{Stop Sign Not Present}) + \log P(\text{Stop Sign Not Present})$$

Here, $P(x \mid \text{Stop Sign Not Present})$ is the likelihood given no stop sign is present (the "background" model's output), and $P(\text{Stop Sign Not Present})$ is the prior (1 minus the prior for "Stop Sign Present").

The decision rule based on these discriminant functions is:

Classify as **"Stop Sign Present"** if $g_1(x) > g_2(x)$
Classify as **"Stop Sign Not Present"** if $g_1(x) < g_2(x)$
If $g_1(x) = g_2(x)$, the decision boundary is met.

Using Bayesian decision theory with these discriminant functions allows the self-driving car to make decisions that are not solely based on what it *sees* right now, but also take into account what is *expected* to be there based on its prior knowledge, leading to more reliable detection and safer driving.

# Question

In pattern recognition, Maximum A Posteriori (MAP) estimation is often preferred over MLE. Justify this preference using a real-world example (e.g., speech recognition, image classification).

---

# Answer

Both Maximum Likelihood Estimation (MLE) and Maximum A Posteriori (MAP) estimation are methods for estimating parameters of a model based on data.

- **MLE** finds the parameter values $\theta$ that **maximize the likelihood** of observing the given data $x$:

$$\hat{\theta}_{\text{MLE}} = \arg\max_{\theta} P(x \mid \theta)$$

- **MAP** finds the parameter values $\theta$ that **maximize the posterior probability** of the parameters given the data:

$$\hat{\theta}_{\text{MAP}} = \arg\max_{\theta} P(\theta \mid x)$$

Using Bayes' theorem, $P(\theta \mid x) = \frac{P(x|\theta)P(\theta)}{P(x)}$. Since $P(x)$ is a constant with respect to $\theta$, maximizing $P(\theta \mid x)$ is equivalent to maximizing $P(x \mid \theta)P(\theta)$:

$$\hat{\theta}_{\text{MAP}} = \arg\max_{\theta} P(x \mid \theta)P(\theta)$$

The key difference is the term $P(\theta)$, which is the **prior distribution** over the parameters $\theta$. MAP incorporates this prior belief about the parameters before the data is even observed, while MLE does not.

## Justification using a Real-World Example (Image Classification)

Consider training a classifier (e.g., estimating the parameters of probability distributions for different classes, or the weights of a logistic regression classifier) to distinguish between images of **dogs** and **wolves** based on certain visual features.

- **MLE Approach:** MLE would estimate the model parameters (e.g., typical feature values for dogs vs. wolves, or classifier weights) solely based on which parameters make the *training images* most probable under the model for their respective classes.
  - **Problem:** If your training dataset is small or not perfectly representative (e.g., mostly pictures of fluffy poodles and large gray wolves), MLE might learn that "fluffiness" is a very strong indicator of "dog" and "grayness" is a very strong indicator of "wolf". This estimate is overly influenced by the limited data's specific characteristics. If it then sees a large, gray poodle or a fluffy husky (which is a dog), the MLE-trained classifier might misclassify it because the estimated parameters are too extreme or specific to the training examples. MLE is prone to **overfitting** to the noise or specific quirks of limited data.
- **MAP Approach:** MAP estimates the parameters by maximizing the likelihood of the data ( $P(x \mid \theta)$) multiplied by a **prior distribution over the parameters** ($P(\theta)$). We can choose $P(\theta)$ to incorporate beneficial prior knowledge. For example, we might use a prior that favors parameter values that are not too large or too far from zero (a common form of **regularization**).
  - **Benefit:** By incorporating a prior, MAP pulls the parameter estimates towards values that are considered more probable *a priori*. If the limited training data suggests an extreme parameter value (like MLE would choose), the prior can "pull" that estimate back towards a more 'reasonable' or less extreme value.
  - In the dog/wolf example, a prior might encourage the estimated feature distributions for dogs and wolves to have some overlap or keep classifier weights from becoming excessively large. This makes the model less sensitive to small variations or non-representative examples in the training data. It helps to **regularize** the model, preventing it from being overly confident about patterns that might just be noise.

**Conclusion:**

MAP estimation is often preferred in pattern recognition, especially with limited or noisy data, because the inclusion of a **prior distribution** $(P(\theta))$ acts as a form of **regularization**. It allows us to incorporate prior knowledge or simply favor less complex/extreme models, which helps to combat **overfitting** and leads to parameter estimates that generalize better to unseen data compared to MLE.

# Question

Discuss the role of decision surfaces in classification. How does a linear decision boundary differ from a quadratic decision boundary?

---

# Answer

In classification, the goal is to divide the **feature space** (the multi-dimensional space where data points live, defined by their features) into regions, with each region corresponding to a specific class. The **decision surface** (or decision boundary) is what separates these regions.

## Role of Decision Surfaces in Classification

The decision surface is critical because:

1. **It defines the classifier's rule:** For any new data point $x$, its location in the feature space relative to the decision surface determines which class it is assigned. If $x$ falls on one side of the boundary, it's classified into Class A; if it falls on the other side, it's classified into Class B (for a two-class problem).
2. **It represents the learned model:** The decision surface is the geometric manifestation of the trained classifier model. The entire training process aims to find the optimal decision surface that best separates the classes based on the training data.
3. **It visualizes class separation:** In 2D or 3D feature spaces, the decision boundary can be visualized, providing insight into how the classifier distinguishes between classes.

## Linear vs. Quadratic Decision Boundaries

Decision boundaries can have different shapes depending on the underlying classifier model. Two common types are linear and quadratic boundaries:

1. **Linear Decision Boundary:**
   - **Shape:** In a 2D feature space $(x_1, x_2)$, a linear decision boundary is a **straight line**. In higher dimensions $(d > 2)$, it's a flat surface called a **hyperplane**.

- **Equation:** It is defined by a linear equation in the features, like
  $w_1x_1 + w_2x_2 + \cdots + w_dx_d + b = 0.$
- **Classifier Logic:** Classifiers that produce linear boundaries make decisions based on a simple weighted sum of the input features. If the weighted sum exceeds a threshold, the output is one class; otherwise, it's another.
- **Suitability:** Linear boundaries are suitable when the classes are **linearly separable** or nearly so – meaning you can draw a single straight line (or flat hyperplane) to effectively separate the points of different classes.
- **Characteristics:** They are simpler models, require fewer parameters, and are less prone to **overfitting** when data is limited. Examples include the Perceptron, Linear SVM, Logistic Regression, and Linear Discriminant Analysis (LDA) with equal covariance matrices.

## 2. Quadratic Decision Boundary:

- **Shape:** In a 2D feature space, a quadratic decision boundary is a **curve**, such as an ellipse, circle, parabola, or hyperbola. In higher dimensions, it's a **quadratic surface** (a hyperquadric).
- **Equation:** It is defined by a quadratic equation in the features, involving squared terms ($x_i^2$) and cross-product terms ($x_ix_j$) in addition to linear terms and a constant, like
  $ax_1^2 + bx_2^2 + cx_1x_2 + dx_1 + ex_2 + f = 0$ (in 2D).
- **Classifier Logic:** Classifiers that produce quadratic boundaries make decisions based on more complex combinations of features, including their interactions and squared values.
- **Suitability:** Quadratic boundaries are needed when the classes are **not linearly separable** but can be separated by a curved shape.
- **Characteristics:** They are more complex models, require more parameters than linear models, and are more prone to **overfitting** if there isn't enough training data to estimate all the parameters reliably. Examples include Quadratic Discriminant Analysis (QDA) and LDA with unequal covariance matrices.

In summary, the decision surface is the visual representation of a classifier's learned rule. A linear boundary is a flat separation suitable for linearly separable data and less prone to overfitting, while a quadratic boundary is a curved separation needed for non-linearly separable data but requires more complexity and data.

# Question

A robot navigates an unknown environment using Bayesian decision theory to classify obstacles. Given prior probabilities and sensor data likelihoods, formulate the decision rule for safe navigation and explain how it adapts over time.

# Answer

In autonomous robot navigation within an unknown environment, classifying potential obstacles is a critical task. Bayesian decision theory provides a robust framework for making decisions under uncertainty, which is inherent in sensor data and environmental unpredictability.

## Problem Formulation:

- **Classes:** Let's define two classes relevant to navigation safety:
    - $C_1$: Threat (e.g., Obstacle Present)
    - $C_2$: No Threat (e.g., Path is Clear)
- **Sensor Data ($x$):** This is the input to the system, derived from the robot's sensors (e.g., lidar scans, sonar readings, camera images processed into features).
- **Priors ($P(C_k)$):** The robot has a prior belief about the likelihood of a threat in a certain direction or location *before* processing the current sensor data. This prior might come from maps (if available), previous sensor readings in that area, or general assumptions about the environment (e.g., obstacles might be more likely near walls).
- **Likelihoods ($P(x \mid C_k)$):** These models represent how likely it is to observe the specific sensor data $x$ given that the true state is a Threat ($C_1$) or No Threat ($C_2$). These come from the robot's trained sensor models (e.g., "this lidar pattern is highly probable if there's an object at 1 meter").

## Role of Bayesian Decision Theory and Decision Rule

The robot needs to make a decision (e.g., "Move Forward", "Stop", "Turn"). This decision is based on classifying the perceived state (Threat or No Threat). Bayesian decision theory makes this decision by calculating the **posterior probability** $P(C_k \mid x) = \frac{P(x|C_k)P(C_k)}{P(x)}$ and then making a decision that minimizes the expected risk.

In navigation, misclassification costs are often **unequal**. A **False Negative** (classifying a Threat as No Threat) typically has a very high cost (collision), while a **False Positive** (classifying No Threat as Threat) has a lower cost (stopping or taking a suboptimal detour).

Let $\lambda(i \mid j)$ be the cost of deciding class $C_i$ when the true class is $C_j$.

- $\lambda(1 \mid 2)$: Cost of False Positive (decide Threat, true state No Threat)
- $\lambda(2 \mid 1)$: Cost of False Negative (decide No Threat, true state Threat)
- Assume $\lambda(1 \mid 1) = \lambda(2 \mid 2) = 0$ (cost of correct decision is zero).

The decision rule is to choose the class $C_i$ that minimizes the expected risk $R(\alpha_i \mid x)$. For our two classes:

- Expected Risk of deciding Threat ($\alpha_1$):
  $$R(\alpha_1 \mid x) = \lambda(1 \mid 1)P(C_1 \mid x) + \lambda(1 \mid 2)P(C_2 \mid x) = 0 \cdot P(C_1 \mid x) + \lambda(1 \mid 2)P(C_2 \mid x) = \lambda(1 \mid 2)P(C$$
  .

- Expected Risk of deciding No Threat ($\alpha_2$):
  $$R(\alpha_2 \mid x) = \lambda(2 \mid 1)P(C_1 \mid x) + \lambda(2 \mid 2)P(C_2 \mid x) = \lambda(2 \mid 1)P(C_1 \mid x) + 0 \cdot P(C_2 \mid x) = \lambda(2 \mid 1)P(C$$
  .

The **decision rule** is to decide **Threat** if $R(\alpha_1 \mid x) < R(\alpha_2 \mid x)$:

$$\lambda(1 \mid 2)P(C_2 \mid x) < \lambda(2 \mid 1)P(C_1 \mid x)$$

This inequality can be rewritten in terms of the posterior ratio:

$$\frac{P(C_1 \mid x)}{P(C_2 \mid x)} > \frac{\lambda(1 \mid 2)}{\lambda(2 \mid 1)}$$

Let $\lambda_{FP} = \lambda(1 \mid 2)$ and $\lambda_{FN} = \lambda(2 \mid 1)$. The rule is:

$$\text{Decide Threat if } \frac{P(C_1 \mid x)}{P(C_2 \mid x)} > \frac{\lambda_{FP}}{\lambda_{FN}}$$

In navigation, $\lambda_{FN}$ (collision cost) is typically much greater than $\lambda_{FP}$ (detour cost). This means the threshold $\frac{\lambda_{FP}}{\lambda_{FN}}$ is a small value (often much less than 1). This makes the decision boundary **conservative**, biasing the robot towards classifying something as a "Threat" even if the evidence is slightly ambiguous, to minimize the severe cost of missing a real obstacle.

Using Bayes' Theorem $P(C_k \mid x) = \frac{P(x \mid C_k)P(C_k)}{P(x)}$, the rule can also be expressed in terms of likelihoods and priors:

$$\text{Decide Threat if } \frac{P(x \mid C_1)P(C_1)}{P(x \mid C_2)P(C_2)} > \frac{\lambda_{FP}}{\lambda_{FN}}$$

## Adaptation Over Time

Bayesian decision theory naturally supports adaptation as the robot explores the environment:

1. **Updating Priors ($P(C_k)$):** As the robot navigates and potentially builds a map, it learns the spatial distribution of obstacles. If it finds an area is reliably clear, the prior $P(\text{Threat})$ for that specific location or context can decrease over time for future visits. Conversely, if an area is found to be cluttered, the prior can increase. This forms the basis of learned environmental models or occupancy grids.

2. **Refining Likelihoods ($P(x \mid C_k)$ Models):** With more experience, the robot gathers more paired data of sensor readings $x$ and the true outcome (whether an obstacle was actually there or not). This data can be used to refine or re-train the sensor models ($P(x \mid C_1)$ and $P(x \mid C_2)$), making them more accurate in distinguishing threats from clear space based on sensor patterns.

3. **Adjusting Threshold (Implicit Cost/Risk):** While costs $\lambda$ are usually fixed, the threshold $\frac{\lambda_{FP}}{\lambda_{FN}}$ might be dynamically adjusted based on the robot's task or situation (e.g., higher speed might warrant an even more conservative threshold).

By continuously updating its prior beliefs based on experience and refining its sensor interpretation models, the Bayesian approach allows the robot's obstacle classification and navigation decisions to become more informed, accurate, and safer over time in an unknown environment.