# Data Science Project Proposal in:

# Harnessing Data Science to Optimize Service Dog Training in Healthcare and Public Service: A Data-Driven Approach for Enhanced Performance and Outcomes

## 1. Description and Background

Service dogs have expanded their roles to support physical and psychological disabilities, including warning of seizures and alleviating PTSD. Global demand is rising, with the market projected to reach $2.5 billion by 2027.

However, traditional training, which takes 18-24 months, has a low 50% success rate, straining organizations. This project uses data science to optimize training, potentially reducing time by 20% and boosting success rates by 15%. By identifying early training success indicators, organizations can focus on the most promising dogs, improving overall efficiency.

## 2. Project Objective

The primary objective of this project is to improve the efficiency and success rate of service dog training programs through a data-driven approach. By applying data science methods, we aim to:

- **Reduce training time by 20%** through personalized and adaptive training plans based on real-time behavioural data.
- **Increase success rates by 15%** via predicting and focusing resources on dogs with a higher likelihood of successful training.
- Provide **real-time monitoring** of dogs' health and performance, allowing for immediate adjustments in training methods.
- Enhance overall training outcomes and ensure the growing demand for service dogs is met.

## 3. Justification

Service dogs provide vital support for medical and mental health needs, with demand rapidly increasing. However, traditional training methods are slow and resource-intensive, with only a 50% success rate. There is a pressing need for innovative approaches to streamline training and improve outcomes

Data science offers a powerful solution to these challenges by enabling organizations to:

- **Analyse behavioural patterns and progress in real time** using machine learning models.
- **Develop customized training programs** based on each dog's unique capabilities and responses to training.
- **Optimize resource allocation** by identifying dogs with higher potential for successful training early on.
- **Enhance the monitoring of dogs' physical and emotional well-being**, reducing health-related training delays and improving overall success rates.

## *4. Data Science Roles and Responsibilities*

Several key data science roles will be involved in this project, each contributing critical expertise to different aspects of the training optimization process:

- **Data Scientist:**
  - **Responsibilities:**
    - Develop machine learning models to predict training outcomes.
    - Analyse behavioural data and wearable device metrics.
    - Identify traits correlated with training success.
    - Create personalized training recommendations.
  - **Contribution:**
    - Use machine learning to optimize training focus and reduce time.

- **Data Engineer:**
  - **Responsibilities:**
    - Build and maintain data pipelines for training logs and wearable data.
    - Ensure data accessibility and integration of real-time data.
  - **Contribution:**
    - Provide a seamless data foundation for analysis and real-time decisions.

- **System Architect:**
  - **Responsibilities:**
    - Design system architecture for data flow and analysis.
    - Ensure scalability and data security
  - **Contribution:**
    - Enable real-time processing and data-driven decision-making.

- **Data Analyst:**
  - **Responsibilities:**
    - Create reports and visualizations on training progress.
    - Perform exploratory data analysis to identify trends.
    - Provide actionable insights to trainers.

- o **Contribution:**
  - Simplify data for trainers, enabling informed decisions based on real-time feedback.

# 5. *Business Model: Subscription-Based Model*

## 5.1 Business/Application Area

This project focuses on the service dog training industry, crucial for healthcare, disability support, and public safety. The industry is growing rapidly due to rising demand for service dogs assisting with physical disabilities and medical conditions like seizures and heart attacks.

- **Market Growth**: The global service dog market is projected to reach **$2.5 billion** by **2027**.
- **Current Challenges**: Traditional training methods take **18-24 months** with only **50%** success rates.
- **Data-Driven Impact**: By integrating **data science**, training time can be reduced by **20%** and success rates increased by **15%**, improving efficiency and scalability to meet growing demand.

## 5.2 Value Proposition:

This data-driven project creates several key benefits and values for the service dog training industry:

- **Optimized Training Efficiency:**
  - Data science techniques reduce training time by 20% and improve success rates by 15%, allowing more dogs to be trained and available faster, meeting increased demand.
- **Cost Savings for Training Organizations:**
  - Training costs per dog range from $20,000 to $30,000. Improving success rates by 15% could save centers up to $675,000 annually for a 100-dog program.
- **Scalable Solutions for Training Centers:**
  - Smaller centers can access affordable subscription-based models, while larger institutions benefit from enterprise solutions, supporting higher volumes.
- **Improved Service Dog Availability:**
  - Reducing training time by 20% can result in 10,000 additional trained dogs being deployed annually, supporting more individuals with disabilities.

## 5.3 Who Can Benefit:

- **Service Dog Training Centers:**
  - Train more dogs with fewer resources and improved success rates. The subscription model offers flexible pricing for both small and large centers.
- **Healthcare Providers and Organizations:**
  - Hospitals and disability organizations will have quicker access to better-trained service dogs, improving service deployment.
- **Individuals with Disabilities:**

- Benefit from faster, more effective service dog training, improving quality of life for those with physical disabilities, PTSD, and mental health challenges.
  - **Insurance and Government Agencies:**
    - Reduced training costs lead to more efficient allocation of resources, enabling support for a greater number of individuals.

## 5.4 Challenges of the Project:

| | Cause | Effect | Solution/Action | Case Examples |
|---|---|---|---|---|
| **Data Privacy and Security** | Collection of sensitive data from wearable devices | Risk of data breaches, non-compliance with GDPR, HIPAA | Ensure data encryption, compliance with regulations | 1. A breach in a healthcare facility led to exposure of health metrics from wearables; 2. Non-compliance with GDPR resulted in fines for mishandling biometric data; |
| **Technological Integration** | Lack of infrastructure in training centers | Difficulty adopting advanced data and machine learning tech | Provide technical support, offer training and smooth integration | A small training center struggled to integrate wearable data with their outdated database; |
| **Resistance to Change** | Trainers accustomed to traditional methods | Resistance to adopting data science approaches | Demonstrate benefits with strong evidence, offer training | 1. Trainers at a facility were hesitant to use data dashboards, preferring manual tracking; 2. Staff expressed concern over relying on AI predictions for training decisions; |

| Data Quality and Reliability | Inconsistent or poor-quality data from devices and logs | Inaccurate predictions, suboptimal training plans | Implement reliable data collection and monitoring systems | 1. Wearable devices frequently lost connection, resulting in gaps in the data; 2. Training logs recorded inconsistently by trainers led to skewed results; |
|---|---|---|---|---|
| **Scalability in Different Geographies** | Varying regulatory standards and certification globally | Difficulty in applying uniform technology and business model | Build flexible, region-specific solutions | A training center in Europe faced stricter GDPR compliance than those in the U.S; |

## 6. *Characterizing and Analysing Data*

### 6.1 Data Sources

Since I couldn't find a suitable dataset, I created a mock dataset by identifying key features needed for the project, setting reasonable ranges for each, and generating 500 entries using Python. The features closely reflect real-world conditions.

- **Dog_ID**: Unique identifier combining a breed-specific prefix and a sequential number.
- **Dog_Name**: Unique name of the dog with a numeric suffix to ensure uniqueness.
- **Dog_Gender**: The gender of the dog( Male/Female).
- **Dog_Breed**: The breed of the dog, selected from common service dog breeds.
- **Date_of_Birth**: The birth date of the dog(YYYY-MM-DD).
- **Dog_Age**: The age of the dog in years, calculated from 'Date_of_Birth'
- **Training_Start_Date**: The date when the dog's training began(YYYY-MM-DD).
- **Last_Training_Date**: The most recent date of the dog's training session(YYYY-MM-DD).
- **Training_Hours**: The total number of hours the dog has spent in training.
- **Complete_Task**: The number of tasks completed by the dog, ranging from 1 to 15.
- **Health_Score**: A numerical score representing the dog's health status(60.00 - 100.00).
- **Stress_Level**: The stress level of the dog (1.0 - 10.0).
- **Task_Completion_Rate**: The rate of tasks completed by the dog, calculated as Complete_Task divided by total tasks(15 tasks).

o **Time_Index**: The date of the last training session, used as a reference point for time-based analyses.
o **Behavioral_Notes**: Notes describing the dog's behaviour during training sessions.

| Dog_ID | Dog_Name | Dog_Gender | Dog_Breed | Date_of_Birth | Dog_Age | Training_Start_Date | Last_Training_Date | Training_Hours | Complete_Task | Health_Score | Stress_Level | Task_Completion_Rate | Time_Index | Behavioral_Notes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GR001 | Max1 | Female | Golden Retriever | 2020-01-12 | 4.7 | 2020-06-04 | 2021-04-15 | 285.71 | 5 | 78.27 | 9.0 | 0.33 | 2021-04-15 | Max was distracted. |
| BC002 | Daisy1 | Male | Border Collie | 2022-06-11 | 2.3 | 2022-08-27 | 2023-05-07 | 175.56 | 3 | 89.52 | 6.1 | 0.2 | 2023-05-07 | Daisy needs more focus. |
| PO003 | Tucker1 | Female | Poodle | 2019-08-29 | 5.1 | 2019-10-01 | 2024-02-01 | 303.48 | 5 | 76.48 | 4.1 | 0.33 | 2024-02-01 | Tucker executed commands flawlessly. |
| BE004 | Bella1 | Male | Beagle | 2021-09-27 | 3.0 | 2022-05-09 | 2022-09-15 | 206.86 | 4 | 83.44 | 3.4 | 0.27 | 2022-09-15 | Bella showed improvement. |
| DO005 | Bailey1 | Male | Doberman | 2018-07-14 | 6.2 | 2018-09-22 | 2022-02-03 | 345.68 | 6 | 70.19 | 8.8 | 0.4 | 2022-02-03 | Bailey showed good scent tracking skills. |
| BE006 | Cooper1 | Male | Beagle | 2020-03-22 | 4.6 | 2020-09-16 | 2021-03-28 | 286.71 | 5 | 80.66 | 4.4 | 0.33 | 2021-03-28 | Cooper showed great agility and focus. |
| HU007 | Duke1 | Female | Husky | 2018-10-17 | 6.0 | 2019-03-03 | 2023-01-03 | 315.68 | 5 | 74.37 | 6.8 | 0.33 | 2023-01-03 | Duke showed good scent tracking skills. |

## 6.2 Characteristics of the Data
**Volume**:
Generates substantial data over time from multiple dogs, with various metrics (health, stress, task completion) collected across training sessions.
**Velocity**:
Data is updated frequently, potentially in real-time during training, requiring systems that can handle continuous updates.
**Variety**:
Involves diverse data types, including numerical (health, stress), categorical (breed, gender), and textual (behavioural notes).
**Veracity**:
Accuracy and reliability of data are crucial for effective analysis, especially for sensors and trainer inputs.

## 6.3 Platforms, Software, and Tools
- **Software:**
  - **R/Python:** For data visualisation, model building and model predictions. **SAS** and **Julia** can do as well.
  - **SQL:** To communicate with databases for tasks such as updating, retrieving, and managing data. Common databases include **SAS**, **Amazon Redshift**, **Teradata**, **BigQuery**, **Google Cloud**, and **Alibaba Cloud.**
  - **Hadoop/Spark:** used to process large datasets, offering a complete view of risk by analyzing transactional data and identifying patterns and risks in distributed environments.
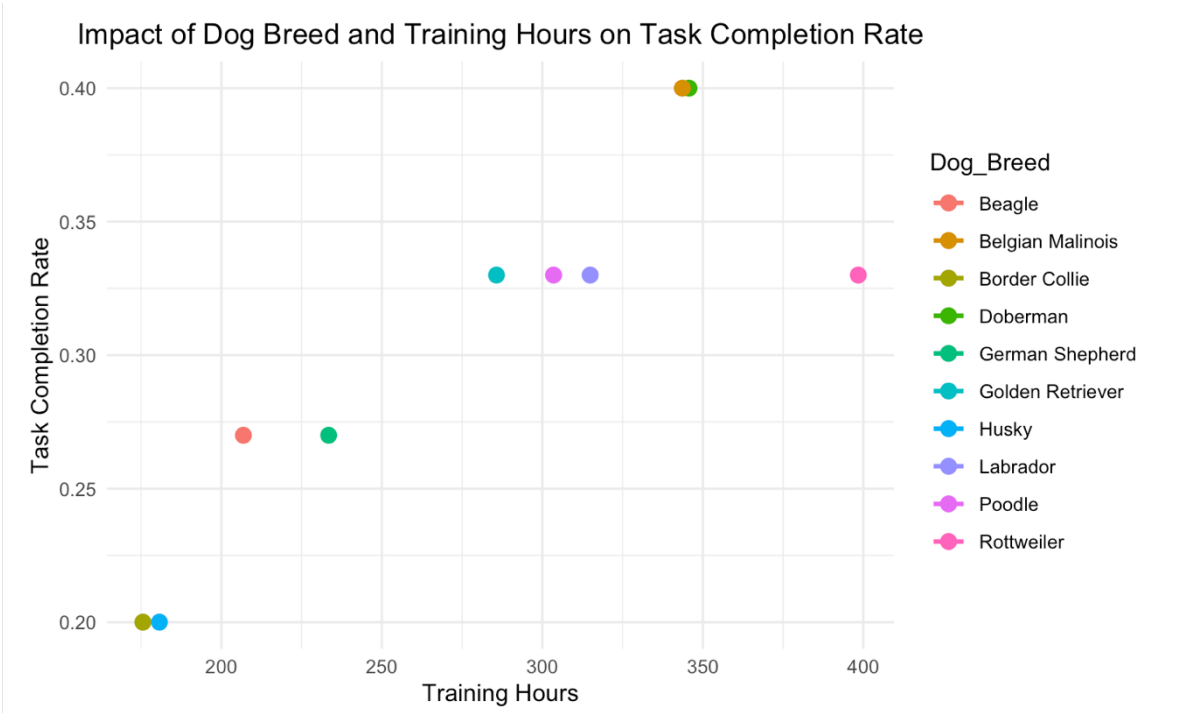
- **Hardware:**
  - **Virtual machines:** To enhance security by preventing tracking of sensitive information.
  - **Large disk space**: For storing large datasets.
  - **High-performance computers**: For running complex algorithms like KNN and K-means.
  - **Good network:** For downloading and frequently updating data.

## 6.4 Data Exploration and Visualisation
Data exploration and visualizations are used to present data in a way that is easily understood by humans. During this process, we can observe how the data is distributed, identify any erroneous values, and evaluate which models might be appropriate for use in later stages.

First, through linear regression analysis, I identified training hours and dog breed as key factors in dog training speed. The coefficient for training hours is 0.0003, showing slight improvements with more hours, while Belgian Malinois has a positive impact (0.087) on learning speed compared to Husky and Border Collie (-0.06). These findings suggest that training strategies should be tailored based on breed and training duration.
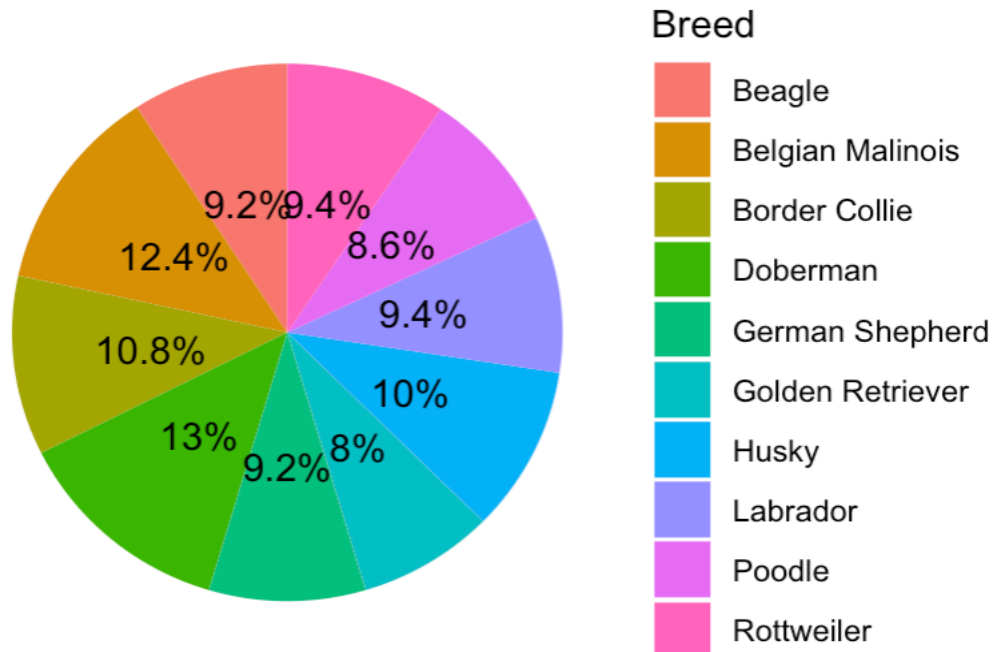


Impact of Dog Breed and Training Hours on Task Completion Rate

```
Call:
lm(formula = Task_Completion_Rate ~ Training_Hours + Dog_Breed,
    data = data)

Coefficients:
              (Intercept)          Training_Hours  Dog_BreedBelgian Malinois
                0.2052044               0.0003132                  0.0871715
    Dog_BreedBorder Collie        Dog_BreedDoberman   Dog_BreedGerman Shepherd
               -0.0601958               0.0865168                 -0.0083226
 Dog_BreedGolden Retriever          Dog_BreedHusky          Dog_BreedLabrador
                0.0353015              -0.0618089                  0.0261613
          Dog_BreedPoodle     Dog_BreedRottweiler
                0.0297353                      NA
```
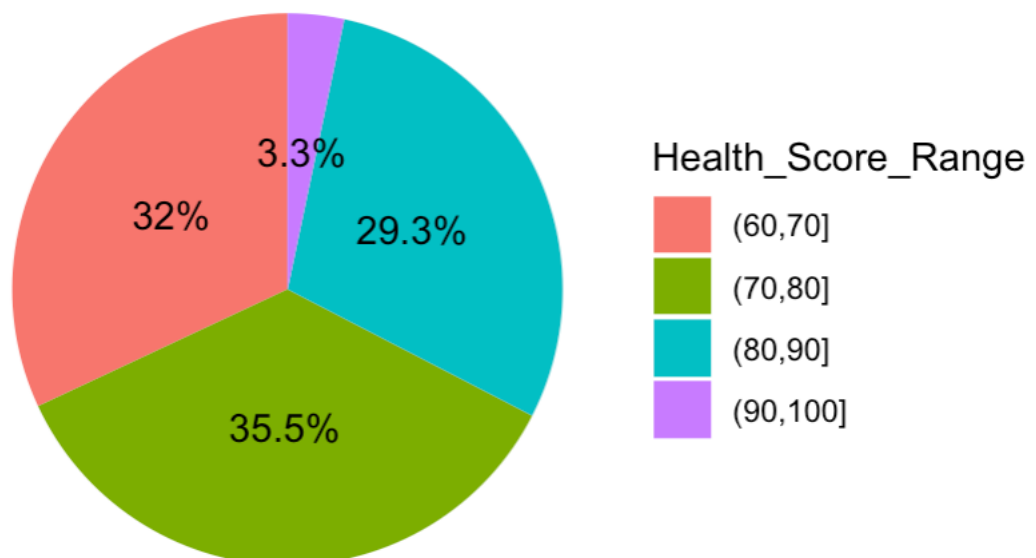
We can know that most service dogs are Golden Retrievers (13%), German Shepherds (12.4%), and Labradors (10.8%), valued for their intelligence and trainability. Other breeds, such as Belgian Malinois, Dobermans, and Rottweilers, are also selected based on task-specific needs rather than popularity.

## Dog Breed Distribution



**Breed**
- Beagle
- Belgian Malinois
- Border Collie
- Doberman
- German Shepherd
- Golden Retriever
- Husky
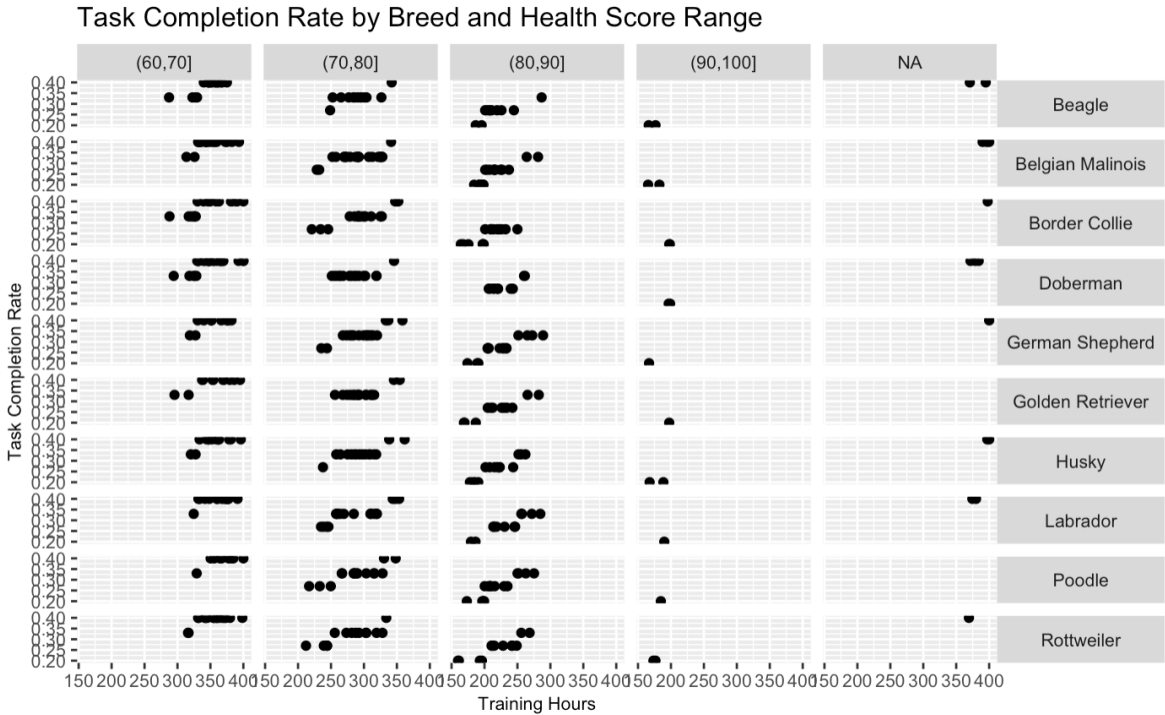- Labrador
- Poodle
- Rottweiler

Most dogs fall within the 70-80 (35.5%) and 60-70 (32%) health score ranges, while only 3.3% of dogs are in the highest range (90-100), indicating that very few dogs are in peak health. This highlights the importance of health monitoring and potential interventions to improve training outcomes for dogs with lower health scores.

## Health Score Distribution



**Health_Score_Range**
- (60,70]
- (70,80]
- (80,90]
- (90,100]

To explore the relationship between breed, health score, and task completion rate, I created a facet grid. The data suggests that both breed and health score are significant

factors influencing training success. Breeds like Golden Retrievers and German Shepherds perform consistently well across different health conditions, while other breeds may require more specialized training approaches. Moreover, dogs with higher health scores tend to have better task completion rates, emphasizing the importance of maintaining good health throughout the training process.


Task Completion Rate by Breed and Health Score Range

## 6.5 Statistical Approach

From this summary, the dataset appears well-balanced in terms of dog age, training hours, and health scores. The distributions of most variables suggest normality, with a few showing potential right skew (e.g., Task Completion Rate). The consistency in training outcomes (Task Completion Rate) and the central tendencies across health and stress scores suggest a relatively uniform population of dogs in training. This balanced distribution will allow for meaningful insights when analysing factors affecting training success.

```
        Dog_Breed         Date_of_Birth           Dog_Age
     Length:500         Length:500          Min.   :2.000
     Class :character   Class :character    1st Qu.:3.500
     Mode  :character   Mode  :character    Median :5.100
                                            Mean   :5.027
                                            3rd Qu.:6.500
                                            Max.   :8.000
     Training_Start_Date Last_Training_Date Training_Hours
     Length:500         Length:500          Min.   :159.8
     Class :character   Class :character    1st Qu.:233.3
     Mode  :character   Mode  :character    Median :291.8
                                            Mean   :288.8
                                            3rd Qu.:343.8
                                            Max.   :400.0
     Complete_Task      Health_Score     Stress_Level
     Min.   :3.000   Min.   :60.00   Min.   : 1.000
     1st Qu.:4.000   1st Qu.:67.51   1st Qu.: 3.100
     Median :5.000   Median :74.44   Median : 5.500
     Mean   :4.908   Mean   :74.81   Mean   : 5.486
     3rd Qu.:6.000   3rd Qu.:82.16   3rd Qu.: 7.825
     Max.   :6.000   Max.   :94.46   Max.   :10.000
     Task_Completion_Rate  Time_Index         Behavioral_Notes
     Min.   :0.2000        Length:500         Length:500
     1st Qu.:0.2700        Class :character   Class :character
     Median :0.3300        Mode  :character   Mode  :character
     Mean   :0.3266
     3rd Qu.:0.4000
     Max.   :0.4000
     Health_Score_Cut
     (60,70] :155
     (70,80] :172
     (80,90] :142
     (90,100]: 16
     NA's    : 15
```

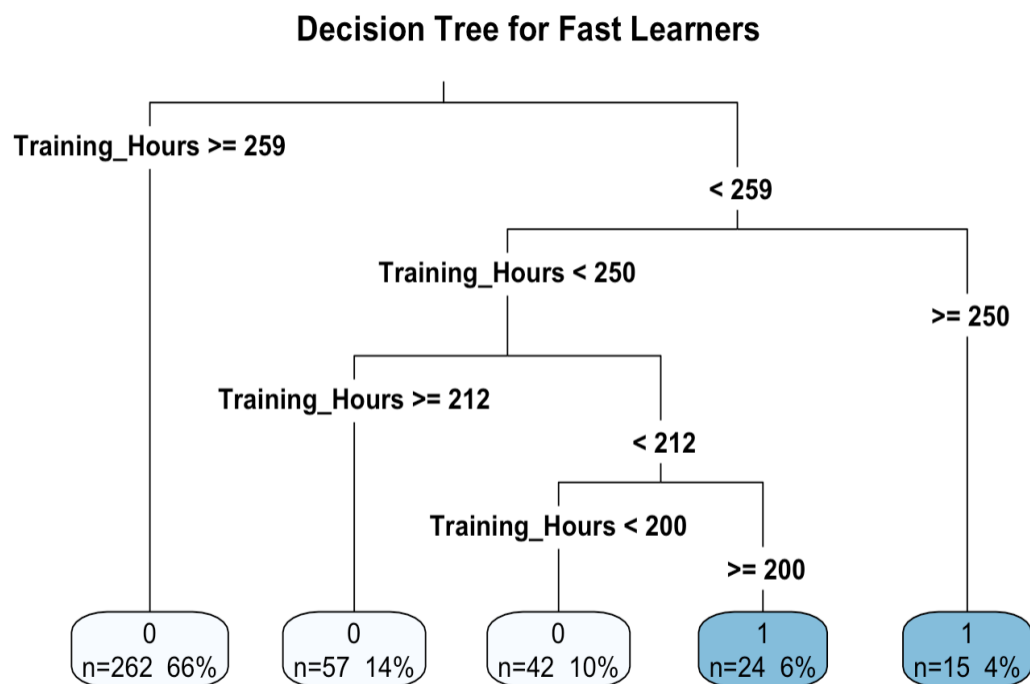### 6.6 Missing Values and Data Processing

Missing values occur due to various reasons, and I plan to address them by using the median for numerical data and the mode for categorical data. Abnormal values, such as health scores outside the 60-100 range or negative training hours, will be corrected or removed to ensure data accuracy. Proper handling of these issues is essential for producing reliable insights into dog training outcomes.

## 7. Data Analysis and Inference

### 7.1 Model Building - Decision Tree

I chose the Decision Tree model for this project because it's easy to interpret and handles non-linear relationships between variables like training hours and health score. It's ideal for binary classification (fast learner vs. not fast learner) and is well-

suited to handle imbalanced data, making it a good fit for identifying the top 10% of fast learners.

## Decision Tree for Fast Learners



### 7.2 Model Inference

From the decision tree result above, we observe that the model mainly relies on the number of training hours to classify dogs as fast learners (label "1") or not (label "0"). The model tends to classify dogs with fewer training hours (less than 259) as potential fast learners. Specifically, dogs that have undergone more than 200 hours but fewer than 250 hours of training are predicted as fast learners.

However, the dataset is imbalanced, with a much higher proportion of non-fast learners (class "0"). This imbalance can make the model biased towards predicting the outcome as "0." Although this model provides a clear separation based on training hours, it may struggle in real-world scenarios due to this imbalance.

In real-world cases, balancing class weights or exploring additional features might help improve the model's performance, as it could better capture patterns related to fast learners beyond just training hours. Additionally, further tuning might be necessary to handle the small number of positive cases in the data.

### 7.3 Model Evaluation

About how about this model, I created a confusion matrix to show.

It demonstrate the decision tree's perfect performance, with an accuracy of 1. Both sensitivity (true positive rate) and specificity (true negative rate) are 1, showing that the model accurately classified all fast and non-fast learners. Furthermore, the positive predictive value (PPV) and negative predictive value (NPV) are 1, indicating no false positives or false negatives. A kappa score of 1 suggests complete agreement between

the predicted and actual values. Despite this excellent performance, there is a risk of overfitting, and the model requires further validation on unseen data.

```
Confusion Matrix and Statistics

          Reference
Prediction  0  1
         0 89  0
         1  0 11

               Accuracy : 1
                 95% CI : (0.9638, 1)
    No Information Rate : 0.89
    P-Value [Acc > NIR] : 8.69e-06

                  Kappa : 1

 Mcnemar's Test P-Value : NA

            Sensitivity : 1.00
            Specificity : 1.00
         Pos Pred Value : 1.00
         Neg Pred Value : 1.00
             Prevalence : 0.89
         Detection Rate : 0.89
   Detection Prevalence : 0.89
      Balanced Accuracy : 1.00

       'Positive' Class : 0
```

## 8. Data Governance and Management

Effective data governance is critical for this service dog training project, ensuring data is handled responsibly and in compliance with relevant regulations. The following principles must be followed:

- Privacy and Confidentiality
  Health scores, performance metrics, and other sensitive data collected from the dogs must be anonymized to protect privacy.

- Data Security
  All data, including behavioral metrics and training outcomes, must be securely

stored using encryption. Both cloud and local storage solutions should be regularly updated and backed up to ensure no data is lost or compromised.

- Ethical Data Use
  The data collected must be used solely for improving service dog training outcomes, with consent from all parties involved. Furthermore, the data collection process must not interfere with the well-being of the dogs or impose unnecessary stress on them.

## 9. Conclusion

In conclusion, while no model is perfect, the data scientist's role is to improve accuracy and optimize service dog training outcomes. Proper data governance and management are crucial for ethical and secure data handling. As technology advances, data science will continue to enhance the training process, benefiting those who rely on service dogs.