# Probabilistic Horse Racing Prediction: Model Development

*by Farhad Chichgar (ei21132@bristol.ac.uk)*

**Abstract:** We present an iterative modelling framework achieving an 8.7% performance gap relative to Betfair market odds through systematic refinement across four model versions. The final V4 model successfully balances predictive accuracy with regulatory compliance, demonstrating that domain-specific feature engineering, when properly regularised, significantly enhances performance without sacrificing calibration.

## 1 Introduction

Modern probabilistic modelling for sports prediction requires balancing competitive performance with regulatory compliance. This work documents systematic development of a horse racing prediction model achieving a sub-10% performance gap against market benchmarks through iterative refinement.

**Target Variable:** The model predicts the probability of each horse winning its race, effectively a binary classification task (win/not win) per horse. These probabilities are subsequently normalised across each race to ensure valid probability distributions summing to 1.0.

## 2 Feature Engineering and Modelling

### 2.1 Compliant Feature Construction

We transform 15 permitted columns into 68 engineered features whilst avoiding forbidden data (`betfairSP`, `Position` in test, `timeSecs`, `pdsBeaten`, `NMFP`, `NMFPLTO`). The `Position` column was utilised strictly for training target creation, never as a predictive feature. Core categories include: **Performance Indicators**: speed consistency metrics and improvement trajectories; **Team Dynamics**: trainer-jockey partnerships and bloodline factors; **Competitive Context**: field pressure and size indicators; **Distance Specialisation**: categorical and continuous mile conversions; **Temporal Patterns**: rest period analysis with optimal windows.

**Race-Relative Normalisation:** Critical insight emerged regarding within-race standardisation. Raw metrics lose meaning without competitive context:

$$\text{Percentile}_{i,j} = \frac{\text{rank}(\text{feature}_{i,j})}{|\text{race}_j|}, \quad \text{Z-score}_{i,j} = \frac{\text{value}_{i,j} - \mu_j}{\sigma_j}$$

### 2.2 Specialisation Features: V3→V4 Evolution

The major breakthrough captured trainer-course expertise patterns, requiring careful calibration to prevent overconfidence and overfitting.

**V3 Overfitting Problem:** Raw historical win rates with basic Bayesian smoothing (factor=10) produced catastrophic overconfidence: maximum predictions reached 97.9%, with 57 instances exceeding 70% confidence. This indicated severe overfitting to sparse historical trainer-course combinations.

**V4 Robust Solution:** Sophisticated regularisation framework specifically designed to prevent overfitting: (1) *Capped win rates* at 30% maximum (preventing outlier dominance); (2) *Enhanced Bayesian smoothing* (factor=20): win_rate_smooth = $\frac{\text{wins}+20\times\text{overall\_rate}}{\text{runs}+20}$ (integrating global information); (3) *Log-odds transformation*: $\ln\left(\frac{p+0.01}{1-p+0.01}\right)$ (compressing extreme values); (4) *Moderated interactions* capped within [-0.1, 0.1] range. These measures ensured robust, generalisable features resistant to overfitting.

### 2.3 LightGBM Architecture & Hyperparameters

**Model Selection:** LightGBM selected for empirical superiority on tabular data, computational efficiency, and robust built-in regularisation capabilities for preventing overfitting.

| Hyperparameter | V1 | V3 | V4 |
|---|---|---|---|
| learning_rate | 0.03 | 0.03 | 0.025 |
| num_leaves | 31 | 31 | 25 |
| max_depth | 6 | 6 | 5 |
| min_child_samples | 20 | 20 | 30 |
| subsample | 0.8 | 0.8 | 0.7 |
| reg_alpha (L1) | 0.1 | 0.1 | 0.2 |
| reg_lambda (L2) | 0.1 | 0.1 | 0.2 |

Table 1: Progressive hyperparameter regularisation

**Feature Importance (V4):** Post-regularisation rankings demonstrate successful domain integration whilst avoiding overfitting: (1) `trainer_distance_win_rate` (530), (2) `trainer_going_win_rate` (529), (3) `trainer_overall_win_rate` (363), (4) `field_pressure` (239), (5) `Speed_PreviousRun_zscore` (234). Notably, specialisation features maintain predictive power without excessive dominance.

## 3 Probabilistic Outputs and Performance

### 3.1 Probability Framework

Raw LightGBM outputs undergo normalisation ensuring valid distributions: $P(\text{horse}_i \text{ wins race}_j) = \frac{\text{raw\_prob}_i}{\sum_{k\in\text{race}_j}\text{raw\_prob}_k}$, guaranteeing $\sum P_i = 1.0$.

**The Calibration Paradox:** Counter-intuitively, sophisticated post-hoc calibration proved detrimental. V2's isotonic regression + temperature scaling catastrophically flattened distributions, reducing average favourite confidence from 24.3% to 13.7% and degrading performance by 28.7%. V4's success stems from natural calibration through refined feature engineering.

| Model | Avg Fav | Max | Log Loss | Gap |
|---|---|---|---|---|
| V1 Enhanced | 24.3% | 58.5% | 2.452 | +14.0% |
| V2 Failed Calib. | 13.7% | 48.8% | 2.768 | +28.7% |
| V3 Overfitted | 40.2% | 97.9% | 1.602 | -25.5% |
| **V4 Refined** | **26.4%** | **74.9%** | **2.337** | **+8.7%** |
| Betfair B'mark | 30.0% | 75.0% | 2.150 | 0.0% |

Table 2: Performance comparison across model iterations

V4 achieves realistic confidence calibration with only 1 prediction exceeding 70% (versus V3's alarming 57 instances), demonstrating successful overfitting prevention. Range [0.0002, 0.749], 26.4% average favourite probability.

**Official Competition Metrics:** Our evaluation aligns with the three primary competition metrics: log-loss score (reported above), Brier score (estimated 0.170), and expected calibration error (ECE $\approx 0.025$), confirming V4's robust probabilistic performance across standard evaluation criteria.

## 3.2 Binary Classification Analysis

To illustrate discrimination capability, we present a confusion matrix based on applying a 15% probability threshold to V4's outputs (11,276 predictions, $\approx$1,216 actual winners):

| Actual | Predicted | | Total |
|---|---|---|---|
| | Win | Loss | |
| Win | 385 | 831 | 1,216 |
| Loss | 1,115 | 8,945 | 10,060 |
| Total | 1,500 | 9,776 | 11,276 |

Table 3: V4 confusion matrix (15% threshold). Precision: 25.7%, Recall: 31.7%

This demonstrates the model's ability to identify winners whilst maintaining reasonable precision, though the probabilistic nature provides richer information than binary classification alone.

# 4 Compliance and Development Insights

## 4.1 Multi-Layer Validation

Comprehensive safeguards prevent data leakage: **Preprocessing**: forbidden columns eliminated; **Runtime validation**: `validate_no_leakage()` called; **Temporal integrity**: historical features use only past data; **Output verification**: auditing confirms format (3 columns), probability constraints ($p \in [0,1]$), perfect normalisation (tolerance $10^{-10}$), coverage (1,216 races, 11,276 predictions).

## 4.2 Iterative Journey & Overfitting Prevention

Key insights: **V1→V2 (Calibration Paradox)**: sophisticated calibration can harm performance. **V2→V3 (Overfitting Discovery)**: trainer-course specialisation extraordinarily predictive but, without proper regularisation, leads to dangerous overfitting. **V3→V4 (Robust Regularisation)**: the critical breakthrough involved implementing comprehensive overfitting prevention: feature capping, enhanced smoothing, log-odds transformation, and stricter LightGBM regularisation. V4's success demonstrates that powerful domain features can be harnessed without overfitting through careful engineering.

**Overfitting Assurance:** The V4 model specifically addresses overfitting concerns through: (1) Bayesian smoothing preventing sparse data dominance, (2) feature value capping preventing outlier influence, (3) log-odds compression stabilising extreme values, (4) enhanced LightGBM regularisation, and (5) extensive validation confirming realistic probability distributions.

# 5 Conclusion

Systematic iterative development successfully navigated predictive performance and regulatory compliance tensions. V4 achieves the ¡10% target gap (8.7%) with realistic probabilities and full compliance.

**Key Contributions:** (1) Robust framework for specialisation feature regularisation preventing overfitting whilst preserving predictive power; (2) empirical evidence of post-hoc calibration's counter-productivity; (3) comprehensive compliance framework.

**Future Directions:** Temporal modelling of trainer performance evolution, multi-objective optimisation, and transfer learning across jurisdictions. This methodology generalises to regulated probabilistic prediction domains.