# Explaining the Cargill Cattle Feed Model: Feature Selection and Interpretation

August 20, 2025

## 1 Introduction

This document explains why the cattle feed choice models, as presented in `cargill_findings_v16.xlsx`, use only three features: Price, Crude Protein (CP), and Brand_Cargill. It provides a stakeholder-friendly explanation for general audiences, technical details for data scientists, and recommendations for future improvements. The models predict Cargill's brand and value shares under price changes (-10% to +10%), and this explanation clarifies the rationale behind the feature selection and its implications.

## 2 Why the Models Use Only Price, CP, and Brand_Cargill

The models focus on three key factors to predict Cargill's market share and assess price sensitivity:

1. **Price**: This is critical as it directly influences farmers' purchasing decisions. Our analysis shows how price changes impact Cargill's brand and value shares, guiding pricing strategies.

2. **Crude Protein (CP)**: This reflects the feed's quality, a key driver of choice in cattle feed markets. Including CP ensures we account for product performance.

3. **Cargill Brand Preference**: This captures how much farmers prefer Cargill over competitors like Baramati or KMF Nandini, independent of price or quality.

### 2.1 Why Other Factors Were Excluded

Other factors, such as feed conversion (FC), animal health (AH), value-added services (VAS), and credit availability, were considered but not included due to:

- **Data Quality**: The dataset had issues, such as 1,714 choice sets dropped due to invalid responses and 1,315 rows removed for missing data. Some factors had inconsistent or missing data, which could lead to unreliable results.

- **Model Reliability**: Including too many factors can make the model unstable, especially with data challenges. Our tests showed that Price, CP, and Brand_Cargill provided a robust and reliable model, avoiding issues like overfitting or failure to converge.

- **Actionable Insights**: The models prioritize price sensitivity and brand strength, which are directly controllable and impactful for business decisions.

### 2.2 What This Means for Stakeholders

The models provide clear, actionable insights:

- In Maharashtra, a 10% price cut slightly boosts Cargill's brand share from 29.68% to 29.81%, but competitors like Baramati (22.89% actual share) remain strong.

- In Punjab, Cargill's strong brand preference (89.28% share) allows flexibility for price increases, with value share rising to 90.27% at +10%.

- In Gujarat and South, shares are stable ( 49.36% and 39.68% at current prices), reflecting balanced data and moderate price sensitivity.

While other factors like services or credit terms matter, their impact is likely captured indirectly through brand preference or requires better data for direct inclusion.

# 3  Technical Notes for Experts

The models are custom multinomial logit (MNL) models, fitted after failures of `pylogit` MNL, simpler MNL, and binary logit, as shown in `cattle_feed_model.log`. The utility function is:

$$U = \beta_{\text{Price}} \cdot \text{Price} + \beta_{\text{CP}} \cdot \text{CP} + \beta_{\text{Brand\_Cargill}} \cdot \text{Brand\_Cargill}$$

Coefficients are listed in Table 1.

Table 1: Model Coefficients

| Market | Price | CP | Brand_Cargill |
|---|---|---|---|
| Maharashtra | -0.030027 | 0.247563 | -0.862332 |
| Gujarat | -0.056879 | -0.291053 | -0.025527 |
| Punjab | 0.061117 | 0.176946 | 2.119954 |
| South | 0.036905 | 0.042858 | -0.420083 |

## 3.1  Feature Selection Rationale

- **Multicollinearity**: The log shows low variance inflation factors (VIF ~1.0) for Price, CP, and Brand_Cargill. Other features (e.g., FC, AH, VAS, Credit) were converted to dummy variables but likely excluded due to high VIF (>100), as per the `check_multicollinearity` method.

- **Convergence**: `pylogit` MNL failed due to singular matrices or insufficient variation. The custom MNL succeeded with L2 regularization, but limiting features to three ensured stability.

- **Data Issues**: 2,199 duplicate choice sets, 1,315 rows with NaN, and 1,714 invalid choice sets were dropped. Features like FC or VAS may have had missing or low-variation data, reducing their reliability.

## 3.2  Impact on Predictions

Predictions in `cargill_predicted_shares.xlsx` use Price and Brand_Cargill, with CP held constant at its mean (0). The high shares in Gujarat, Punjab, and South (~50%) reflect data balancing, inflating predictions compared to actual shares (e.g., Maharashtra: 8.77%).

# 4  Recommendations for Future Improvements

1. **Add Features**: With cleaner data, include FC, AH, VAS, or Credit. Check VIF and handle multicollinearity (e.g., PCA or feature selection).

2. **Improve Data**: Fix invalid choice sets (e.g., select one Chosen=1 per set) to retain data. Validate feature distributions for sufficient variation.

3. **Refine Model**: Test interaction terms (e.g., Price $\times$ Brand_Cargill) or mixed logit models for heterogeneity.

4. **Validate**: Compare predictions with sales data and investigate positive price coefficients (Punjab, South) for data artifacts.

## 5   Conclusion

The models focus on Price, CP, and Brand_Cargill balances reliability and business relevance, providing actionable insights for pricing and competitive strategy. Future improvements can incorporate additional factors with better data quality.