

Mapping the Latent Structure of Economic Bias: A DML-LSM Fusion Framework

Team Sharks

November 22, 2025

Abstract

This report introduces a novel methodological framework that fuses Double Machine Learning (DML) with Latent Space Mapping (LSM) to visualize the topology of economic bias in professional sports. By addressing the limitations of deterministic salary structures through a stratified learning approach, we isolate the true market price of bias factors. We then project these effects into a high-dimensional Attribution Matrix, which is unfolded into a low-dimensional map. This transformation moves beyond static coefficient tables, revealing the geometric structure of bias and identifying distinct player clusters driven by shared structural forces.

1 Introduction

1.1 The Analytical Gap: The “Wall of Coefficients”

Traditional econometric analyses of salary discrimination and market inefficiency predominantly rely on regression-based techniques. While robust for hypothesis testing, these methods typically culminate in a “wall of coefficients”—a static list of penalties and premiums (e.g., “Age reduces salary by β_{age} ”).

While these coefficients successfully identify *what* factors are statistically significant, they fail to reveal the *system* of bias. Standard regression outputs cannot easily visualize collinearity in the error terms or reveal latent archetypes among the population. For example, a coefficient table cannot intuitively demonstrate if “Draft Pedigree” and “Market Size” act as independent forces or if they coalesce to form a singular “Hype Bias” that affects a specific cluster of players.

1.2 The Deterministic Contract Problem

A significant challenge in analyzing professional sports markets (such as the NBA) is the prevalence of non-market compensation. A substantial subset of the population operates under deterministic contracts—specifically, Rookie Scale and Maximum Salary contracts—where compensation is dictated by Collective Bargaining Agreement (CBA) rules rather than statistical performance valuation.

Applying standard statistical models to this mixed population constitutes model misspecification. A high-performing rookie may appear to have a large negative residual not because of market bias, but because of a salary ceiling. Training a bias model on these artificially constrained salaries would bias the learned coefficients, obscuring the true market valuation of contextual factors.

1.3 The Proposed Contribution: A Fusion Framework

This paper introduces a novel two-stage framework that integrates **Double Machine Learning (DML)** with **Latent Space Mapping (LSM)**. This fusion leverages the strengths of both disciplines: DML provides the rigorous, causal isolation of the “price” of bias, while LSM provides the geometric tools to visualize the relationships between these biases.

Our approach shifts the analytical goal from simply *quantifying* bias to *mapping* its topology. By treating the magnitude of salary attribution as a geometric attractor, we generate a visual landscape where:

- **Bias Factors** act as fixed anchors in a latent space.
- **Players** cluster around the specific factors that drive their valuation.
- **Distance** encodes the magnitude of economic sensitivity.

This framework provides stakeholders with a diagnostic tool to identify market inefficiencies, separating players whose salaries are driven by pure on-court performance from those whose valuations are heavily distorted by structural biases.

2 Stage 1: Isolating the Market Signal (The DML Engine)

The first stage of our framework is strictly econometric. Our objective is to isolate the causal effect of contextual bias factors (vector Z) on player salary (Y), independent of on-court performance statistics (vector X). To achieve this in the presence of deterministic contracts, we employ a **Stratified Double Machine Learning (DML)** approach.

2.1 The Statistical Objective

We assume a partially linear structural model for player valuation:

$$Y_i = \alpha + Z_i^T \gamma + f(X_i) + \epsilon_i \quad (1)$$

Where:

- Y_i is the log-salary of player i .
- $Z_i \in \mathbb{R}^m$ is the vector of bias factors (e.g., Age, Draft Number, Market Size).
- $X_i \in \mathbb{R}^p$ is the high-dimensional vector of performance metrics.
- $f(X_i)$ is an unknown, potentially non-linear function representing the “fair” market value of performance.
- $\gamma \in \mathbb{R}^m$ is the vector of coefficients representing the marginal market price of each bias factor.

Standard regression fails here because X is a confounder: performance affects both salary and the bias factors (e.g., better players are often drafted higher). To recover an unbiased estimate of γ , we use the Frisch-Waugh-Lovell (FWL) theorem extended by DML. This involves training flexible machine learning models to “residualize” both salary and bias factors, stripping away the variance explained by performance.

2.2 Methodology: The “Learn vs. Apply” Framework

A standard DML implementation would fail in the NBA context due to the **Deterministic Contract Problem** described in Section 1.2. Training the model on the full population would allow the artificial salary caps of Rookie and Max contracts to contaminate the estimate of market valuation.

To resolve this, we introduce a stratified **Learn-Apply** protocol:

2.2.1 Step 1: Learning the Market Price (Training Phase)

We first isolate the sub-population of players operating under negotiated, free-market contracts, denoted as D_{FM} . We train our nuance parameters solely on this group to learn the true market dynamics.

For each player $i \in D_{FM}$, we estimate:

1. **The Outcome Model:** We train a gradient boosting model $\hat{g}(X)$ to predict salary from stats. The residual $\hat{\epsilon}_{Y,i} = Y_i - \hat{g}(X_i)$ represents the portion of salary unexplained by on-court performance.
2. **The Treatment Models:** For each bias factor $Z^{(j)}$, we train a model $\hat{h}_j(X)$ to predict the factor from stats. The residual $\hat{\epsilon}_{Z^{(j)},i} = Z_i^{(j)} - \hat{h}_j(X_i)$ represents the portion of the bias factor orthogonal to performance.

We then regress the outcome residuals on the treatment residuals via OLS to obtain the vector of market prices, $\hat{\gamma}$:

$$\hat{\epsilon}_Y = \sum_{j=1}^m \hat{\gamma}_j \hat{\epsilon}_{Z^{(j)}} + \nu \quad (2)$$

2.2.2 Step 2: Counterfactual Application (Inference Phase)

Once the market prices ($\hat{\gamma}$) and nuisance models (\hat{g}, \hat{h}) are fixed, we apply them to the **entire** player population D_{ALL} , including those on fixed contracts.

For a Rookie player $k \in D_{Rookie}$, we calculate their **Counterfactual Residuals**:

$$\hat{\epsilon}_{Z^{(j)},k} = Z_k^{(j)} - \hat{h}_j(X_k) \quad (3)$$

This quantity represents how much “excess bias” the rookie possesses relative to their performance (e.g., how much earlier they were drafted than their stats would justify). By applying the free-market price $\hat{\gamma}$ to these residuals in the next stage, we can estimate the latent market pressure acting on the rookie, even if their actual salary is censored by the CBA.

3 Stage 2: The Attribution Matrix (The Bridge)

The second stage of our framework acts as the bridge between the econometric analysis (DML) and the geometric visualization (LSM). The raw output from Stage 1 consists of a vector of bias prices, $\hat{\gamma} \in \mathbb{R}^m$, and a matrix of residuals, $\hat{\epsilon}_Z \in \mathbb{R}^{n \times m}$.

While statistically rigorous, these raw outputs are geometrically incompatible. The bias factors have incommensurable units (e.g., Age in years vs. Market Size in millions of dollars), making a direct Euclidean embedding impossible. To resolve this, we construct the **Attribution Matrix** (L).

3.1 Constructing the Matrix

We define the Attribution Matrix $L \in \mathbb{R}^{n \times m}$ such that each entry L_{ij} quantifies the total economic impact of bias factor j on the salary of player i . The transformation is defined as:

$$L_{ij} = |\hat{\gamma}_j \cdot \hat{\epsilon}_{Z^{(j)}, i}| \quad (4)$$

This formulation achieves two critical objectives:

1. **Unit Unification:** By multiplying the residual (in units of Z) by the price (in log-dollars per unit of Z), we convert all factors into a single, unified currency: *log-dollars of unexplained salary*.
2. **Focus on Relevance:** By taking the absolute value, we shift the analytical focus from “direction” (overpaid vs. underpaid) to “magnitude” (relevance). A high value of L_{ij} indicates that factor j is a primary driver of player i ’s valuation variance, regardless of whether it inflates or suppresses their salary.

3.2 Geometric Interpretation

This matrix L serves as the input “liking” or “proximity” matrix for the Latent Space Mapping algorithm. In the context of our map:

- **High Attribution ($L_{ij} \gg 0$):** Implies a strong structural link. The player’s salary is heavily determined by this factor. Geometrically, the player will be pulled *close* to the factor’s anchor point.
- **Low Attribution ($L_{ij} \approx 0$):** Implies independence. The player’s salary is unaffected by this factor (either they have no residual bias, or the market price for that bias is zero). Geometrically, the player will be pushed *away* from the factor.

This transformation allows us to visualize the *drivers of valuation*. Clusters on the resulting map represent groups of players whose economic fate is governed by the same set of structural biases.

4 Stage 3: Latent Space Mapping (The Visualization Engine)

Once the Attribution Matrix L is constructed, our task shifts from econometrics to geometry. We seek to visualize this high-dimensional matrix in a low-dimensional space (typically \mathbb{R}^3) to reveal its latent structure.

4.1 The Geometric Model: A Probabilistic Formulation

We model the relationship between players and bias factors not as fixed points, but as probabilistic interactions. We represent each entity as a Gaussian distribution to capture inherent uncertainty.

Let the position of player i be a random variable $c_i \sim \mathcal{N}(\mu_{c,i}, V_{c,i})$, and the position of bias factor j be $p_j \sim \mathcal{N}(\mu_{p,j}, V_{p,j})$.

We define the *difference distribution* for each player-factor pair (i, j) . Since c_i and p_j are independent, their difference $\delta = c_i - p_j$ is also Gaussian:

$$\delta \sim \mathcal{N}(\mu_{ij}, V_{ij}), \quad \text{where } \mu_{ij} = \mu_{c,i} - \mu_{p,j} \quad \text{and} \quad V_{ij} = V_{c,i} + V_{p,j}. \quad (5)$$

Our similarity function, $g(\delta)$, maps this difference vector to a predicted attribution score. We use the Gaussian kernel, $g(\delta) = \exp(-\delta^\top \delta)$, which corresponds to $\exp(-d^2)$ where $d = \|\delta\|_2$.

4.2 Derivation of the Expected Attribution

The predicted attribution, \hat{L}_{ij} , is the expected value of this similarity function, taken over the entire distribution of possible difference vectors δ . Mathematically, this is defined by the multi-dimensional integral:

$$\begin{aligned}\hat{L}_{ij} &= \mathbb{E}_{\delta \sim \mathcal{N}(\mu_{ij}, V_{ij})}[g(\delta)] \\ &= \underbrace{\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \frac{1}{\sqrt{(2\pi)^d |V_{ij}|}} \exp\left(-\frac{1}{2}(\delta - \mu_{ij})^\top V_{ij}^{-1}(\delta - \mu_{ij})\right)}_{\text{Probability of difference } \delta} \underbrace{\exp(-\delta^\top \delta)}_{\text{Similarity kernel}} d\delta.\end{aligned}$$

This integral represents the theoretical heart of the model. It is a weighted average over every possible perceptual difference δ , summing the probability-weighted similarity scores to yield a single expected value.

4.2.1 The Closed-Form Solution via Moment-Generating Functions

While the integral above is the correct definition, it is computationally intractable to solve numerically within an optimization loop. Fortunately, for our specific choice of a Gaussian kernel, this integral has an elegant, closed-form solution.

The integral is an application of the *moment-generating function (MGF)* for a quadratic form of a Gaussian random variable. Specifically, we are computing the MGF of the random variable $Q = \delta^\top I \delta$ evaluated at $t = -1$. This standard result provides the exact solution:

$$\hat{L}_{ij} = \underbrace{|I + 2V_{ij}|^{-1/2}}_{\text{Variance Penalty}} \underbrace{\exp\left(-\mu_{ij}^\top (I + 2V_{ij})^{-1} \mu_{ij}\right)}_{\text{Adjusted Distance Penalty}}, \quad (6)$$

where I is the identity matrix and $|\cdot|$ denotes the determinant.

This algebraic formula is the computational engine of our model, allowing us to calculate the exact expected attribution without numerical integration. It naturally penalizes uncertainty: as the variance V_{ij} increases, the expected attribution \hat{L}_{ij} decreases, reflecting lower confidence in the link.

4.3 Objective Function: Global Maximum Likelihood

We solve for the optimal map coordinates $\theta = \{\mu_c, \mu_p, V_c, V_p\}$ by minimizing the error between the observed attribution matrix L (from Stage 2) and our model's prediction $\hat{L}(\theta)$.

To ensure robustness against noise, we employ a **Maximum Likelihood Estimation (MLE)** framework with a single, global noise parameter σ_{err} . We assume the observed values are drawn from a Gaussian around the prediction:

$$L_{ij} \sim \mathcal{N}(\hat{L}_{ij}(\theta), \sigma_{err}^2) \quad (7)$$

This leads to the Negative Log-Likelihood (NLL) objective function:

$$\mathcal{J}(\theta, \sigma_{err}) = N \log(\sigma_{err}) + \frac{1}{2\sigma_{err}^2} \sum_{i,j} (L_{ij} - \hat{L}_{ij}(\theta))^2 \quad (8)$$

Minimizing this objective is mathematically equivalent to a precision-weighted Sum of Squared Errors (SSE). The global σ_{err} acts as an adaptive regularization term, preventing the model from overfitting to the idiosyncratic noise in the DML residuals.

4.3.1 Justification for Gaussian Likelihood

The choice of a Gaussian likelihood function is not arbitrary; it provides a direct link to the intuitive geometric goal of minimizing squared errors.

Under the assumption of a constant global error variance σ_{err}^2 , maximizing the likelihood is mathematically equivalent to minimizing the Sum of Squared Errors (SSE). The log-likelihood function for our Gaussian model is:

$$\ln \mathcal{L}(\theta, \sigma_{err}) = -\frac{N}{2} \ln(2\pi\sigma_{err}^2) - \frac{1}{2\sigma_{err}^2} \sum_{i,j} (L_{ij} - \hat{L}_{ij}(\theta))^2 \quad (9)$$

Maximizing this expression with respect to the map coordinates θ (while holding σ_{err} constant) requires minimizing the only term that depends on θ :

$$\text{maximize } \ln \mathcal{L} \iff \text{minimize } \sum_{i,j} (L_{ij} - \hat{L}_{ij}(\theta))^2 \quad (10)$$

This equivalence demonstrates that our probabilistic MLE approach is a principled generalization of standard least-squares unfolding. It preserves the intuitive geometric interpretation of minimizing errors while adding the statistical robustness of a learned noise parameter.

4.4 Optimization Strategy

The objective function \mathcal{J} is non-convex. To find a stable global minimum, we implement a three-stage optimization protocol:

1. **Initialization:** We use Principal Component Analysis (PCA) on the Attribution Matrix to generate a rigorous initial guess for the coordinates.
2. **Alternating Optimization:** We iteratively update the player coordinates while fixing factors, and vice-versa, to descend rapidly into the correct basin of attraction.
3. **JAX Acceleration:** The entire optimization loop is Just-In-Time (JIT) compiled using Google's JAX library, allowing for thousands of gradient descent steps to be computed in seconds.

5 Interpretation and Conclusion

The final output is an interactive 3D scatter plot.

- **Bias Anchors:** Fixed points representing the bias factors.
- **Player Clusters:** Groups of players clustering around the factors that drive their valuation.
- **The Void:** The center of the map, representing players with low residual bias (pure performance valuation).

By fusing the causal identification of DML with the rigorous probabilistic geometry of PULS, we provide a new lens for understanding market inefficiencies.

6 Interpretation of the Map

The final output is an interactive 3D scatter plot that maps the topology of market bias. Interpretation relies on analyzing the geometric relationships between the two sets of points:

6.1 The Geometry of Influence

- **Bias Anchors (Diamonds):** The bias factors (e.g., “Age”, “Draft”) act as fixed anchors in the latent space. Their relative positions reveal the correlation structure of the biases themselves. Anchors that appear close together (e.g., “Market Size” and “Team Revenue”) affect the same players in similar ways, suggesting a shared underlying economic force.
- **Player Clusters (Points):** Players cluster around the specific anchors that drive their salary.
 - *Example:* A dense cluster of players forming around the “Draft Number” anchor indicates a segment of the league where draft pedigree is the primary determinant of pay. This cluster is likely to contain Rookies and early-career stars whose contracts are rigidly tied to their entry position.
- **The “Void” (Center):** Players located in the geometric center of the map are those with low attribution scores across *all* bias factors. These represent the “Pure Performance” players—their salaries are almost entirely explained by their on-court statistics (X), with very little residual variance attributable to contextual biases (Z).

6.2 Metadata Overlay and Validation

By coloring the player points according to their Contract Type (Rookie vs. Free Market vs. Max), we can visually validate our structural hypotheses. For instance, if the model is correct, we should observe a distinct spatial separation where Rookies cluster tightly around the “Draft Bias” anchor, while Free Market veterans disperse towards anchors like “Age” or “Experience.” This visual confirmation transforms the abstract statistical output into an intuitive diagnostic tool for market analysis.

7 Conclusion

This report has detailed a rigorous framework for mapping the latent structure of economic bias. By fusing the causal identification power of Stratified Double Machine Learning with the geometric visualization capabilities of Latent Space Mapping, we have created a tool that transcends traditional regression analysis.

We successfully transformed the problem from a static list of coefficients into a dynamic topological map. This framework allows stakeholders—from General Managers to Agents—to identify “market inefficiencies” not just by magnitude, but by type. It reveals not just *if* a player is misvalued, but *why*, and situates them within a peer group facing the same structural economic forces.