

Estimating the Conditional Average Treatment Effect (CATE) of Credit Access Using Causal Forests in Conjunction with Double Machine Learning

Warade Atharv Abhijit - 2022582

Yash Sinha - 2022590



INDRAPRASTHA INSTITUTE *of*
INFORMATION TECHNOLOGY
DELHI



Problem Statement



Credit: Whom Does It Help?

- The core challenge is identifying who benefits from agricultural credit.
- Farmers possess diverse resources, constraints, and productivity levels.
- Access to credit is expected to improve their agricultural output.
- We must estimate the credit effect across different farmer types.
- A single "average effect" obscures the substantial heterogeneity of impact.

The Problem with "Average" Effects



The Flaw of ATE: Traditional econometrics focuses on the **Average Treatment Effect (ATE)**.

- *Question:* "Does credit work on average?"
- *Critique:* This assumes that a loan helps a wealthy landowner in Kapanimbargi exactly the same way it helps a marginal smallholder in Aurepalle .

The Reality of "Problem Soils":

- The return on financial investment in farming is fundamentally dependent on complementary assets, such as having fertile land and necessary resources like nitrogen (fertilizer) and operational equipment
- *Hypothesis:* Financial inclusion and providing credit may be ineffective for farmers if they face critical environmental limitations, such as severely degraded soil or a lack of adequate rainfall.

The ICRISAT Data and Variables



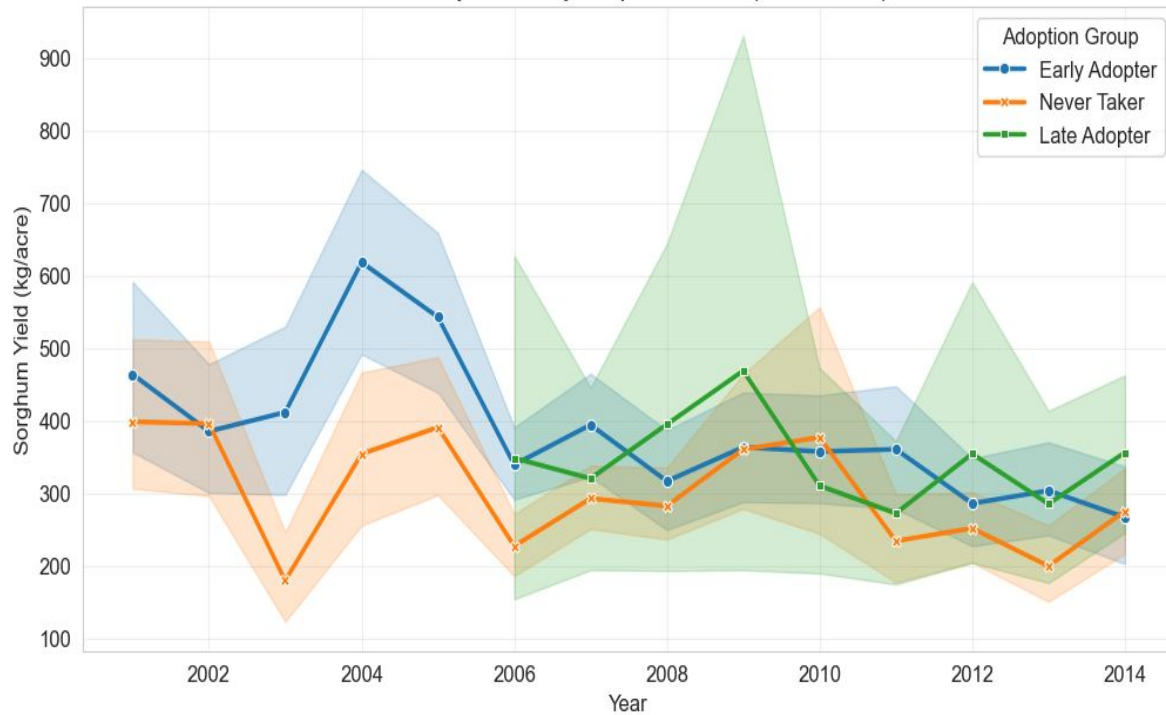
- **Scope:** Panel of 462 households across 14 years (2001–2014).
- **Treatment (D):** `credit_access` (Binary: 1 if accessed formal/informal credit, 0 otherwise).
- **Outcome (Y):** `crop_yield` (Sorghum yield in kg/acre).
- **High-Dimensional Covariates (X):**
 - **Biophysical:** `vdeepsoil_plotcount` (Deep Vertisols), `problemsoil_plotcount` (Saline/Erosive), `irrigation_indicator`.
 - **Economic:** `wealth_index`, `operational land` (Farm size).
 - **Inputs (Lagged):** `lag_nitropa` (Nitrogen), `lag_motorpa` (Mechanization).
Used lags to avoid post-treatment bias.

EDA

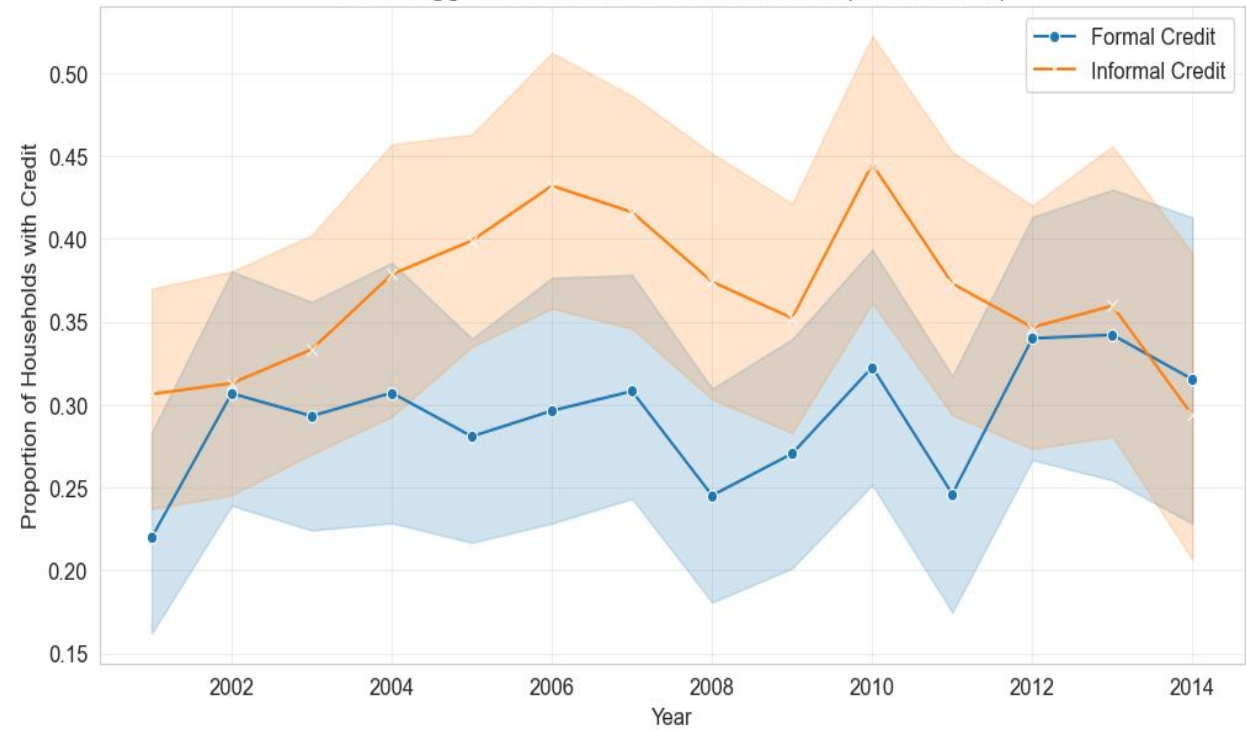


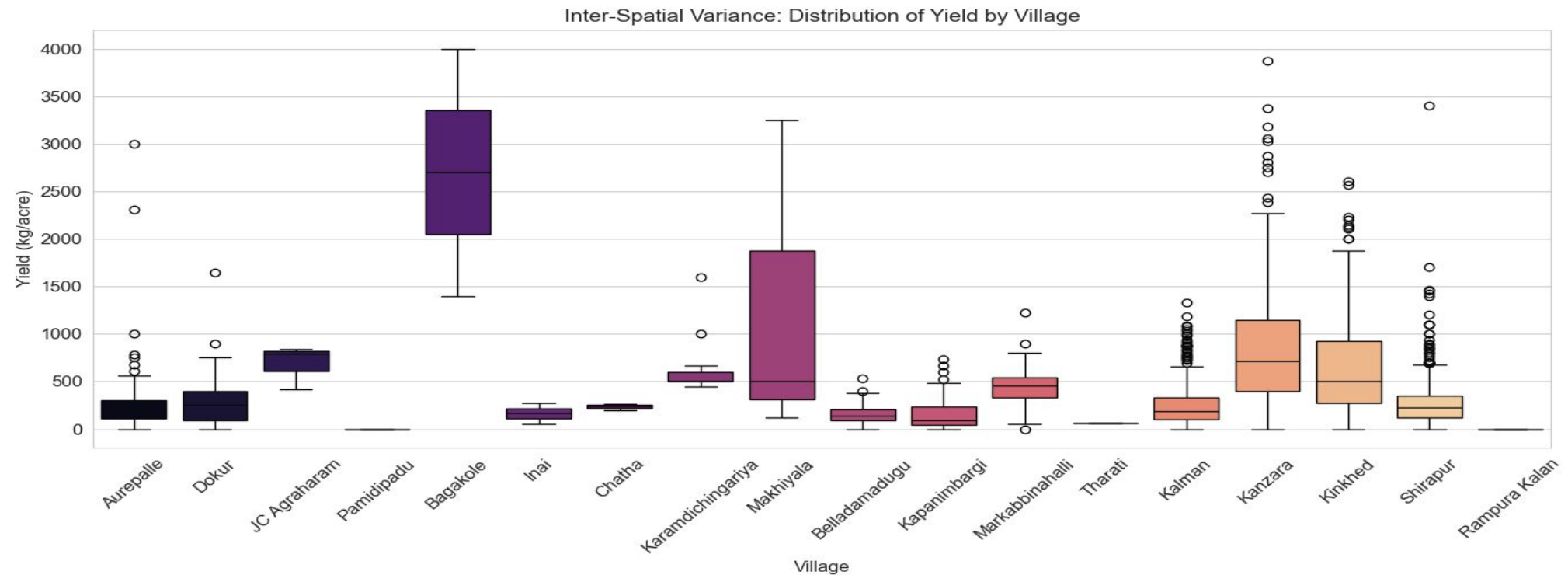
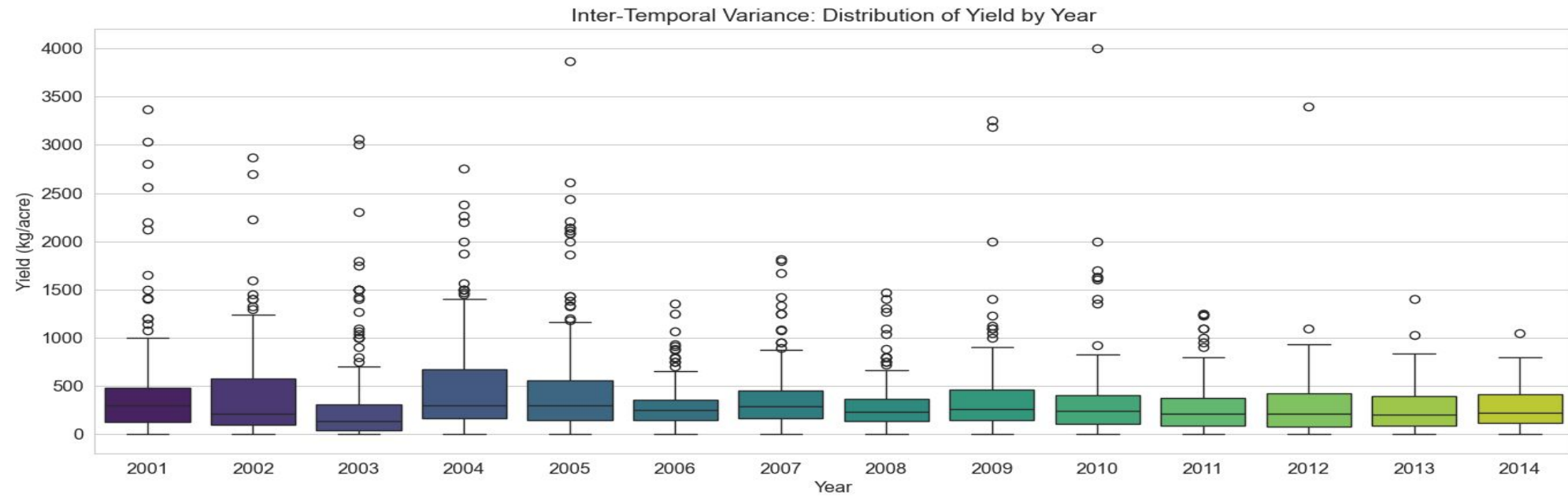
column	NA	NA %	mean	sd	p0	p25	p50	p75	p100
year	0	0	2007	3.851	2001	2004	2007	2010	2014
plotcount	0	0	1.35	0.6966	1	1	1	2	7
problemsoil_	13	0.60437006043	0.5463	0.6674	0	0	0	1	4
plotcount		7006							
alkaline_aci	13	0.60437006043	0.01964	0.1454	0	0	0	0	2
dic_plotcoun		7006							
t									
erosive_plot	18	0.83682008368	0.3638	0.593	0	0	0	1	4
count		20083							
deepsoil_plo	18	0.83682008368	0.2058	0.4913	0	0	0	0	5
tcount		20083							
vdeepsoil_pl	18	0.83682008368	0.2639	0.5641	0	0	0	0	4
otcount		20083							
croparea	0	0	2.855	3.317	0.025	1	2	3.5	43.5
yield	0	0	359	421.9	0	116.7	240	440	4000
fertilizer_f	0	0	2.615	2.992	0	0	2	4	26
requency									
fertilizer_i	0	0	0.5997	0.4668	0	0	1	1	1
ndicator									
irrigation_f	0	0	0.7146	1.705	0	0	0	1	28
requency									
irrigation_i	0	0	0.2223	0.3939	0	0	0	0.3333	1
ndicator									
motorpa	0	0	0.1796	0.5845	0	0	0	0	10
nitropa	0	0	0.2969	0.5022	0	0	0	0.4091	7.667
phospa	0	0	0.2468	0.4061	0	0	0	0.384	3.594
potashpa	0	0	0.03191	0.1365	0	0	0	0	1.857

Yield Trajectories by Adoption Cohort (with 95% CI)

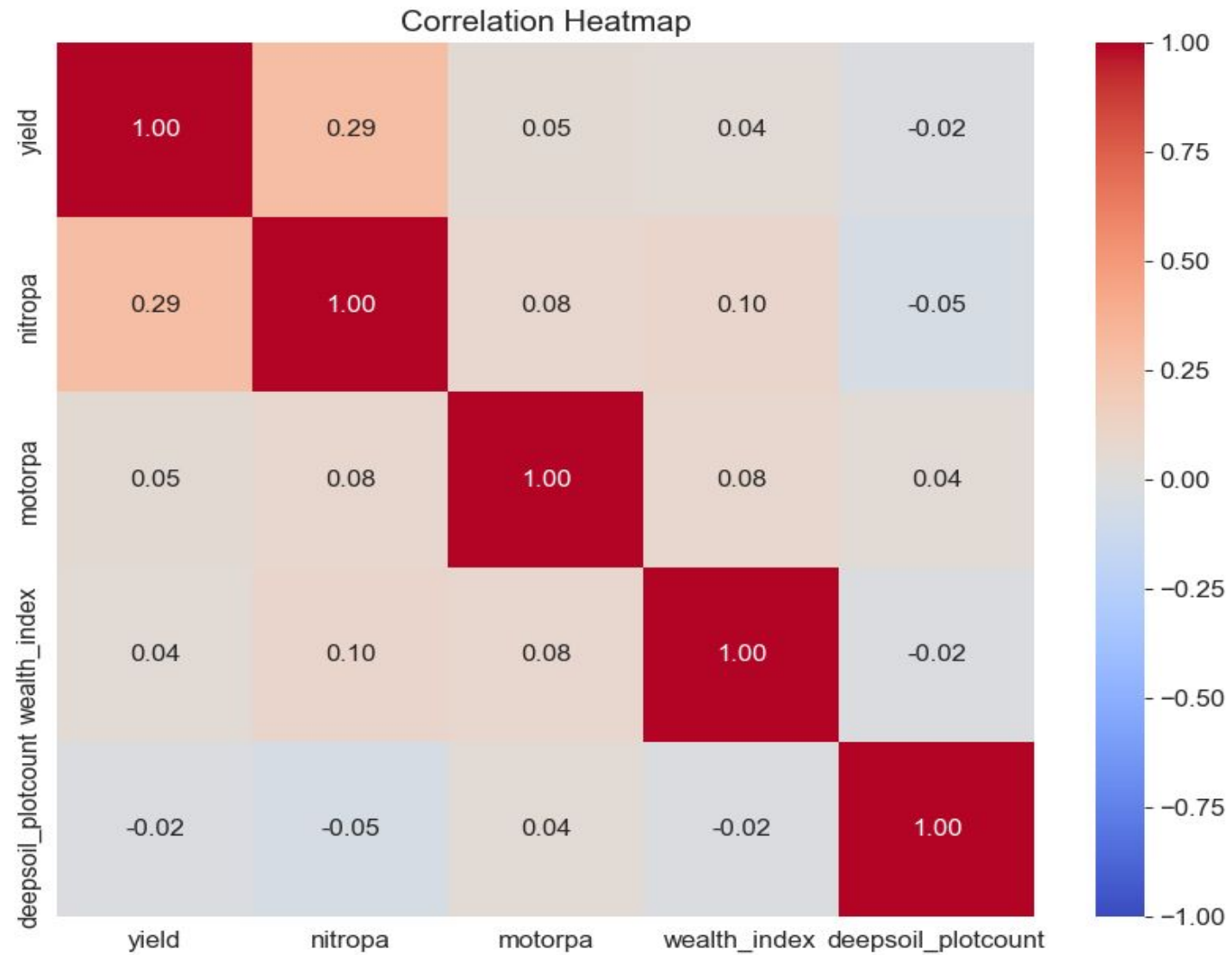


The Staggered Diffusion of Credit Access (with 95% CI)

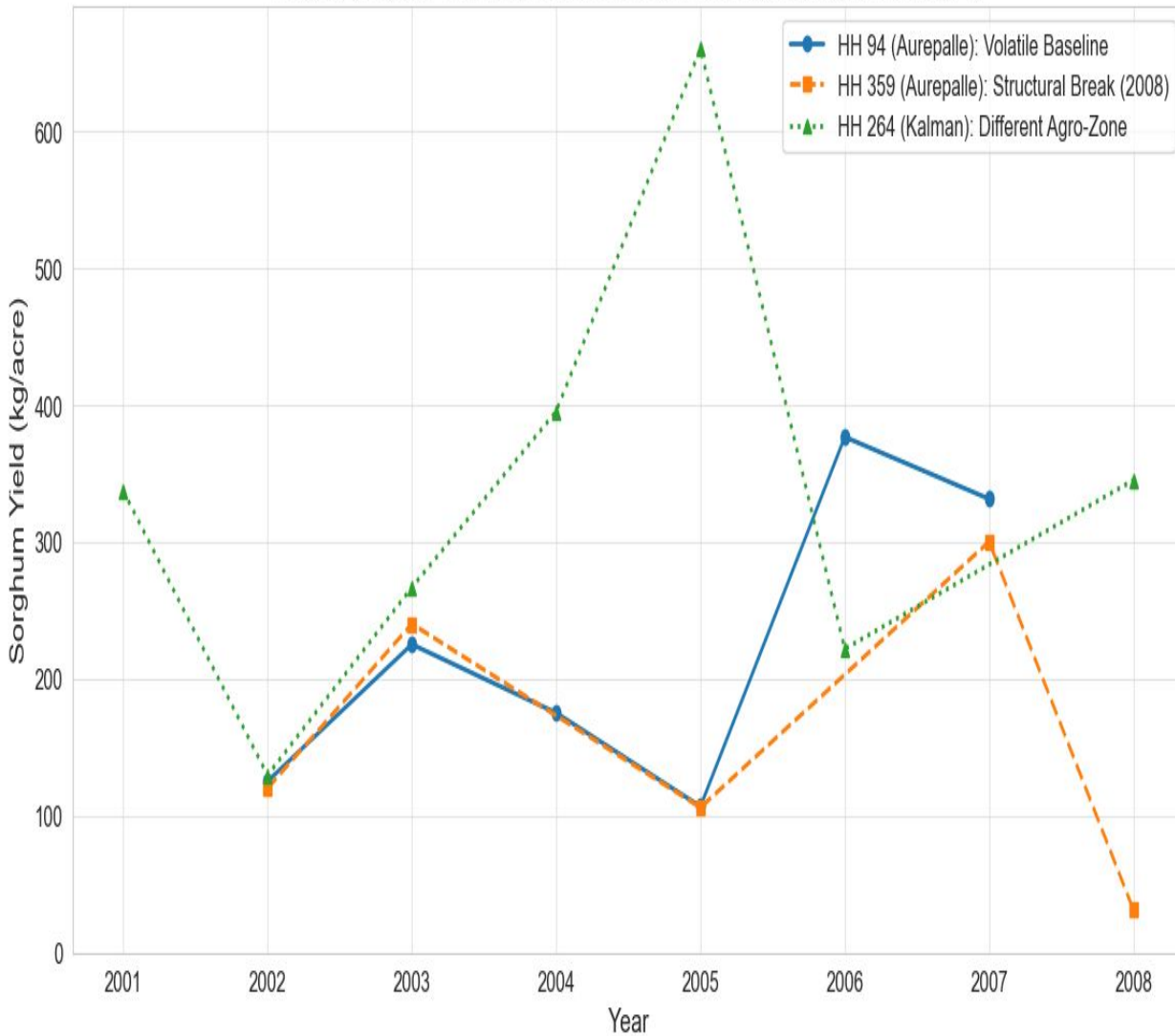




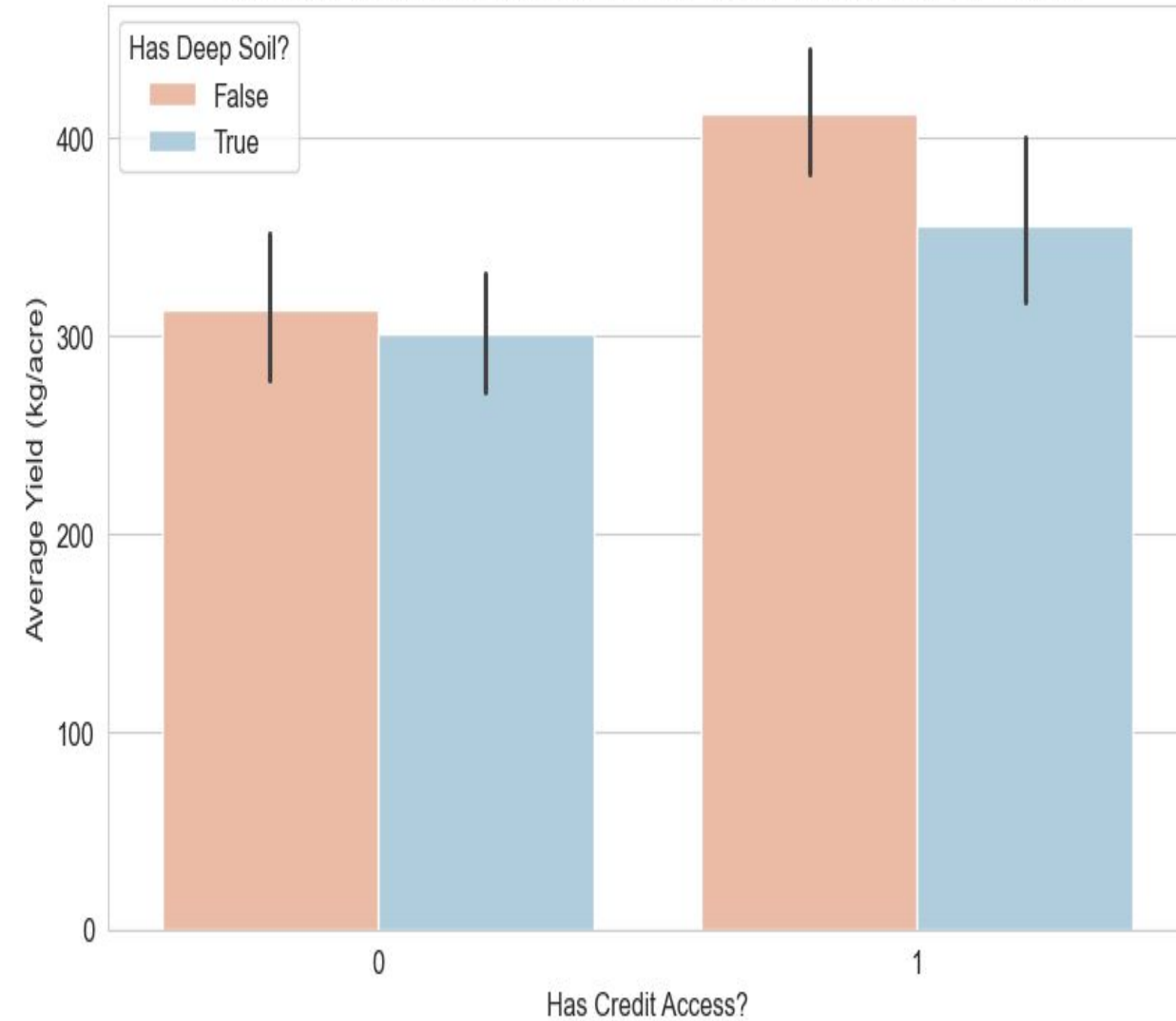
EDA



Micro-Analysis: Yield Trajectories of Specific Households (2001-2014)



Heterogeneous Treatment Effects: Credit Impact by Soil Depth



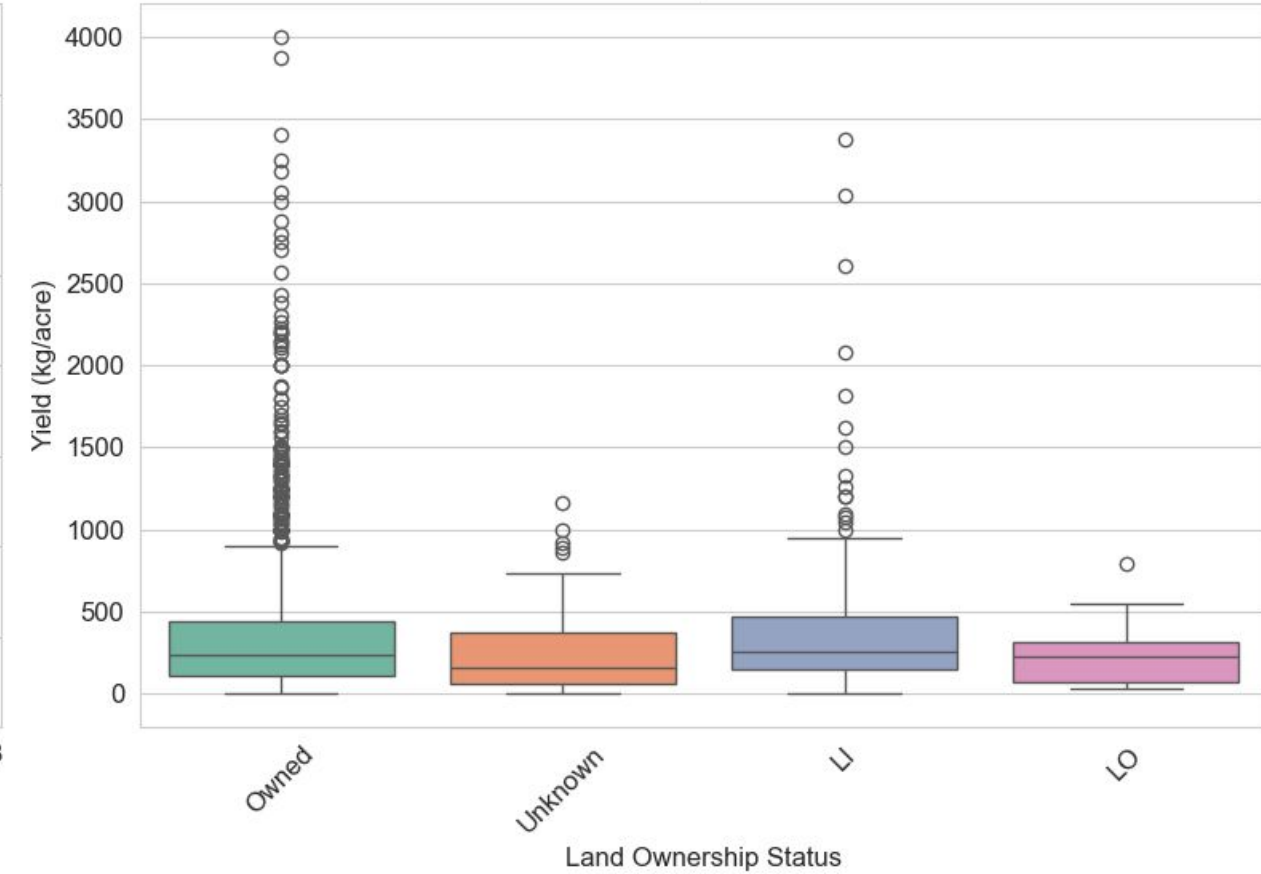
EDA

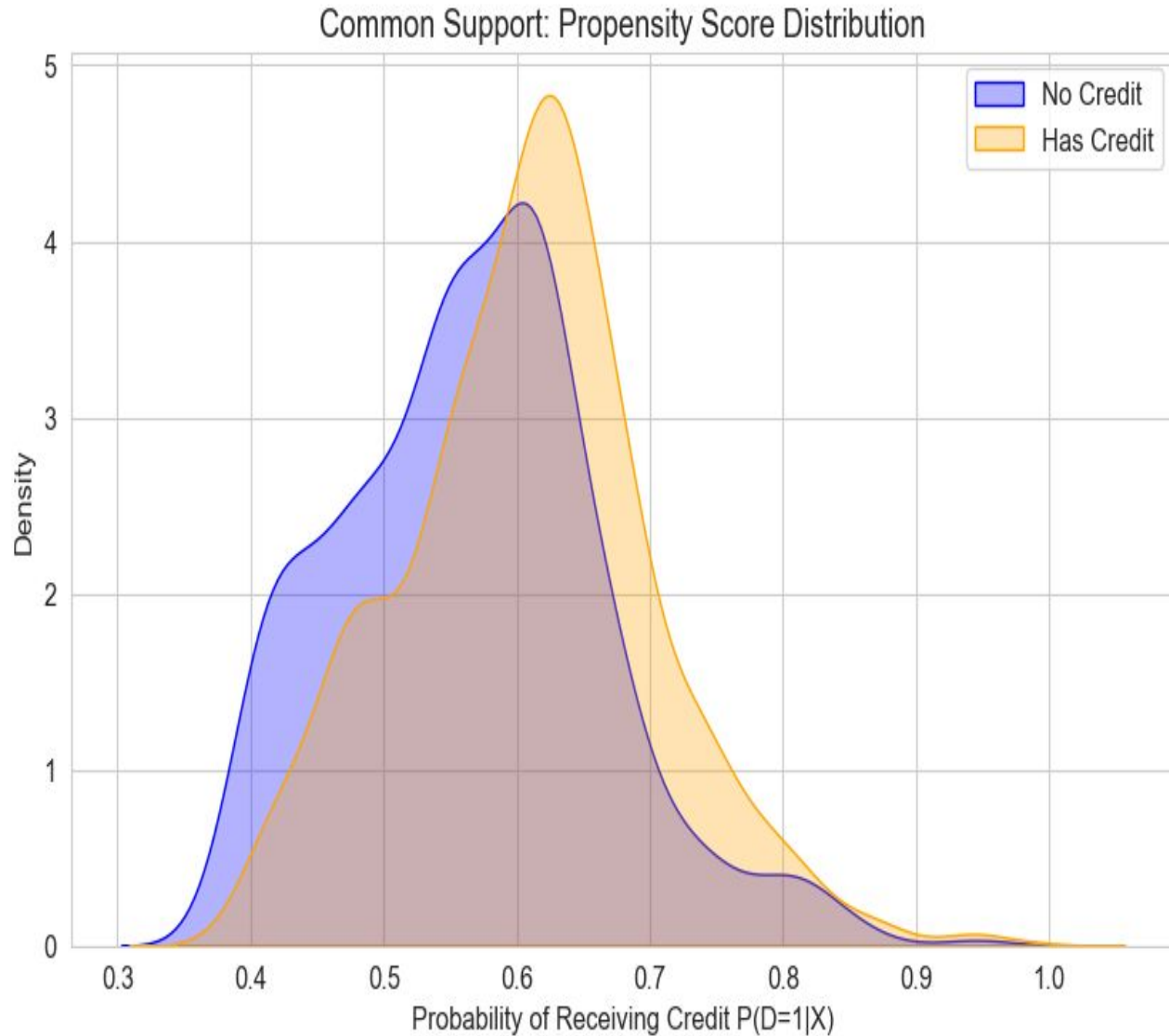
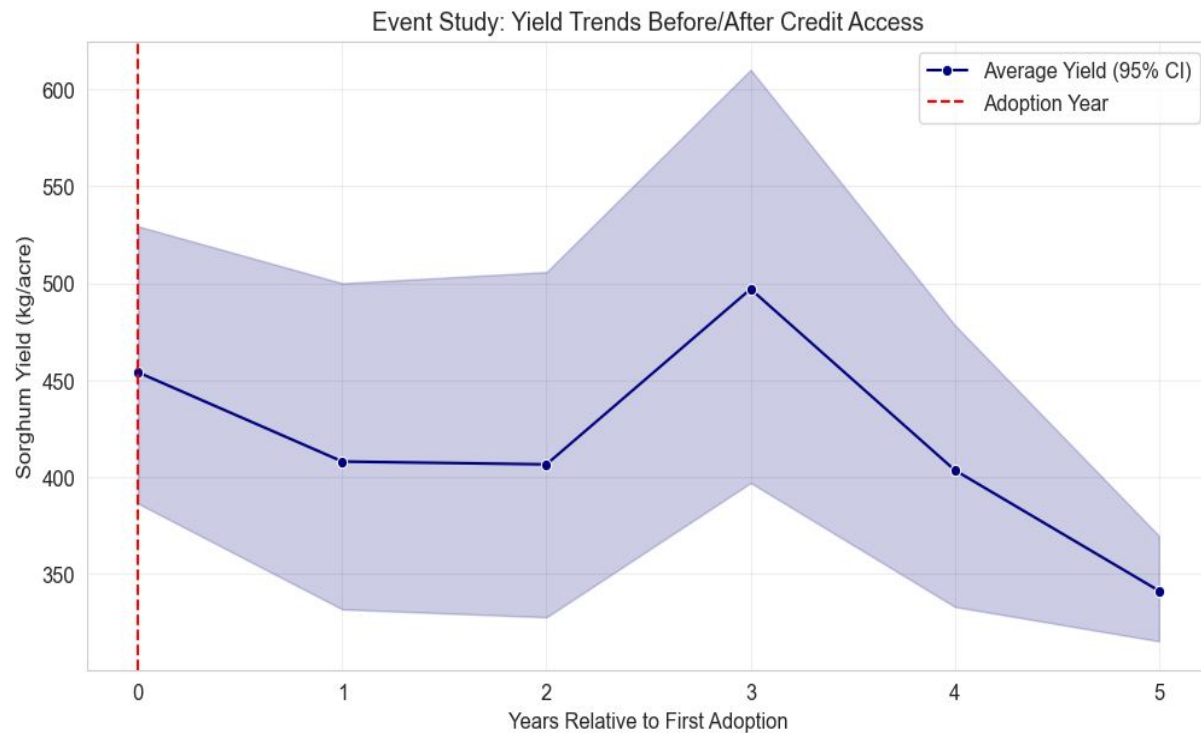


Input Intensity: Nitrogen (kg/acre) vs Yield

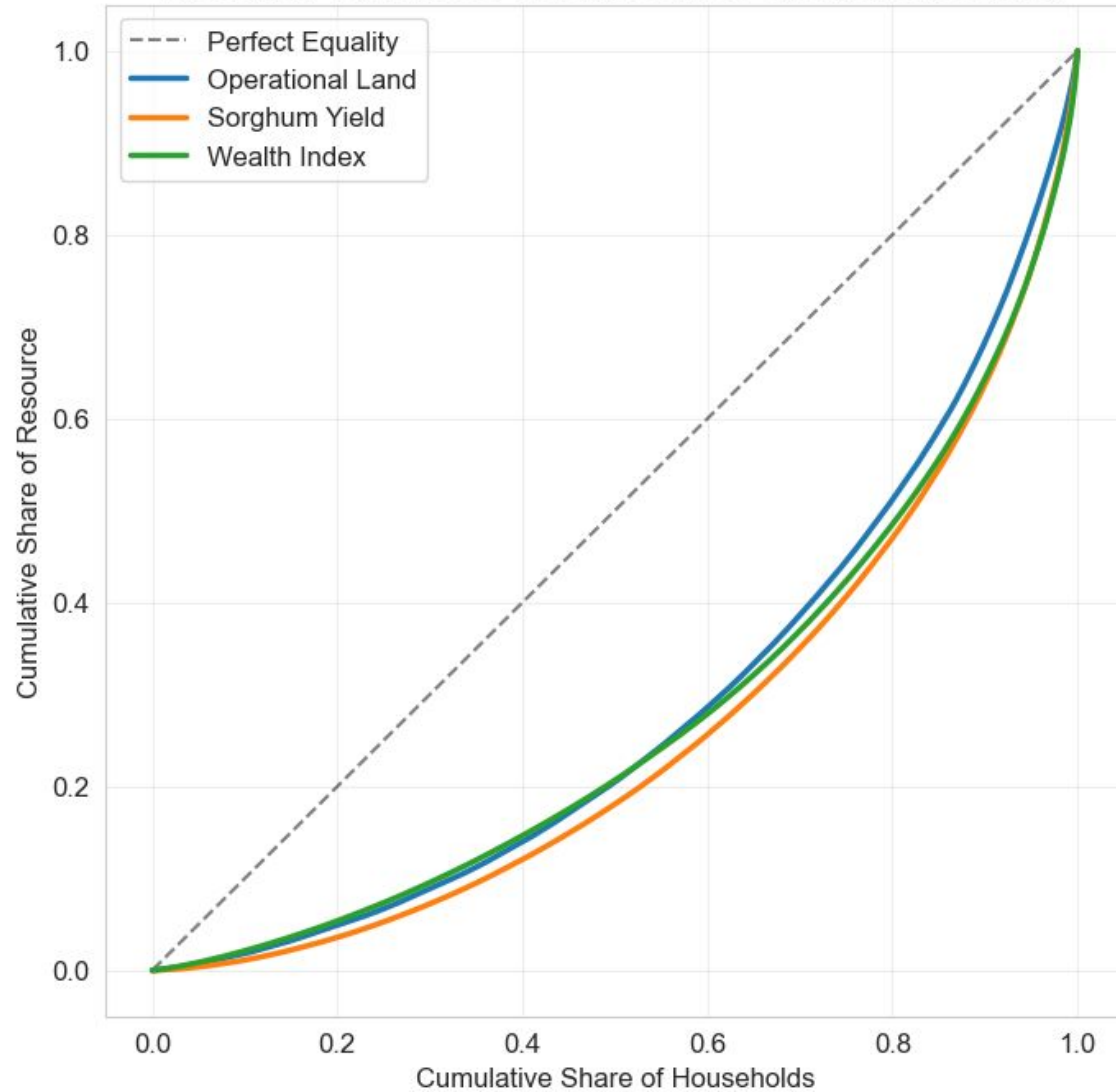


Stratification: Yield by Land Tenure Status

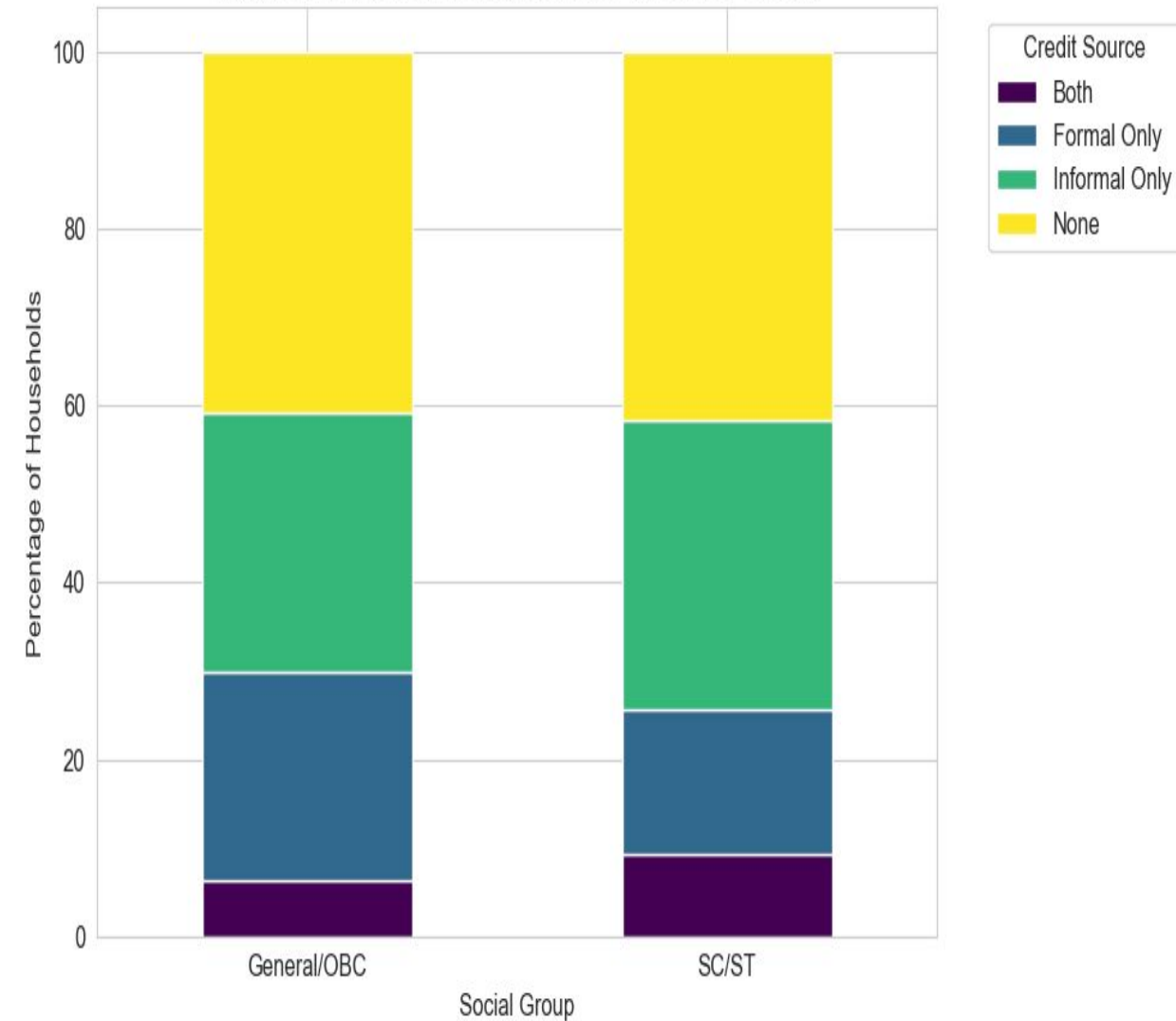




Inequality Dynamics: Lorenz Curves of Agrarian Capital



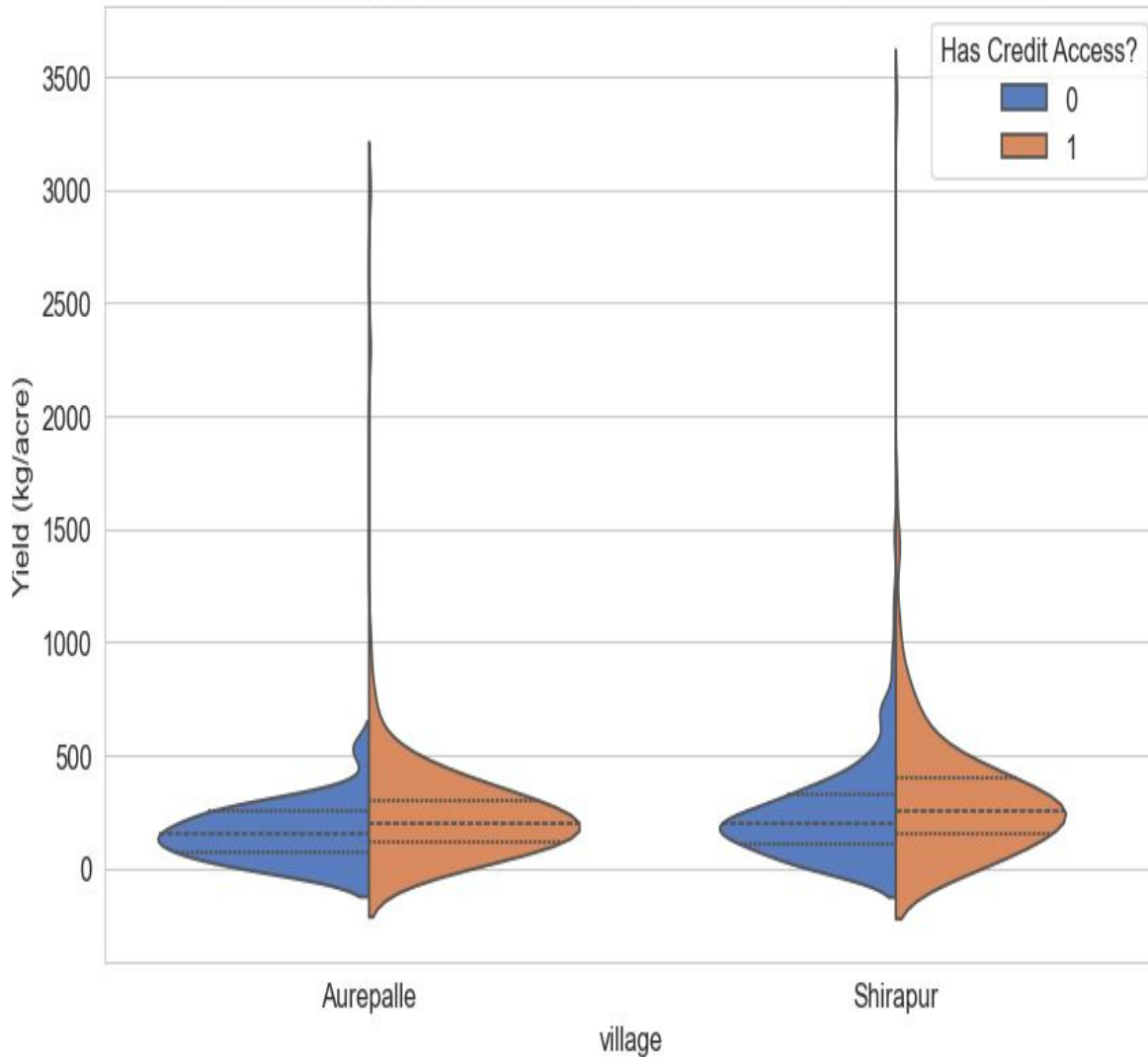
Financial Exclusion: Credit Access by Social Group



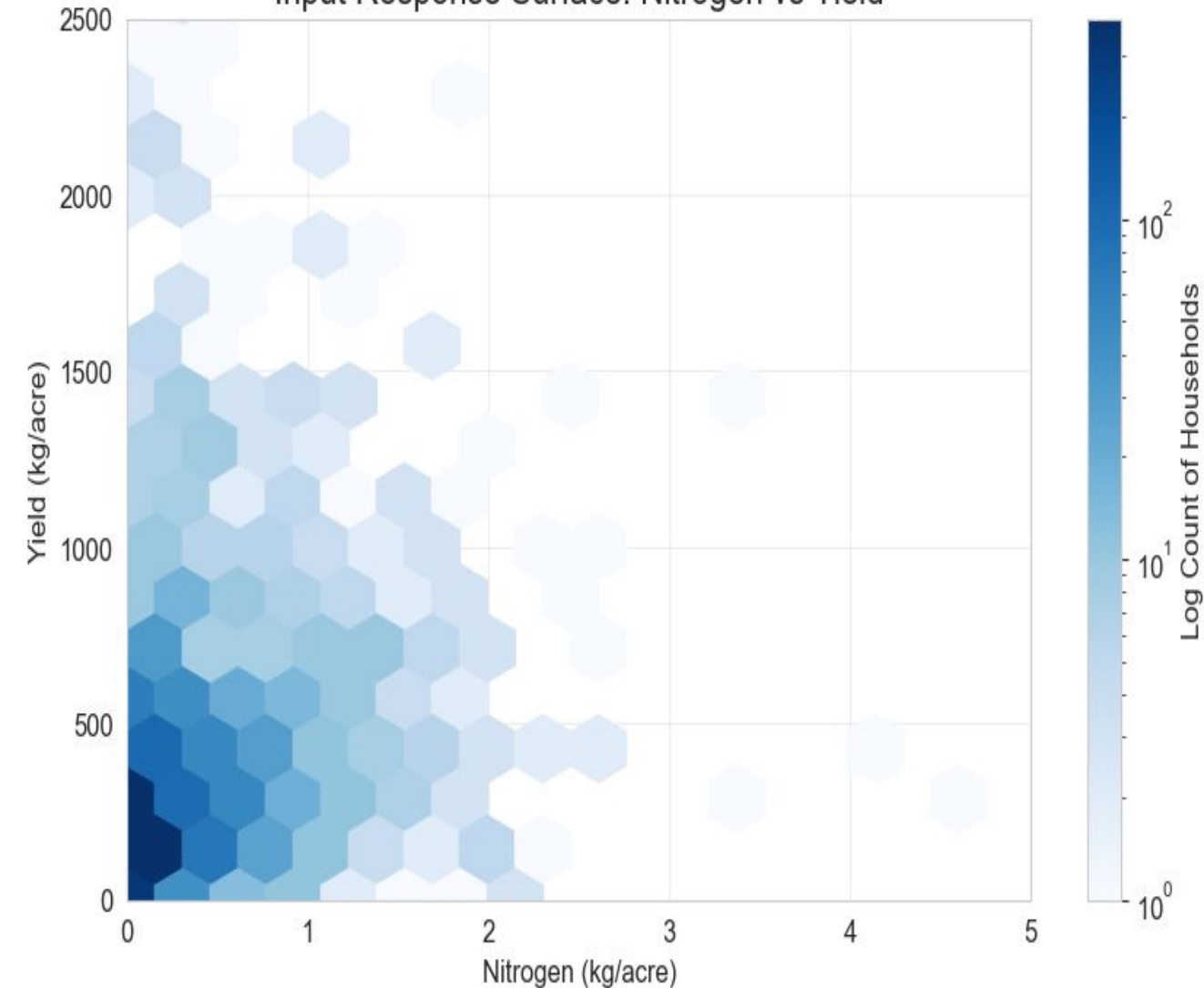
EDA



Distribution of Yield: Aurepalle vs Shirapur (by Credit Status)



Input Response Surface: Nitrogen vs Yield



Baseline – Data Assignment



Method: We estimated the standard “Average Effect” model.

Result: We found a single coefficient, leading to the naive interpretation that "Credit increases yield for everyone."

Critical Flaw: Being linear in nature, the model assumes that every farmer—regardless of soil or wealth—receives the exact same benefit. It masks heterogeneity by failing to account for interacting factors.

Motivation: To propose a new approach that explicitly accounts for these varying factors.

Conceptual and Practical Advancement

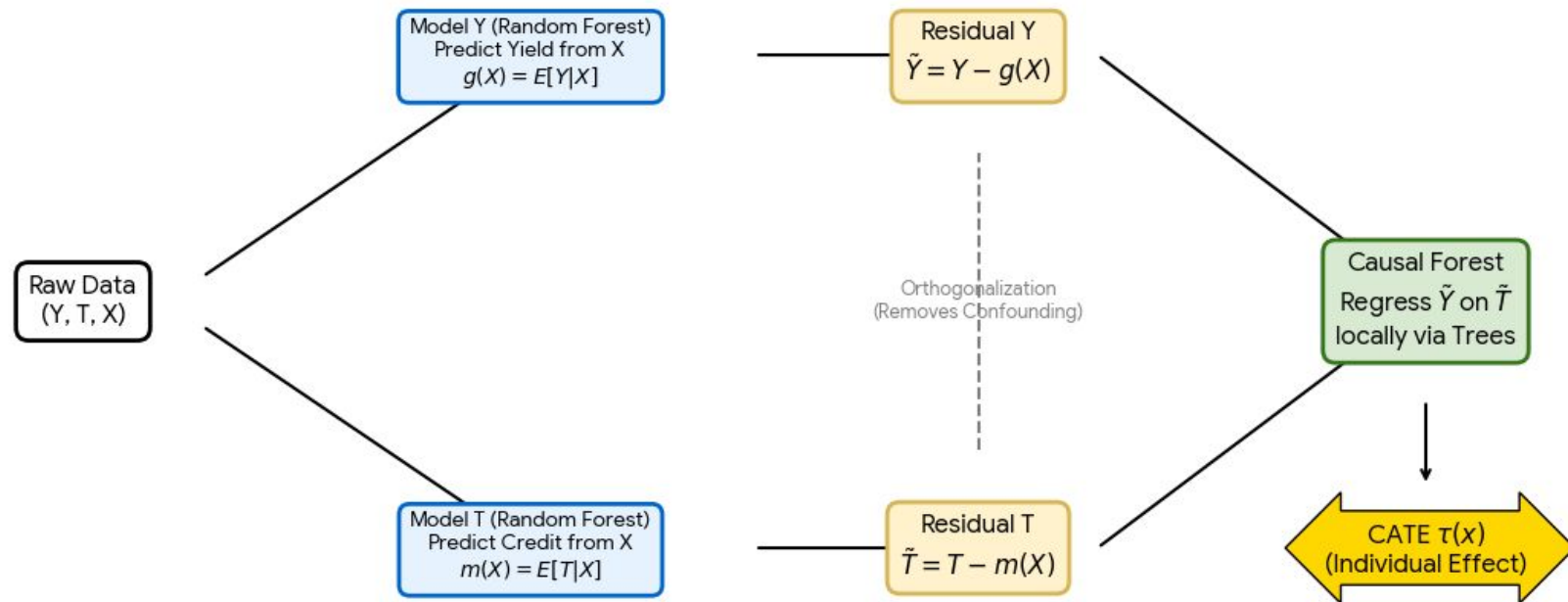


Feature	Data Assignment (Baseline)	EML Project (Advancement)
Assumption	Homogeneity: Homogeneity: Assumed $\tau_i = \tau$ (Constant Effect for all).	Heterogeneity: Heterogeneity: Assumes $\tau_i = \tau(x)$ (Effect depends on soil/wealth).
Method	Naive OLS / Standard DiD: Relies on Linearity and Parallel Trends.	Double Machine Learning (DML): Uses Random Forests to relax linearity assumptions .
Identification	Selection on Observables: Used simple Logit PSM (fails with high-dim data).	Orthogonalization: Removes regularization bias via the Double/Debiased Machine Learning.
Controls	Potential Bad Controls: Potential Bad Controls: Included current inputs (t) which bias results.	Clean Controls: Clean Controls: Uses Lagged inputs (t-1) to control for baseline skill.

Methodology in brief



Methodology: Double Machine Learning (DML) Pipeline



Methodology I – Double Machine Learning



The Goal: Clean the data - We must first isolate the effect of confounding variables on both Credit(D) and Yield(Y)

Step 1 - Train two separate Random Forests to predict expected outcomes:

$$g(X) = E[Y|X] \quad (\text{Predict Yield from Soil/Wealth})$$

$$m(X) = E[D|X] \quad (\text{Predict Credit from Soil/Wealth})$$

Step 2 - Orthogonalization or Residualization - We will subtract these values from their original values to get residuals

$$\tilde{Y} = Y - \hat{g}(X) \quad (\text{Yield unexplained by } X)$$

$$\tilde{D} = D - \hat{m}(X) \quad (\text{Credit unexplained by } X)$$

so that we get true Yield and true Credit which are unexplained or not influenced by the confounding variables

Methodology II – Splitting and Estimation



Goal - **We regress Y on D** by creating a decision tree and splitting such that we get maximum variance (For maximum heterogeneity)

Splitting Criterion - Unlike standard trees which minimize error, Causal Trees maximize Variance of Effects.

At each node, the algorithm tests every variable X_j to find split S that maximizes: where

$$(\tau_L - \tau_R)^2$$

$$\hat{\tau}_L = \frac{\sum_{i \in L} \tilde{Y}_i \tilde{D}_i}{\sum_{i \in L} (\tilde{D}_i)^2}$$

τ_L

is the estimated Conditional Average Treatment for the observations that fall into the **Left Child Node** after a split

τ_R

is the estimated Conditional Average Treatment for the observations that fall into the **Right Child Node** after a split

For discrete values, we split at different points. For continuous values, we create partitions and split at lesser/greater than terms.

Methodology III – Calculating the Score

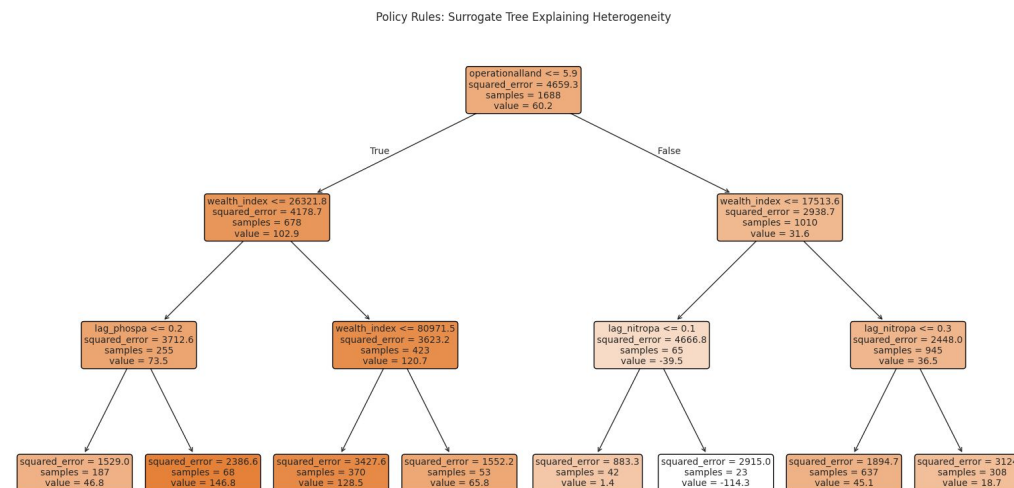


What happens at inference time?

We create 200 trees as explained in methodology II and call it a forest.

In each tree, we drop the farmer on top of the tree and we get the final value at the leaf. We then take the average treatment effect of the leaves they land in

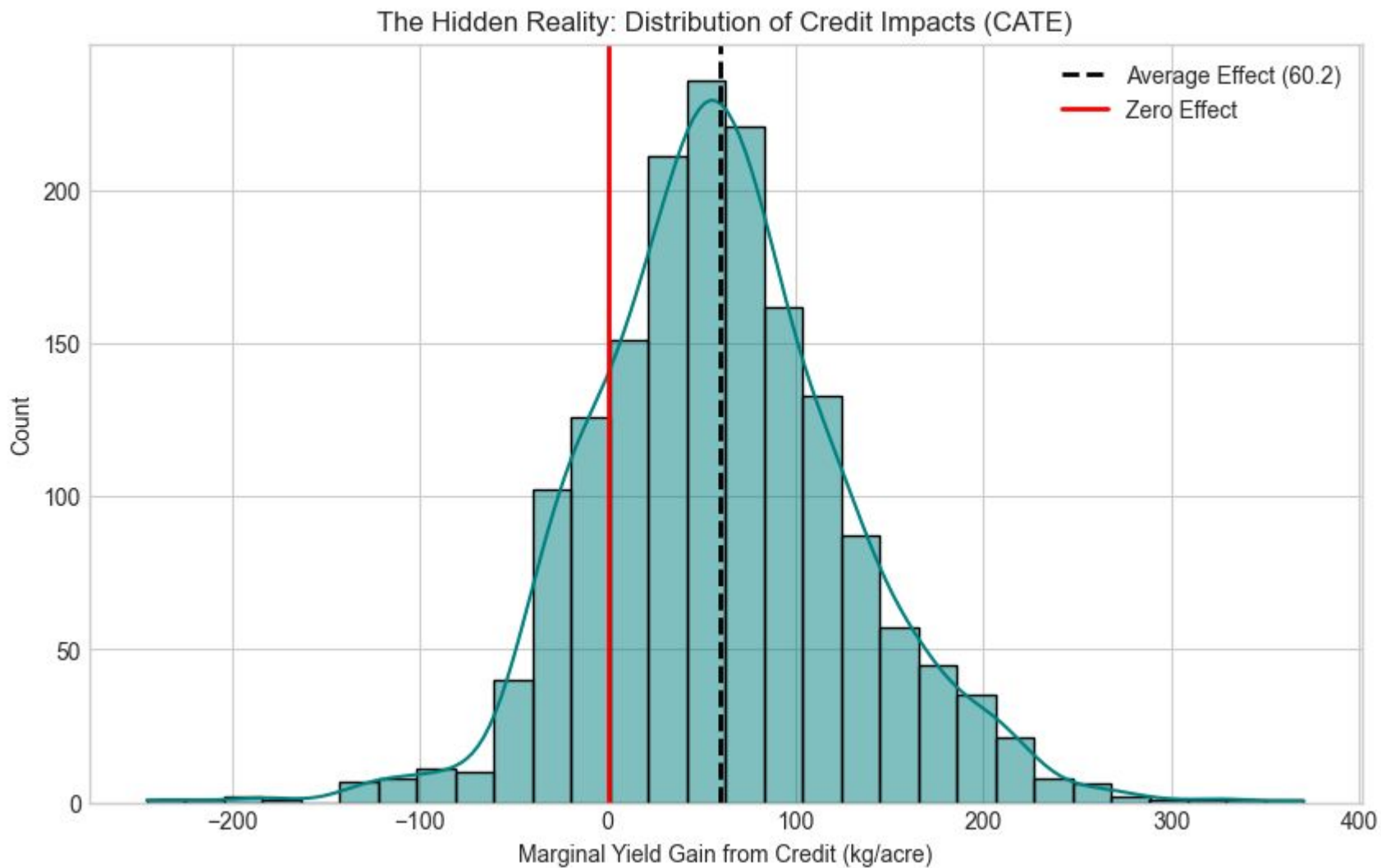
$$\hat{\tau}(x_i) = \frac{1}{B} \sum_{b=1}^B \hat{\tau}_{L_b}(x_i)$$



RESULTS

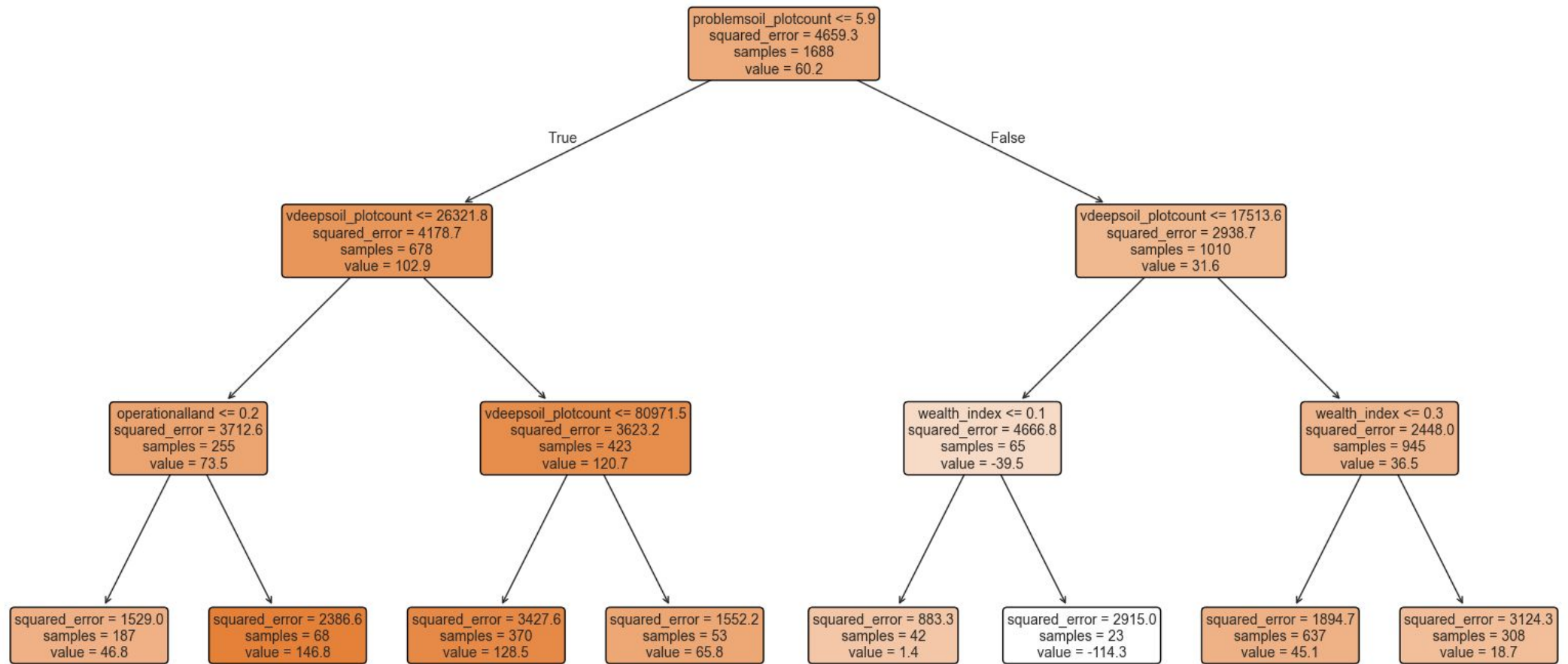
AND

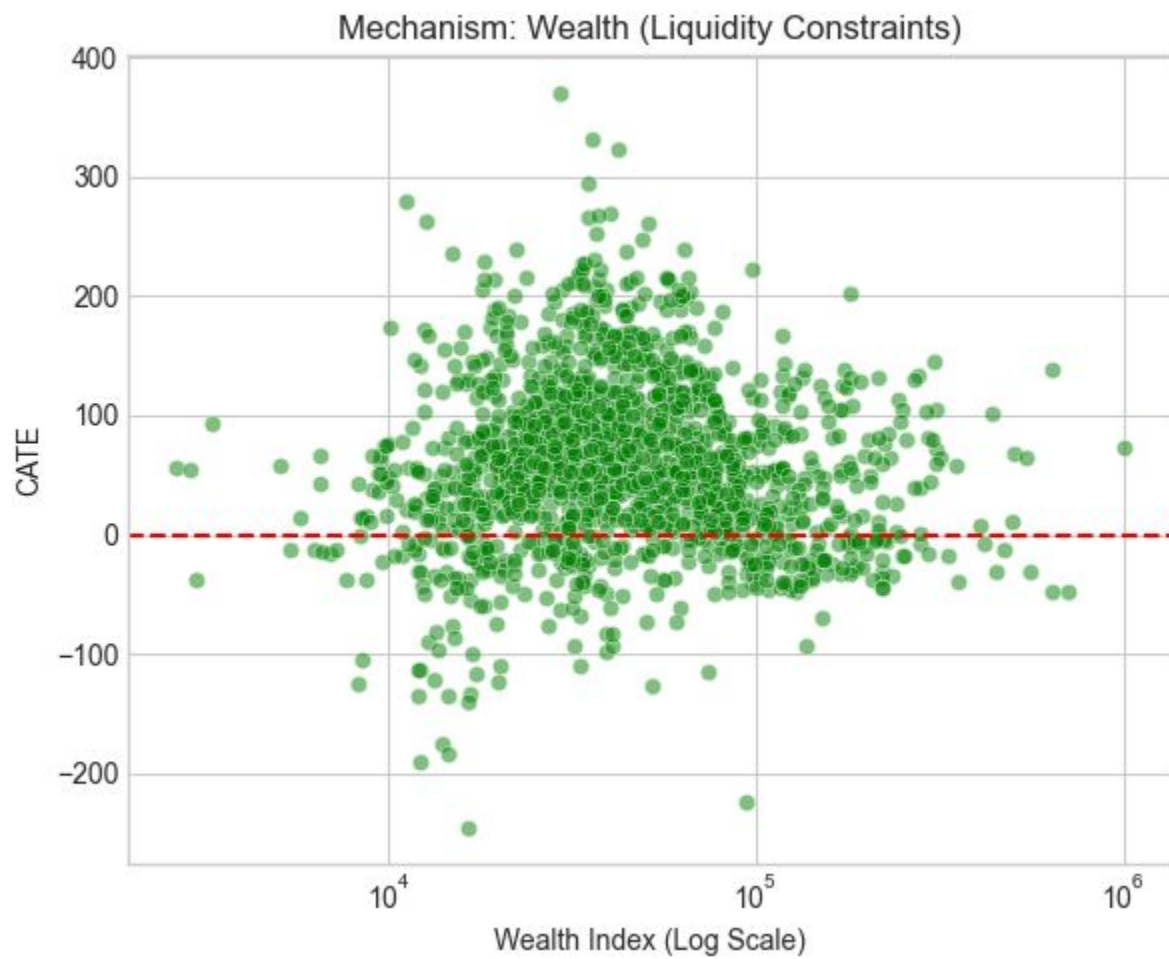
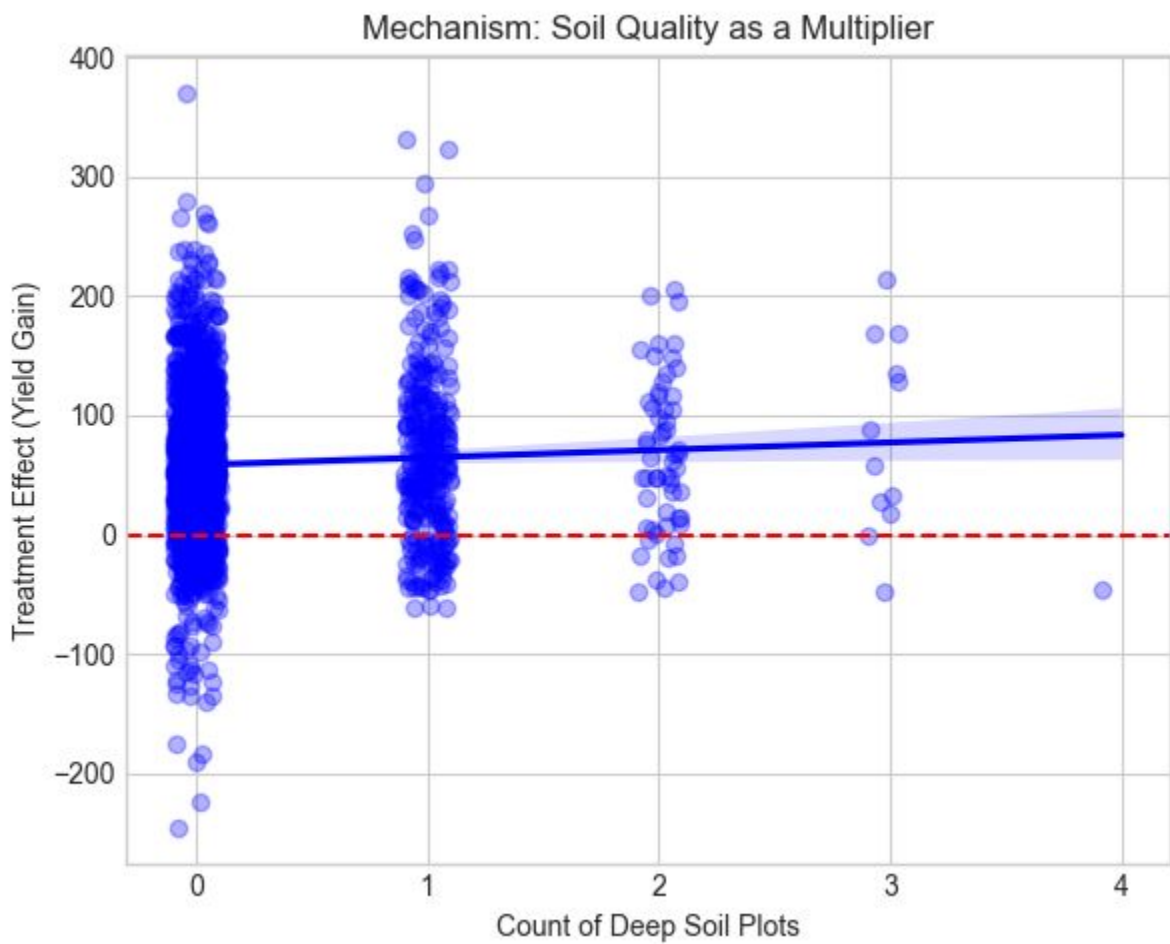
INTERPRETATION

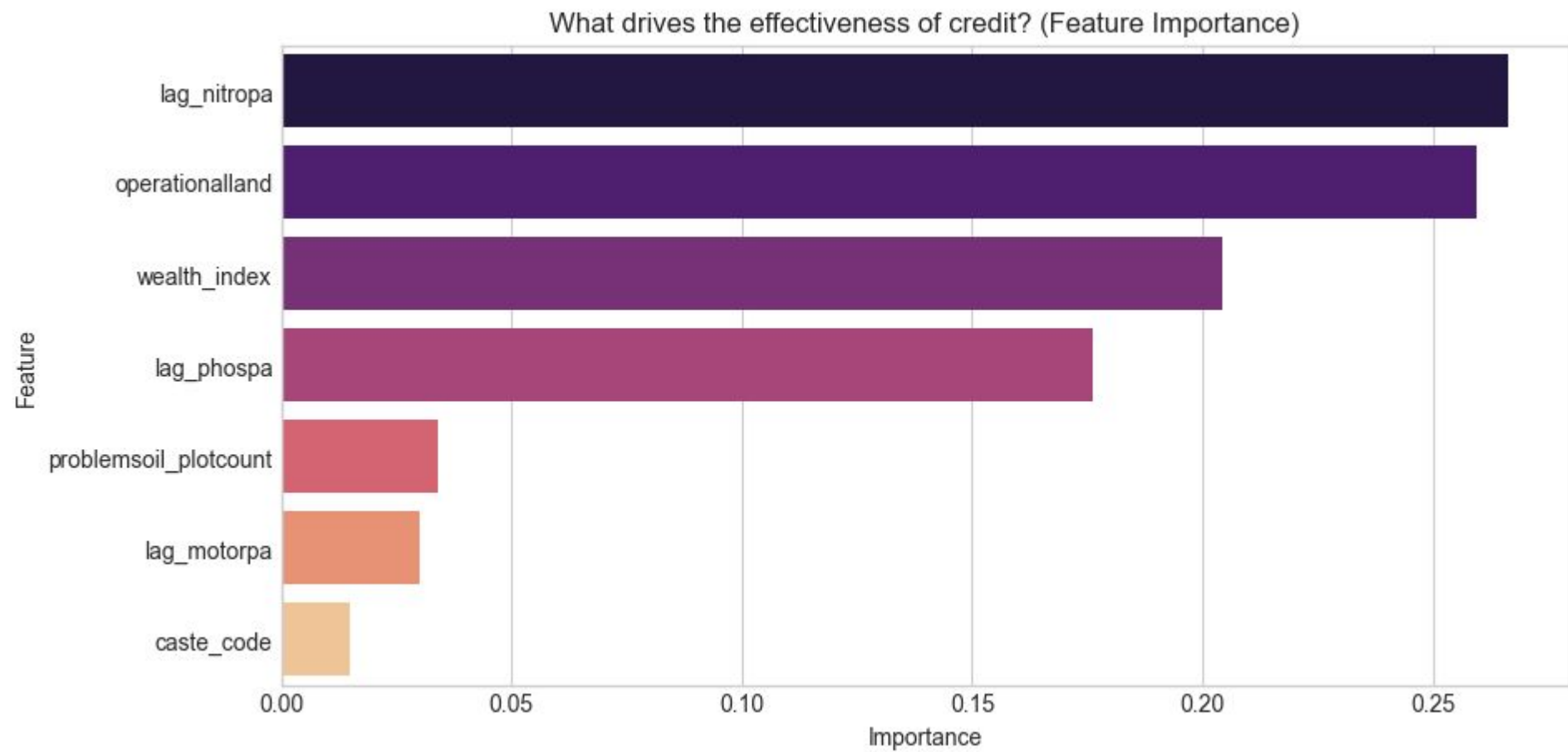


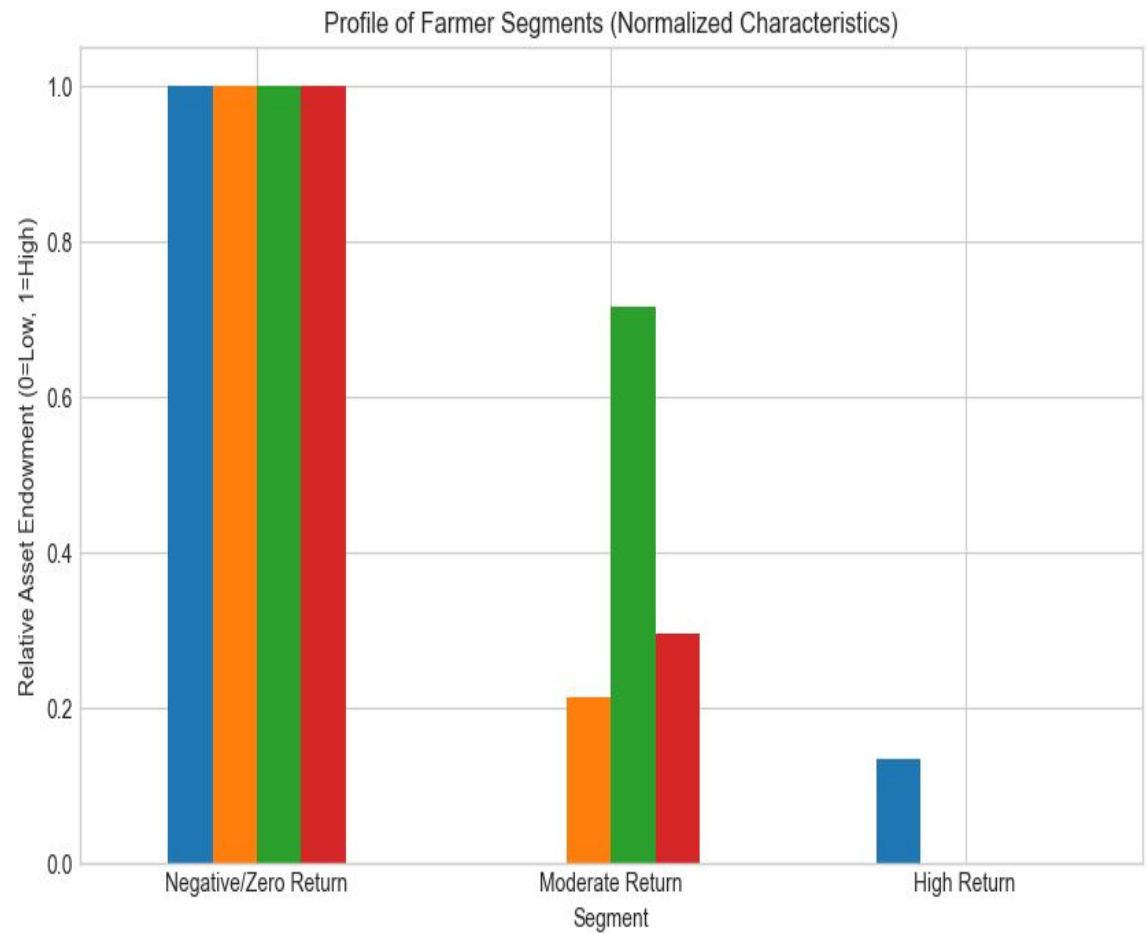
The **Distribution Plot** proves that while the average effect is positive, a significant minority of farmers actually lose out

Surrogate Tree: The Rules of Heterogeneity
(How the model decides who benefits)







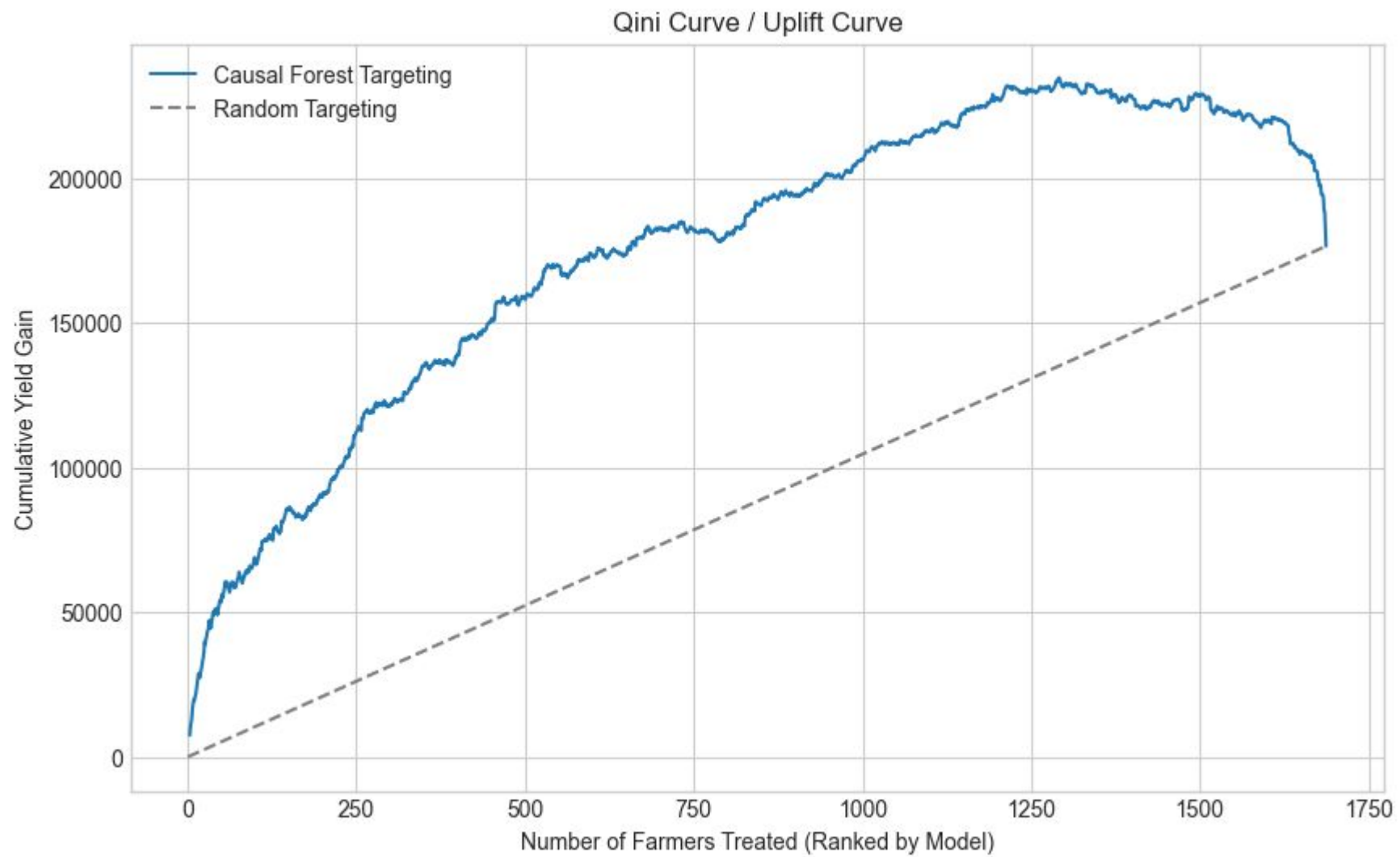


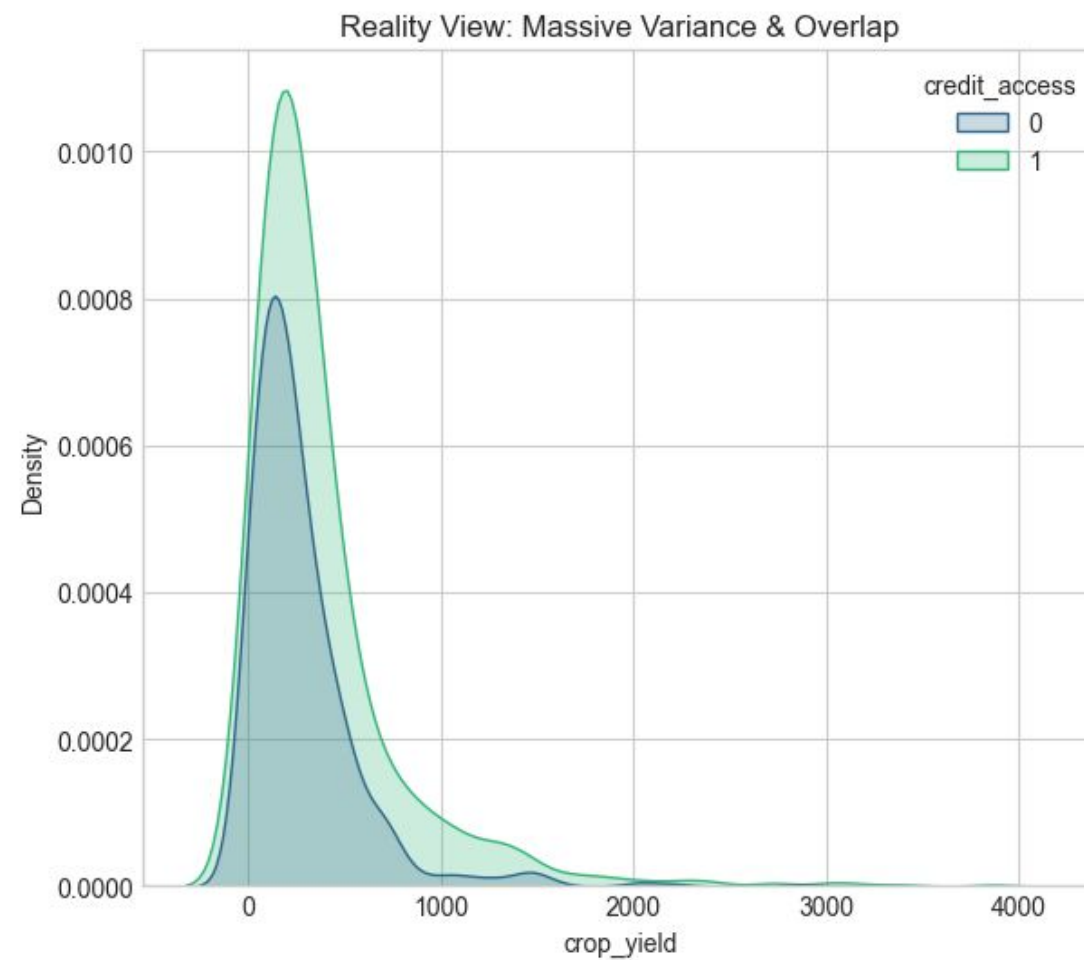
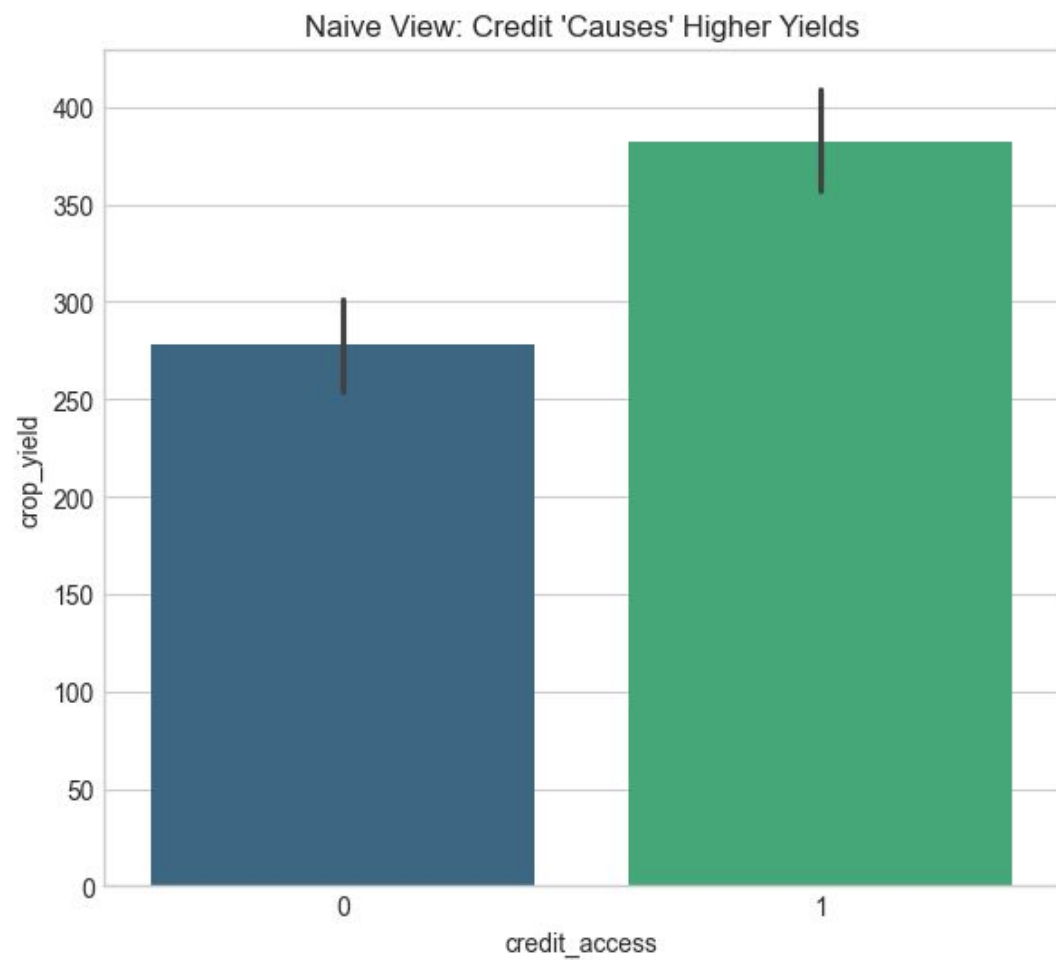
vdeepsoil_plotcount
wealth_index
operationalland
lag_nitropa

--- FARMER PERSONAS: WHO WINS? ---

	vdeepsoil_plotcount	wealth_index	operationalland
Segment			
Negative/Zero Return	0.299342	82310.180597	12.700033
Moderate Return	0.271300	60980.221653	11.167130
High Return	0.275053	55241.715091	7.312228

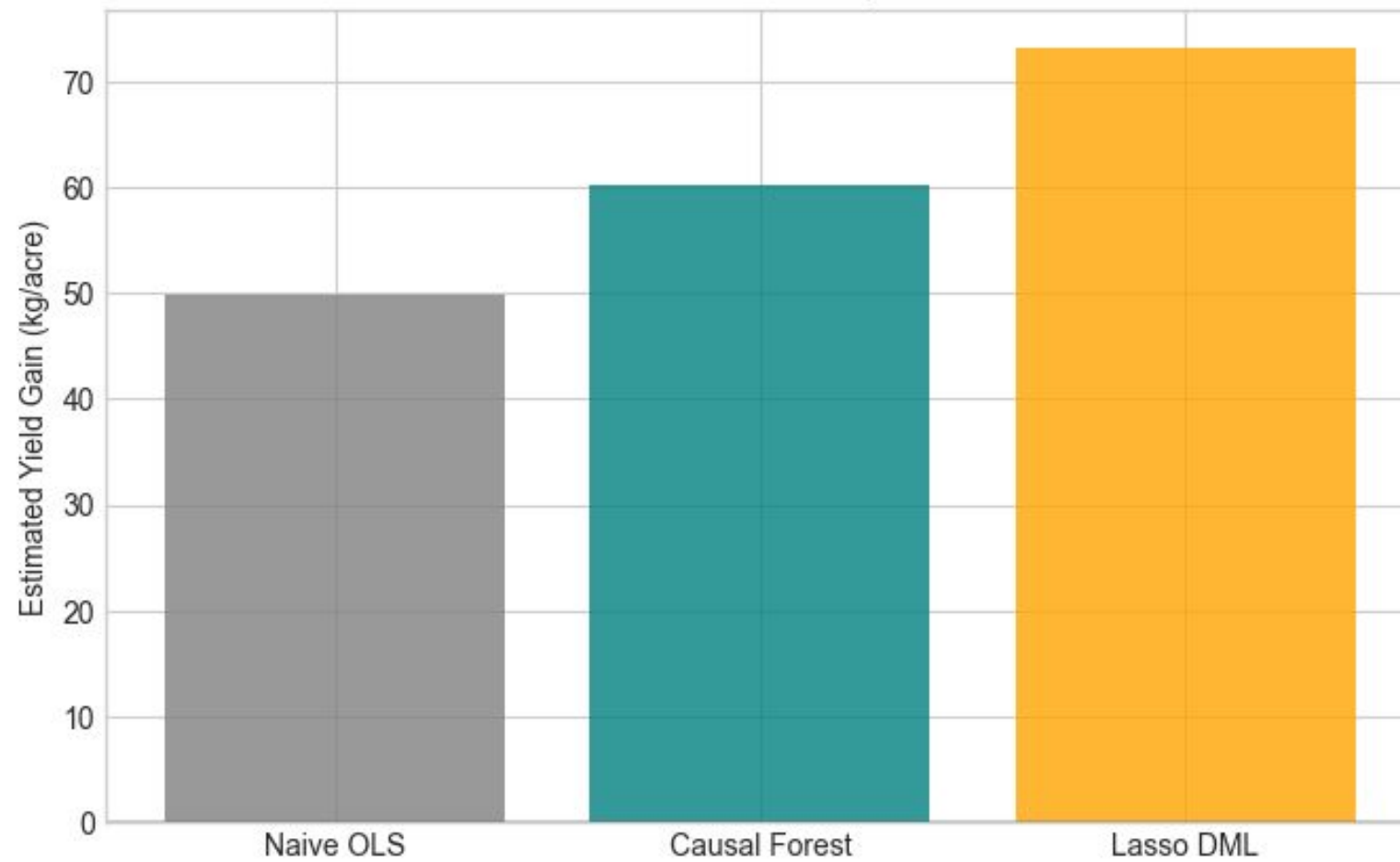
	lag_nitropa
Segment	
Negative/Zero Return	0.381969
Moderate Return	0.291712
High Return	0.254235





Naive OLS Coefficient for Credit: 44.35

Robustness Check: Effect Stability Across Estimators



--- BEST LINEAR PREDICTOR (BLP) TEST ---

OLS Regression Results

```

=====
Dep. Variable:          CATE      R-squared:                0.163
Model:                  OLS      Adj. R-squared:            0.158
Method:                 Least Squares      F-statistic:          29.74
Date:                  Wed, 26 Nov 2025      Prob (F-statistic):      1.05e-57
Time:                  00:51:04      Log-Likelihood:          -9373.6
No. Observations:      1688      AIC:                    1.877e+04
Df Residuals:          1676      BIC:                    1.884e+04
Df Model:              11
Covariance Type:       nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	60.2454	1.525	39.502	0.000	57.254	63.237
crop_yield	9.5202	1.607	5.924	0.000	6.368	12.672
credit_access	-4.0952	1.582	-2.589	0.010	-7.198	-0.992
vdeepsoil_plotcount	4.2489	1.585	2.681	0.007	1.140	7.358
problemsoil_plotcount	-0.1993	1.561	-0.128	0.898	-3.261	2.862
wealth_index	-0.6152	1.708	-0.360	0.719	-3.966	2.735
operationalland	-23.6677	1.677	-14.110	0.000	-26.958	-20.378
lag_nitropa	-8.2401	1.877	-4.391	0.000	-11.921	-4.559
lag_phospa	8.6333	1.851	4.665	0.000	5.004	12.263
lag_motorpa	-0.6667	1.835	-0.363	0.716	-4.267	2.933
lag_irrigation_indicator	-5.0583	1.814	-2.788	0.005	-8.616	-1.500
caste_code	-4.3717	1.542	-2.835	0.005	-7.397	-1.347

```

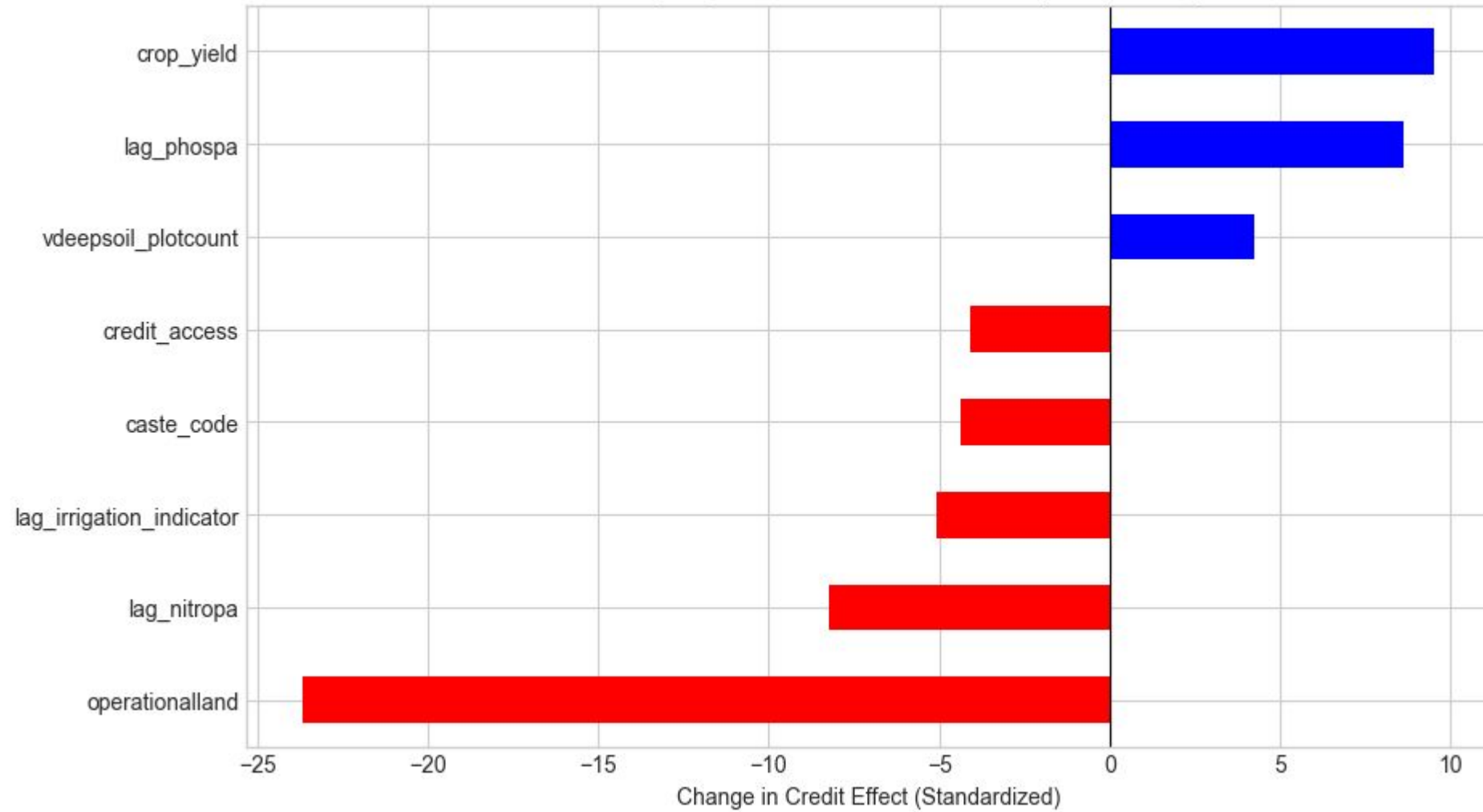
=====
Omnibus:                90.277      Durbin-Watson:          1.277
Prob(Omnibus):          0.000      Jarque-Bera (JB):       228.282
Skew:                   -0.285      Prob(JB):               2.69e-50
Kurtosis:               4.709      Cond. No.                2.28
=====

```

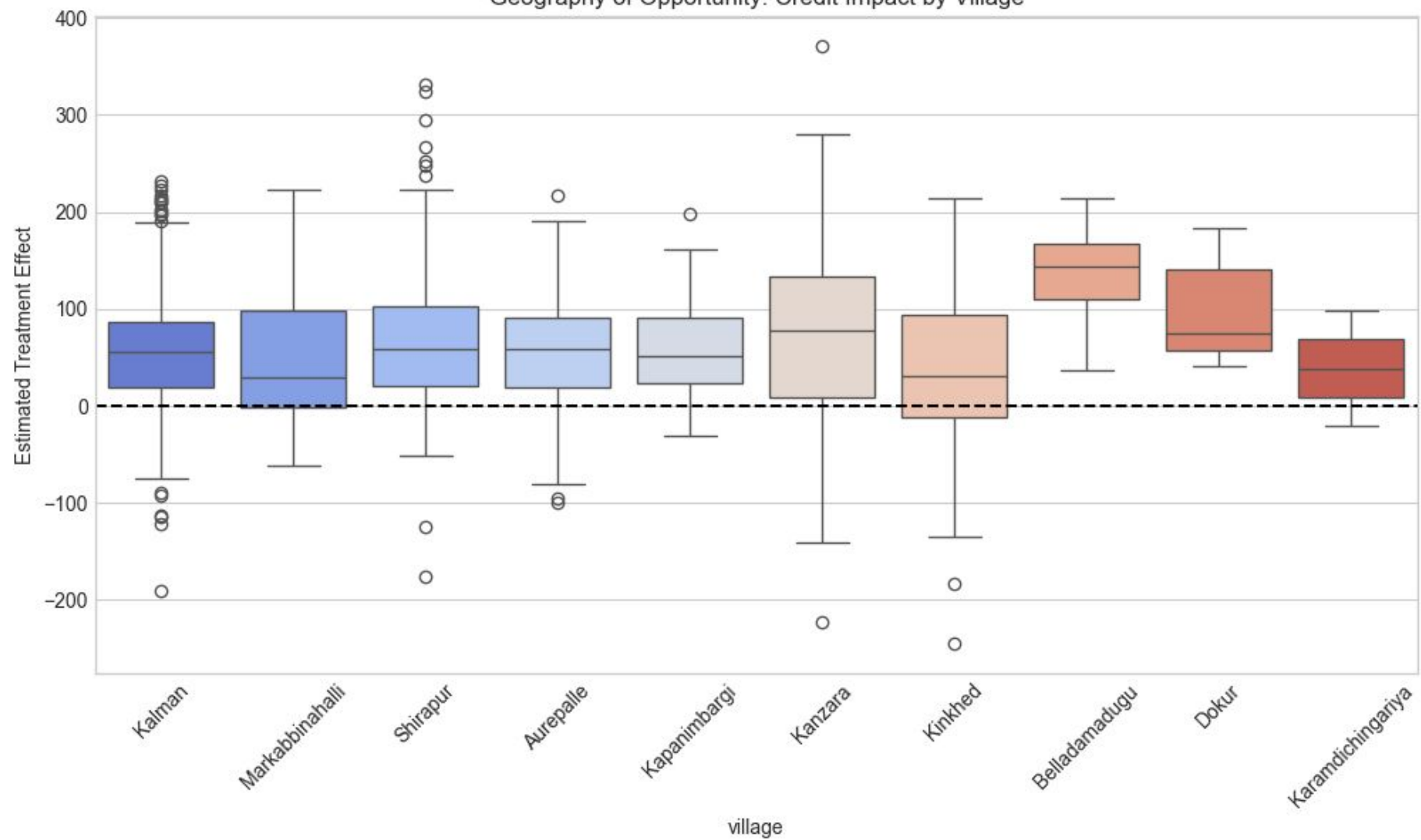
Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

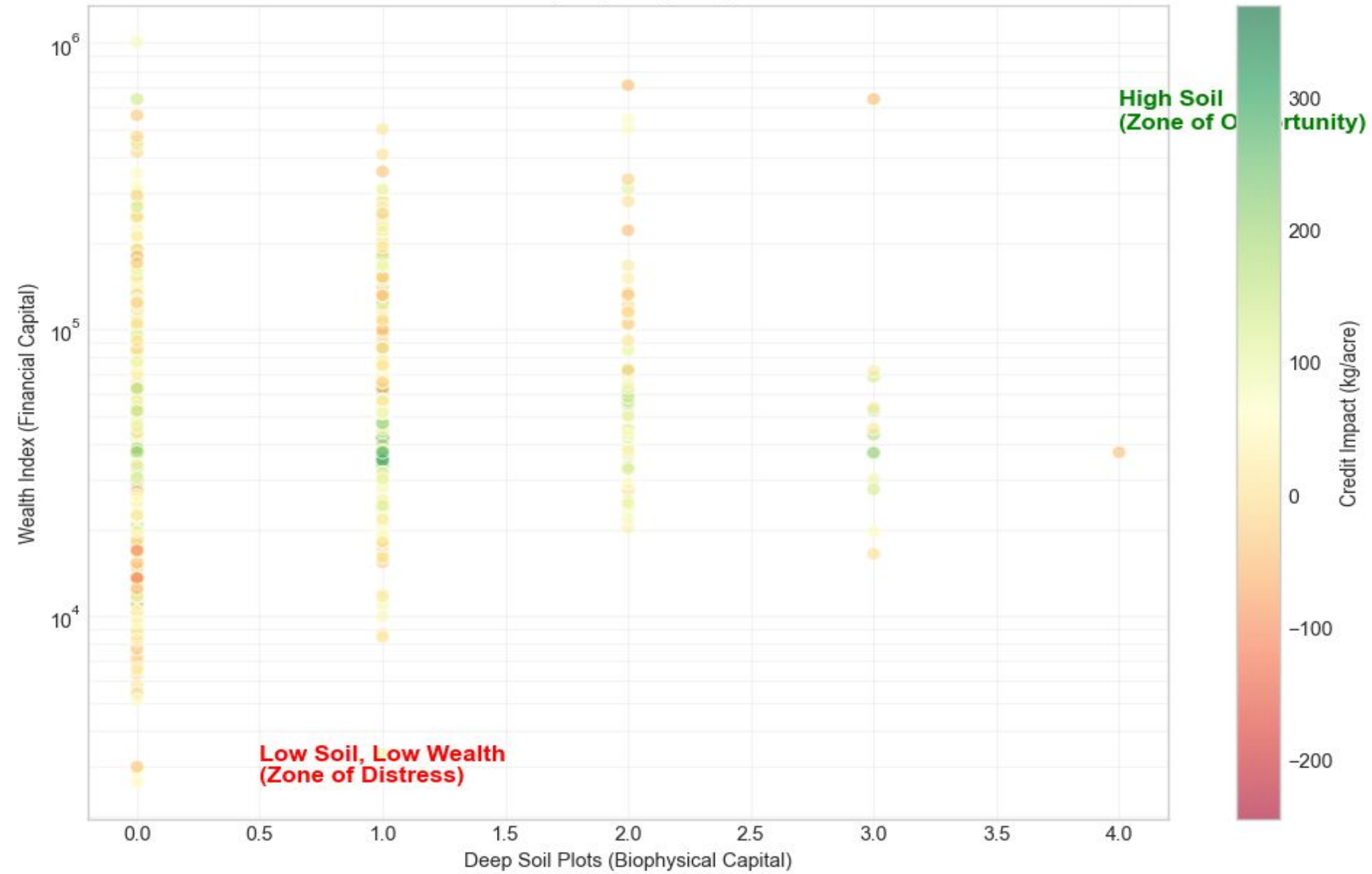
Statistically Significant Drivers of Credit Impact ($p < 0.1$)



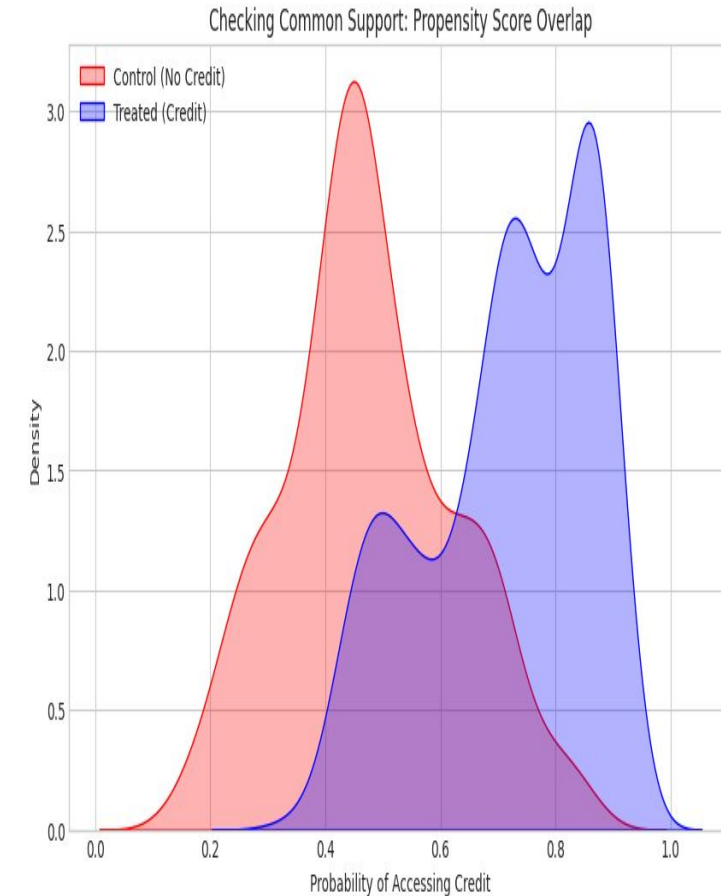
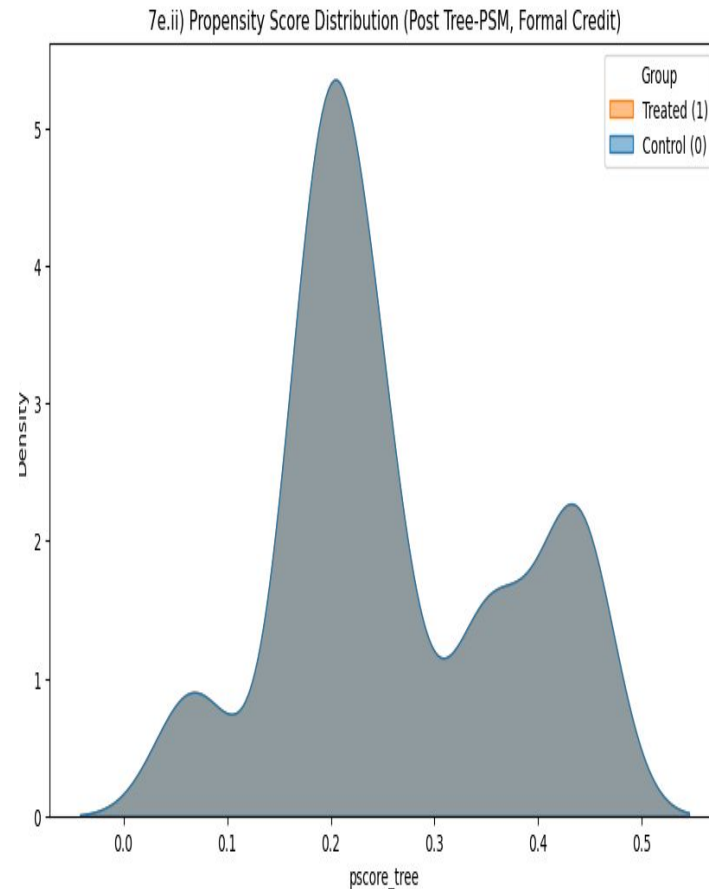
Geography of Opportunity: Credit Impact by Village



The Precision Policy Map: Targeting Credit Allocation



Model Diagnostic and Rigor Analysis



With the simple ATE model, treated and control households show almost identical propensity scores, falsely suggesting random credit access. The CATE/DML model reveals higher scores for treated households, correctly capturing real selection patterns and confirming the need for heterogeneous causal methods.

Placebo Test



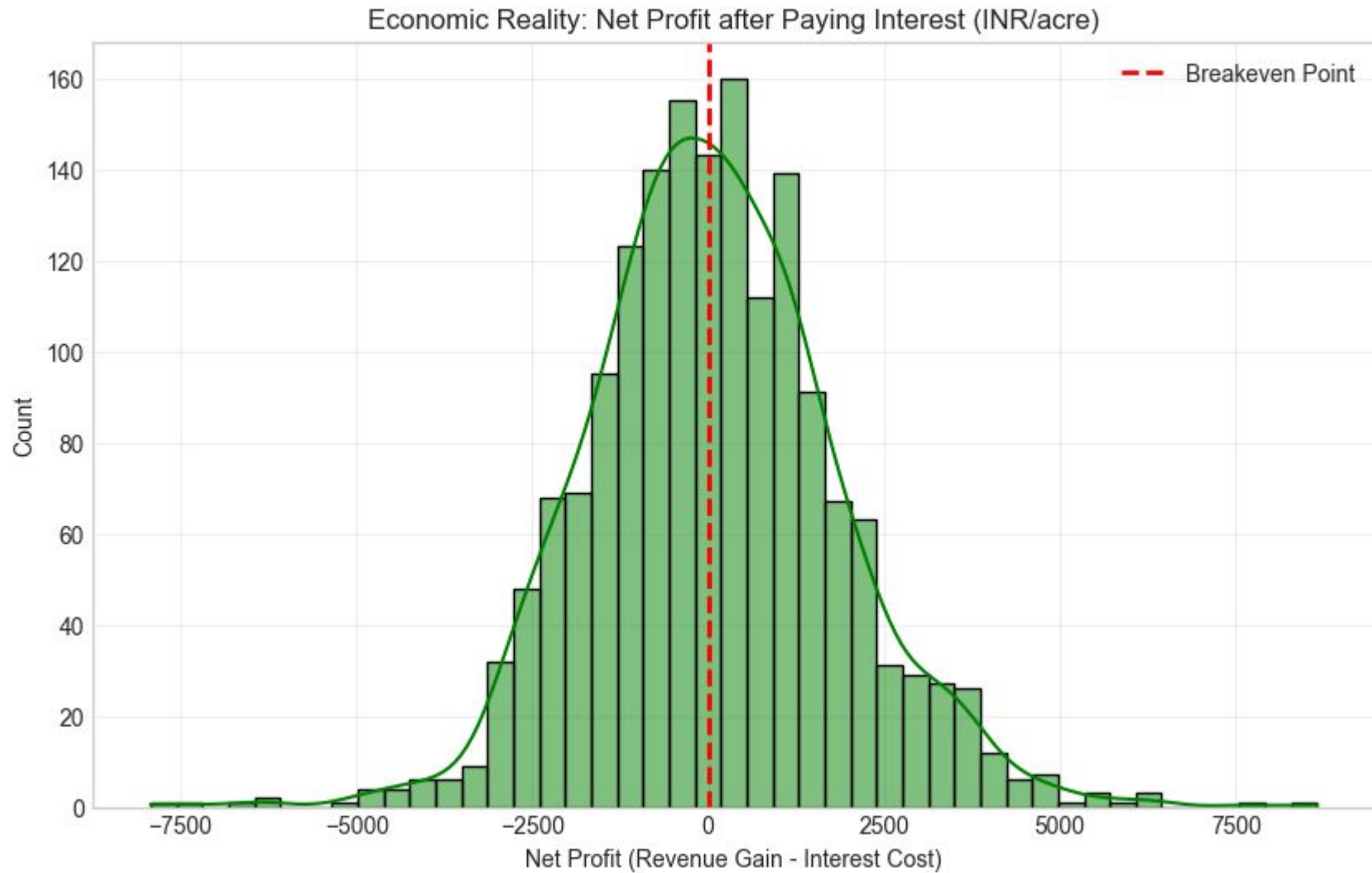
```
T_placebo = np.random.permutation(df['credit_access'])

est_placebo = CausalForestDML(
    model_y=RandomForestRegressor(n_estimators=50, min_samples_leaf=10),
    model_t=RandomForestClassifier(n_estimators=50, min_samples_leaf=10),
    discrete_treatment=True,
    n_estimators=100,
    random_state=42
)
est_placebo.fit(Y, T_placebo, X=X)
placebo_ate = est_placebo.ate(X)

print(f"Placebo ATE (Target ~0): {placebo_ate:.2f}")
if abs(placebo_ate) < abs(ate_val) * 0.2:
    print("Pass: Placebo effect is noise.")
else:
    print("Caution: Placebo effect is significant.")
```

Placebo ATE (Target ~0): -0.12
Pass: Placebo effect is noise.

Robustness Check (Placebo Test) When we randomized the treatment vector, the estimated effect dropped to **-0.12 kg/acre** (effectively zero). This null result confirms that the heterogeneity found in our main model is driven by real economic signals, not statistical noise.



SORGHUM PRICE = 25

CREDIT COST = 5000

FORMAL RATE = 0.12

INFORMAL RATE = 0.36

References



1. **Banerjee et al. (2015):** *The Miracle of Microfinance?* – Demonstrates that credit demand is often inelastic among the poor due to a lack of complementary assets.
2. **Chernozhukov et al. (2018):** *Double Machine Learning.* – Provides the rigorous identification strategy to handle high-dimensional confounding without regularization bias.
3. **Wager & Athey (2018):** *Causal Forests.* – Introduces data-driven estimation of heterogeneous treatment effects, allowing for the discovery of non-linear subgroups