

phData Case Study

(Technical Version)

Michael Rowlands

Tax Company's Challenge

For years, Tax Company has been unable to identify leads that will result in sales of their software.

This has resulted in thousands of dollars in unnecessary labor and advertising costs.

Tax Company has created a dataset containing two years of customer information and if they were successful at selling to each customer.

Our Objective

Investigate the dataset, and determine if a machine learning approach is viable.

If so, create a model to predict if a lead will convert.

Preliminary Assumptions

The dataset is historical so we assume future data is subject to the same data generating process.

Customers only appear once in the dataset.

Client has no prior knowledge on what features are predictive of sale.

There is a fixed profit for sale (true positives) and fixed loss for attempted sales (false positives).

Exploratory Analysis

Basic Dataset Info

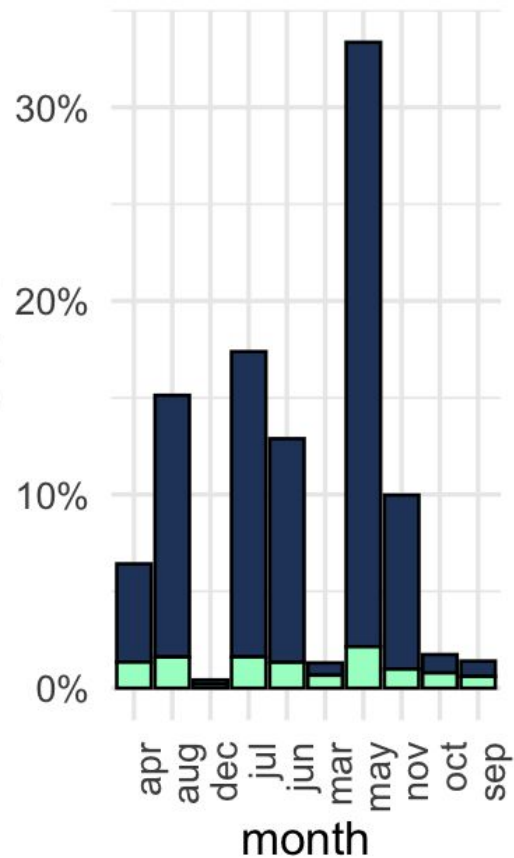
Full dataset contains 41,188 rows and 23 columns

We pulled a test set of 20% and put it aside for later

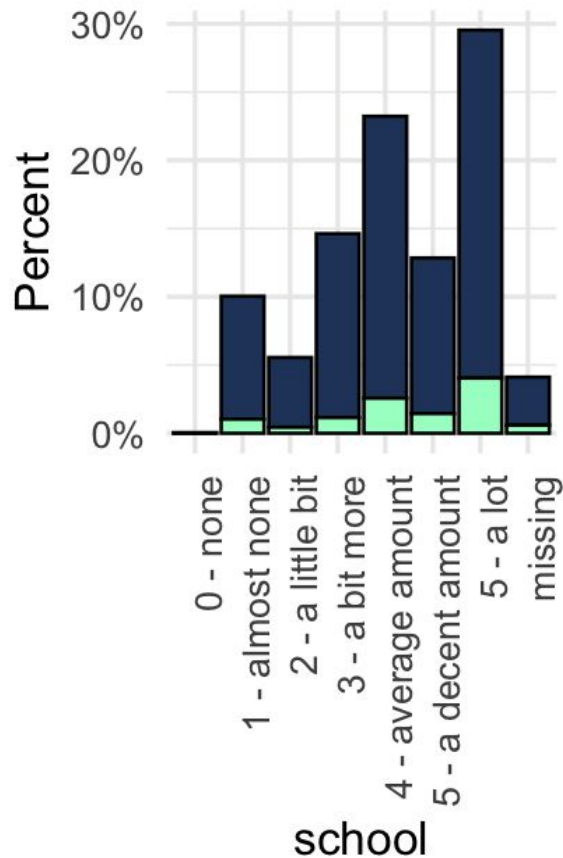
Target variable is “successful_sell”

About half the features are numeric and the other half are categorical

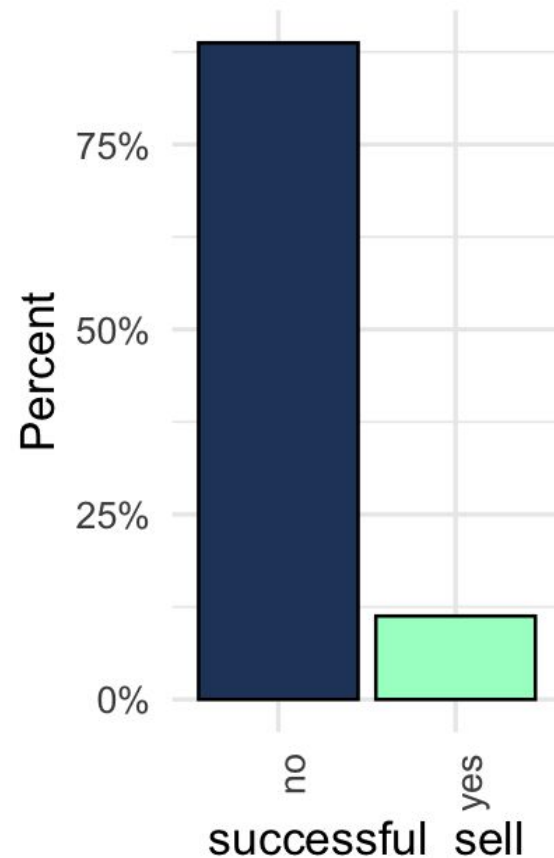
month



school

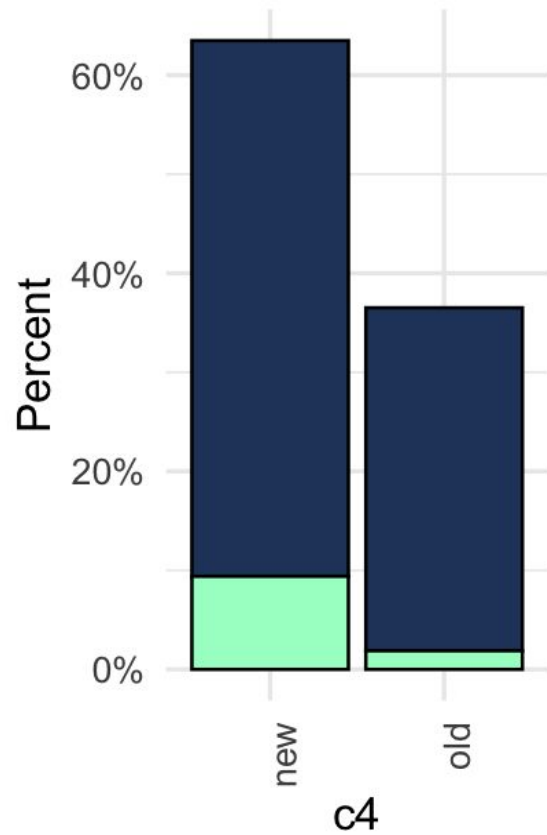


successful_sell

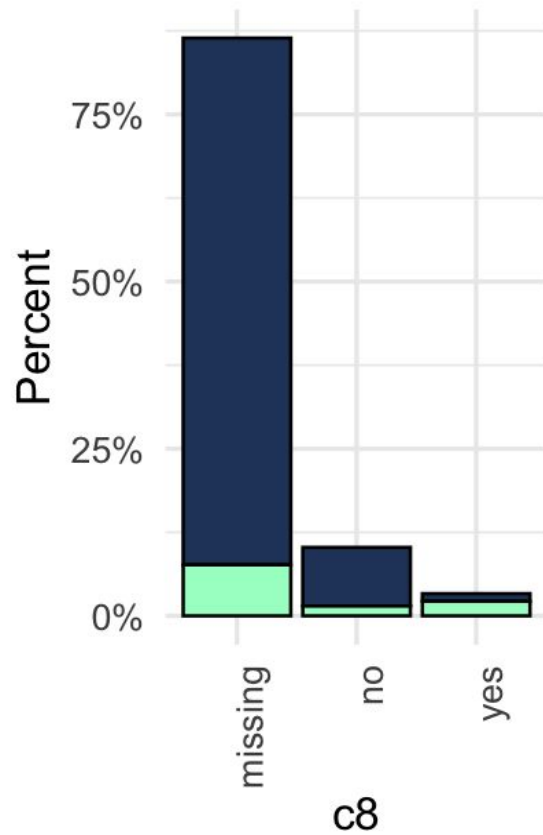


successful_sell no yes

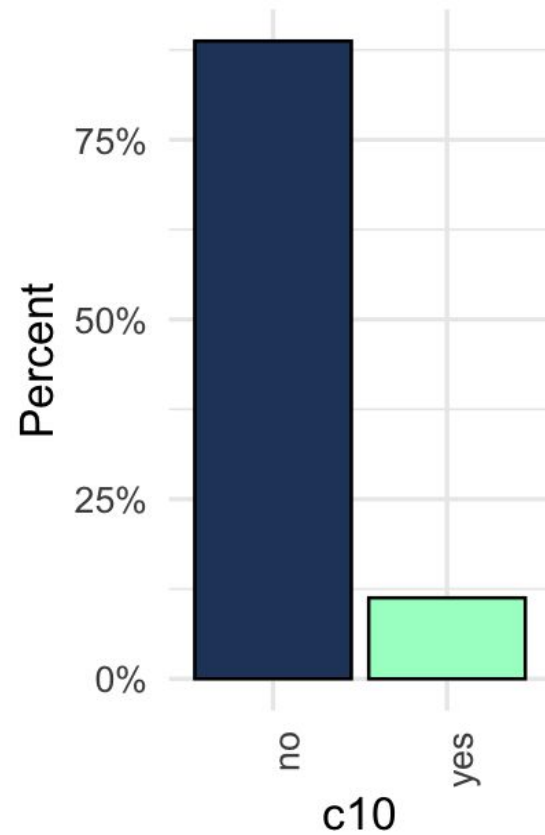
c4



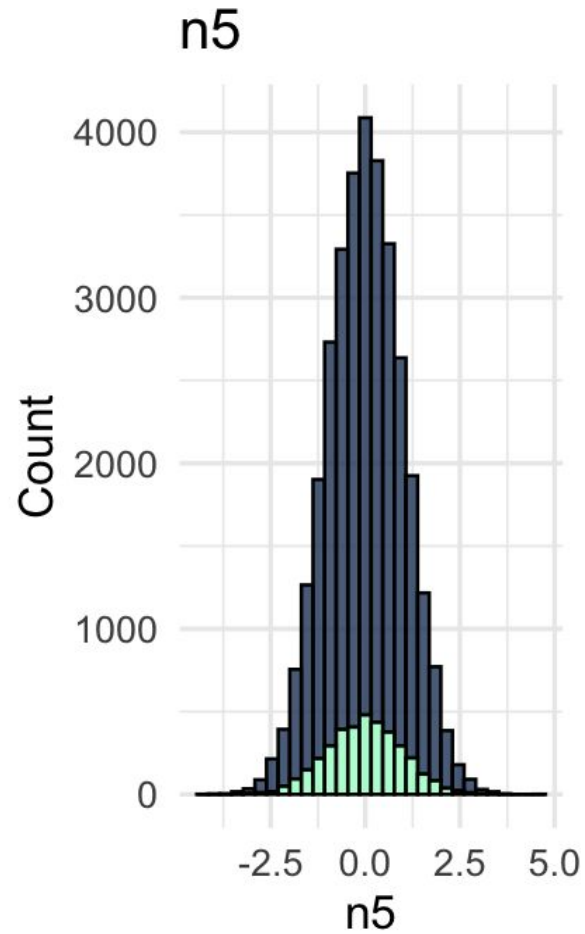
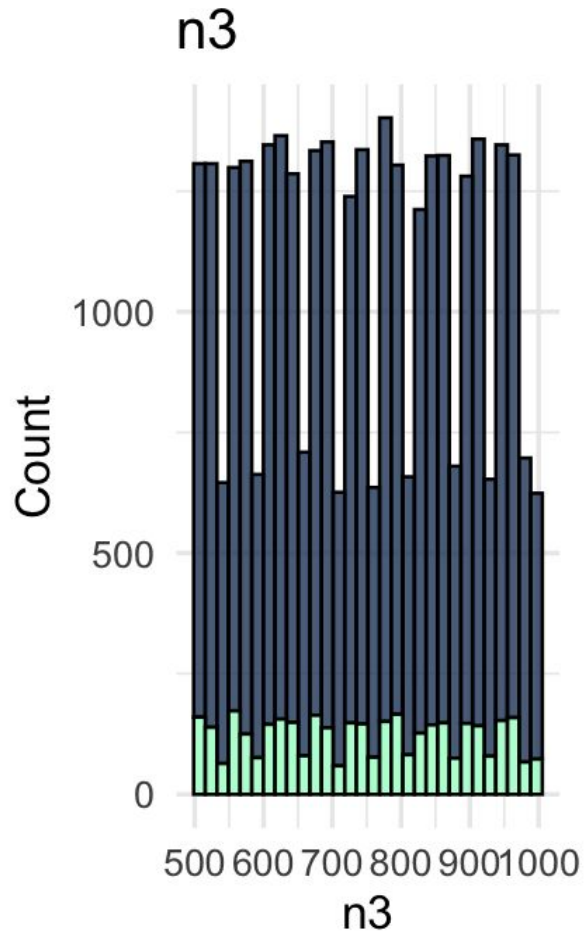
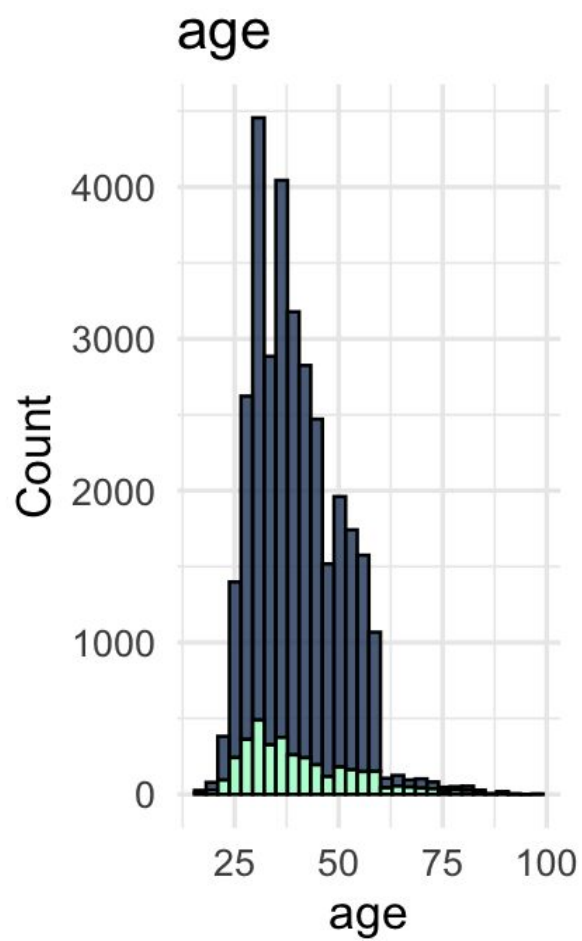
c8



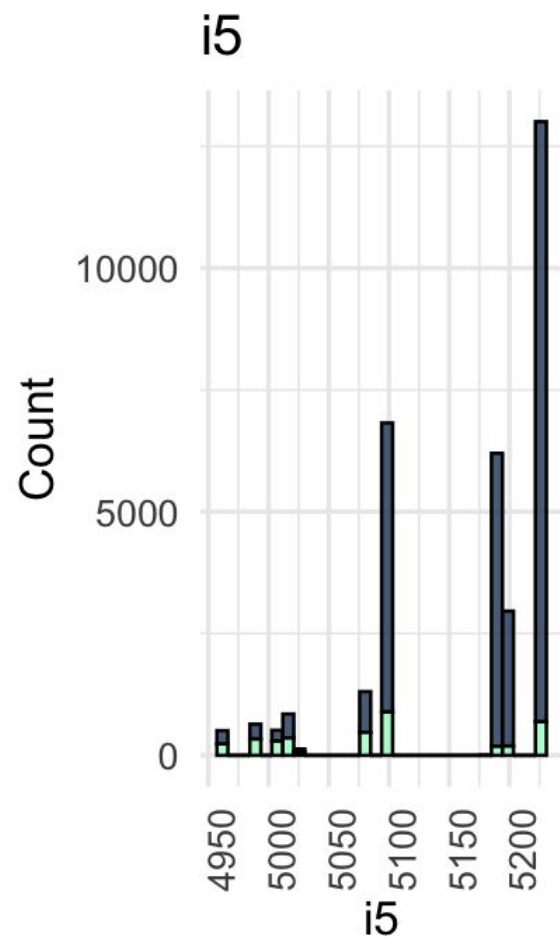
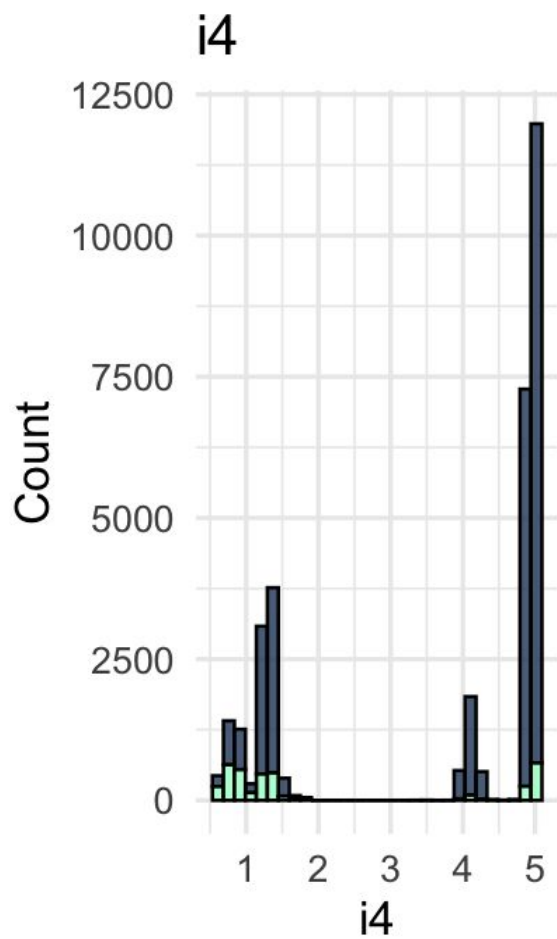
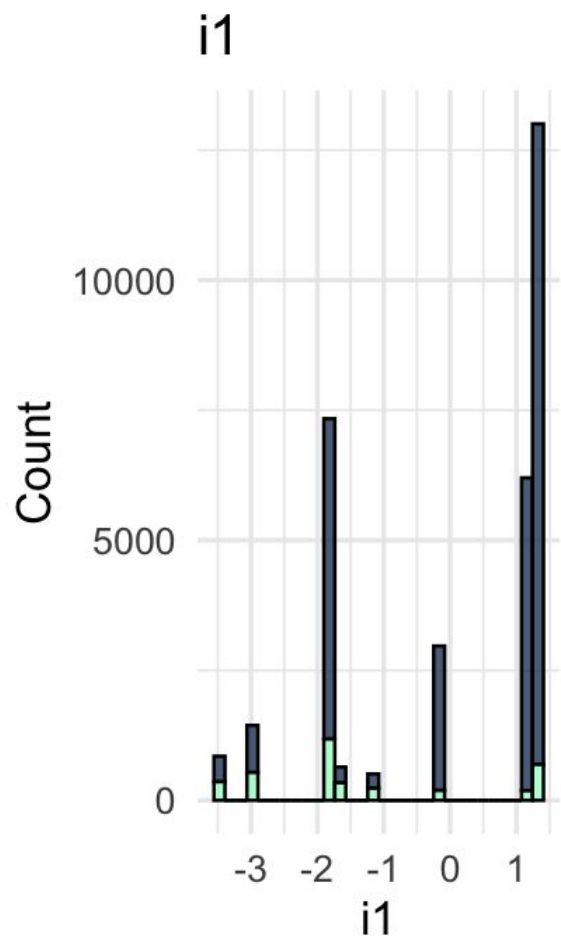
c10



successful_sell no yes

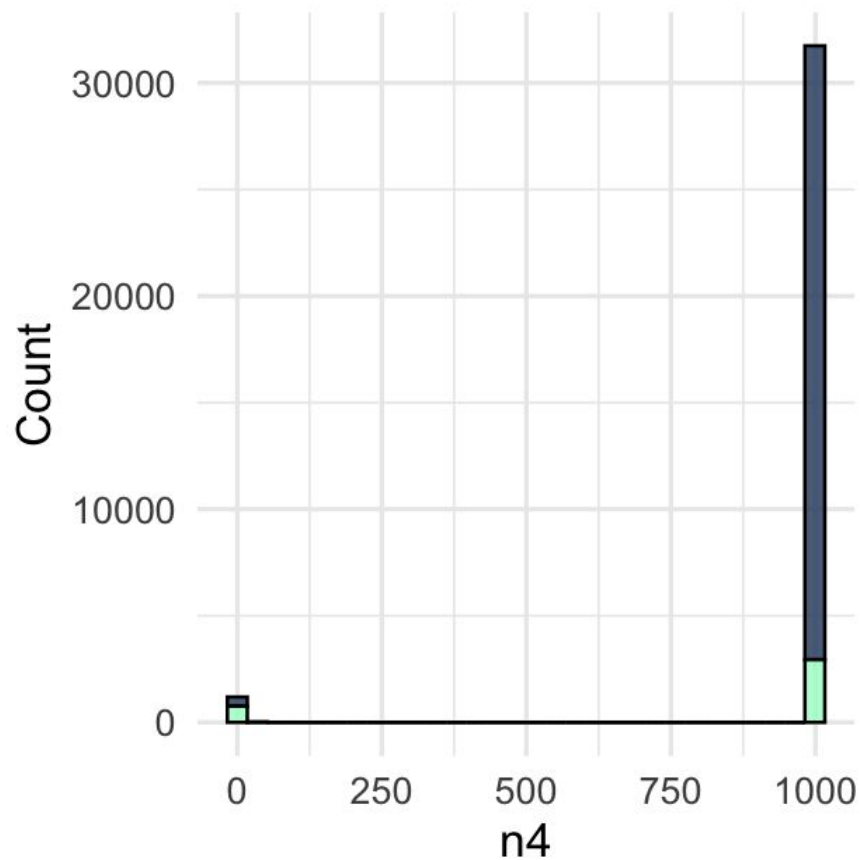


successful_sell no yes

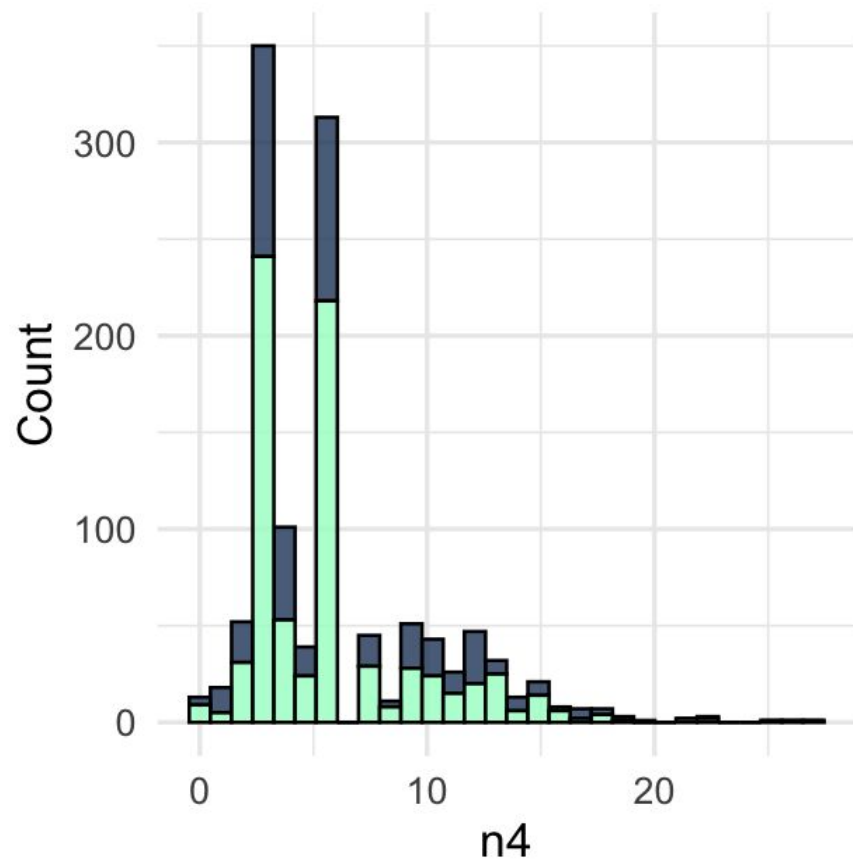


successful_sell no yes

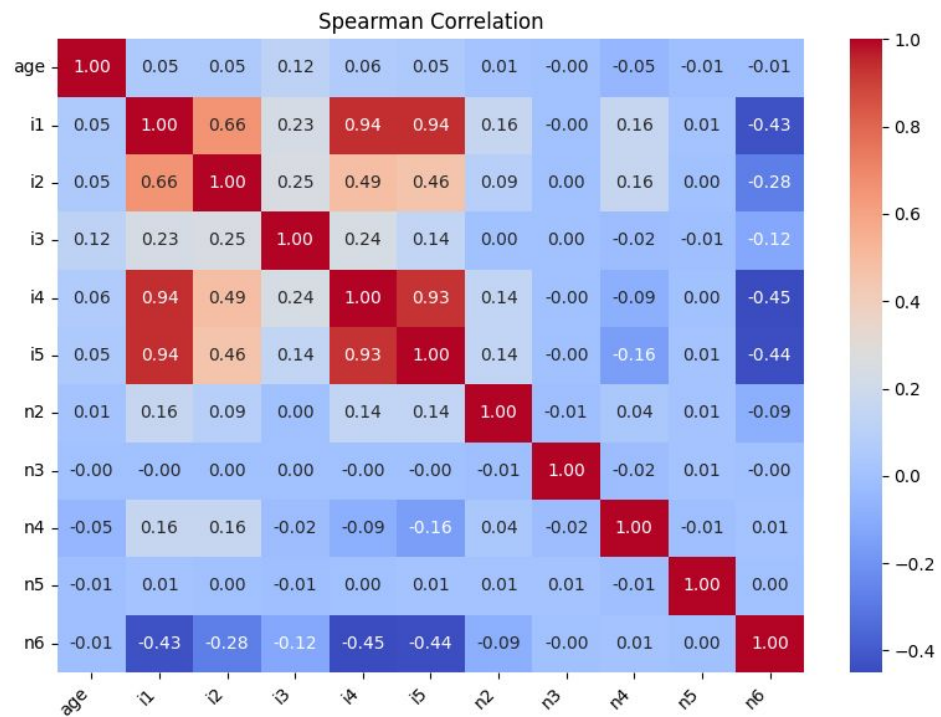
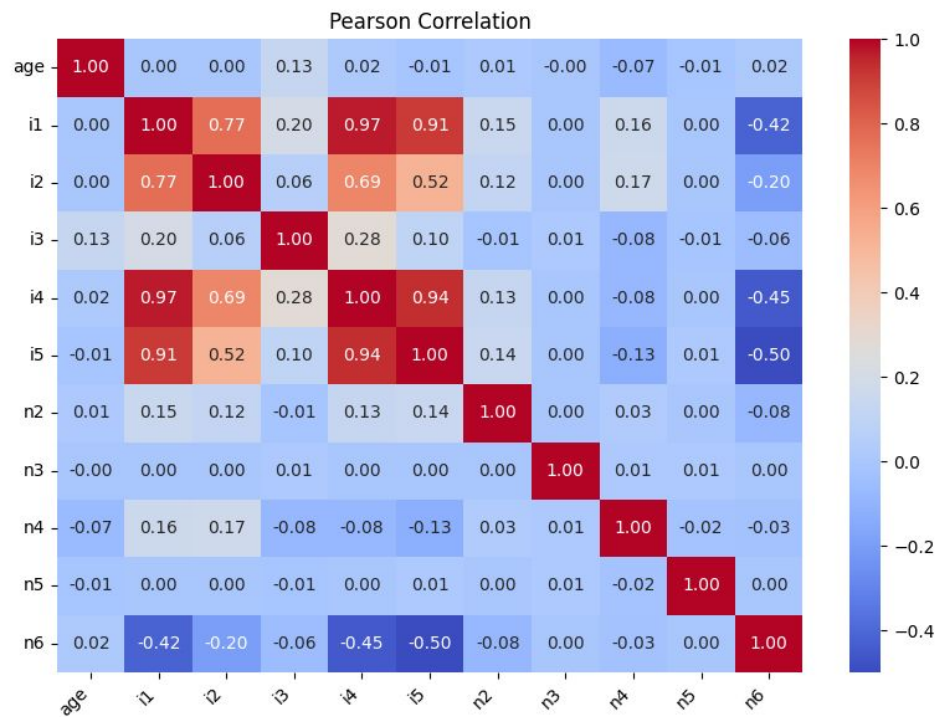
n4

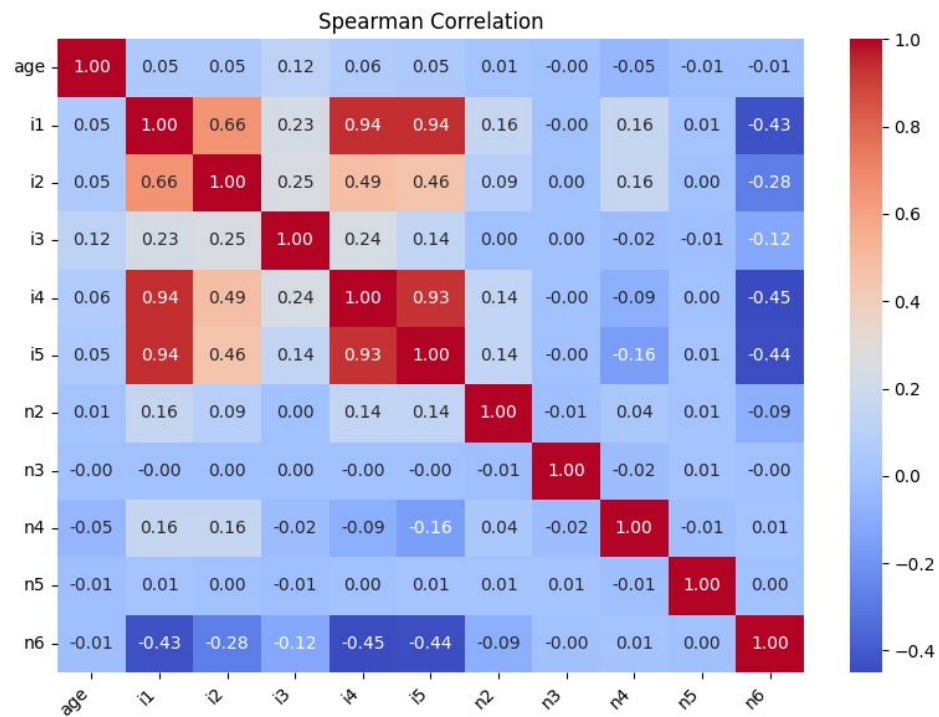
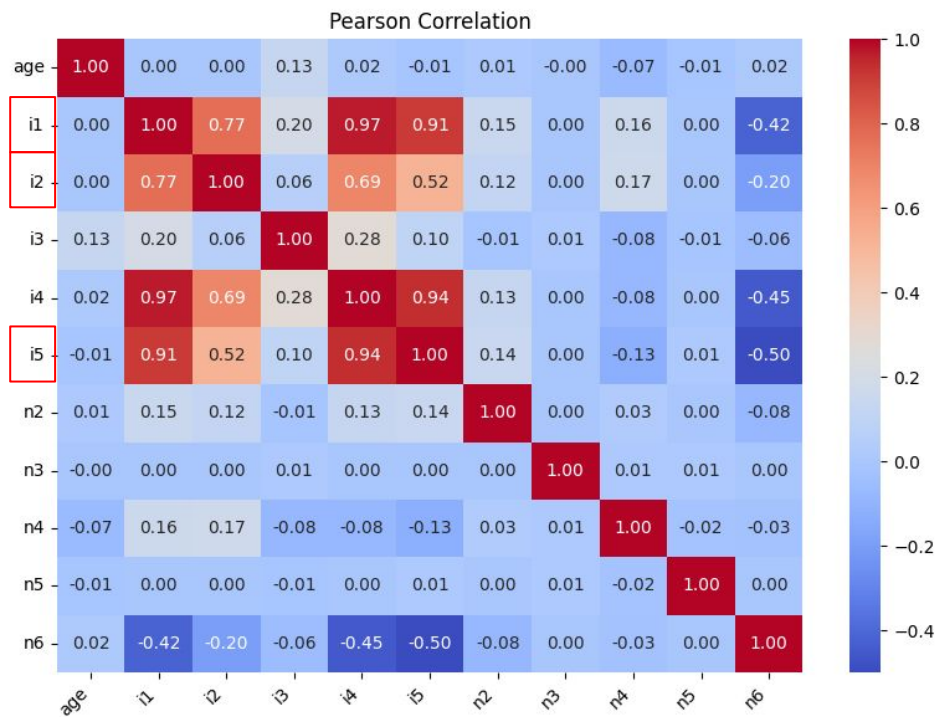


n4 Filtered

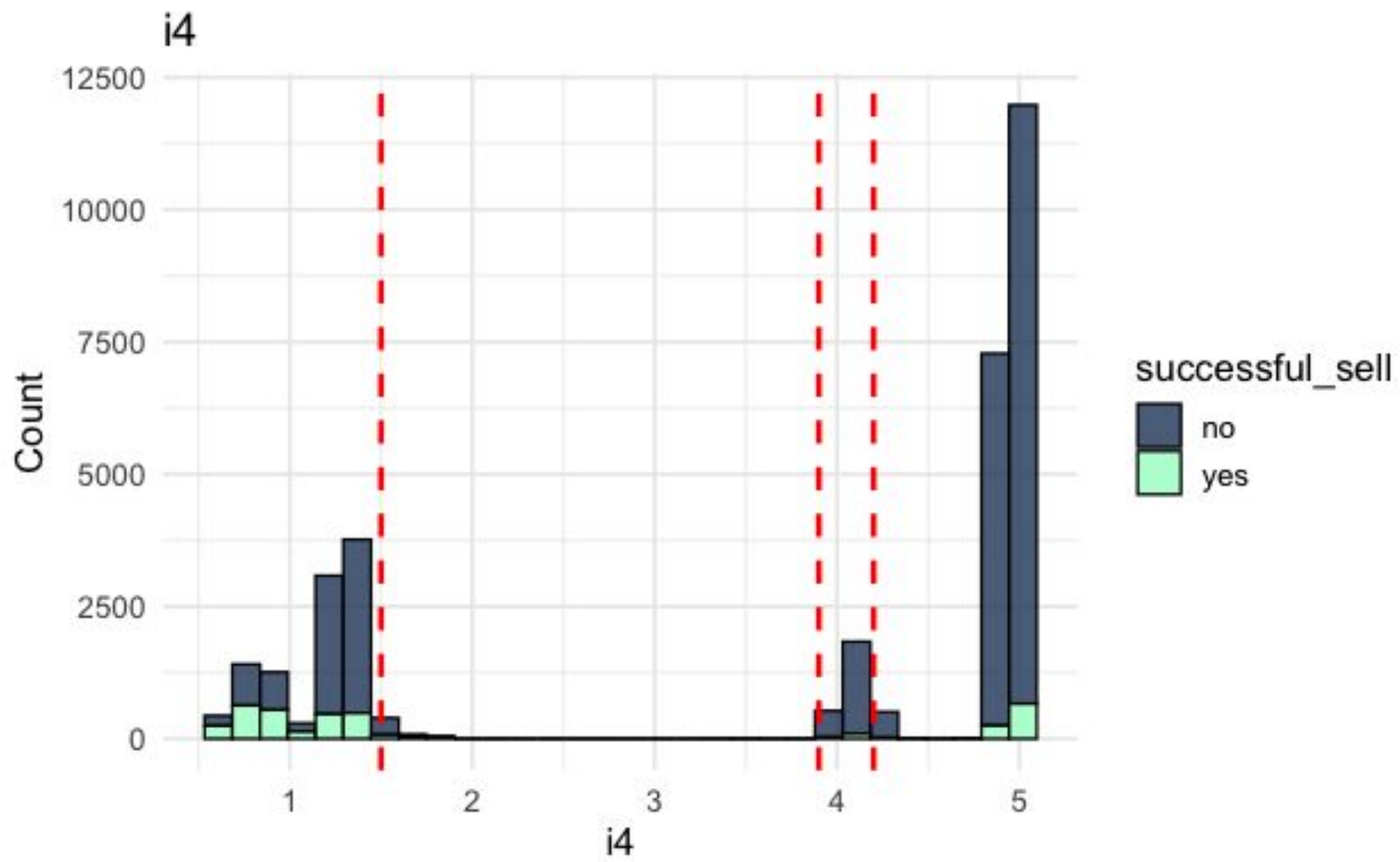


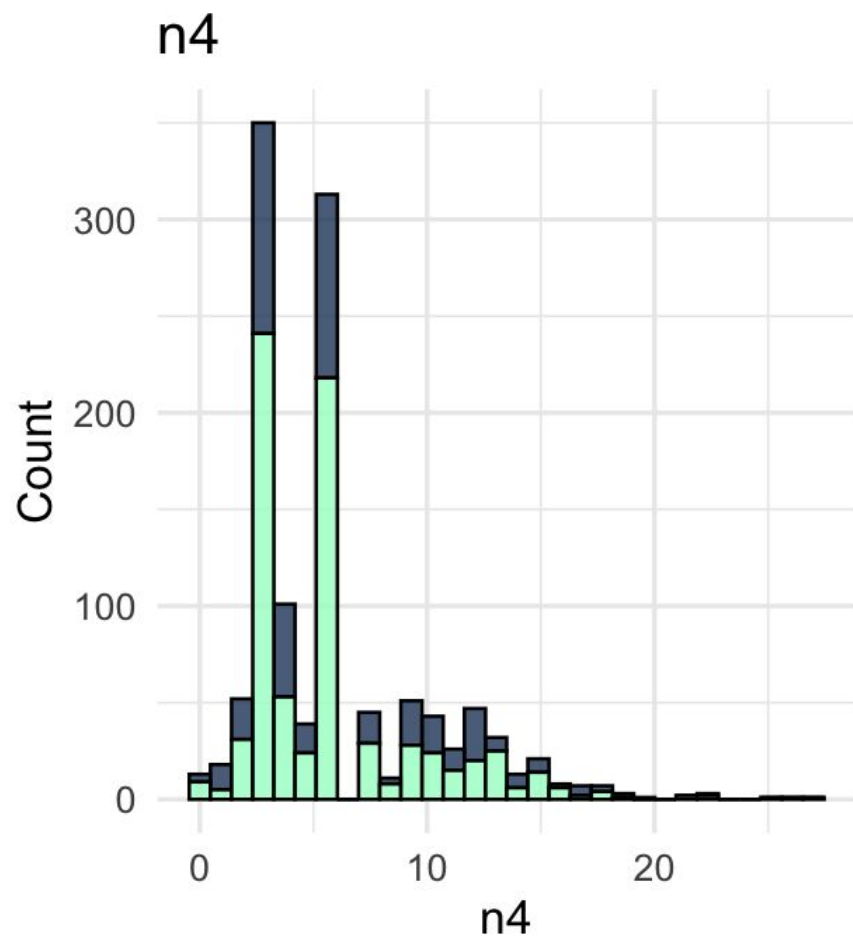
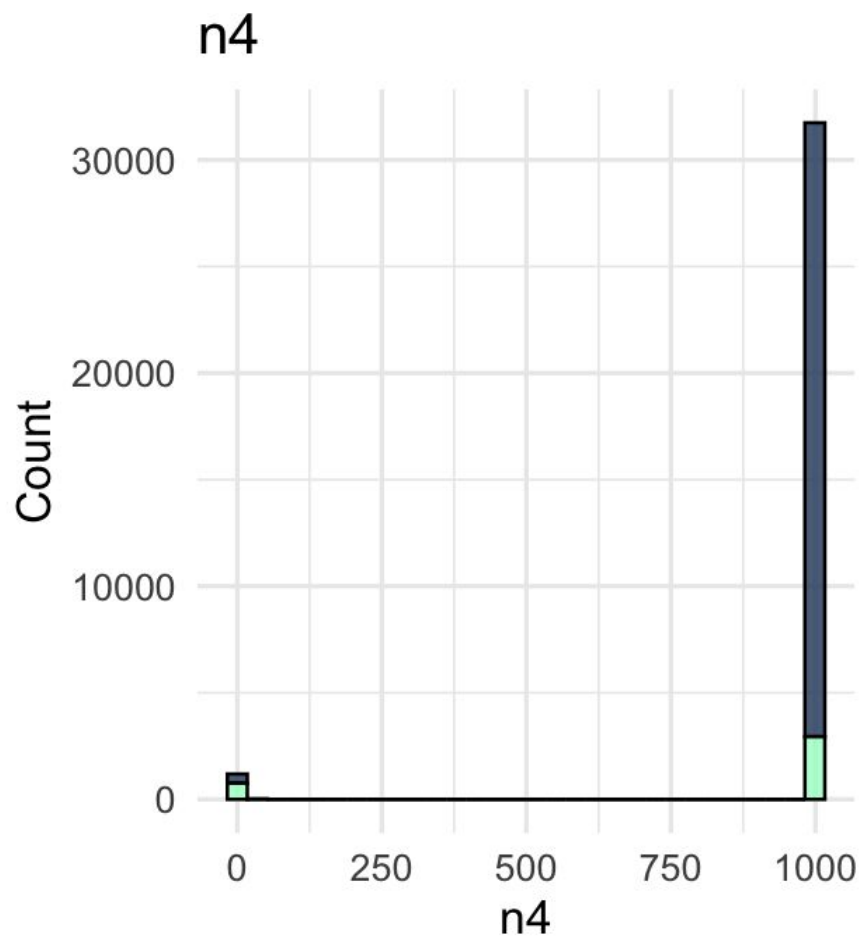
successful_sell no yes





Feature Engineering





successful_sell no yes

Imputation

Almost all of the missing data is in the categorical columns (with the exception of possibly n4)

We will leave these values as a “missing” level

While this can introduce bias (since we are assuming all data is missing for the same reason), we will leave more sophisticated methods for imputation such as MICE for future analysis

Modeling and Results

```

def evaluate_models_with_thresholds(
    models, X, y, preprocessor, n_splits=5, random_state=42,
    cost_params=None, thresholds=np.linspace(0.1, 0.9, 9), sampling_strategies=None
):
    """
    Evaluate multiple models with hyperparameter tuning using cross-validation,
    optimize classification thresholds for custom cost function, and calculate costs.

    Returns:
    - results: dict, evaluation metrics, costs, optimal thresholds, and predictions for each model
    """

    # Set default values of cost parameters and sampling strategies
    if cost_params is None:
        cost_params = {"tp_cost": 100, "tn_cost": 0, "fp_cost": -5, "fn_cost": 0}

    if sampling_strategies is None:
        sampling_strategies = [
            None,
            RandomOverSampler(random_state=random_state),
            RandomUnderSampler(random_state=random_state)
        ]

    # Initialize cross-validation strategy and results dictionary
    cv = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=random_state)
    results = {}

    # Iterate over models and sampling strategies with hyperparameter tuning using GridSearchCV
    for model_name, model, param_grid in models:

        for sampler in sampling_strategies:
            sampler_name = sampler.__class__.__name__ if sampler else "NoSampler"
            pipeline = ImbPipeline([
                ("preprocessor", preprocessor),
                ("sampler", sampler if sampler else "passthrough"),
                ("model", model)
            ])

            grid_search = GridSearchCV(
                pipeline,
                param_grid={"model__" + key: value for key, value in param_grid.items()},
                cv=cv,
                scoring='average_precision', # Use average precision for hyperparameter tuning
                n_jobs=-1,
                return_train_score=True
            )

            # Fit the GridSearchCV object
            grid_search.fit(X, y)

            # Pull best model from GridSearchCV using average precision as metric
            best_model = grid_search.best_estimator_

```

```

        # Calculate cross-validation score variation for the best model
        best_index = grid_search.best_index_
        cv_scores = grid_search.cv_results_['mean_test_score']
        cv_std = grid_search.cv_results_['std_test_score'][best_index]

        # Get cross-validated probabilities
        probs = cross_val_predict(best_model, X, y, cv=cv, method="predict_proba", n_jobs=-1)[ :, 1]

        # Optimize threshold based on custom cost
        best_threshold = None
        best_cost = float('-inf')
        best_preds = None

        for threshold in thresholds:

            # If the probability is greater than the threshold, predict as positive
            thresholded_preds = (probs >= threshold).astype(int)
            cost = custom_cost(y, thresholded_preds, **cost_params)

            # If threshold increases profit, update best cost, threshold, and predictions
            if cost > best_cost:
                best_cost = cost
                best_threshold = threshold
                best_preds = thresholded_preds

        # Pull feature importance from the best model if available
        feature_importance = None
        if hasattr(best_model.named_steps["model"], "feature_importances_"):
            feature_importance = best_model.named_steps["model"].feature_importances_
        elif hasattr(best_model.named_steps["model"], "coef_"):
            feature_importance = np.abs(best_model.named_steps["model"].coef_).flatten()

        # Confusion matrix
        confusion_mat = confusion_matrix(y, best_preds)

        # F1 score
        f1 = f1_score(y, best_preds)

        # Average precision
        average_precision = average_precision_score(y, probs)

        # Store results
        results[f"{model_name}_{sampler_name}"] = {
            "best_params": grid_search.best_params_,
            "cv_score_variation": cv_std, # Standard deviation of CV scores
            "best_threshold": best_threshold,
            "total_profit": best_cost,
            "average_profit_per_sale_attempt": best_cost / (confusion_mat[1,1] + confusion_mat[0,1]),
            "feature_importance": feature_importance,
            "model": best_model,
            "confusion_matrix": confusion_mat,
            "F1": f1,
            "average_precision_score": average_precision,
            "probs": probs, # Store probabilities
            "best_preds": best_preds, # Store thresholded predictions
        }

    return results

```

```

def evaluate_models_with_thresholds(
    models, X, y, preprocessor, n_splits=5, random_state=42,
    cost_params=None, thresholds=np.linspace(0.1, 0.9, 9), sampling_strategies=None
):
    """
    Evaluate multiple models with hyperparameter tuning using cross-validation,
    optimize classification thresholds for custom cost function, and calculate costs.

    Returns:
    - results: dict, evaluation metrics, costs, optimal thresholds, and predictions for each model
    """

    # Set default values of cost parameters and sampling strategies
    if cost_params is None:
        cost_params = {"tp_cost": 100, "tn_cost": 0, "fp_cost": -5, "fn_cost": 0}

    if sampling_strategies is None:
        sampling_strategies = [
            None,
            RandomOverSampler(random_state=random_state),
            RandomUnderSampler(random_state=random_state)
        ]

    # Initialize cross-validation strategy and results dictionary
    cv = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=random_state)
    results = {}

    # Iterate over models and sampling strategies with hyperparameter tuning using GridSearchCV
    for model_name, model, param_grid in models:

        for sampler in sampling_strategies:
            sampler_name = sampler.__class__.__name__ if sampler else "NoSampler"
            pipeline = ImbPipeline([
                ("preprocessor", preprocessor),
                ("sampler", sampler if sampler else "passthrough"),
                ("model", model)
            ])

            grid_search = GridSearchCV(
                pipeline,
                param_grid={"model__" + key: value for key, value in param_grid.items()},
                cv=cv,
                scoring='average_precision', # Use average precision for hyperparameter tuning
                n_jobs=-1,
                return_train_score=True
            )

            # Fit the GridSearchCV object
            grid_search.fit(X, y)

            # Pull best model from GridSearchCV using average precision as metric
            best_model = grid_search.best_estimator_

```

```

            # Calculate cross-validation score variation for the best model
            best_index = grid_search.best_index_
            cv_scores = grid_search.cv_results_['mean_test_score']
            cv_std = grid_search.cv_results_['std_test_score'][best_index]

            # Get cross-validated probabilities
            probs = cross_val_predict(best_model, X, y, cv=cv, method="predict_proba", n_jobs=-1)[ :, 1]

            # Optimize threshold based on custom cost
            best_threshold = None
            best_cost = float('-inf')
            best_preds = None

            for threshold in thresholds:

                # If the probability is greater than the threshold, predict as positive
                thresholded_preds = (probs >= threshold).astype(int)
                cost = custom_cost(y, thresholded_preds, **cost_params)

                # If threshold increases profit, update best cost, threshold, and predictions
                if cost > best_cost:
                    best_cost = cost
                    best_threshold = threshold
                    best_preds = thresholded_preds

            # Pull feature importance from the best model if available
            feature_importance = None
            if hasattr(best_model.named_steps["model"], "feature_importances_"):
                feature_importance = best_model.named_steps["model"].feature_importances_
            elif hasattr(best_model.named_steps["model"], "coef_"):
                feature_importance = np.abs(best_model.named_steps["model"].coef_).flatten()

            # Confusion matrix
            confusion_mat = confusion_matrix(y, best_preds)

            # F1 score
            f1 = f1_score(y, best_preds)

            # Average precision
            average_precision = average_precision_score(y, probs)

            # Store results
            results[f"{model_name}_{sampler_name}"] = {
                "best_params": grid_search.best_params_,
                "cv_score_variation": cv_std, # Standard deviation of CV scores
                "best_threshold": best_threshold,
                "total_profit": best_cost,
                "average_profit_per_sale_attempt": best_cost / (confusion_mat[1,1] + confusion_mat[0,1]),
                "feature_importance": feature_importance,
                "model": best_model,
                "confusion_matrix": confusion_mat,
                "F1": f1,
                "average_precision_score": average_precision,
                "probs": probs, # Store probabilities
                "best_preds": best_preds, # Store thresholded predictions
            }

    return results

```

```

def evaluate_models_with_thresholds(
    models, X, y, preprocessor, n_splits=5, random_state=42,
    cost_params=None, thresholds=np.linspace(0.1, 0.9, 9), sampling_strategies=None
):
    """
    Evaluate multiple models with hyperparameter tuning using cross-validation,
    optimize classification thresholds for custom cost function, and calculate costs.

    Returns:
    - results: dict, evaluation metrics, costs, optimal thresholds, and predictions for each model
    """

    # Set default values of cost parameters and sampling strategies
    if cost_params is None:
        cost_params = {"tp_cost": 100, "tn_cost": 0, "fp_cost": -5, "fn_cost": 0}

    if sampling_strategies is None:
        sampling_strategies = [
            None,
            RandomOverSampler(random_state=random_state),
            RandomUnderSampler(random_state=random_state)
        ]

    # Initialize cross-validation strategy and results dictionary
    cv = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=random_state)
    results = {}

    # Iterate over models and sampling strategies with hyperparameter tuning using GridSearchCV
    for model_name, model, param_grid in models:
        for sampler in sampling_strategies:
            sampler_name = sampler.__class__.__name__ if sampler else "NoSampler"
            pipeline = ImbPipeline([
                ("preprocessor", preprocessor),
                ("sampler", sampler if sampler else "passthrough"),
                ("model", model)
            ])

            grid_search = GridSearchCV(
                pipeline,
                param_grid={"model__" + key: value for key, value in param_grid.items()},
                cv=cv,
                scoring='average_precision', # Use average precision for hyperparameter tuning
                n_jobs=-1,
                return_train_score=True
            )

            # Fit the GridSearchCV object
            grid_search.fit(X, y)

            # Pull best model from GridSearchCV using average precision as metric
            best_model = grid_search.best_estimator_

```

```

            # Calculate cross-validation score variation for the best model
            best_index = grid_search.best_index_
            cv_scores = grid_search.cv_results_['mean_test_score']
            cv_std = grid_search.cv_results_['std_test_score'][best_index]

            # Get cross-validated probabilities
            probs = cross_val_predict(best_model, X, y, cv=cv, method="predict_proba", n_jobs=-1)[ :, 1]

            # Optimize threshold based on custom cost
            best_threshold = None
            best_cost = float('-inf')
            best_preds = None

            for threshold in thresholds:
                # If the probability is greater than the threshold, predict as positive
                thresholded_preds = (probs >= threshold).astype(int)
                cost = custom_cost(y, thresholded_preds, **cost_params)

                # If threshold increases profit, update best cost, threshold, and predictions
                if cost > best_cost:
                    best_cost = cost
                    best_threshold = threshold
                    best_preds = thresholded_preds

            # Pull feature importance from the best model if available
            feature_importance = None
            if hasattr(best_model.named_steps["model"], "feature_importances_"):
                feature_importance = best_model.named_steps["model"].feature_importances_
            elif hasattr(best_model.named_steps["model"], "coef_"):
                feature_importance = np.abs(best_model.named_steps["model"].coef_).flatten()

            # Confusion matrix
            confusion_mat = confusion_matrix(y, best_preds)

            # F1 score
            f1 = f1_score(y, best_preds)

            # Average precision
            average_precision = average_precision_score(y, probs)

            # Store results
            results[f"{model_name}_{sampler_name}"] = {
                "best_params": grid_search.best_params_,
                "cv_score_variation": cv_std, # Standard deviation of CV scores
                "best_threshold": best_threshold,
                "total_profit": best_cost,
                "average_profit_per_sale_attempt": best_cost / (confusion_mat[1,1] + confusion_mat[0,1]),
                "feature_importance": feature_importance,
                "model": best_model,
                "confusion_matrix": confusion_mat,
                "F1": f1,
                "average_precision_score": average_precision,
                "probs": probs, # Store probabilities
                "best_preds": best_preds, # Store thresholded predictions
            }

    return results

```



```

def evaluate_models_with_thresholds(
    models, X, y, preprocessor, n_splits=5, random_state=42,
    cost_params=None, thresholds=np.linspace(0.1, 0.9, 9), sampling_strategies=None
):
    """
    Evaluate multiple models with hyperparameter tuning using cross-validation,
    optimize classification thresholds for custom cost function, and calculate costs.

    Returns:
    - results: dict, evaluation metrics, costs, optimal thresholds, and predictions for each model
    """

    # Set default values of cost parameters and sampling strategies
    if cost_params is None:
        cost_params = {"tp_cost": 100, "tn_cost": 0, "fp_cost": -5, "fn_cost": 0}

    if sampling_strategies is None:
        sampling_strategies = [
            None,
            RandomOverSampler(random_state=random_state),
            RandomUnderSampler(random_state=random_state)
        ]

    # Initialize cross-validation strategy and results dictionary
    cv = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=random_state)
    results = {}

    # Iterate over models and sampling strategies with hyperparameter tuning using GridSearchCV
    for model_name, model, param_grid in models:

        for sampler in sampling_strategies:
            sampler_name = sampler.__class__.__name__ if sampler else "NoSampler"
            pipeline = ImbPipeline([
                ("preprocessor", preprocessor),
                ("sampler", sampler if sampler else "passthrough"),
                ("model", model)
            ])

            grid_search = GridSearchCV(
                pipeline,
                param_grid={"model__" + key: value for key, value in param_grid.items()},
                cv=cv,
                scoring='average_precision', # Use average precision for hyperparameter tuning
                n_jobs=-1,
                return_train_score=True
            )

            # Fit the GridSearchCV object
            grid_search.fit(X, y)

            # Pull best model from GridSearchCV using average precision as metric
            best_model = grid_search.best_estimator_

```

```

            # Calculate cross-validation score variation for the best model
            best_index = grid_search.best_index_
            cv_scores = grid_search.cv_results_['mean_test_score']
            cv_std = grid_search.cv_results_['std_test_score'][best_index]

            # Get cross-validated probabilities
            probs = cross_val_predict(best_model, X, y, cv=cv, method="predict_proba", n_jobs=-1)[ :, 1]

            # Optimize threshold based on custom cost
            best_threshold = None
            best_cost = float('-inf')
            best_preds = None

            for threshold in thresholds:

                # If the probability is greater than the threshold, predict as positive
                thresholded_preds = (probs >= threshold).astype(int)
                cost = custom_cost(y, thresholded_preds, **cost_params)

                # If threshold increases profit, update best cost, threshold, and predictions
                if cost > best_cost:
                    best_cost = cost
                    best_threshold = threshold
                    best_preds = thresholded_preds

            # Pull feature importance from the best model if available
            feature_importance = None
            if hasattr(best_model.named_steps["model"], "feature_importances_"):
                feature_importance = best_model.named_steps["model"].feature_importances_
            elif hasattr(best_model.named_steps["model"], "coef_"):
                feature_importance = np.abs(best_model.named_steps["model"].coef_).flatten()

            # Confusion matrix
            confusion_mat = confusion_matrix(y, best_preds)

            # F1 score
            f1 = f1_score(y, best_preds)

            # Average precision
            average_precision = average_precision_score(y, probs)

            # Store results
            results[f"{model_name}_{sampler_name}"] = {
                "best_params": grid_search.best_params_,
                "cv_score_variation": cv_std, # Standard deviation of CV scores
                "best_threshold": best_threshold,
                "total_profit": best_cost,
                "average_profit_per_sale_attempt": best_cost / (confusion_mat[1,1] + confusion_mat[0,1]),
                "feature_importance": feature_importance,
                "model": best_model,
                "confusion_matrix": confusion_mat,
                "F1": f1,
                "average_precision_score": average_precision,
                "probs": probs, # Store probabilities
                "best_preds": best_preds, # Store thresholded predictions
            }

    return results

```

```

def evaluate_models_with_thresholds(
    models, X, y, preprocessor, n_splits=5, random_state=42,
    cost_params=None, thresholds=np.linspace(0.1, 0.9, 9), sampling_strategies=None
):
    """
    Evaluate multiple models with hyperparameter tuning using cross-validation,
    optimize classification thresholds for custom cost function, and calculate costs.

    Returns:
    - results: dict, evaluation metrics, costs, optimal thresholds, and predictions for each model
    """

    # Set default values of cost parameters and sampling strategies
    if cost_params is None:
        cost_params = {"tp_cost": 100, "tn_cost": 0, "fp_cost": -5, "fn_cost": 0}

    if sampling_strategies is None:
        sampling_strategies = [
            None,
            RandomOverSampler(random_state=random_state),
            RandomUnderSampler(random_state=random_state)
        ]

    # Initialize cross-validation strategy and results dictionary
    cv = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=random_state)
    results = {}

    # Iterate over models and sampling strategies with hyperparameter tuning using GridSearchCV
    for model_name, model, param_grid in models:

        for sampler in sampling_strategies:
            sampler_name = sampler.__class__.__name__ if sampler else "NoSampler"
            pipeline = ImbPipeline([
                ("preprocessor", preprocessor),
                ("sampler", sampler if sampler else "passthrough"),
                ("model", model)
            ])

            grid_search = GridSearchCV(
                pipeline,
                param_grid={"model__" + key: value for key, value in param_grid.items()},
                cv=cv,
                scoring='average_precision', # Use average precision for hyperparameter tuning
                n_jobs=-1,
                return_train_score=True
            )

            # Fit the GridSearchCV object
            grid_search.fit(X, y)

            # Pull best model from GridSearchCV using average precision as metric
            best_model = grid_search.best_estimator_

```

```

        # Calculate cross-validation score variation for the best model
        best_index = grid_search.best_index_
        cv_scores = grid_search.cv_results_['mean_test_score']
        cv_std = grid_search.cv_results_['std_test_score'][best_index]

        # Get cross-validated probabilities
        probs = cross_val_predict(best_model, X, y, cv=cv, method="predict_proba", n_jobs=-1)[ :, 1]

        # Optimize threshold based on custom cost
        best_threshold = None
        best_cost = float('-inf')
        best_preds = None

        for threshold in thresholds:

            # If the probability is greater than the threshold, predict as positive
            thresholded_preds = (probs >= threshold).astype(int)
            cost = custom_cost(y, thresholded_preds, **cost_params)

            # If threshold increases profit, update best cost, threshold, and predictions
            if cost > best_cost:
                best_cost = cost
                best_threshold = threshold
                best_preds = thresholded_preds

        # Pull feature importance from the best model if available
        feature_importance = None
        if hasattr(best_model.named_steps["model"], "feature_importances_"):
            feature_importance = best_model.named_steps["model"].feature_importances_
        elif hasattr(best_model.named_steps["model"], "coef_"):
            feature_importance = np.abs(best_model.named_steps["model"].coef_).flatten()

        # Confusion matrix
        confusion_mat = confusion_matrix(y, best_preds)

        # F1 score
        f1 = f1_score(y, best_preds)

        # Average precision
        average_precision = average_precision_score(y, probs)

        # Store results
        results[f"{model_name}_{sampler_name}"] = {
            "best_params": grid_search.best_params_,
            "cv_score_variation": cv_std, # Standard deviation of CV scores
            "best_threshold": best_threshold,
            "total_profit": best_cost,
            "average_profit_per_sale_attempt": best_cost / (confusion_mat[1,1] + confusion_mat[0,1]),
            "feature_importance": feature_importance,
            "model": best_model,
            "confusion_matrix": confusion_mat,
            "F1": f1,
            "average_precision_score": average_precision,
            "probs": probs, # Store probabilities
            "best_preds": best_preds, # Store thresholded predictions
        }

```

return results

```

def evaluate_models_with_thresholds(
    models, X, y, preprocessor, n_splits=5, random_state=42,
    cost_params=None, thresholds=np.linspace(0.1, 0.9, 9), sampling_strategies=None
):
    """
    Evaluate multiple models with hyperparameter tuning using cross-validation,
    optimize classification thresholds for custom cost function, and calculate costs.

    Returns:
    - results: dict, evaluation metrics, costs, optimal thresholds, and predictions for each model
    """

    # Set default values of cost parameters and sampling strategies
    if cost_params is None:
        cost_params = {"tp_cost": 100, "tn_cost": 0, "fp_cost": -5, "fn_cost": 0}

    if sampling_strategies is None:
        sampling_strategies = [
            None,
            RandomOverSampler(random_state=random_state),
            RandomUnderSampler(random_state=random_state)
        ]

    # Initialize cross-validation strategy and results dictionary
    cv = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=random_state)
    results = {}

    # Iterate over models and sampling strategies with hyperparameter tuning using GridSearchCV
    for model_name, model, param_grid in models:

        for sampler in sampling_strategies:
            sampler_name = sampler.__class__.__name__ if sampler else "NoSampler"
            pipeline = ImbPipeline([
                ("preprocessor", preprocessor),
                ("sampler", sampler if sampler else "passthrough"),
                ("model", model)
            ])

            grid_search = GridSearchCV(
                pipeline,
                param_grid={"model__" + key: value for key, value in param_grid.items()},
                cv=cv,
                scoring='average_precision', # Use average precision for hyperparameter tuning
                n_jobs=-1,
                return_train_score=True
            )

            # Fit the GridSearchCV object
            grid_search.fit(X, y)

            # Pull best model from GridSearchCV using average precision as metric
            best_model = grid_search.best_estimator_

```

```

        # Calculate cross-validation score variation for the best model
        best_index = grid_search.best_index_
        cv_scores = grid_search.cv_results_['mean_test_score']
        cv_std = grid_search.cv_results_['std_test_score'][best_index]

        # Get cross-validated probabilities
        probs = cross_val_predict(best_model, X, y, cv=cv, method="predict_proba", n_jobs=-1)[ :, 1]

        # Optimize threshold based on custom cost
        best_threshold = None
        best_cost = float('-inf')
        best_preds = None

        for threshold in thresholds:

            # If the probability is greater than the threshold, predict as positive
            thresholded_preds = (probs >= threshold).astype(int)
            cost = custom_cost(y, thresholded_preds, **cost_params)

            # If threshold increases profit, update best cost, threshold, and predictions
            if cost > best_cost:
                best_cost = cost
                best_threshold = threshold
                best_preds = thresholded_preds

        # Pull feature importance from the best model if available
        feature_importance = None
        if hasattr(best_model.named_steps["model"], "feature_importances_"):
            feature_importance = best_model.named_steps["model"].feature_importances_
        elif hasattr(best_model.named_steps["model"], "coef_"):
            feature_importance = np.abs(best_model.named_steps["model"].coef_).flatten()

        # Confusion matrix
        confusion_mat = confusion_matrix(y, best_preds)

        # F1 score
        f1 = f1_score(y, best_preds)

        # Average precision
        average_precision = average_precision_score(y, probs)

        # Store results
        results[f"{model_name}_{sampler_name}"] = {
            "best_params": grid_search.best_params_,
            "cv_score_variation": cv_std, # Standard deviation of CV scores
            "best_threshold": best_threshold,
            "total_profit": best_cost,
            "average_profit_per_sale_attempt": best_cost / (confusion_mat[1,1] + confusion_mat[0,1]),
            "feature_importance": feature_importance,
            "model": best_model,
            "confusion_matrix": confusion_mat,
            "F1": f1,
            "average_precision_score": average_precision,
            "probs": probs, # Store probabilities
            "best_preds": best_preds, # Store thresholded predictions
        }

```

return results


```

models = [
    ("Dummy",
     DummyClassifier(strategy="constant", constant=1, random_state=42),
     {}),
    (
        "Random Forest",
        RandomForestClassifier(random_state=42),
        {"class_weight": ["balanced", None], "n_estimators": [50, 100, 500], "max_depth": [5, 10, None]}
    ),
    (
        "Penalized Logistic Regression",
        LogisticRegression(random_state=42, solver='liblinear'),
        {"C": [0.1, 1, 10], "penalty": ["l1", "l2"]}
    ),
    (
        "XGBoost",
        XGBClassifier(random_state=42, eval_metric="logloss"),
        {"n_estimators": [50, 100, 500], "max_depth": [3, 6, None], "learning_rate": [0.01, 0.1, 0.2], "reg_alpha": [0, .1, .5, 1, 10]}
    )
]

# Preprocessing pipeline
preprocessor = ColumnTransformer(
    transformers=[
        ("num", StandardScaler(), num_cols_X),      # Scale numeric columns
        ("cat", OneHotEncoder(drop="first"), cat_cols_X) # One-hot encode categorical columns
    ]
)

# Define profit for each type of prediction
cost_params = {"tp_cost": 100, "tn_cost": 0, "fp_cost": -15, "fn_cost": 0}

# Evaluate models with threshold optimization
results = evaluate_models_with_thresholds([models=models, X=X_train, y=y_train, preprocessor=preprocessor, cost_params=cost_params,
                                           thresholds=np.linspace(0.05, .95, 19),
                                           sampling_strategies=[None, RandomUnderSampler(random_state=42), RandomOverSampler(random_state=42)],
                                           h_splits=5, random_state=42])

# Display model diagnostics
model_diagnostics(results)

```



```
Model: XGBoost_RandomOverSampler
Best Parameters: {'model__learning_rate': 0.1, 'model__max_depth': 6, 'model__n_estimators': 50, 'model__reg_alpha': 10}
Total Profit: 175475.00
Average Profit per Sale Attempts: 31.39
Best Threshold: 0.55
Confusion Matrix:
[[25903  3335]
 [ 1457  2255]]
F1 Score: 0.48
Average Precision Score: 0.46
CV Average Precision Variation (std): 0.013
-----
Model: Random Forest_NoSampler
Best Parameters: {'model__class_weight': None, 'model__max_depth': 10, 'model__n_estimators': 500}
Total Profit: 175395.00
Average Profit per Sale Attempts: 30.67
Best Threshold: 0.15
Confusion Matrix:
[[25791  3447]
 [ 1441  2271]]
F1 Score: 0.48
Average Precision Score: 0.46
CV Average Precision Variation (std): 0.015
-----
Model: Dummy_RandomOverSampler
Best Parameters: {}
Total Profit: -67370.00
Average Profit per Sale Attempts: -2.04
Best Threshold: 0.05
Confusion Matrix:
[[    0 29238]
 [    0  3712]]
F1 Score: 0.20
Average Precision Score: 0.11
CV Average Precision Variation (std): 0.000
-----
```

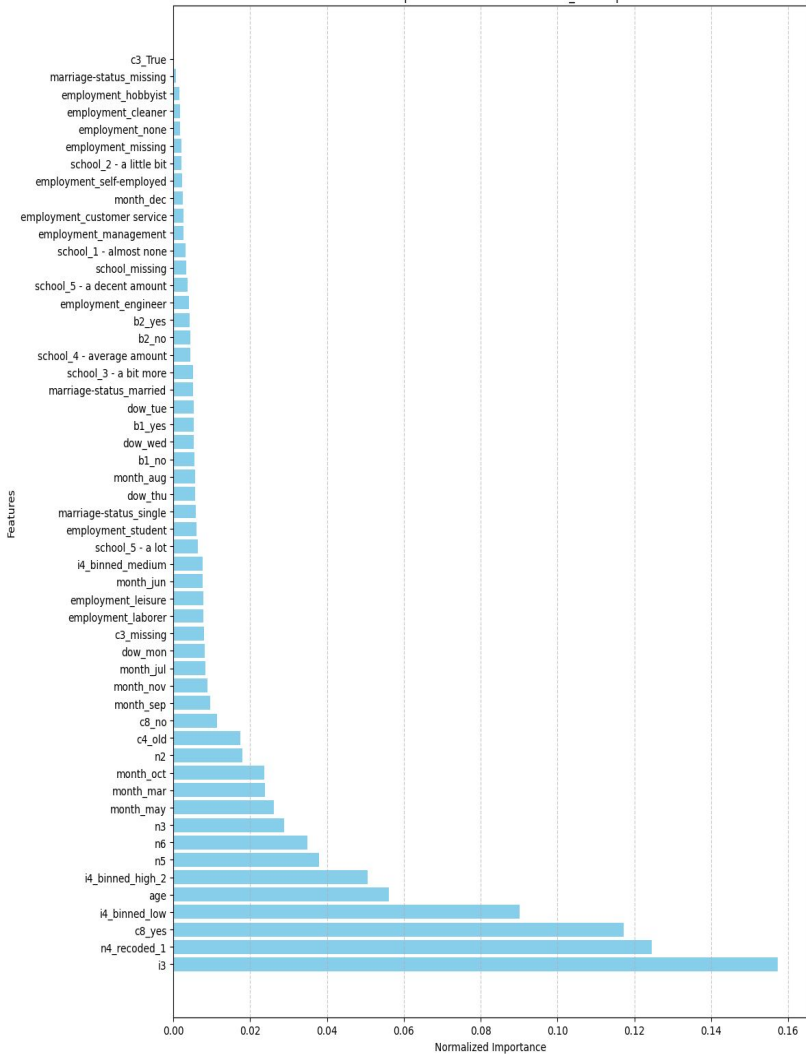
```
Model: XGBoost_RandomOverSampler
Best Parameters: {'model__learning_rate': 0.1, 'model__max_depth': 6, 'model__n_estimators': 50, 'model__reg_alpha': 10}
Total Profit: 175475.00
Average Profit per Sale Attempts: 31.39
Best Threshold: 0.55
Confusion Matrix:
[[25903  3335]
 [ 1457  2255]]
F1 Score: 0.48
Average Precision Score: 0.46
CV Average Precision Variation (std): 0.013
-----
Model: Random Forest_NoSampler
Best Parameters: {'model__class_weight': None, 'model__max_depth': 10, 'model__n_estimators': 500}
Total Profit: 175395.00
Average Profit per Sale Attempts: 30.67
Best Threshold: 0.15
Confusion Matrix:
[[25791  3447]
 [ 1441  2271]]
F1 Score: 0.48
Average Precision Score: 0.46
CV Average Precision Variation (std): 0.015
-----
Model: Dummy_RandomOverSampler
Best Parameters: {}
Total Profit: -67370.00
Average Profit per Sale Attempts: -2.04
Best Threshold: 0.05
Confusion Matrix:
[[    0  29238]
 [    0   3712]]
F1 Score: 0.20
Average Precision Score: 0.11
CV Average Precision Variation (std): 0.000
-----
```

```
Model: XGBoost_RandomOverSampler
Best Parameters: {'model__learning_rate': 0.1, 'model__max_depth': 6, 'model__n_estimators': 50, 'model__reg_alpha': 10}
Total Profit: 175475.00
Average Profit per Sale Attempts: 31.39
Best Threshold: 0.55
Confusion Matrix:
[[25903  3335]
 [ 1457  2255]]
F1 Score: 0.48
Average Precision Score: 0.46
CV Average Precision Variation (std): 0.013
-----
```

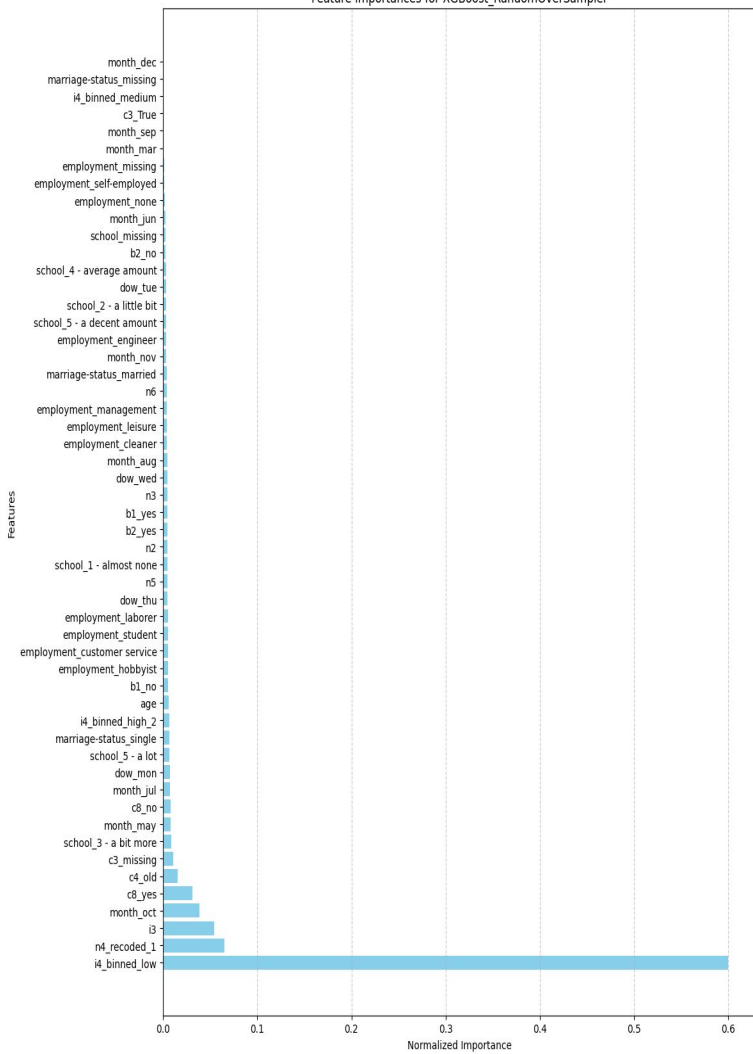
```
Model: Random Forest_NoSampler
Best Parameters: {'model__class_weight': None, 'model__max_depth': 10, 'model__n_estimators': 500}
Total Profit: 175395.00
Average Profit per Sale Attempts: 30.67
Best Threshold: 0.15
Confusion Matrix:
[[25791  3447]
 [ 1441  2271]]
F1 Score: 0.48
Average Precision Score: 0.46
CV Average Precision Variation (std): 0.015
-----
```

```
Model: Dummy_RandomOverSampler
Best Parameters: {}
Total Profit: -67370.00
Average Profit per Sale Attempts: -2.04
Best Threshold: 0.05
Confusion Matrix:
[[    0 29238]
 [    0  3712]]
F1 Score: 0.20
Average Precision Score: 0.11
CV Average Precision Variation (std): 0.000
-----
```

Feature Importances for Random Forest_NoSampler



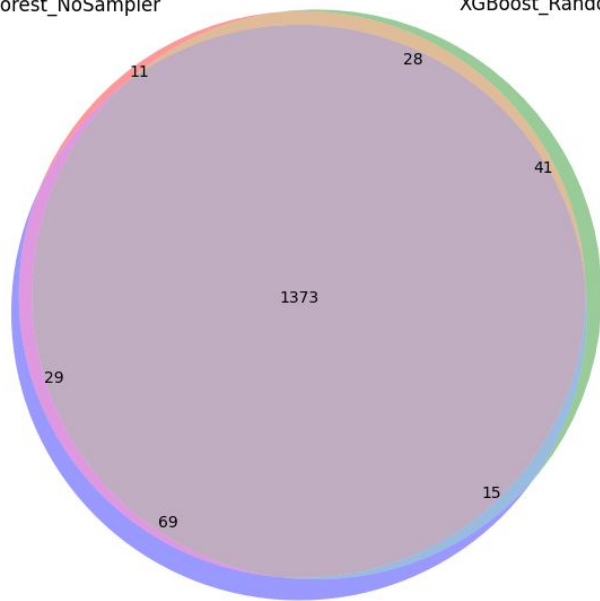
Feature Importances for XGBoost_RandomOverSampler



Overlap of Misclassified Instances for True Class 1

Random Forest_NoSampler

XGBoost_RandomOverSampler

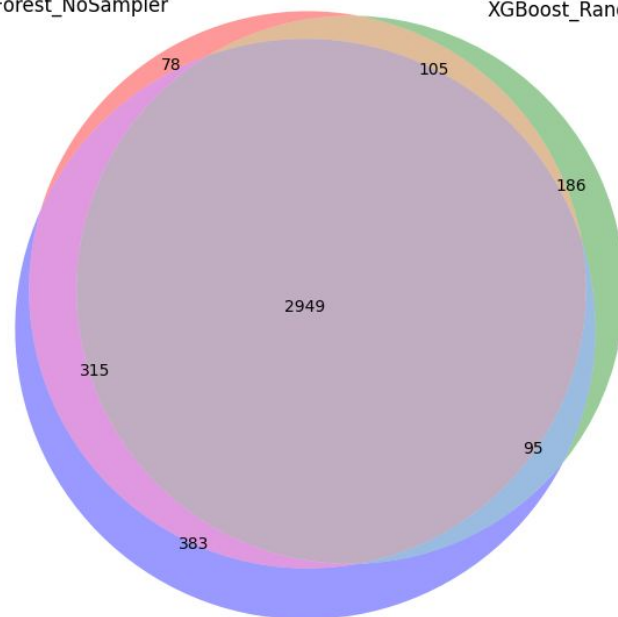


Penalized Logistic Regression_NoSampler

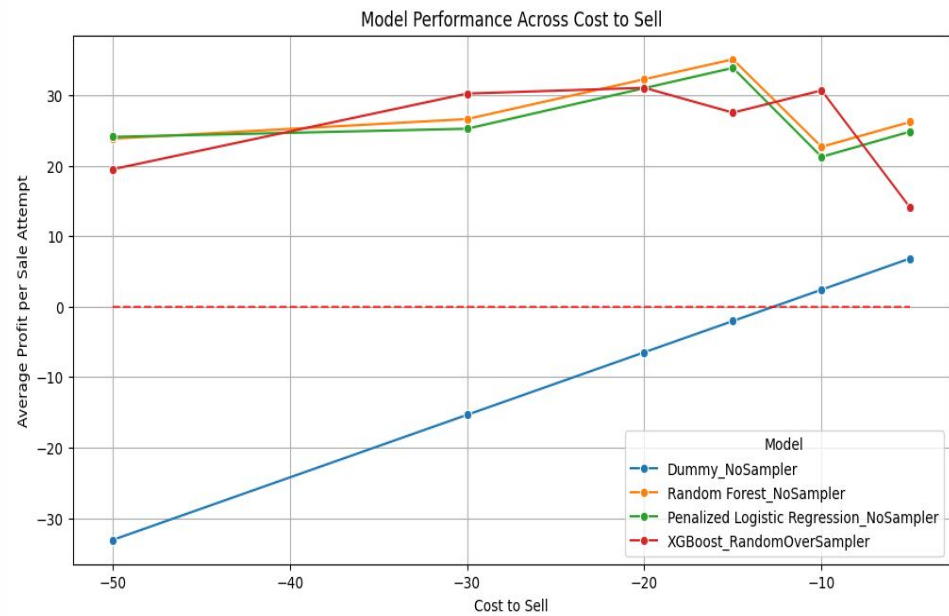
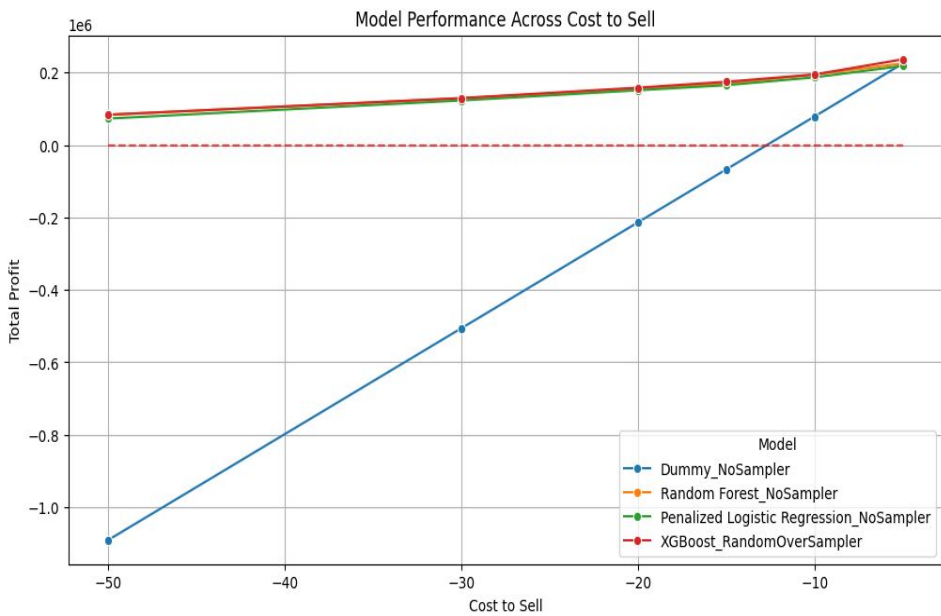
Overlap of Misclassified Instances for True Class 0

Random Forest_NoSampler

XGBoost_RandomOverSampler



Penalized Logistic Regression_NoSampler



```
X_test, y_test = preprocess_tax_df(df_test)
```

```
cost_params = {'tp_cost': 100, 'tn_cost': 0, 'fp_cost': -15, 'fn_cost': 0}
```

```
# Preprocessing pipeline
preprocessor = ColumnTransformer(
    transformers=[
        ("num", StandardScaler(), num_cols_X), # Scale numeric columns
        ("cat", OneHotEncoder(drop="first"), cat_cols_X) # One-hot encode categorical columns
    ]
)
```

```
best_model_name = 'Random Forest_NoSampler'
best_result = results[best_model_name]
best_params = best_result['best_params']
best_threshold = best_result['best_threshold']
```

```
# Extract model hyperparameters
model_params = {key.replace('model_', ''): value for key, value in best_params.items() if key.startswith('model_')}
```

```
# 2. Reconstruct the Pipeline
# Initialize the model with the best hyperparameters
best_model = RandomForestClassifier(random_state=42, **model_params)
```

```
# Build the pipeline steps
pipeline_steps = [
    ('preprocessor', preprocessor)
]
```

```
# Add the model to the pipeline
pipeline_steps.append(('model', best_model))
```

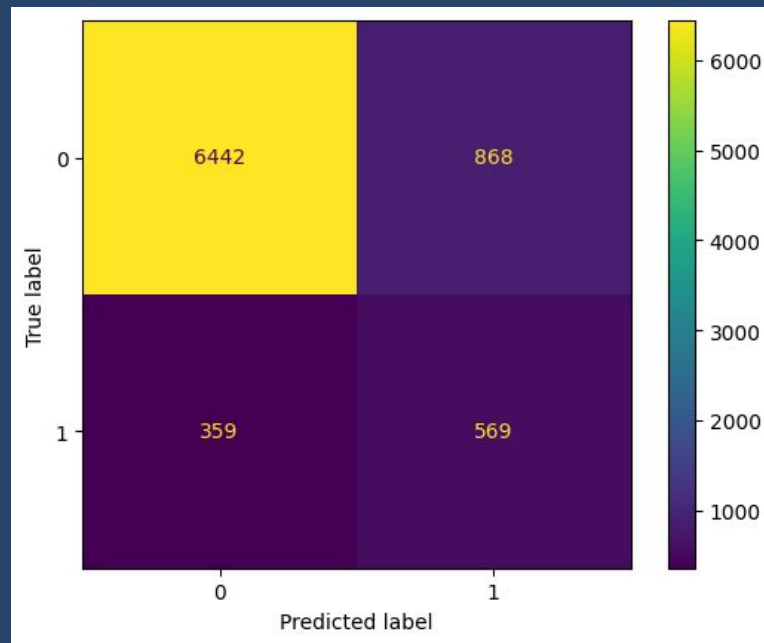
```
# Create the pipeline
pipeline = ImbPipeline(steps=pipeline_steps)
```

```
# 3. Train on the Full Training Set
pipeline.fit(X_train, y_train)
```

```
# 4. Evaluate on the Test Set
# Get predicted probabilities for the positive class
probs_test = pipeline.predict_proba(X_test)[:, 1]
```

```
# Apply the optimal threshold to get binary predictions
y_pred_test = (probs_test >= best_threshold).astype(int)
```

```
# Calculate the confusion matrix
confusion_mat = confusion_matrix(y_test, y_pred_test)
ConfusionMatrixDisplay(confusion_mat).plot()
```



Random Forest:

Profit: \$43,880

Profit per Attempted Sale: \$30.54

Average Precision: 0.45

Current Business Practice:

Profit: -\$16,850

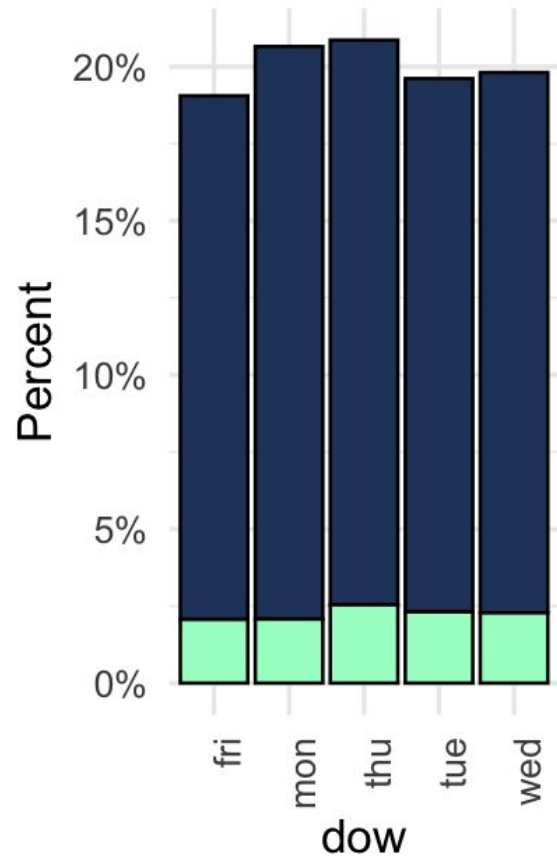
Profit per Attempted Sale: -\$2.05

Next Steps and Places for Improvement

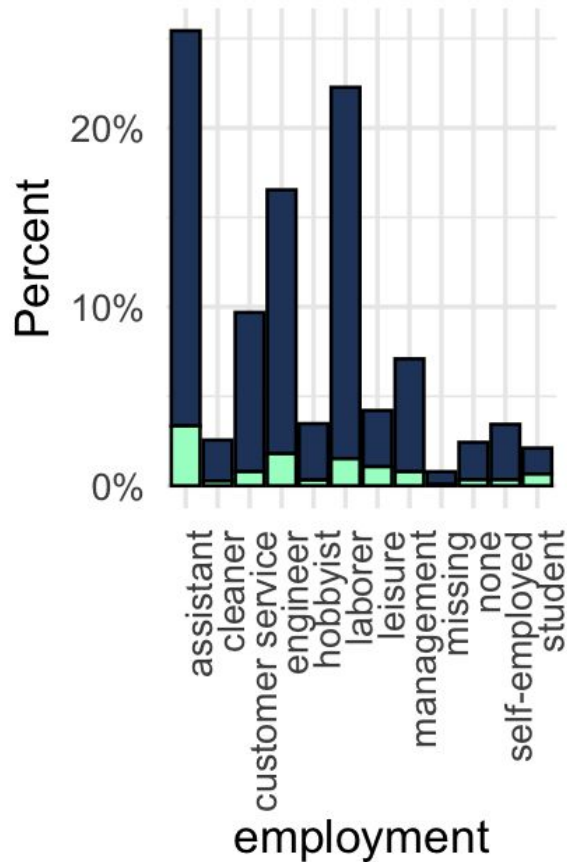
- Deploy the model so the customer can use it to generate likely leads
 - Watch features to ensure data generating process remains the same
 - Monitor model performance
- Look to utilize domain knowledge and other data sources if available in future modeling
 - If we can leverage strong domain knowledge, consider a Bayesian approach which naturally lends itself to decision analysis
- Consider different choices for feature engineering and imputation methods and perform sensitivity analysis on those choices

Questions???

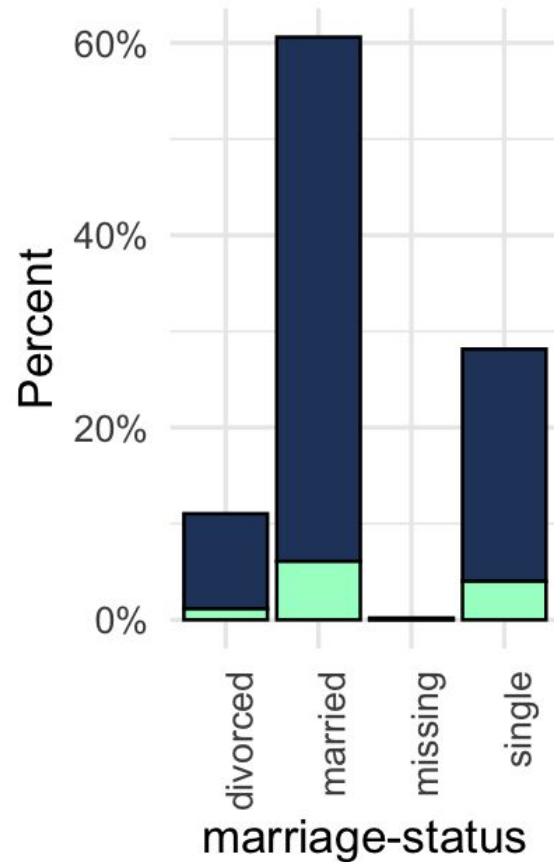
dow



employment



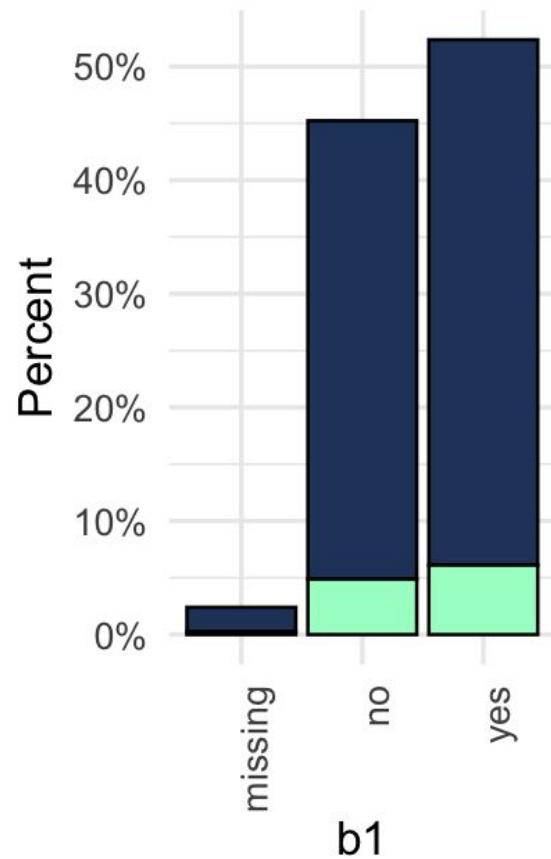
marriage-status



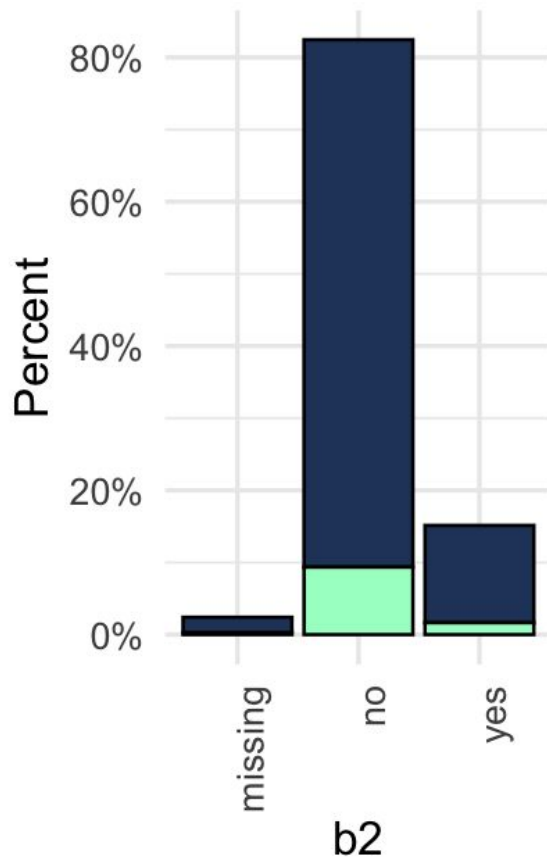
successful_sell

no	yes
----	-----

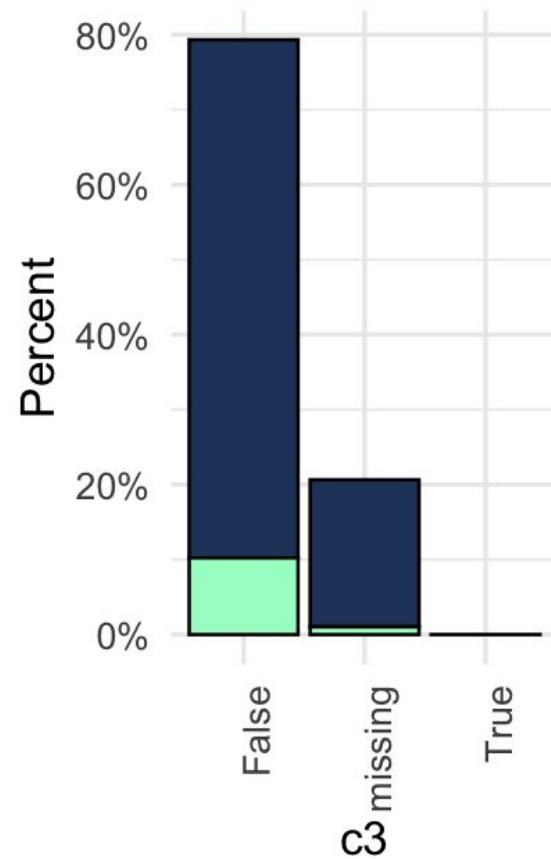
b1



b2



c3



successful_sell



no



yes