

Decision Tree Classifier (Дерево решений)

Дорофеев Демид, Клюкин Павел
М80-307Б-23

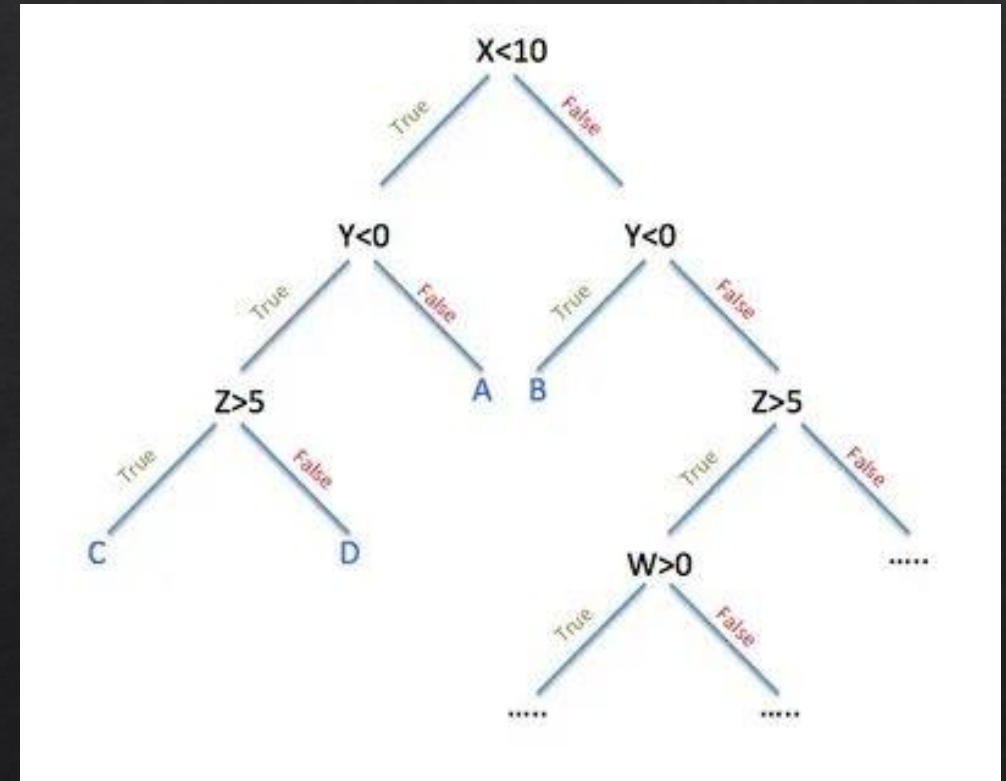
Дерево решений

Преимущества:

- ♦ Простая интерпретация (можно визуализировать дерево).
- ♦ Не требует масштабирования данных.

Недостатки:

- ♦ Склонно к **переобучению**, особенно при большой глубине дерева.
- ♦ Может быть нестабильным (небольшие изменения данных → другое дерево).

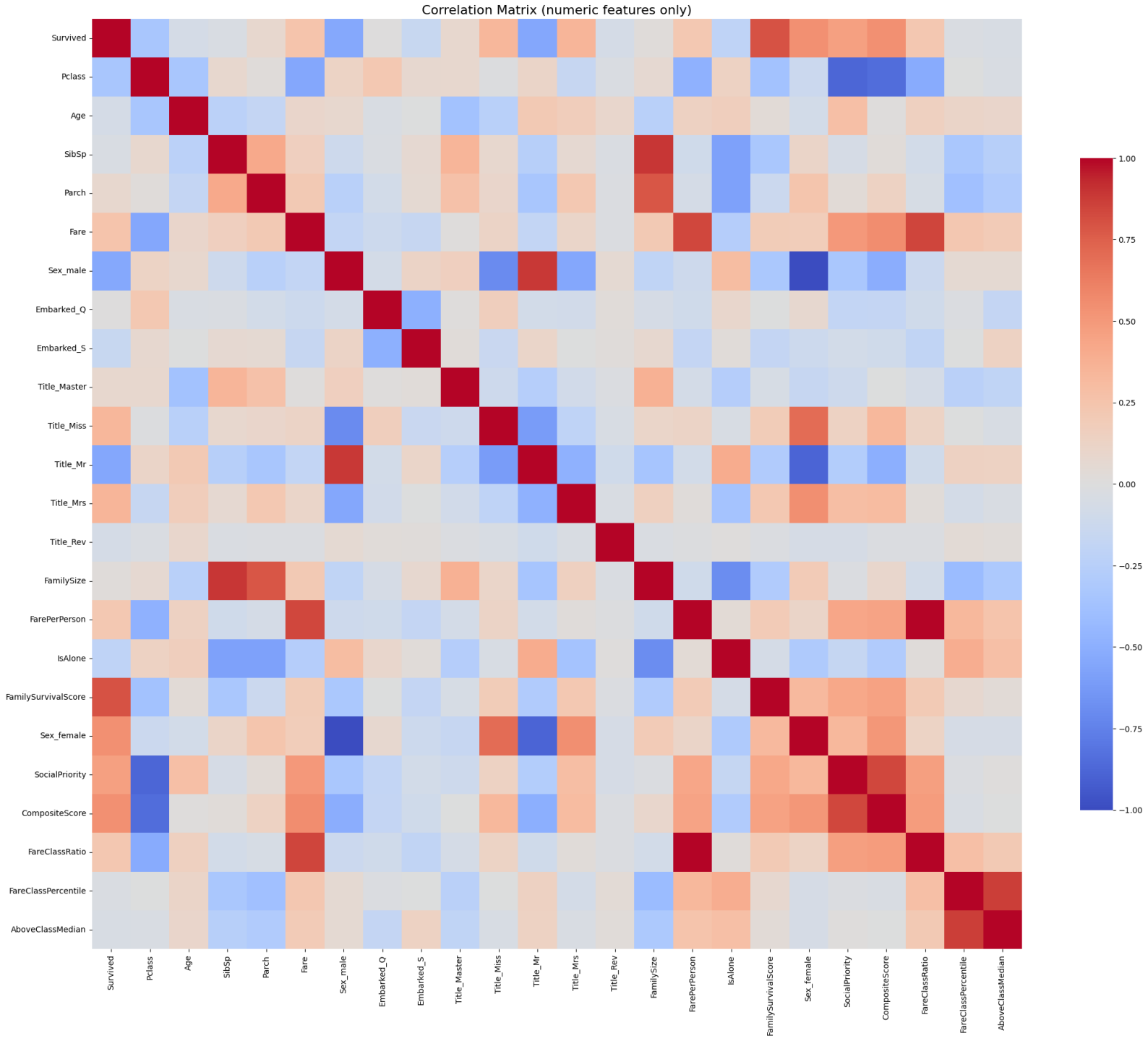


Подготовка

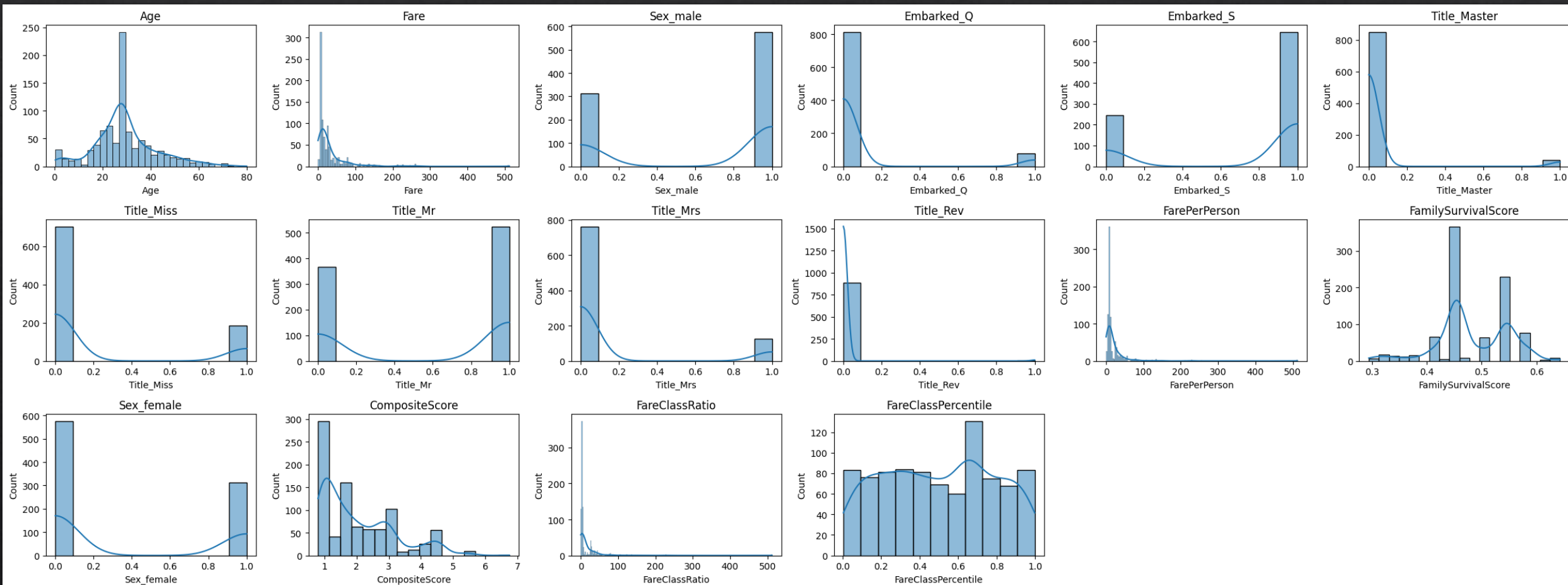
Добавление новых features

PassengerId,	→	Survived, Pclass, Age, SibSp, Parch, Fare, Sex_male,
Survived		Embarked_Q, Embarked_S, Title_Master, Title_Miss, Title_Mr,
Pclass,		Title_Mrs, Title_Rev, FamilySize, FarePerPerson, IsAlone,
Name, Sex,		FamilySurvivalScore, Sex_female, SocialPriority,
Age, SibSp,		CompositeScore, FareClassRatio, FareClassPercentile,
Parch,		AboveClassMedian, Archetype_AveragePassenger,
Ticket,		Archetype_Businessman, Archetype_PoorChild,
Fare, Cabin,		Archetype_RichChild, Archetype_WealthyMother,
Embarked		Archetype_WealthySingleWoman, Archetype_WorkingClassMan

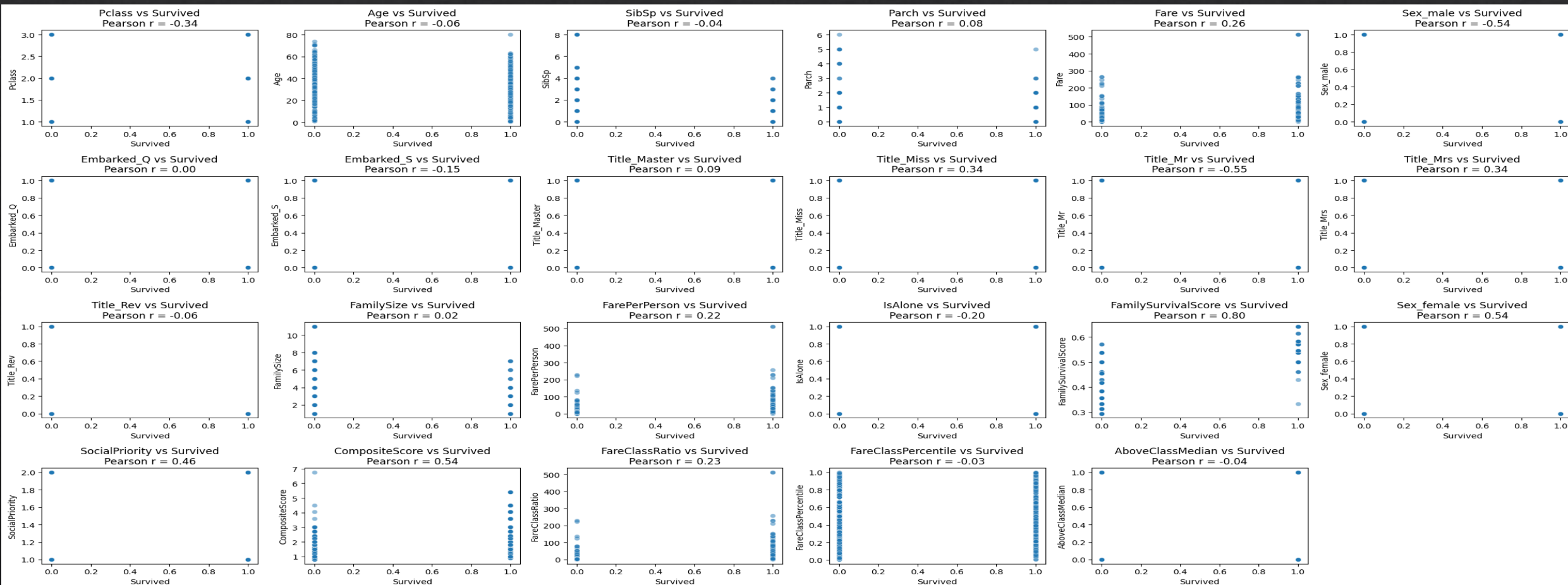
Корреляционная матрица признаков



Распределение данных



Соотношение с целевой переменной + коэффициент пирсона

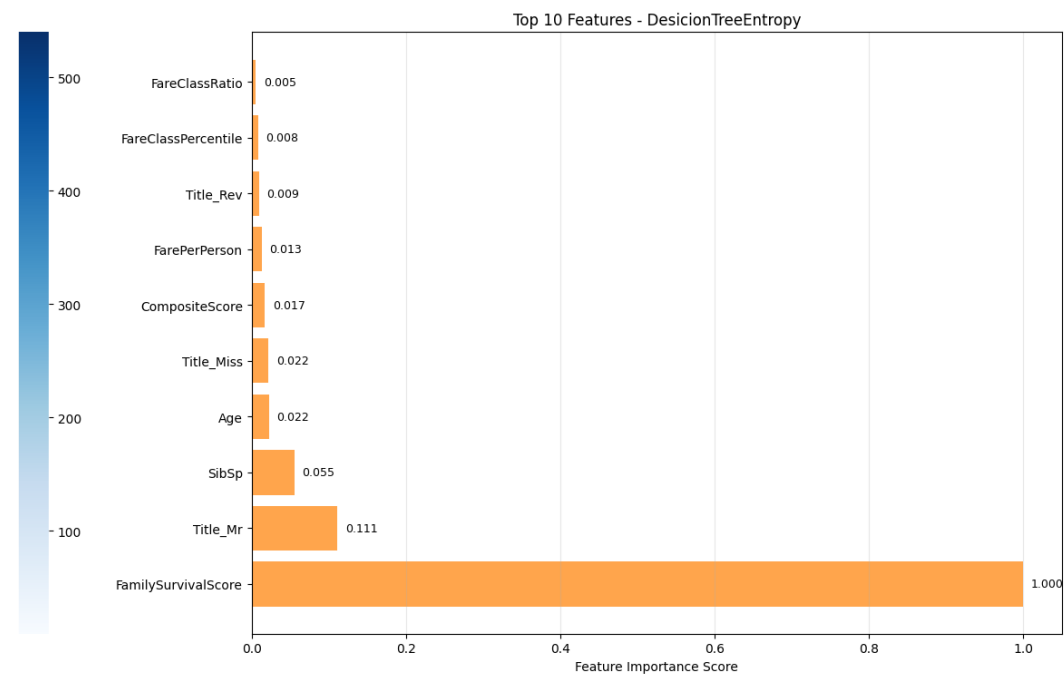
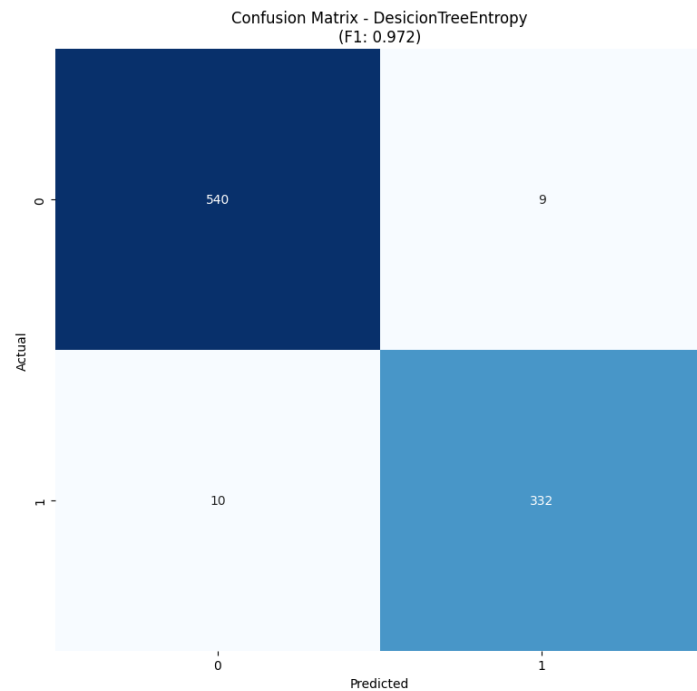
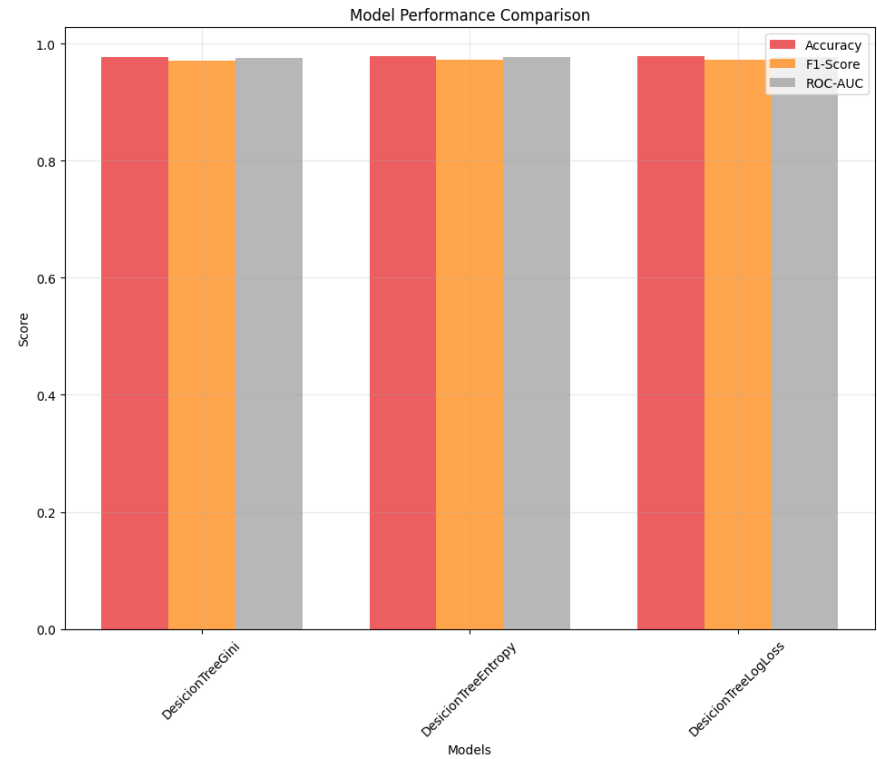
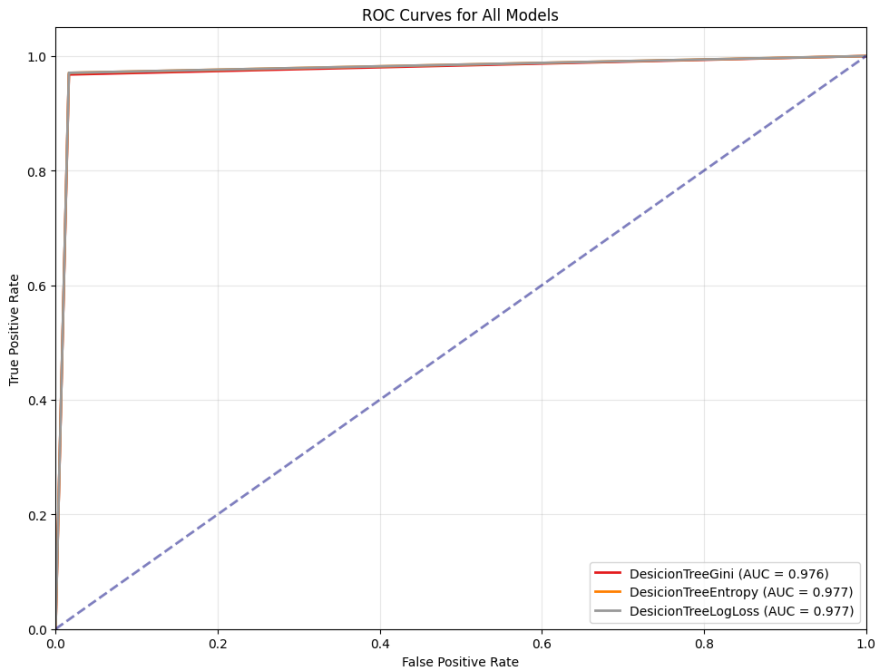


Критерий качества разбиения.

gini — индекс Джини

entropy —энтропия

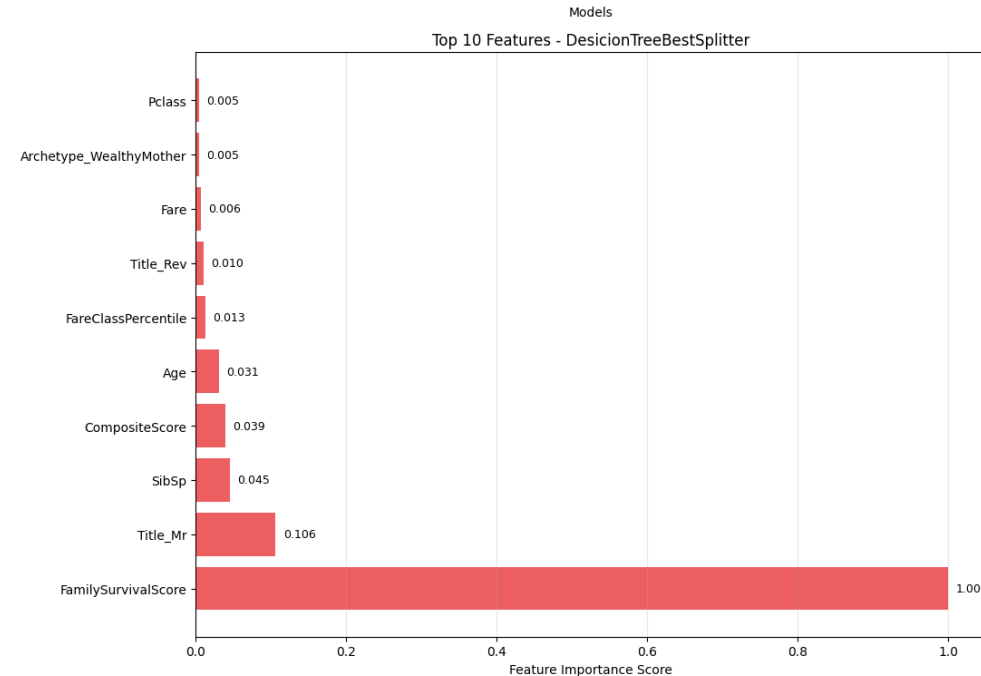
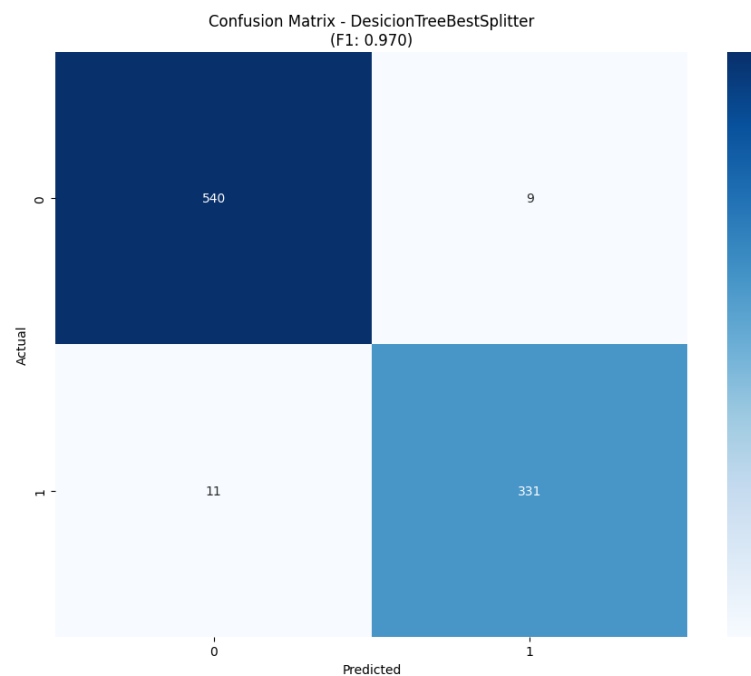
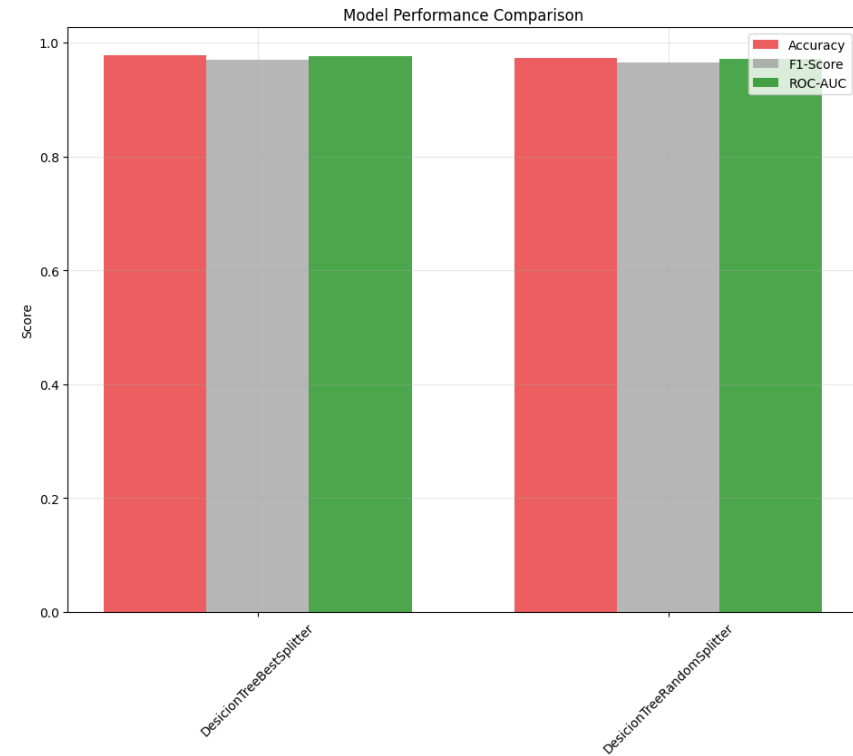
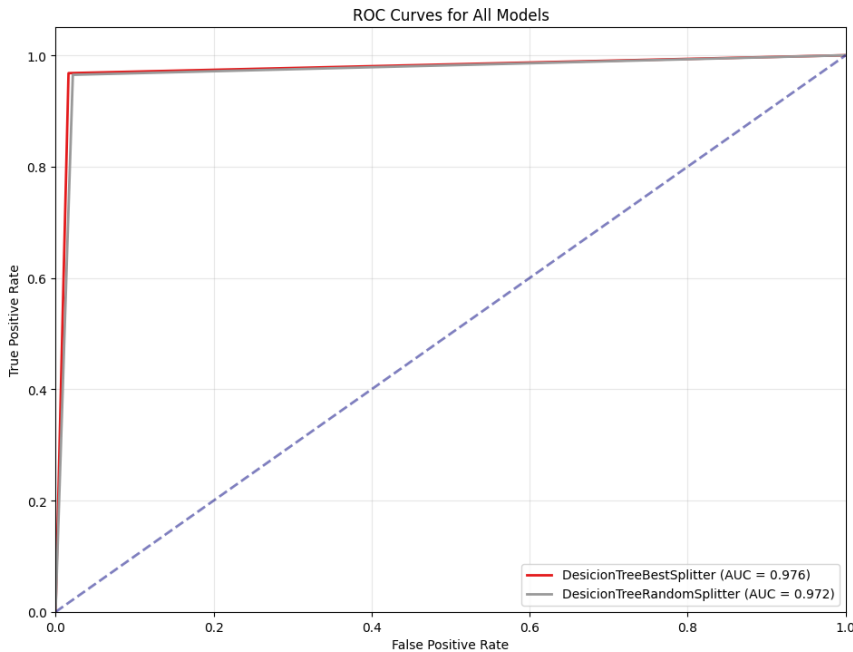
log_loss — основан на
логарифмической потере.



Стратегия выбора признака для разбиения.

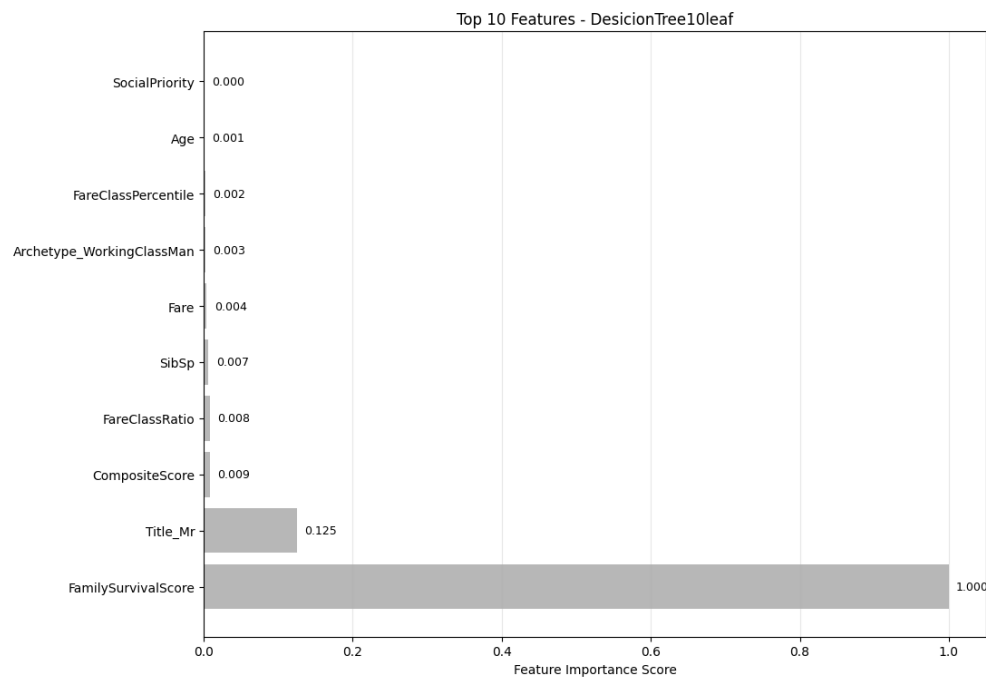
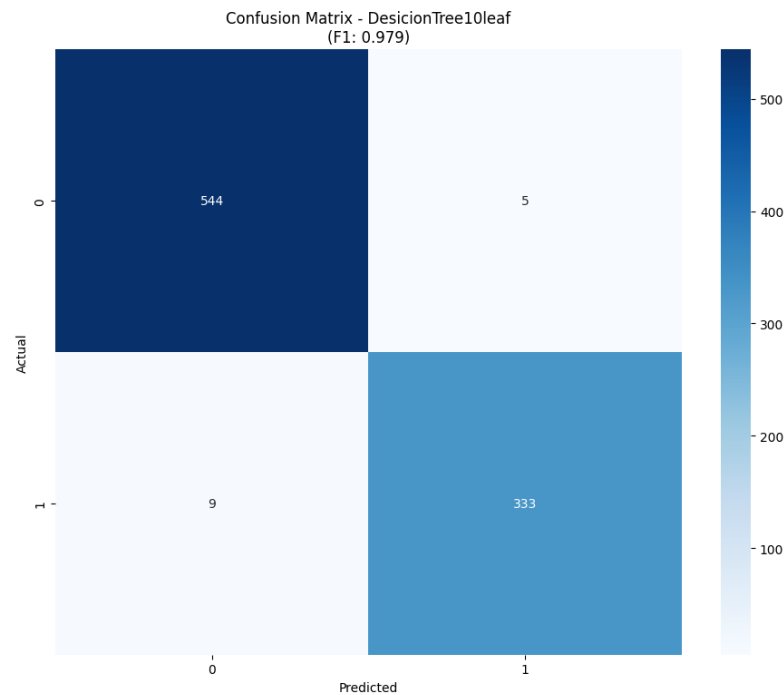
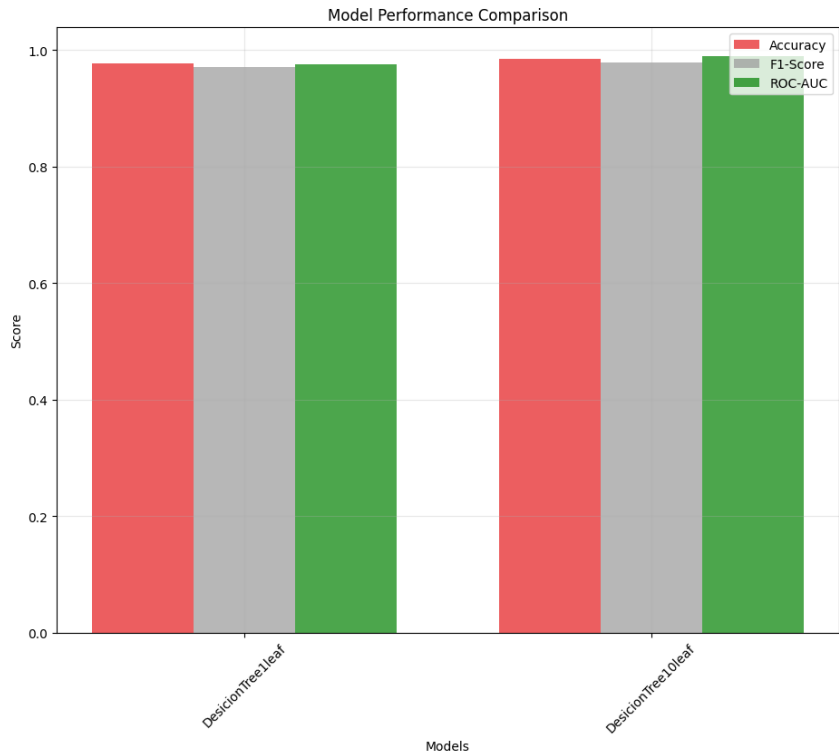
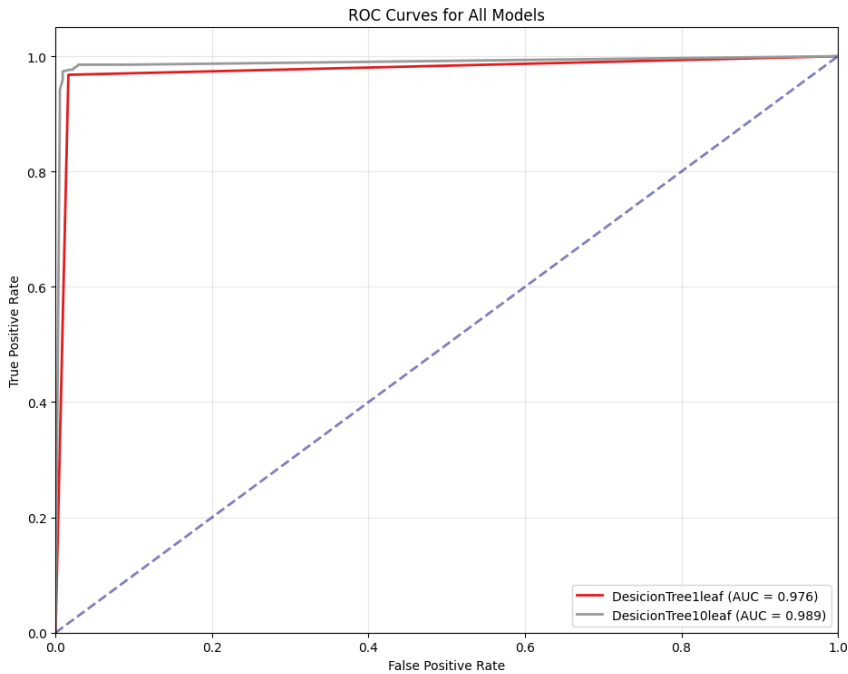
best — выбирает наилучшее разбиение

random' — выбирает случайный признак (для ускорения и разнообразия).



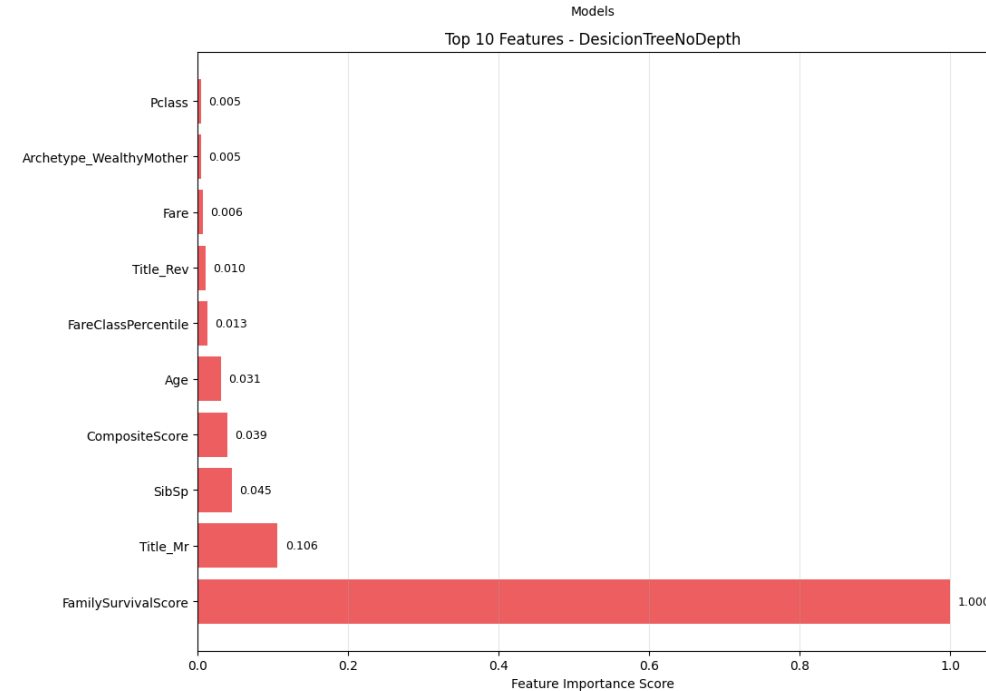
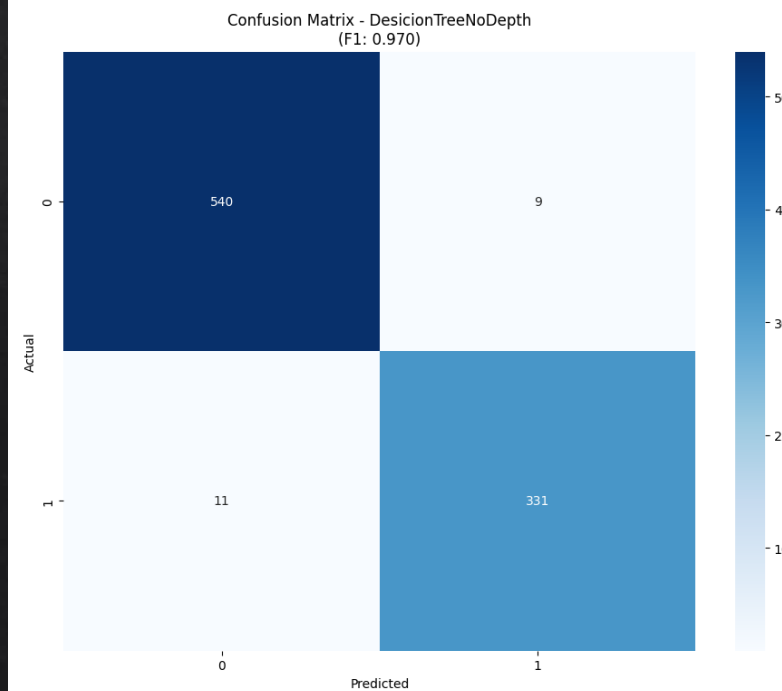
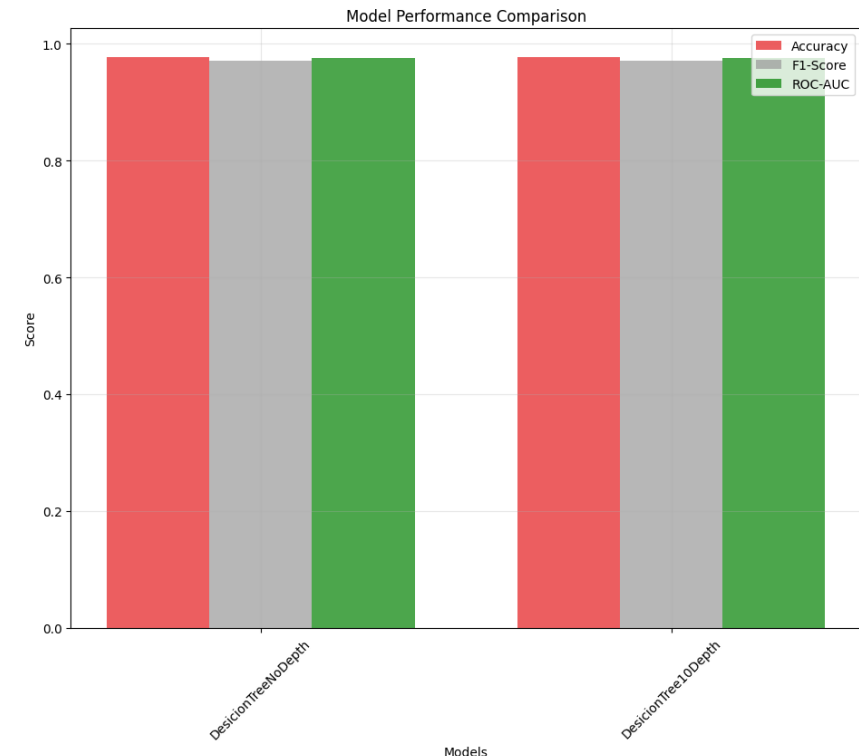
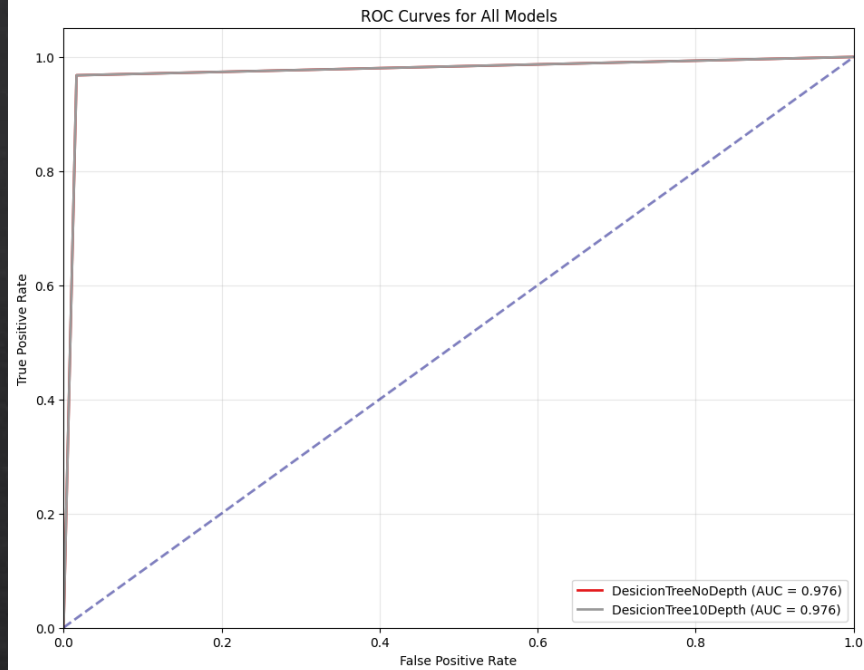
Минимальное количество образцов в листе.

Увеличение значения делает дерево более “гладким”.



Максимальная глубина дерева.

None — без ограничения
(может привести к
переобучению).



- ◇ criterion=gini
- ◇ splitter=best
- ◇ max_depth=None
- ◇ min_samples_split=2
- ◇ min_samples_leaf=15
- ◇ random_state=RANDOM_STATE
- ◇ max_features=None
- ◇ class_weight=None

