

Месторождения ошибок

Анализ пространства ошибок моделей машинного обучения

Дмитрий Колодезев
ООО Промсофт, Новосибирск
DataConf Barnaul 18.06.2021

План

- Модели и ошибки
- Пространство ошибок
- История вопроса
- Месторождения ошибок
- Источники ошибок
- Инструменты
- Что делать

Ошибки

- Когда модель отдает не то, что мы хотим
- Оценка «ошибочности» модели
 - Как оцениваем
 - Метрики регрессии MAE MSE ...
 - Метрики классификации Accuracy, F-мера, ROC-AU ...
 - Точечные оценки
 - На чем оцениваем
 - Метрики на тесте
 - Метрики на OOF
 - На всем датасете

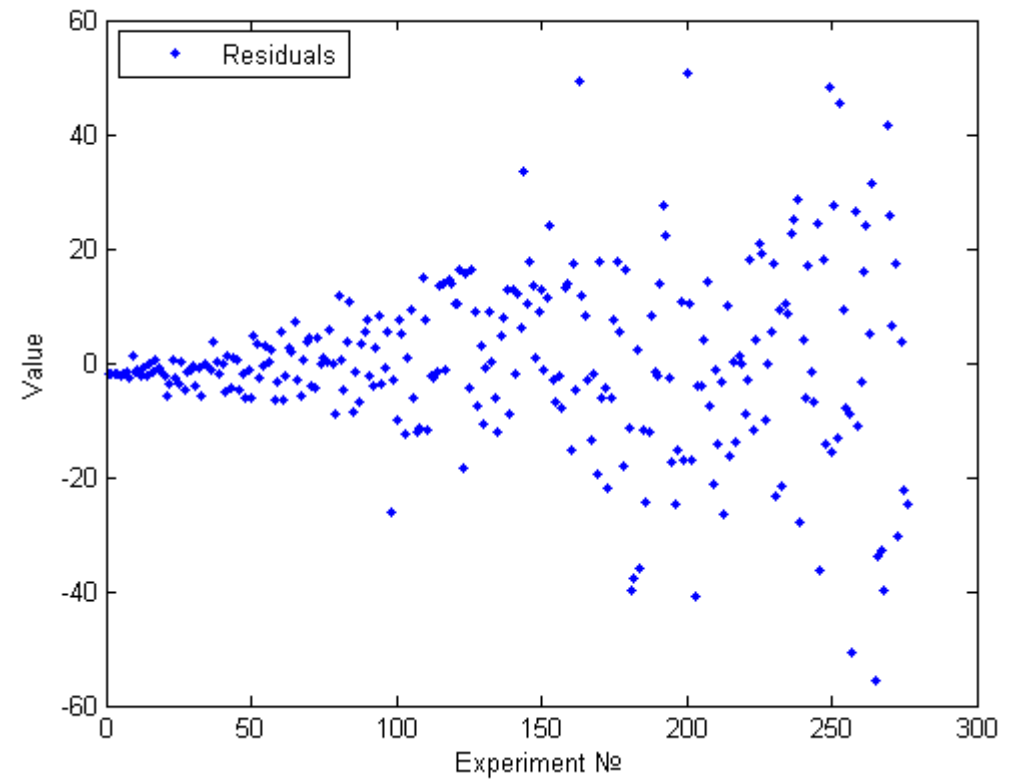
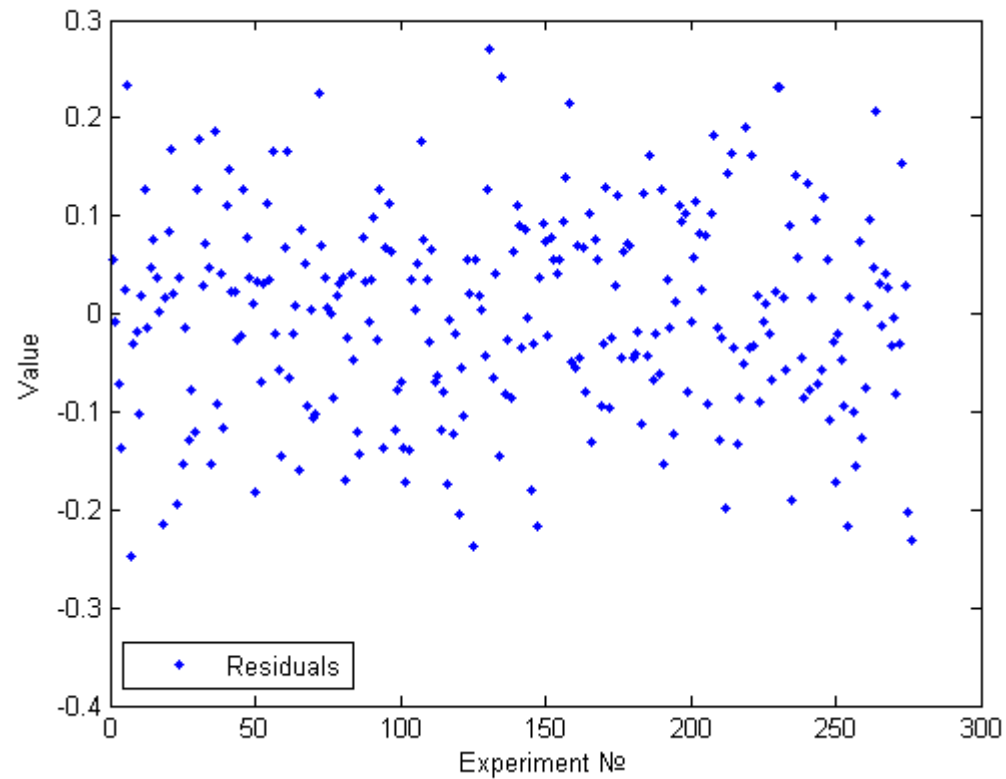
Пространство ошибок

- Берем пространство признаков
- Добавляем точечную оценку ошибки
- Получаем пространство ошибок
- Например, **данные медосмотра**

Age (days)
Height (cm)
Weight (kg)
Gender
Systolic blood pressure
Diastolic blood pressure
Cholesterol
Glucose
Smoking
Alcohol intake
Physical activity
Cardiovascular disease

	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio	bmi
id													
0	18393	2	168	62.0	110	80	1	1	0	0	1	0	21.967120
1	20228	1	156	85.0	140	90	3	1	0	0	1	1	34.927679
2	18857	1	165	64.0	130	70	3	1	0	0	0	1	23.507805
3	17623	2	169	82.0	150	100	1	1	0	0	1	1	28.710479
4	17474	1	156	56.0	100	60	1	1	0	0	0	0	23.011177

Анализ регрессионных остатков



ПОЧИСТИМ, ПОДЕЛИМ, ОБУЧИМ

```
df = pd.read_csv(DATA_FILE, sep=';', index_col='id')
df['bmi'] = df.weight/(df.height/100)**2

idx = df.query('(ap_hi < 50) | (ap_lo < 50) | (ap_lo > 240) | (ap_hi > 300) | (weight < 40) |
              |(height < 120) | (bmi > 100) | (ap_lo > (ap_hi-10))').index
df_clean = df.drop(idx, axis=0)

X = df_clean.drop('cardio', axis=1)
y = df_clean.cardio
features = list(X.columns)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=SEED)

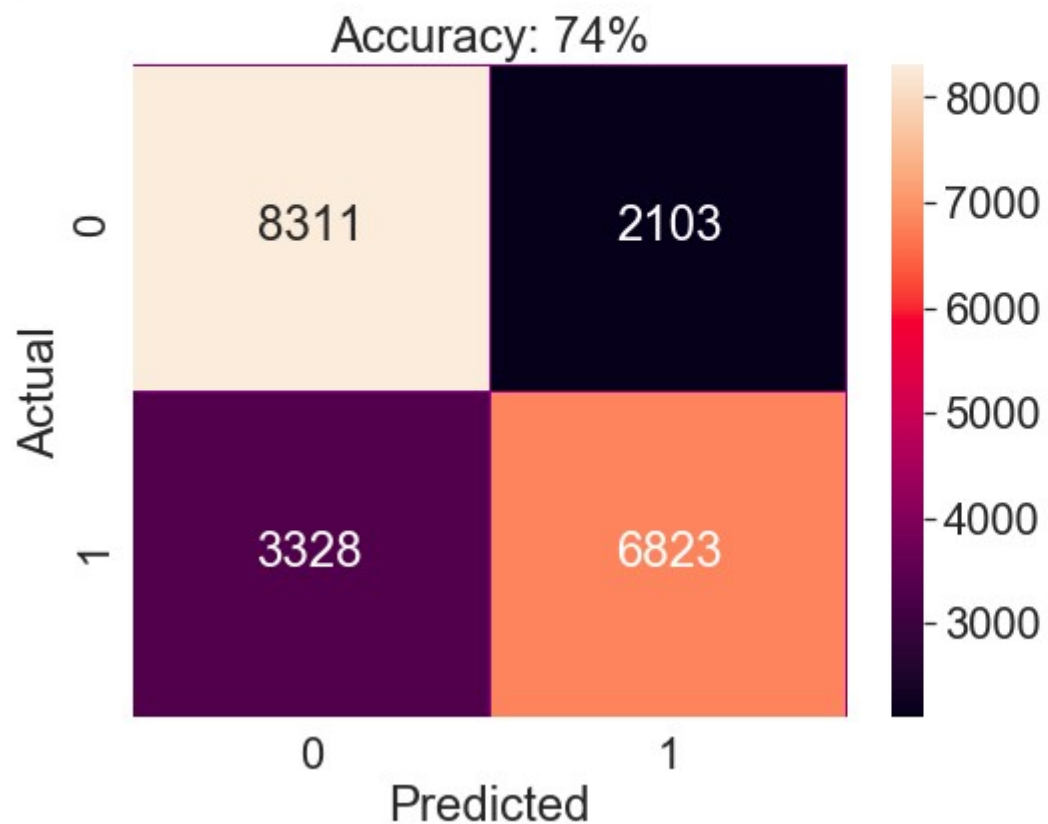
model = RandomForestClassifier(**model_params)
model.fit(X_train.values, y_train.values)

pred = model.predict(X_test.values)
pred_proba = model.predict_proba(X_test.values)[: , 1]
```

Датасет <https://www.kaggle.com/sulianova/cardiovascular-disease-dataset>

ОЦЕНИМ

```
df_cm = pd.DataFrame(confusion_matrix(y_test, pred), columns=np.unique(y_test), index = np.unique(y_test))
df_cm.index.name = 'Actual'
df_cm.columns.name = 'Predicted'
plt.figure(figsize=(8, 6))
score = accuracy_score(y_test.values, pred)
plt.title('Accuracy: {:.0%}'.format(score))
sns.heatmap(df_cm, fmt=".0f", annot=True,linewidths=0.2, linecolor="purple", );
```



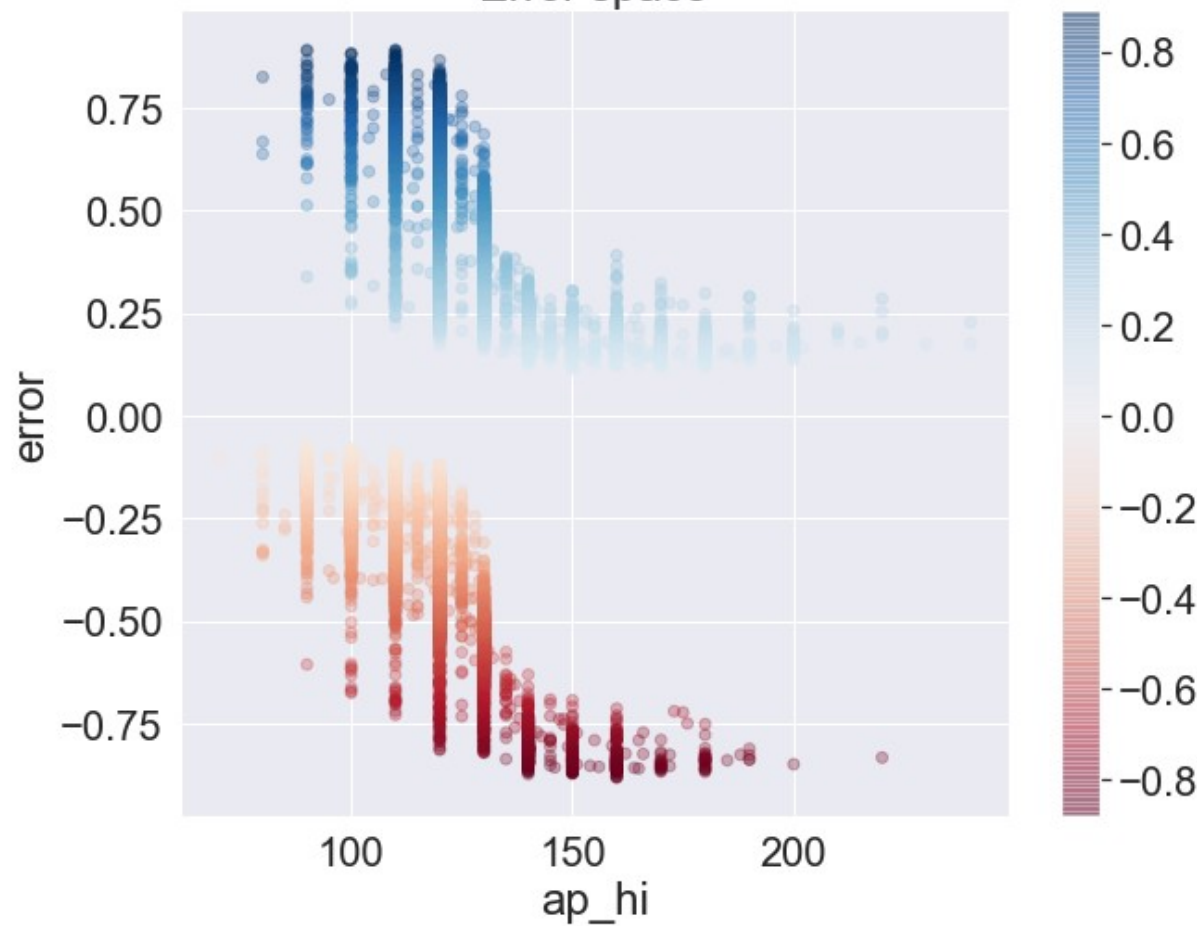
Соберем ошибки

```
df_error = X_test.copy()
df_error['error'] = (y_test - pred_proba)
df_error['abs_error'] = (y_test - pred_proba).abs()
df_error.head()
```

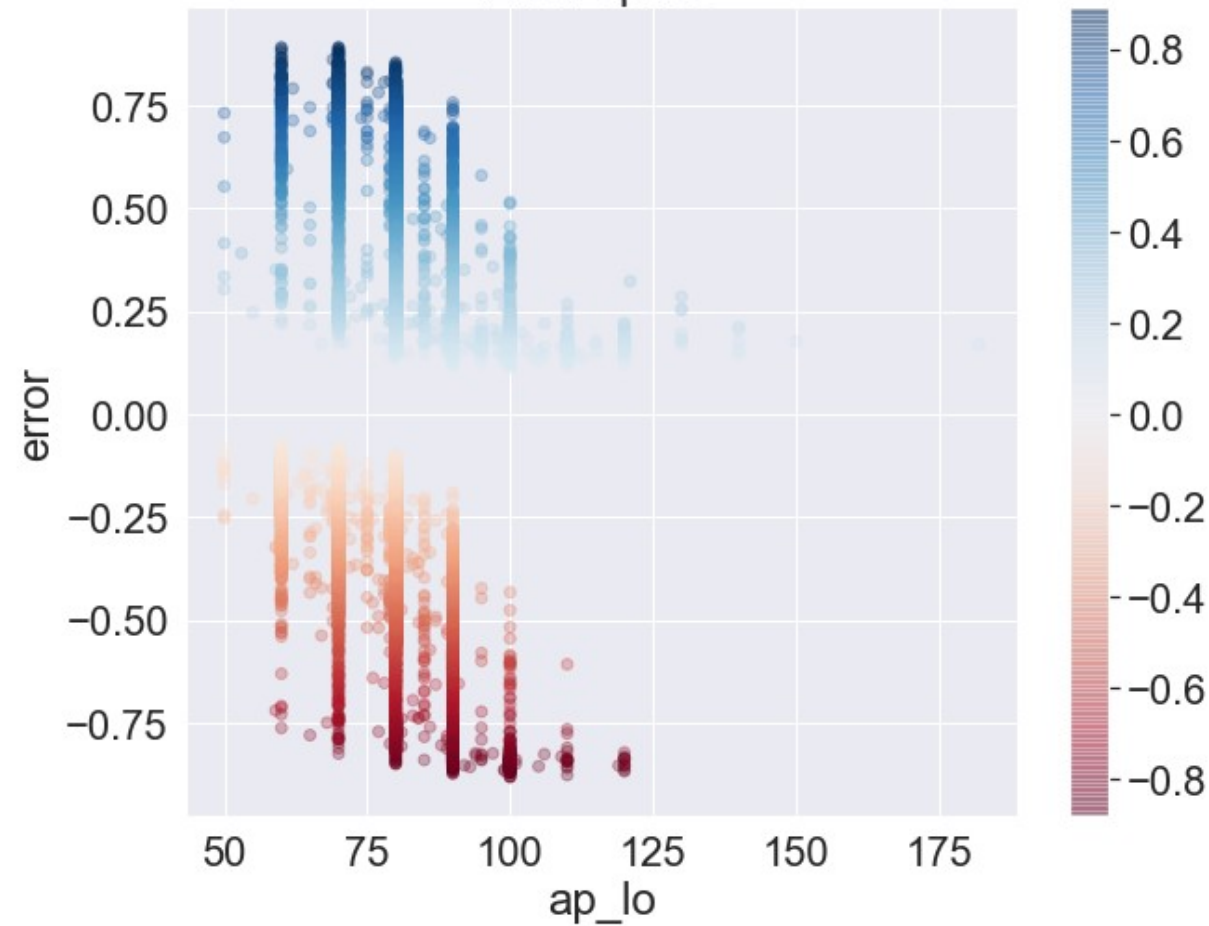
	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	bmi	error	abs_error
id														
10289	23636	1	160	63.0	120	80	1	1	0	0	0	24.609375	-0.528284	0.528284
17843	19676	1	169	75.0	120	80	1	1	0	0	1	26.259585	-0.276148	0.276148
9901	21032	1	168	67.0	120	70	1	1	0	0	1	23.738662	0.629541	0.629541
93721	20214	2	165	90.0	130	80	1	1	1	0	0	33.057851	-0.558215	0.558215
37683	20438	1	160	64.0	100	60	1	1	0	0	1	25.000000	-0.281052	0.281052

Ошибка зависит

Error space



Error space



Microsoft raiWidget

- <https://github.com/microsoft/responsible-ai-widgets>
- <https://github.com/interpretml/interpret-community>
- <https://github.com/fairlearn/fairlearn>

```
from raiwidgets import ErrorAnalysisDashboard
```

```
ErrorAnalysisDashboard(dataset=X_test.values, true_y=y_test.values, features=features, pred_y=pred);
```

```
ErrorAnalysis started at http://localhost:5001
```

Карта ошибок

X-Axis: Feature 1

Y-Axis: Feature 2

ap_hi

ap_lo

Clear all

Select all

ap_lo

(69.83, 91.25]	15%	12%	0%	0%	0%	0%	0%	0%
(91.25, 112.5]	24%	23%	32%	0%	0%	0%	0%	0%
(112.5, 133.75]	28%	32%	34%	36%	0%	0%	0%	0%
(133.75, 155.0]	21%	20%	16%	16%	0%	0%	0%	0%
(155.0, 176.25]	36%	15%	13%	16%	21%	0%	0%	0%
(176.25, 197.5]	0%	13%	10%	12%	26%	0%	0%	0%
(197.5, 218.75]	0%	0%	0%	4%	0%	0%	0%	0%
(218.75, 240.0]	0%	0%	0%	0%	20%	0%	0%	0%
	(49.868, 66.5]	(66.5, 83.0]	(83.0, 99.5]	(99.5, 116.0]	(116.0, 132.5]	(132.5, 149.0]	(149.0, 165.5]	(165.5, 182.0]

ap_hi

Дерево ошибок

Error explorer: Tree map ▾

Fullscreen

Feature list

Cohort settings ▾

Cohort info

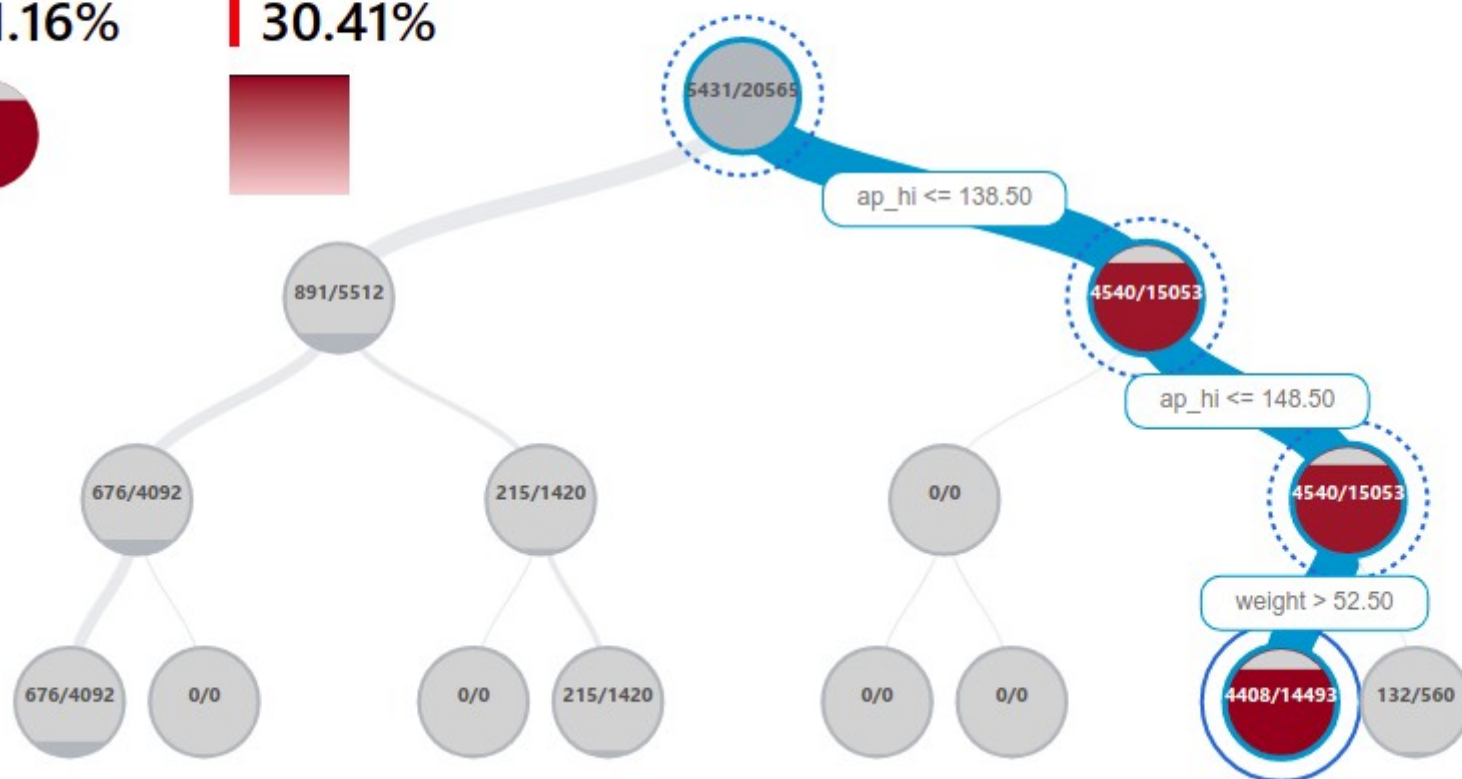
Explanation

ⓘ The tree visualization uses the mutual information between each feature and the error to best separate error instances from success instances hierarchically in the data. [See more](#)

Cohort: All data

Error coverage ⓘ
81.16%

Error rate ⓘ
30.41%



Error coverage ⓘ

81.16%



Error rate ⓘ

30.41%



ap_hi <= 138.50

Save as a new cohort



Save the current cohort to the cohort list. You can revisit the saved cohort via the cohort list.

Cohort name

Cohort info

Error coverage	Error rate	Correct/Total	Incorrect/Total
81.16%	30.41%	10085 / 14493	4408 / 14493

Base cohort and filters

Base cohort All data

Error explorer Tree map

Filters weight > 52.50, ap_hi <= 148.50, ap_hi <= 138.50

Save

Cancel

Работа с когортами

Error explorer

Error explorer: Tree map ▾

Fullscreen

Feature list

Cohort settings ▾

Cohort info

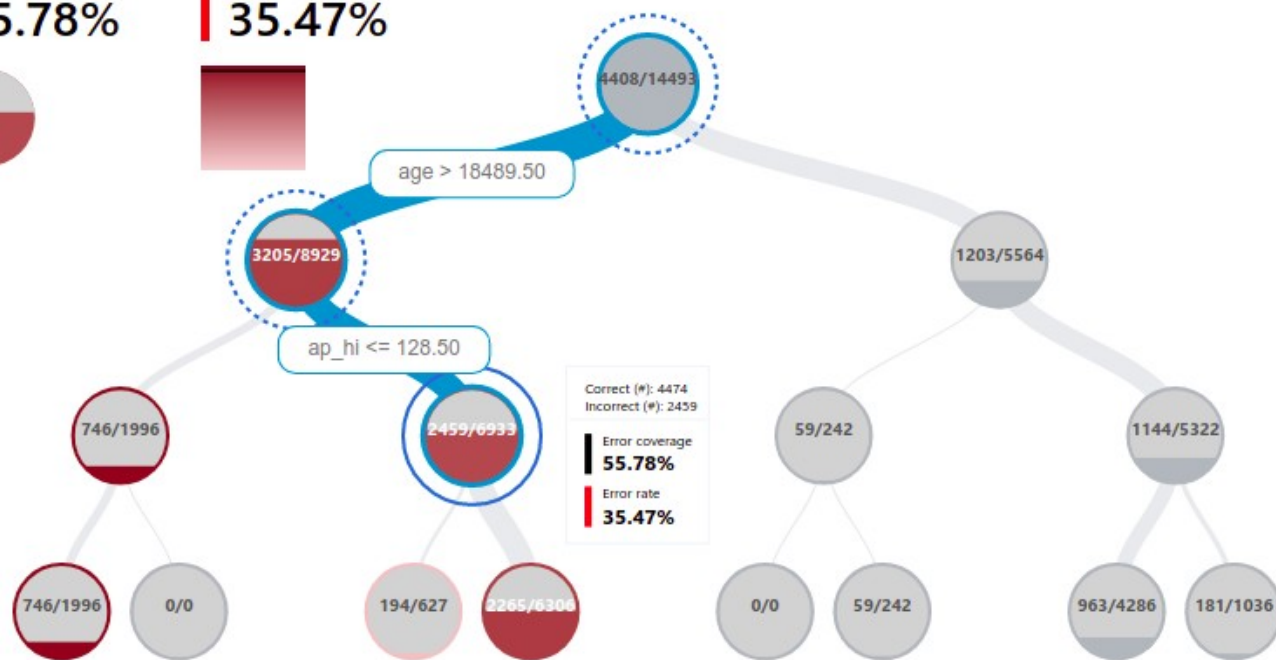
Explanation

The tree visualization uses the mutual information between each feature and the error to best separate error instances from success instances hierarchically in the data. [See more](#)

Cohort: Error pocket

Error coverage
55.78%

Error rate
35.47%



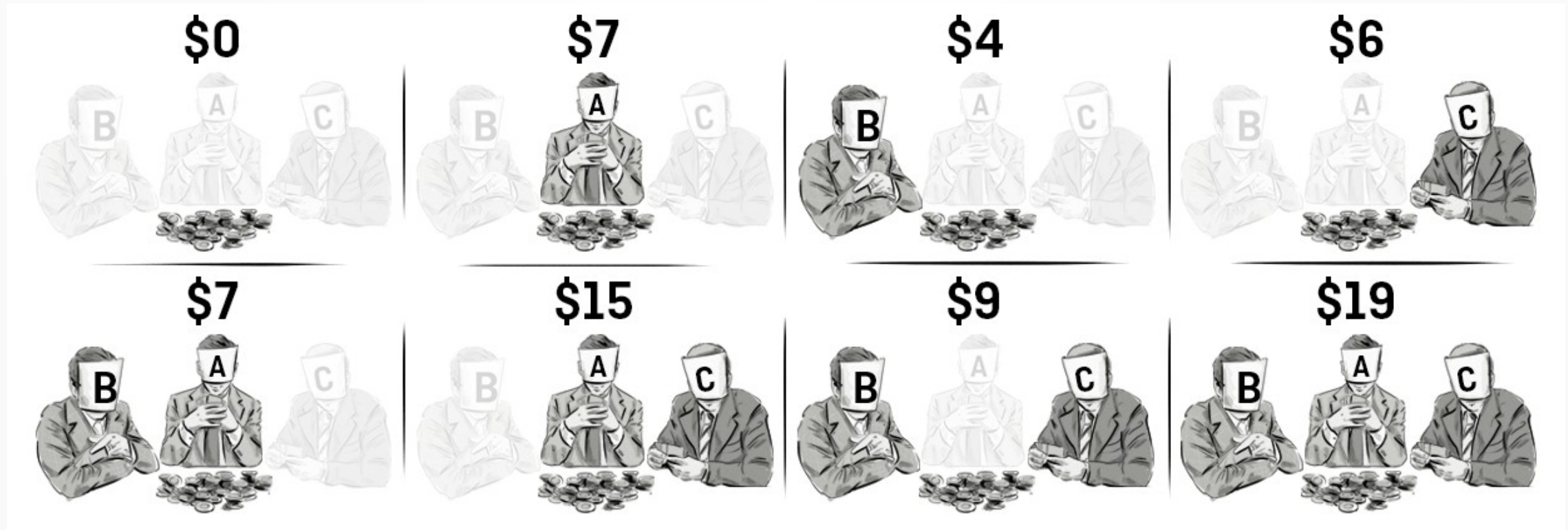
ИСТОЧНИКИ ОШИБОК

- Исходные данные
- Генерация признаков
- Генерация признаков внутри модели
- Особенности моделей
- Параметры моделей

Что делать

- Идентифицировать и рассмотреть месторождения ошибок
- Сравнить месторождения для разных типов моделей
- Добавить новые признаки
- Не верить ответу модели в этих случаях
- Стекать!

Shapley Values



<https://clearcode.cc/blog/game-theory-attribution/>

А еще

- Обычно мы с помощью Shapley Values оцениваем вклад каждого признака в предсказание, для объяснения решений модели
- Можно строить Shapley Values-объяснения для ошибок. Например, [Explaining the Loss of a Tree Model](#) и [Use SHAP loss values to debug your model](#)

Вопросы

Слайды тут



dkolodezev



promsoft



dkolodezev



d_key



dmitry_kolodezev

<https://kolodezev.ru/download/slides-dataconf-2021.pdf>