

# ЛАБОРАТОРНАЯ РАБОТА #2

Линейная классификация

Ивченко Матвей Сергеевич  
Пермяков Герман Алексеевич  
Крушинин Никита Игоревич  
М80-307Б-23

# Цели

- Построить, оптимизировать и оценить классификаторы.
- Понять интерпретацию признаков (веса, permutation importance).
- Освоить расширение и отбор признаков (feature tuning).
- Изучить влияние регуляризации, калибровки и гиперпараметров

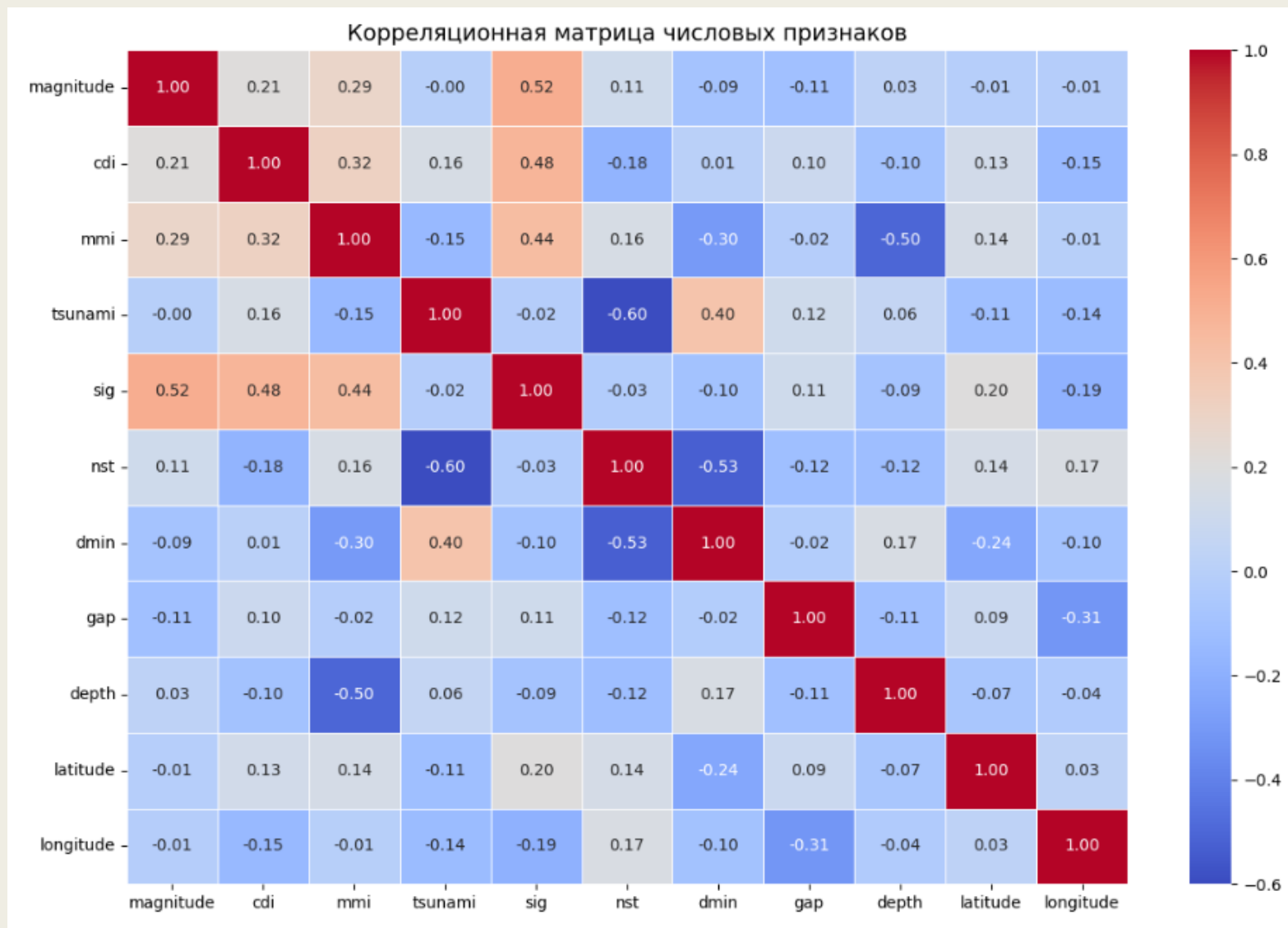
# УСЛОВИЯ

- Numpy
- Pandas
- scikit-learn
- XGBoost

# EDA

- В датасете в основном нет пропусков, кроме continent, country и alert (таргет)
- Есть категориальные и вещественные, а также геоданные (в общем, тоже категориальные)
- Распределение классов неравномерное (наибольший класс – 325 примеров, наименьший – 12)
- Вещественные признаки имеют преимущественно нормальное и логнормальное распределение

# EDA

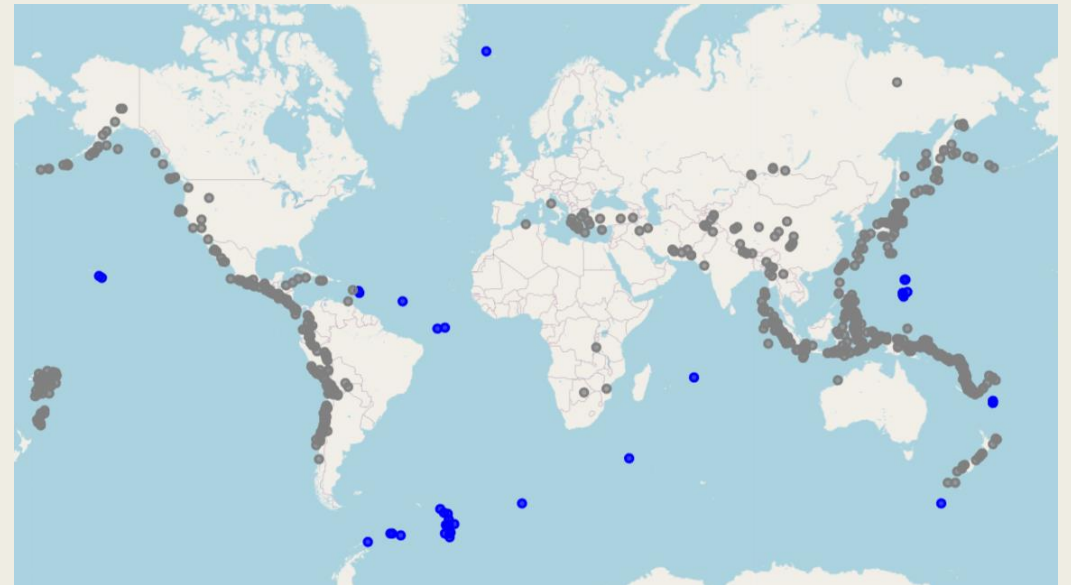


# Feature Engineering

- Признаки longitude и latitude были дискретизированы и закодированы модой в каждом дискретном участке.
- Пропущенные значения в continent и country были заполнены с помощью RadiusNearestNeighbor





До заполнения. Серые точки - пропуски



После заполнения. Серые точки – континент, синие - океан

# Pseudolabeling

- Так как таргет содержит пропуски, попробуем применить трюк под названием псевдолейблинг (хотя он скорее для нейронок, но всё же)
- Суть в том, чтобы использовать модель, обученную на размеченных данных, для разметки тех объектов, на которых она очень уверена, а потом обучить модель на доразмеченных данных, повторяя итерации, пока качество на отложенной выборке растёт.
- Удалось увеличить precision и recall почти на 0.1



```
--- Iteration 1 ---  
Добавлено псевдометок: 97  
Accuracy: 0.840  
Weighted Precision: 0.721  
Weighted Recall: 0.812  
--- Iteration 2 ---  
Добавлено псевдометок: 12  
Accuracy: 0.850  
Weighted Precision: 0.808  
Weighted Recall: 0.851  
--- Iteration 3 ---  
Добавлено псевдометок: 4  
Accuracy: 0.845  
Weighted Precision: 0.795  
Weighted Recall: 0.902  
--- Iteration 4 ---  
Добавлено псевдометок: 2  
Accuracy: 0.855  
Weighted Precision: 0.798  
Weighted Recall: 0.890  
--- Iteration 5 ---  
Добавлено псевдометок: 2  
Accuracy: 0.843  
Weighted Precision: 0.806  
Weighted Recall: 0.896
```

# Обучение моделей

	model	accuracy	f1-score	roc_auc
Best	CatBoost	0.902	0.767	0.971
	CatBoost	0.896	0.758	0.970
	XGBoost	0.906	0.757	0.972
	Logistic Regression	0.896	0.730	0.973
	RandomForest	0.896	0.713	0.972
	Boosting	0.895	0.704	0.961
	DesicionTree	0.866	0.695	0.810
	SVM	0.866	0.511	0.958
	KNN	0.851	0.477	0.887