

# Smote

Synthetic Minority  
Over-sampling  
Technique

# Техники работы с несбалансированными данными

## **I. *Препроцессинг данных:***

- 1) Smote
- 2) Random Under-sampling
- 3) Random Over-sampling
- 4) (на примере классификаторов Random Forest и Logistic Regression)

## **II. *Модификация работы классификаторов:***

Задание априорного распределения классов  
(на примере классификатора Naïve Bayes)



# Алгоритм Smote( $T, N, k$ )

## 1) **Параметры алгоритма:**

- 2) а)  $T$  — число образцов меньшего класса
- 3) б)  $N\%$  от  $T$  — число синтетических образцов, которые хотим получить
- 4) в)  $k$  — число ближайших соседей

5)

## 6) **Суть алгоритма:**

- 7) Для каждого образца  $i$  меньшего класса находим  $k$  ближайших соседей из обоих классов и генерируем  $N/100$  искусственных образцов, повторяя  $N/100$  раз для образца  $i$ :
- 8)
- 9) Среди  $k$  соседей образца  $i$  произвольно выбираем одного —  $np$ . Прибавляем к каждому из атрибутов  $i$  разницу между соответствующими атрибутами  $np$  и  $i$ , умноженную на произвольное число из отрезка  $[0,1]$ .
- 10) Получен новый искусственный меньшего класса.



# Метрики качества

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

В матрице неточностей меньший (более важный) класс назовем положительным, больший — отрицательным.



1) Accuracy — для сбалансированных данных

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

---

**2) Для несбалансированных данных:**

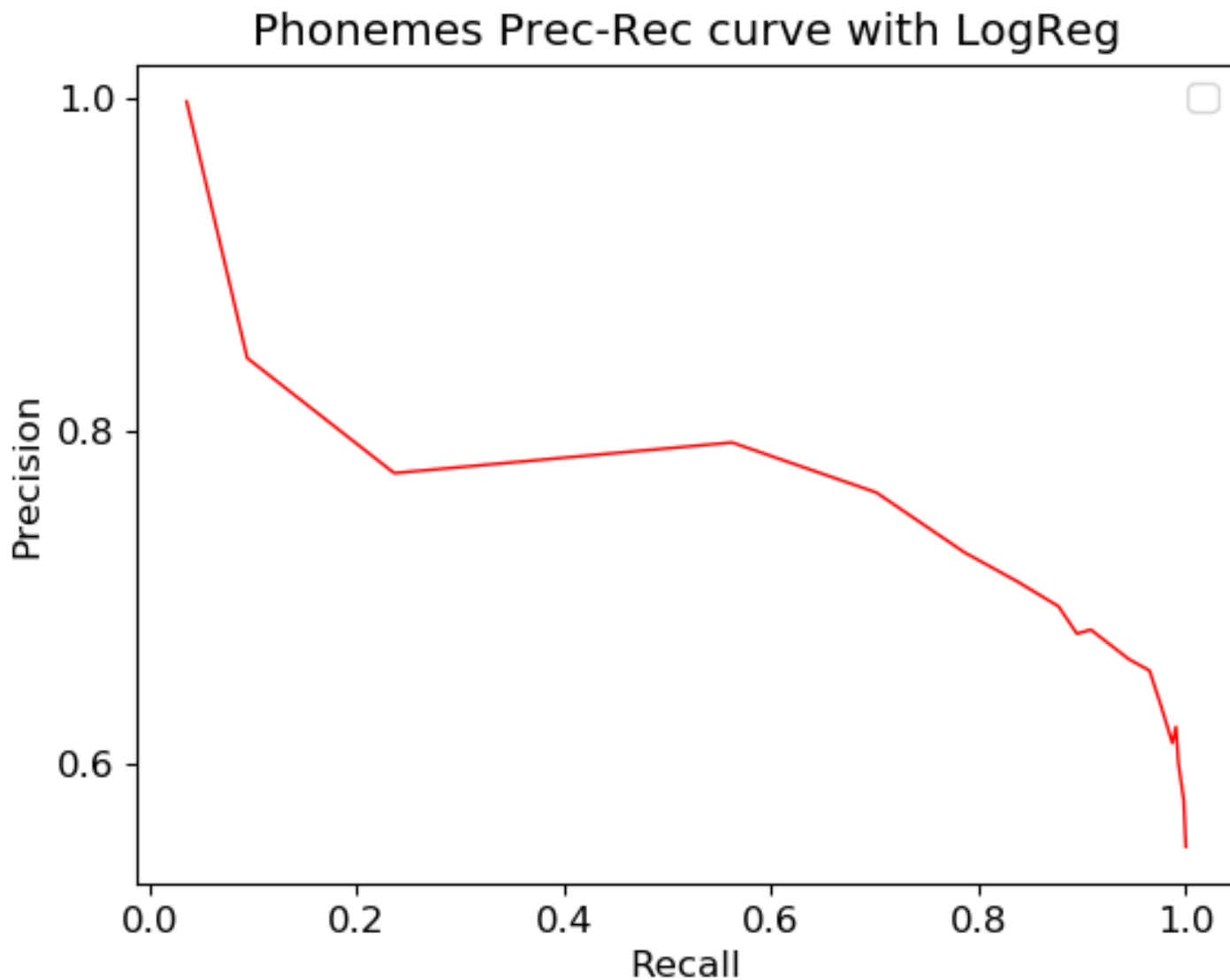
$$recall = \frac{TP}{TP + FN} \quad precision = \frac{TP}{TP + FP}$$

Полнота

Точность



# Recall — Precision кривая



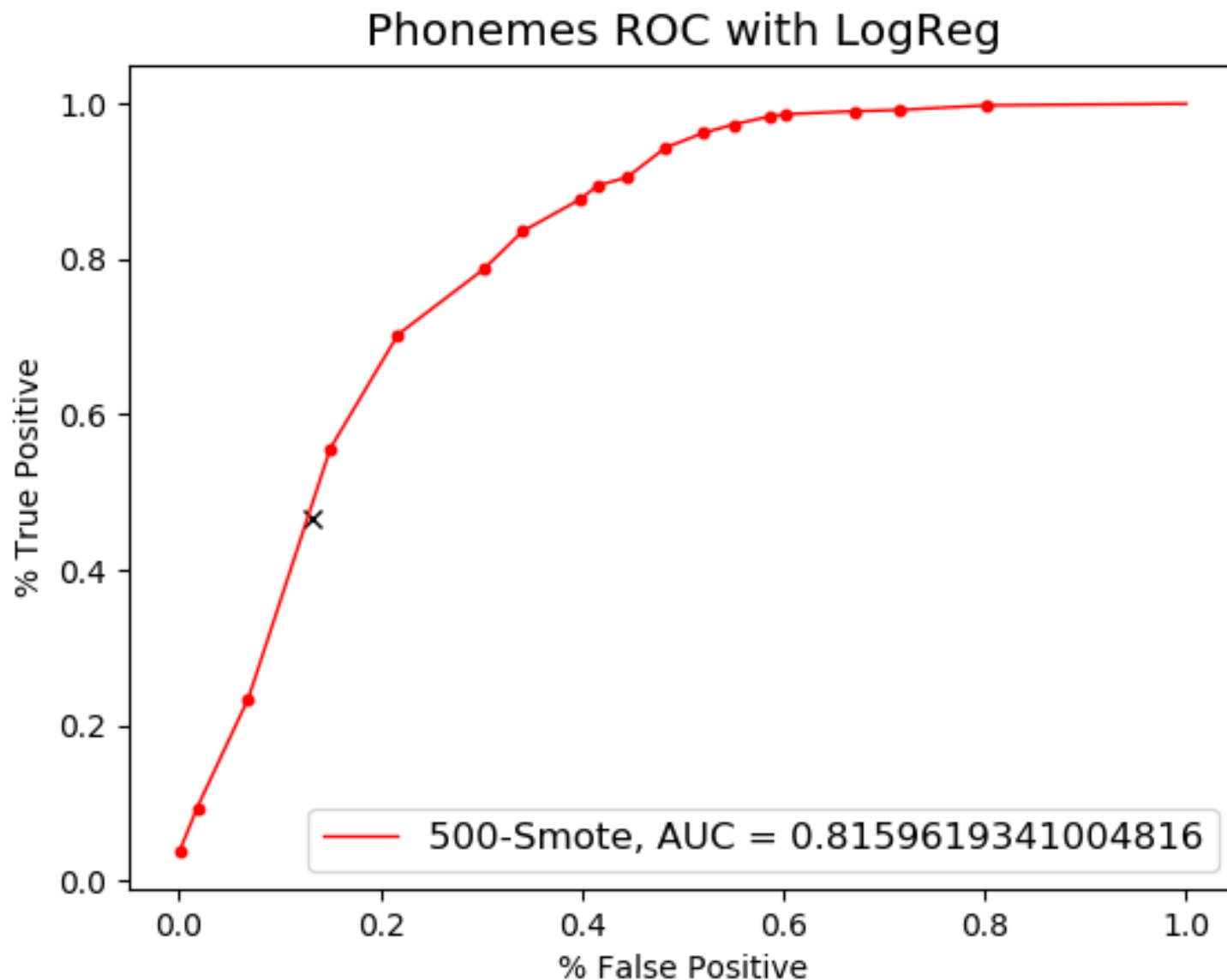
### 3) ROC кривая – для несбалансированных данных

- Кривая отображает соотношение правильно классифицированных объектов положительного класса и неверно классифицированных объектов отрицательного класса.
- 
- Каждая точка с координатами FPR по оси X и TPR по оси Y — результат классификации.
- 

False Positive Rate  $FPR (1\text{-specificity}) = \frac{FP}{TN + FP}$

True Positive Rate  $TPR (sensitivity) = \frac{TP}{TP + FN}$

# Метрика AUC - площадь под ROC кривой





# Smote для дискретных атрибутов объектов

## 1) SMOTE-NC

Пример двух объектов: F1 = 1 2 3 A B C F2 = 4 6 5 A D E

Расстояние между ними:

$$\text{Eucl} = \text{sqrt}[(4-1)^2 + (6-2)^2 + (5-3)^2 + \text{Med}^2 + \text{Med}^2]$$

---

## 2) SMOTE-N

Расстояние между двумя значениями атрибутов  $V(1)$  и  $V(2)$

$$\delta(V_1, V_2) = \sum_{i=1}^n \left| \frac{C_{1i}}{C_1} - \frac{C_{2i}}{C_2} \right|^k$$

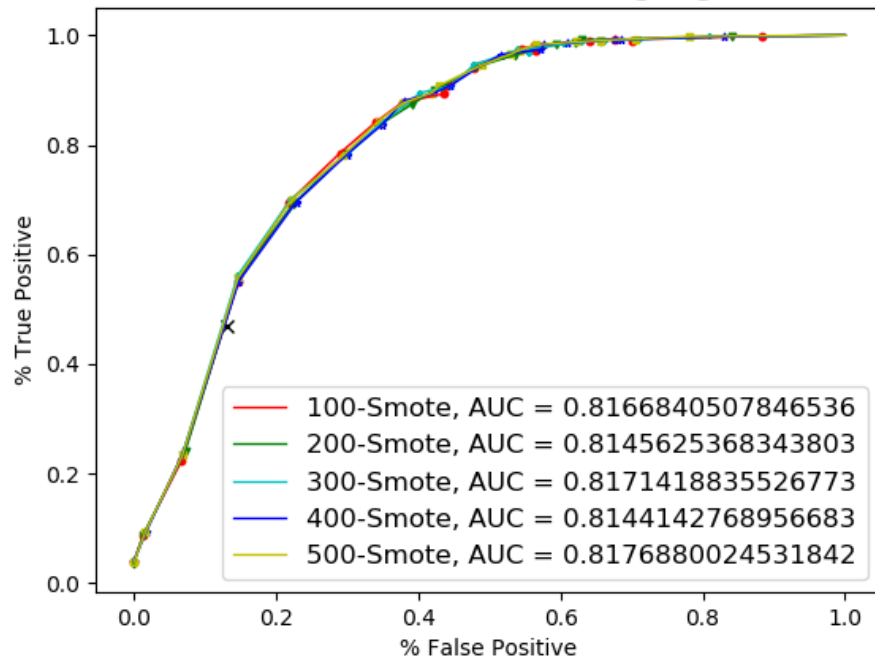
Расстояние между двумя объектами  $X$  и  $Y$

$$\Delta(X, Y) = \sum_{i=1}^N \delta(x_i, y_i)^r$$

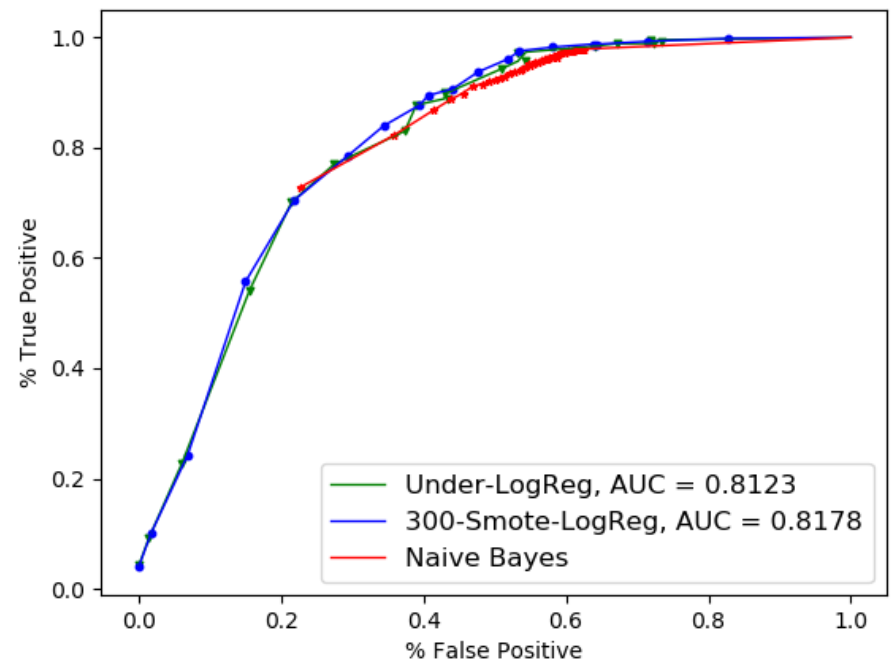


# Эксперимент

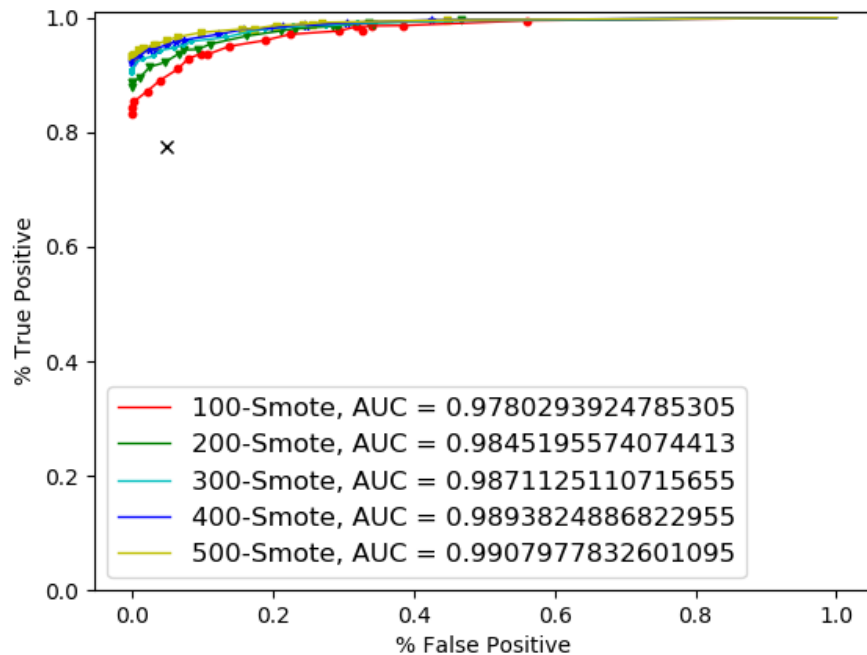
Phoneme ROCs with LogReg



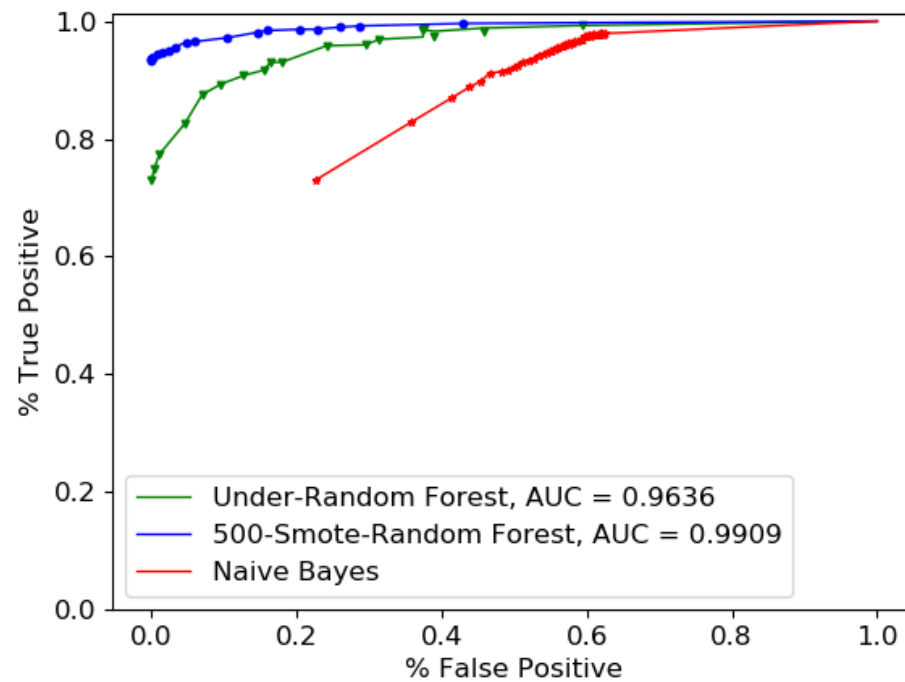
Phoneme ROC



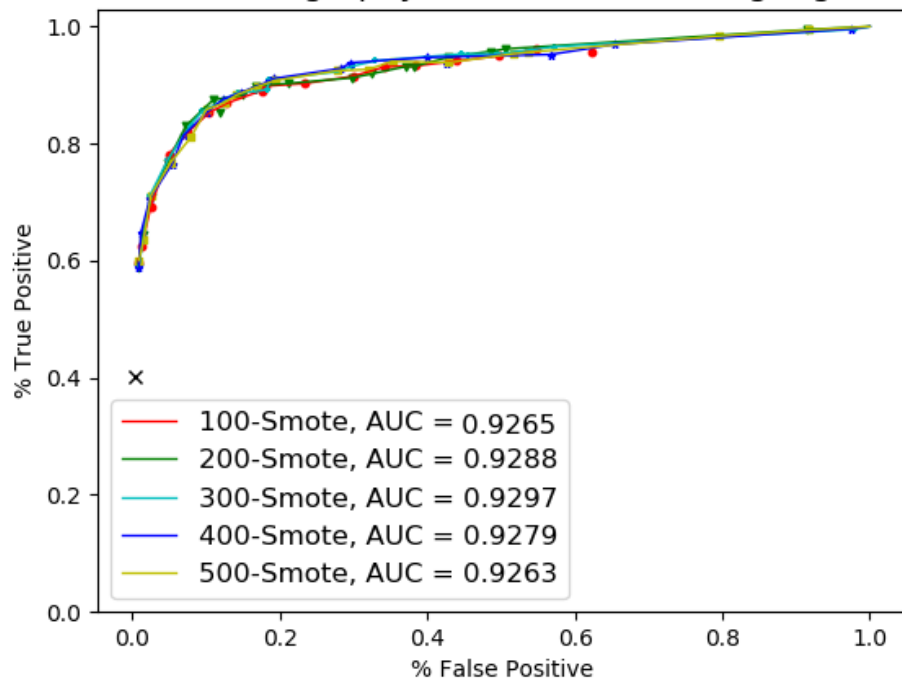
Phoneme ROC curves with Random Forest



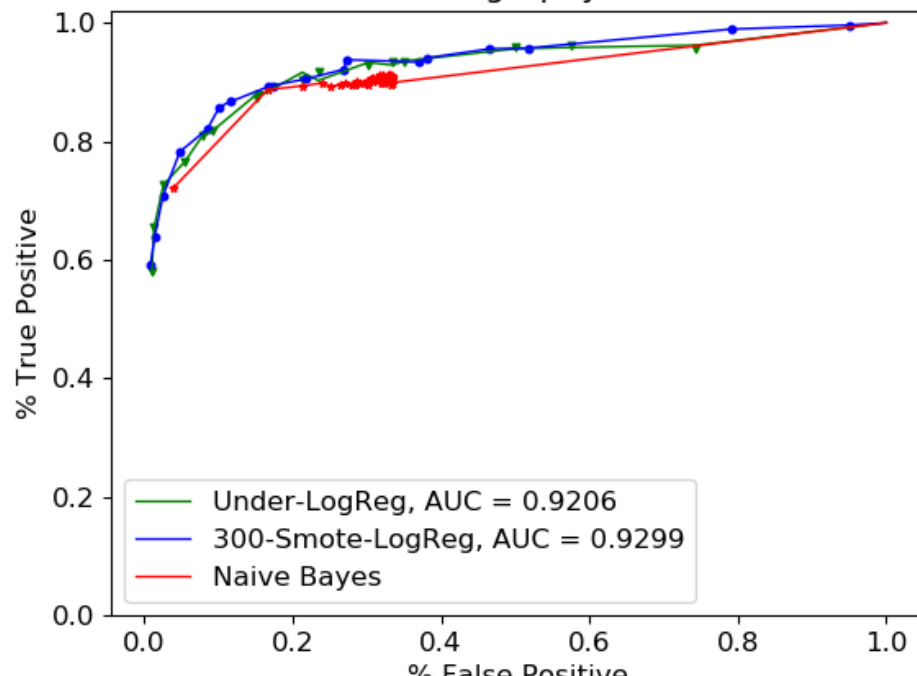
Phoneme ROC



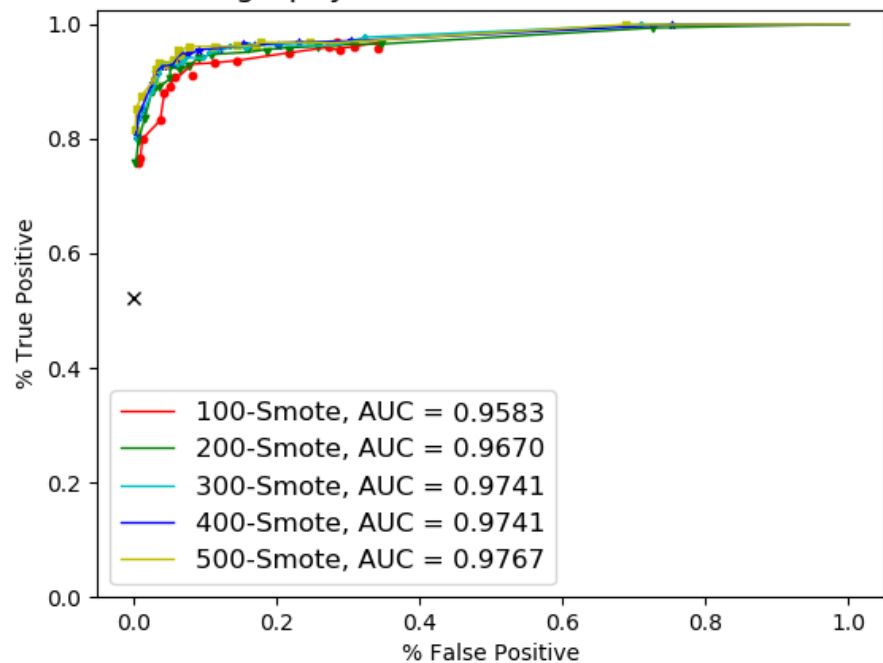
Mammography ROC curves with LogReg



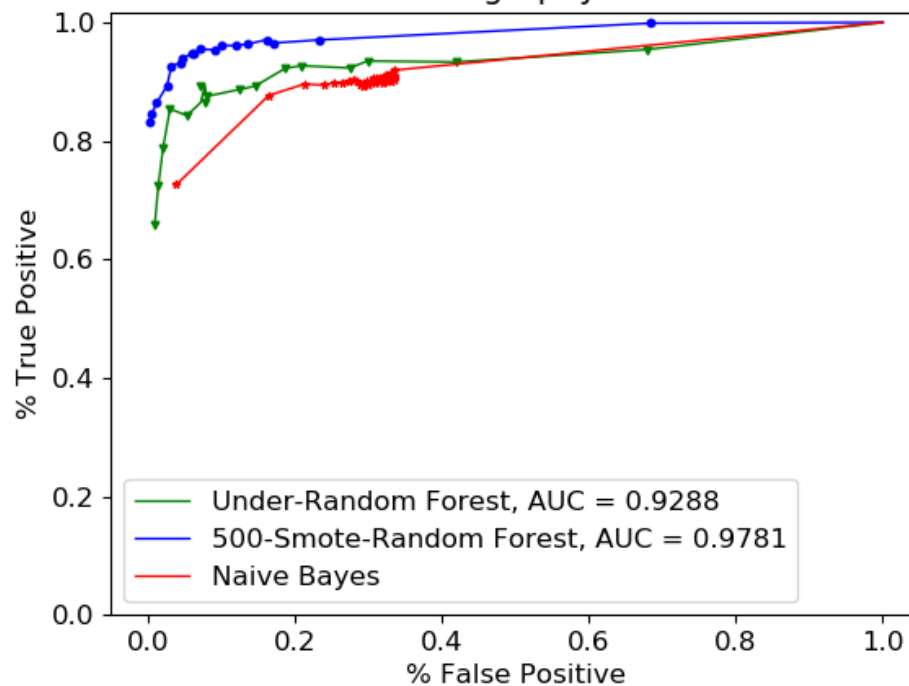
Mammography ROC



Mammography ROC curves with Random Forest



Mammography ROC



# Результаты

Dataset	Majority Class	Minority Class	Ratio	Best result
Phoneme	3818	1586	5:2	Smote + UnderSampling
Mammography	10923	260	41:1	Smote + UnderSampling

