

# ФОРМИРОВАНИЕ МОДЕЛИ



# Балансировка

- Кошка1
- Кошка2
- Кошка3
- НЕкошка1
- НЕкошка2
- ....
- НЕкошка97

	кот	неКот
кот	0	0
неКот	3	97

Precision (кот) = 0  
Recall(кот) = 0  
Accuracy = 0.97

# Балансировка

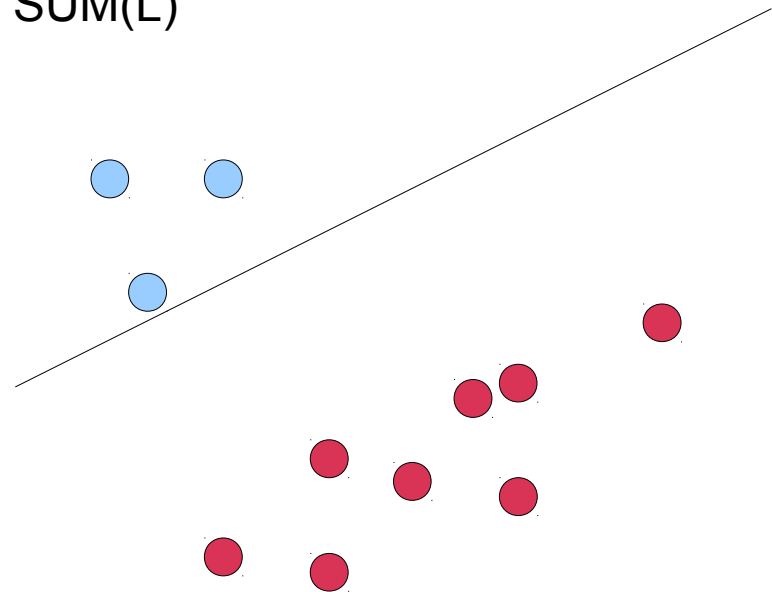
Кошка1

- Кошка2
- Кошка3
- НЕкошка1
- НЕкошка2
- ....
- НЕкошка97

$$L(y(X_i, W)) = [y(X_i, W) \neq y_{truei}]$$

$$y_i = X_i W$$

$$Q = \text{SUM}(L)$$

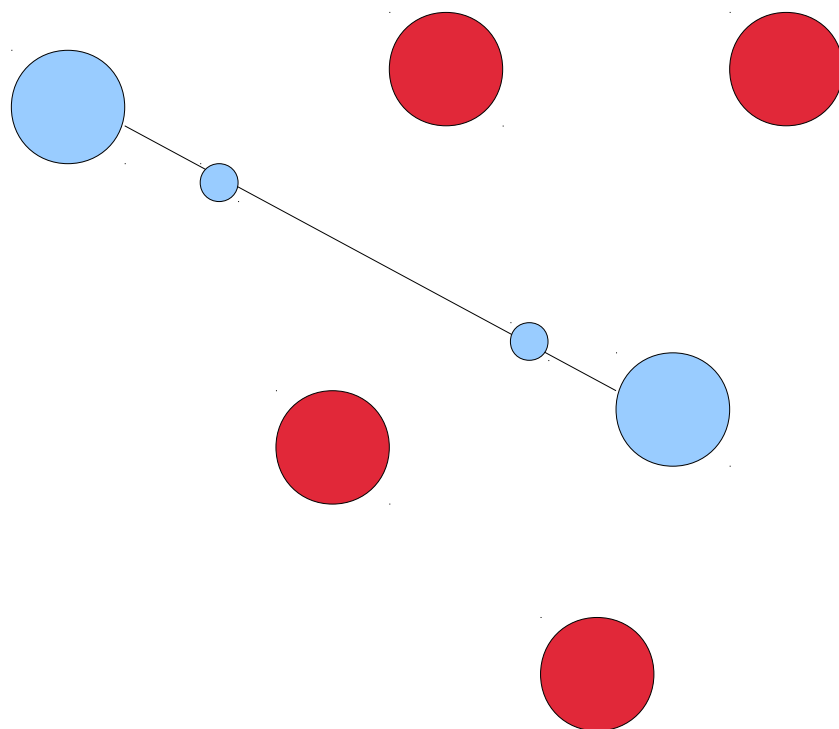


# Балансировка

- Собрать больше данных
- Выбрать подходящую метрику качества
- Попробовать разные модели, одни модели более устойчивы к несбалансированным данным, чем другие
- Штраф за ошибки при прогнозе меньшего класса
- Undersampling и Oversampling
- Создание синтетических примеров для меньшего класса

# Балансировка

SMOTE



# Балансировка

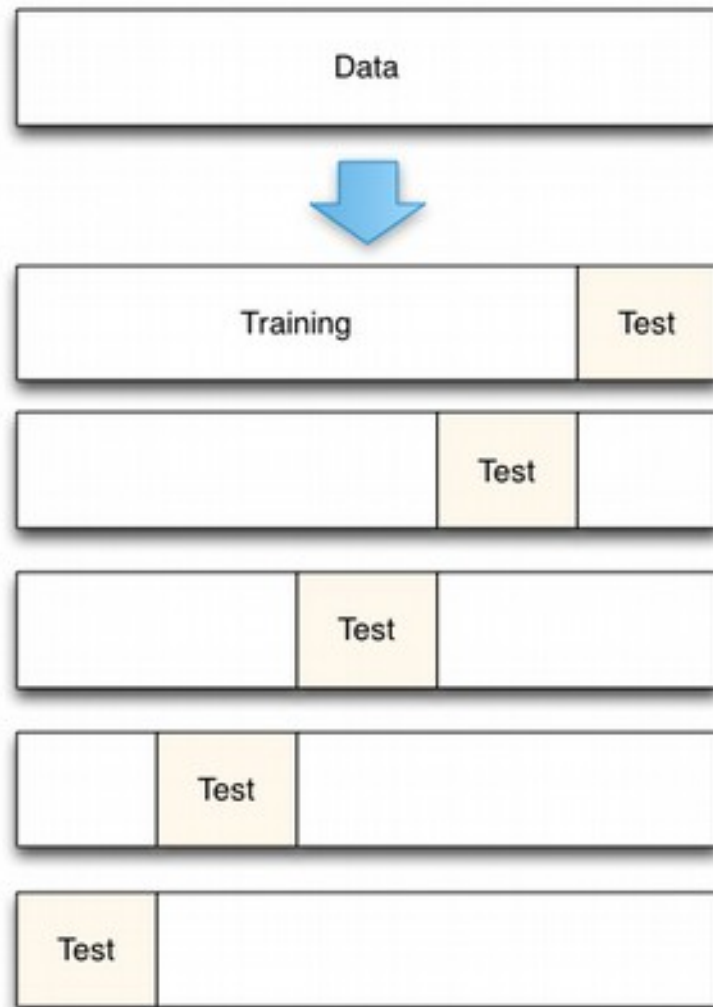
Undersampling



Oversampling



# Кросс-Валидация



Кошка2

- Кошка3

- НЕкошка1

- НЕкошка2

- ....

- НЕкошка97

- Кошка1

- Кошка3

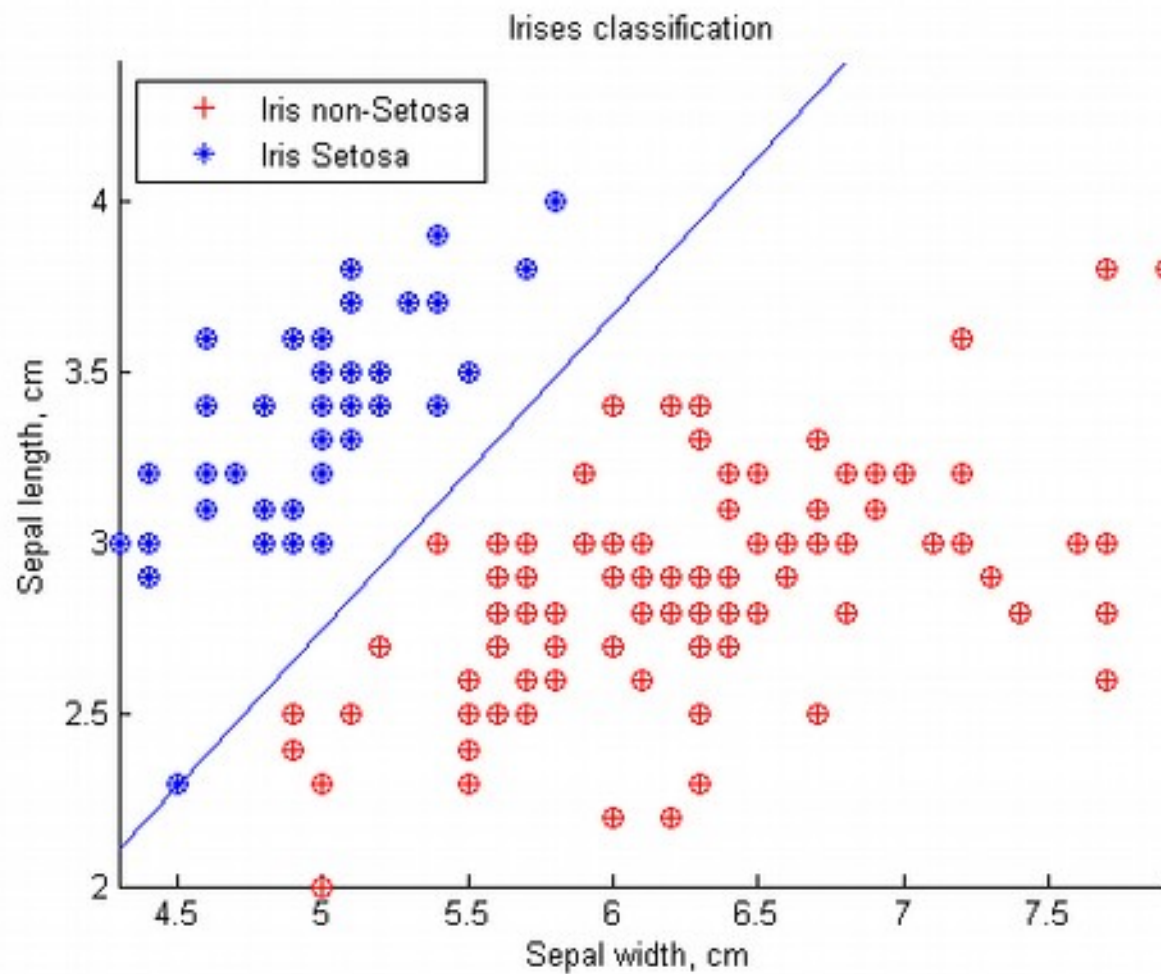
- НЕкошка1

- НЕкошка2

- ....

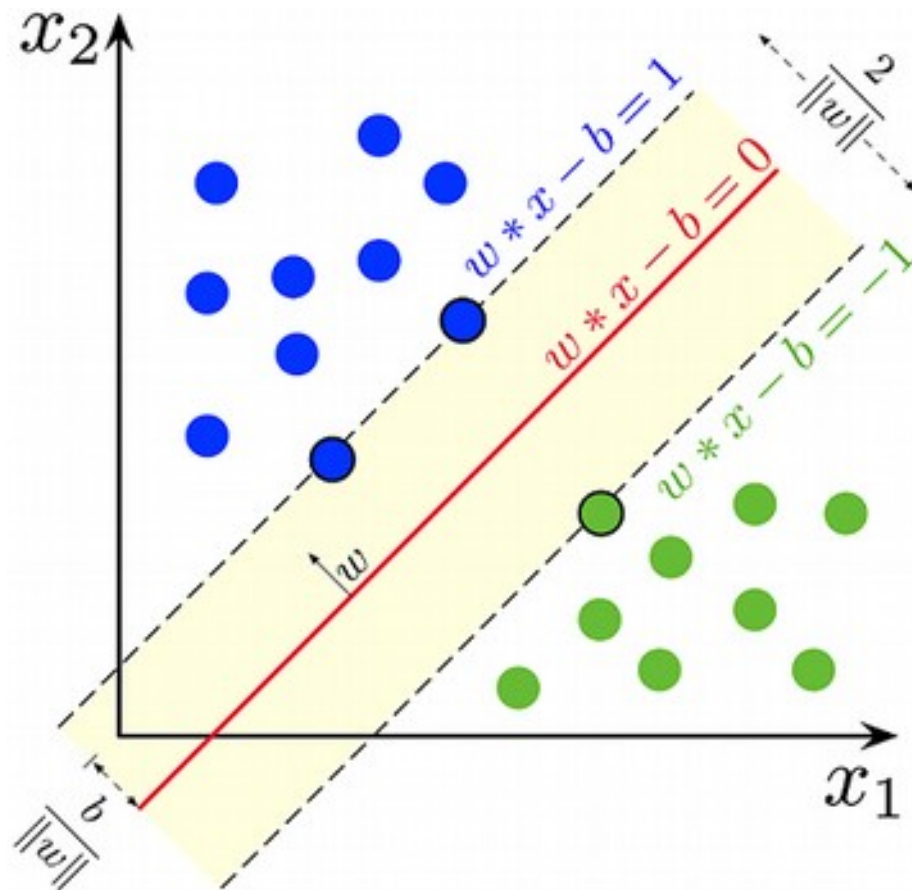
- НЕкошка97

# Классификаторы

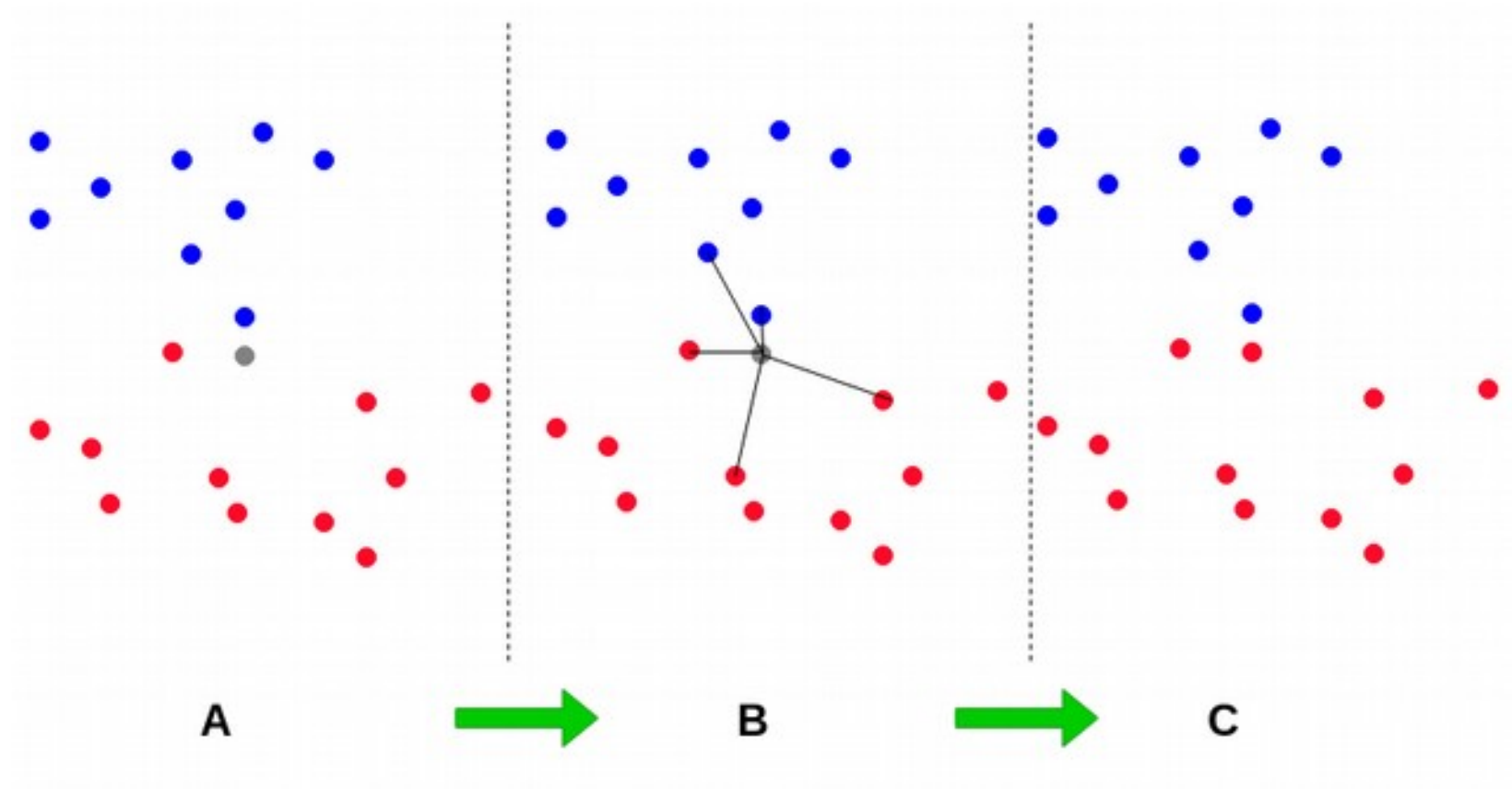




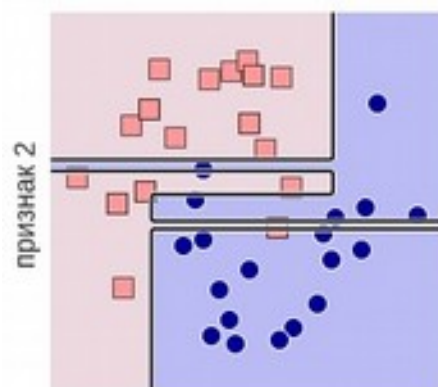
# Классификаторы - SVM



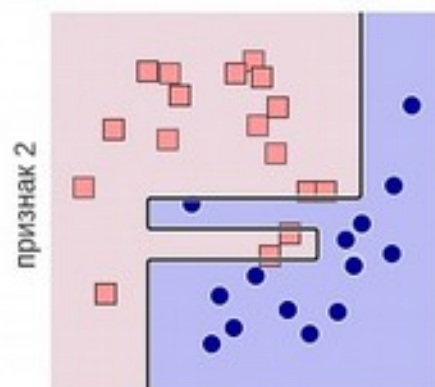
# Классификаторы - KNN



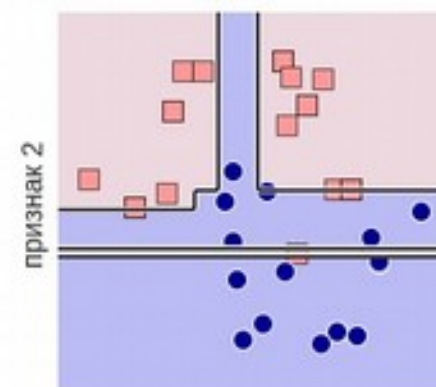
# Классификаторы — Случайный лес



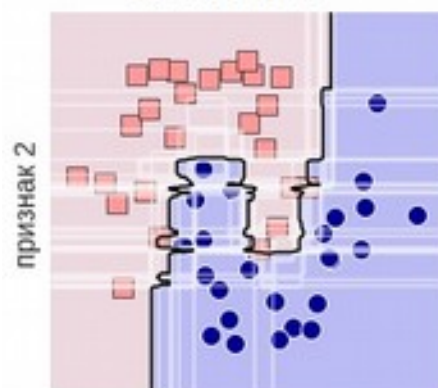
признак 1  
дерево № 1



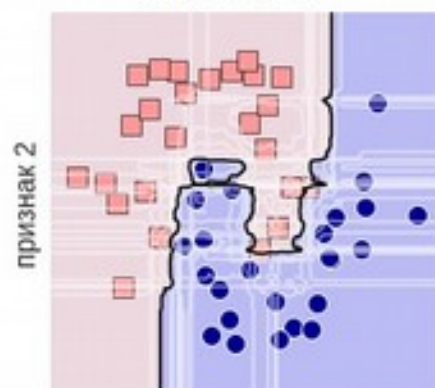
признак 1  
дерево № 2



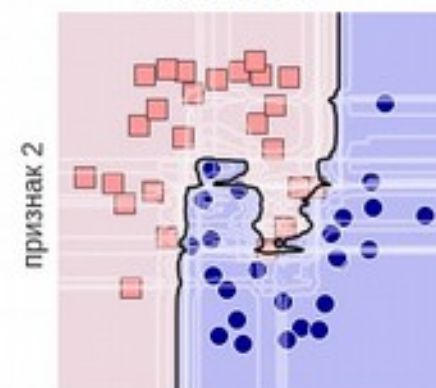
признак 1  
дерево № 3



признак 1  
RF, число деревьев=10



признак 1  
RF, число деревьев=100



признак 1  
RF, число деревьев=1000