



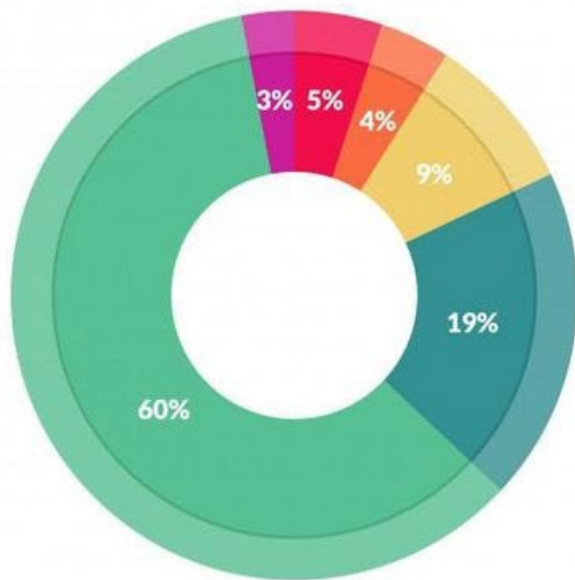
# Introduction to Machine Learning

**Practical Considerations:**  
**Features**  
**and**  
**Handling Class Imbalance**

# Features

“garbage in, garbage out”

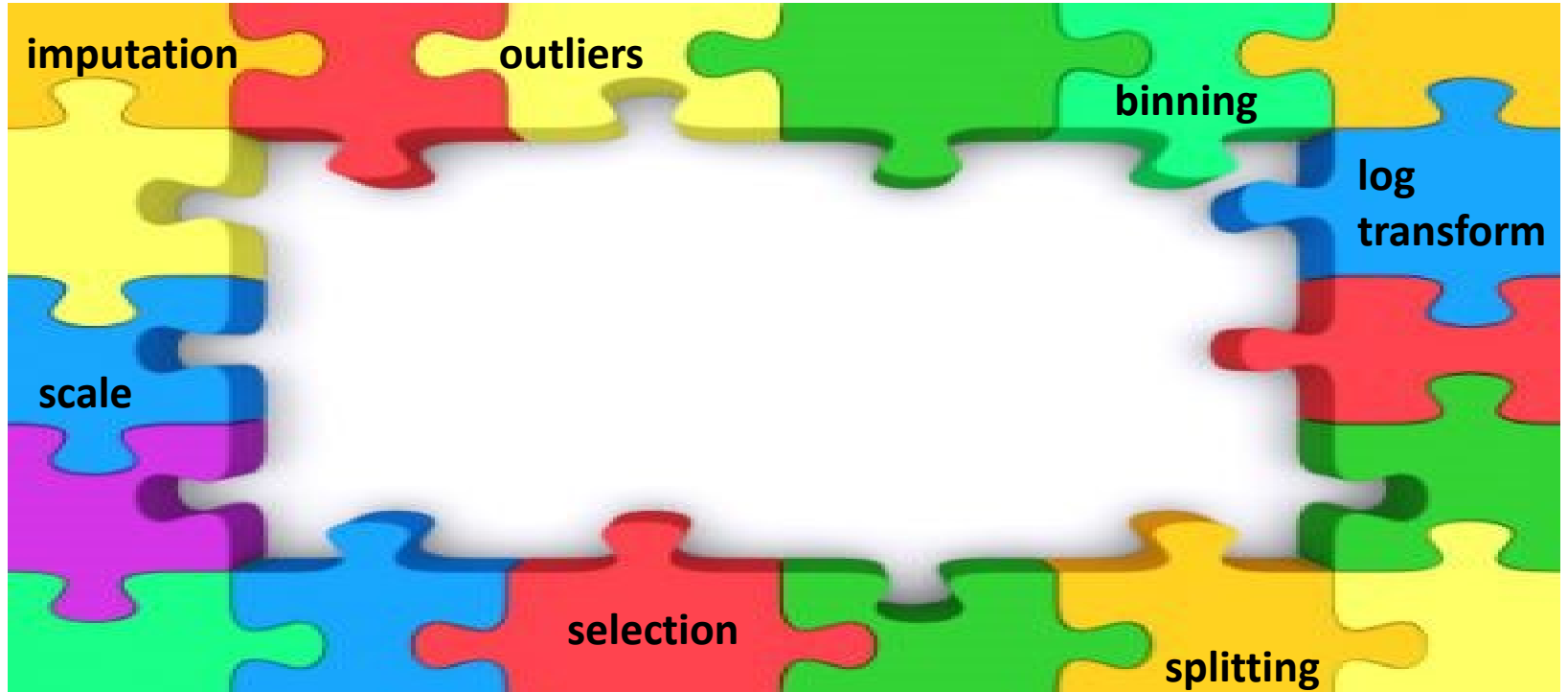
# Features



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

# Feature Engineering



# Imputation

- Why are values missing?
- What do we do?



## Feature Engineering



# Outliers

- Standard deviation
- Percentiles



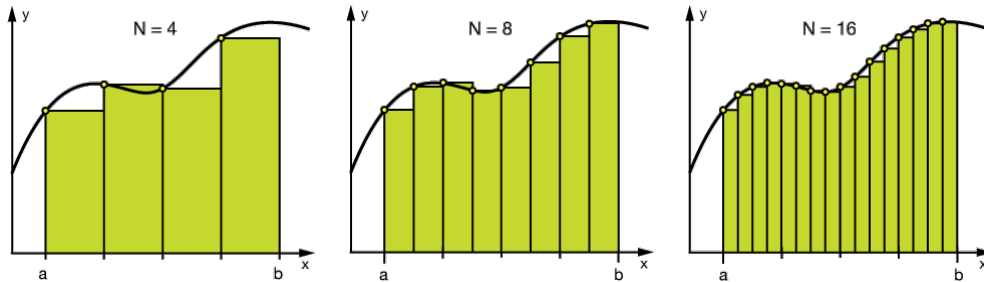
- What do we do?

## Feature Engineering



# Binning

- Equal width
- Equal frequency
- Bin value



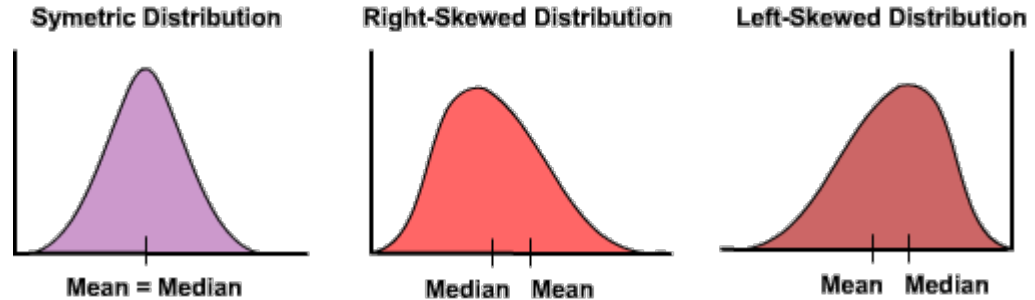
## Feature Engineering



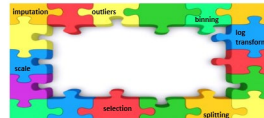
# Log Transform

- Another way to normalize

- $\text{Log}(x+1)$



Feature Engineering





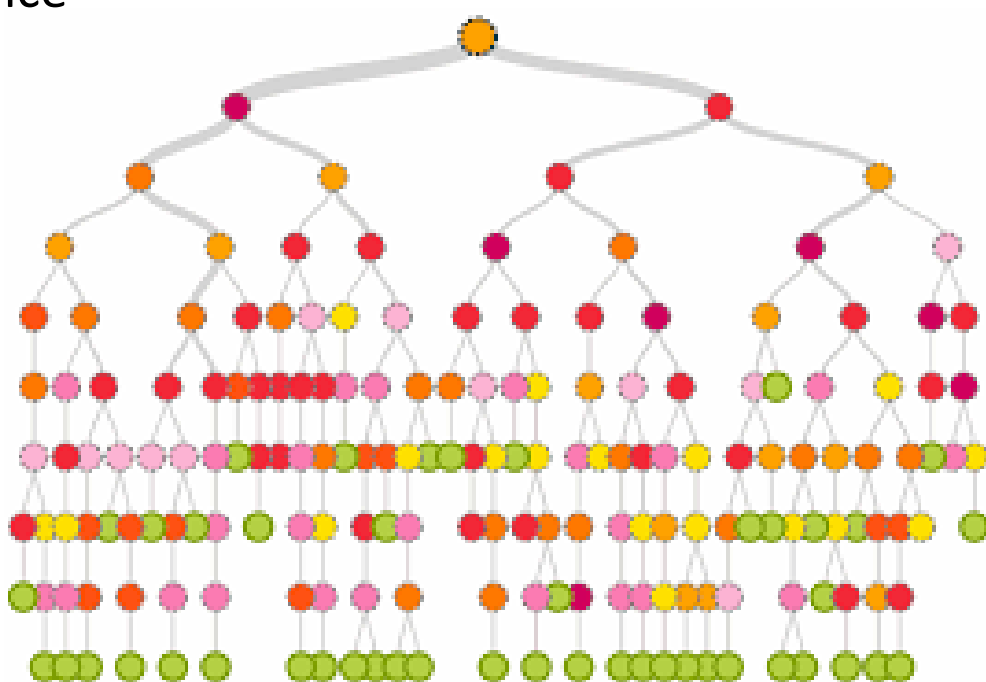
# Splitting

Feature Engineering



# Selection

- Prune features with low variance
- Prune redundant features
  - Use decision tree
  - Wrapper
  - Correlation



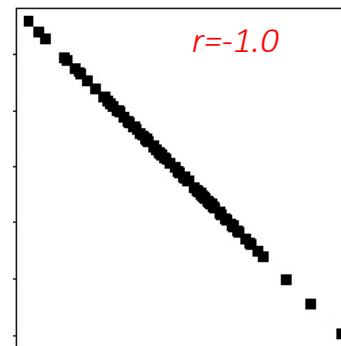
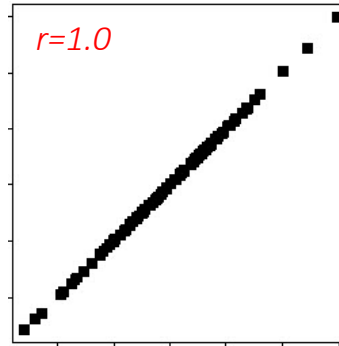
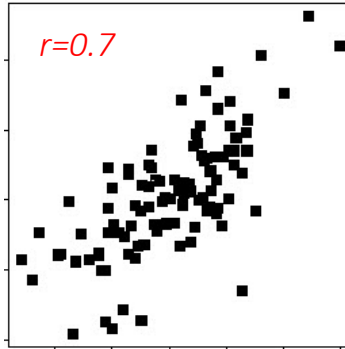
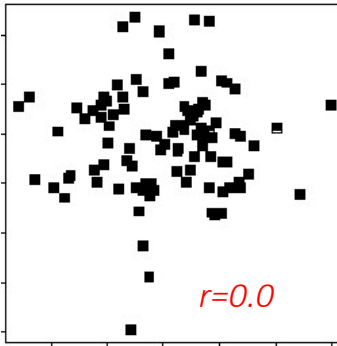
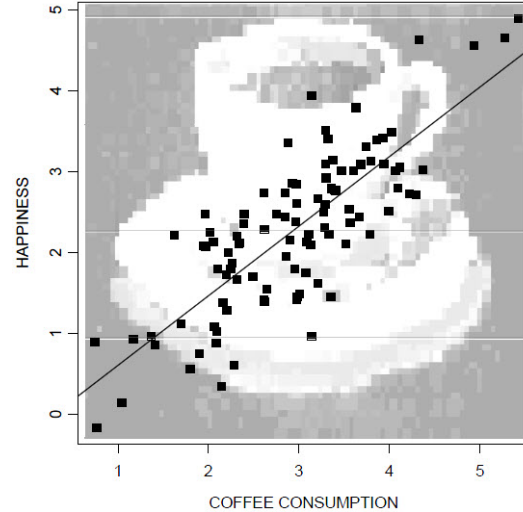
# Correlation

- Covariance

$$\text{Cov}(X, Y) = \sum_{i=1}^N \frac{(x_i - \mu_x)(y_i - \mu_y)}{N}$$

- Correlation

$$\text{Cor}(X, Y) = r = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

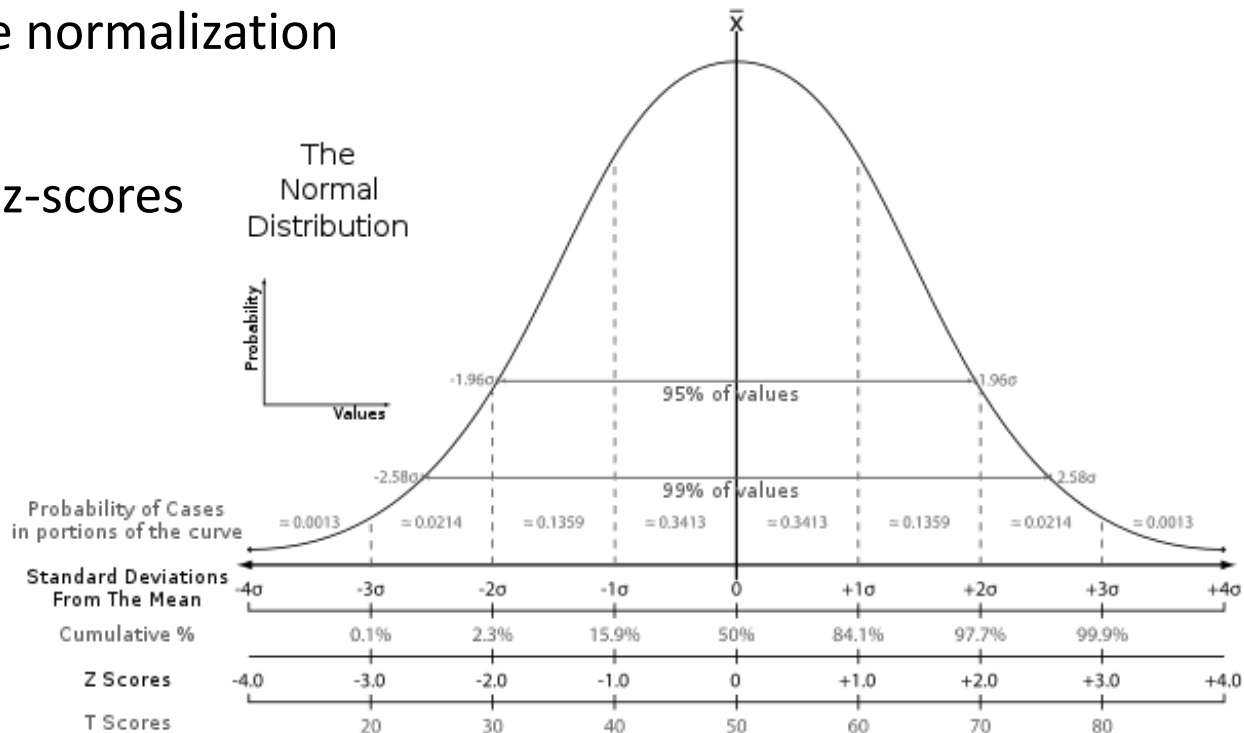


Feature Engineering



# Scale

- We discussed feature normalization
- Another approach is z-scores (standard scores)



THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG  
PILE OF LINEAR ALGEBRA, THEN COLLECT  
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL  
THEY START LOOKING RIGHT.



# Handling Class Imbalance

© 2013 Ted Goff



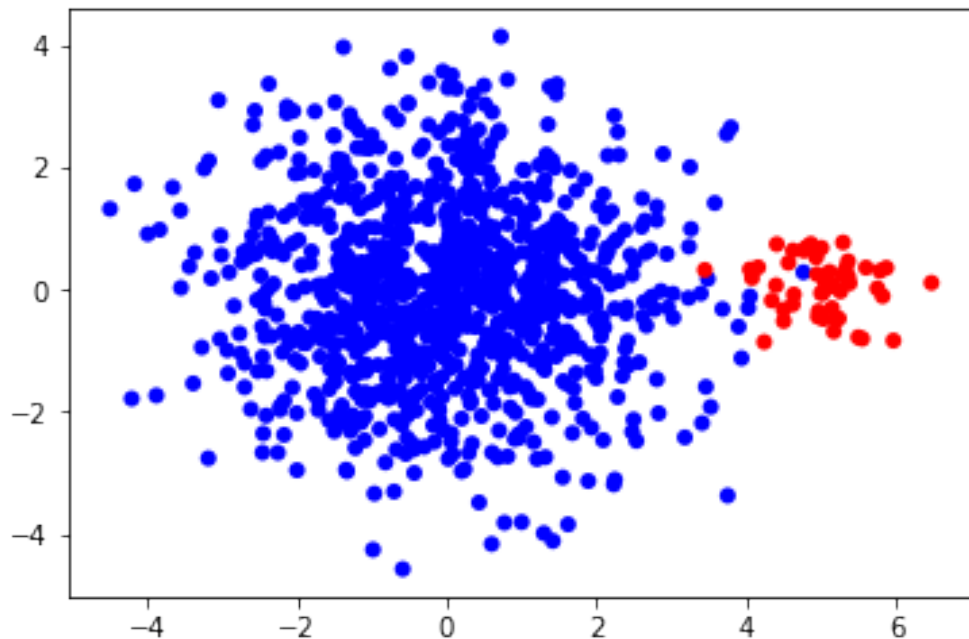
“I achieved that 99.9%  
fraud detection accuracy  
you requested.”

The catch?

# Imbalanced class distributions - misleading results



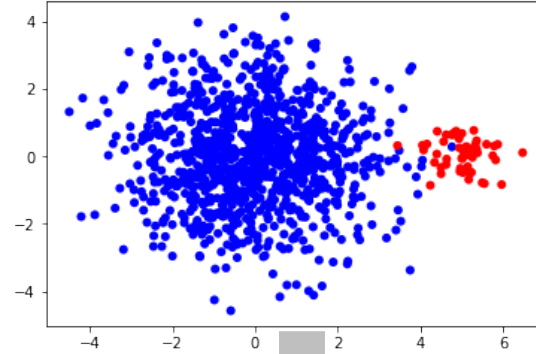
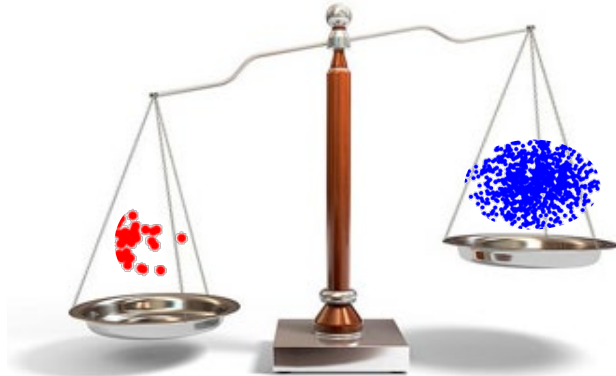
# Imbalanced class distributions





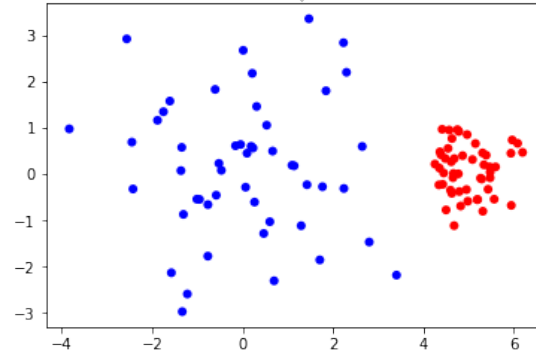
# Imbalanced class distributions - misleading results

Weight

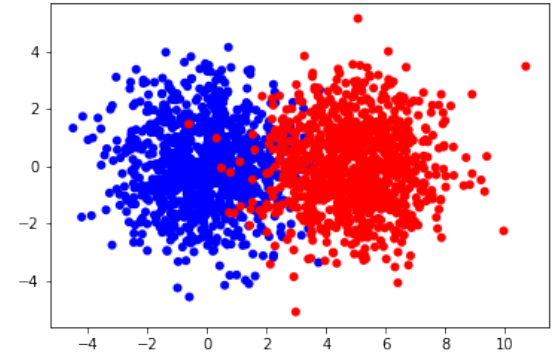


Imbalanced data

Under-sample



Over-sample



# Standard binary classification

## TASK: BINARY CLASSIFICATION

*Given:*

1. An input space  $\mathcal{X}$
2. An unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times \{-1, +1\}$
3. A training set  $D$  sampled from  $\mathcal{D}$

*Compute:* A function  $f$  minimizing:  $\mathbb{E}_{(x,y) \sim \mathcal{D}} [f(x) \neq y]$

# Standard binary classification

## TASK: $\alpha$ -WEIGHTED BINARY CLASSIFICATION

*Given:*

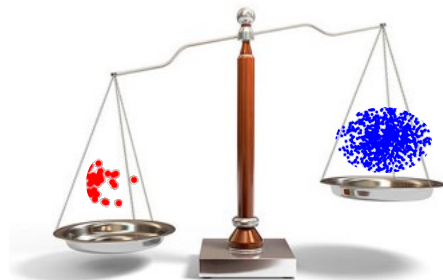
1. An input space  $\mathcal{X}$
2. An unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times \{-1, +1\}$
3. A training set  $D$  sampled from  $\mathcal{D}$

*Compute:* A function  $f$  minimizing:  $\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \alpha^{y=1} [f(x) \neq y] \right]$

# How use this for imbalanced class distribution?

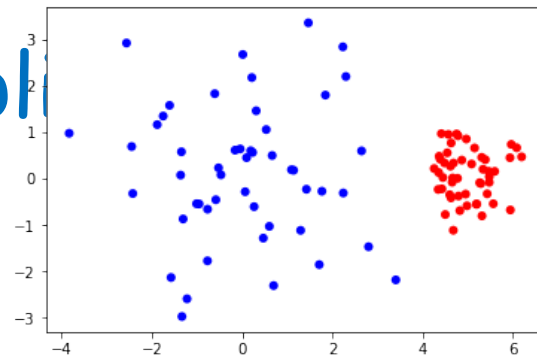
- Fraudulent transactions, 0.01% data
  - Weight = 0.0001
- Normal transactions
  - Weight = 0.9999

# Importance weight



- $\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$
- $\text{WeightedEntropy}(S) = (\alpha) -p_+ \log_2 p_+ - (1-\alpha) p_- \log_2 p_-$
- $\text{Gain}(S,A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$
- $\text{WeightedGain}(S,A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{\alpha |S_v^+| + (1-\alpha) |S_v^-|}{\alpha |S^+| + (1-\alpha) |S^-|}$

# Sampling (undersampling)



---

**Algorithm 11** SUBSAMPLEMAP( $\mathcal{D}^{\text{weighted}}, \alpha$ )

---

```
1: while true do
2:    $(x, y) \sim \mathcal{D}^{\text{weighted}}$  // draw an example from the weighted distribution
3:    $u \sim$  uniform random variable in  $[0, 1]$ 
4:   if  $y = +1$  or  $u < \frac{1}{\alpha}$  then
5:     return  $(x, y)$ 
6:   end if
7: end while
```

---

---

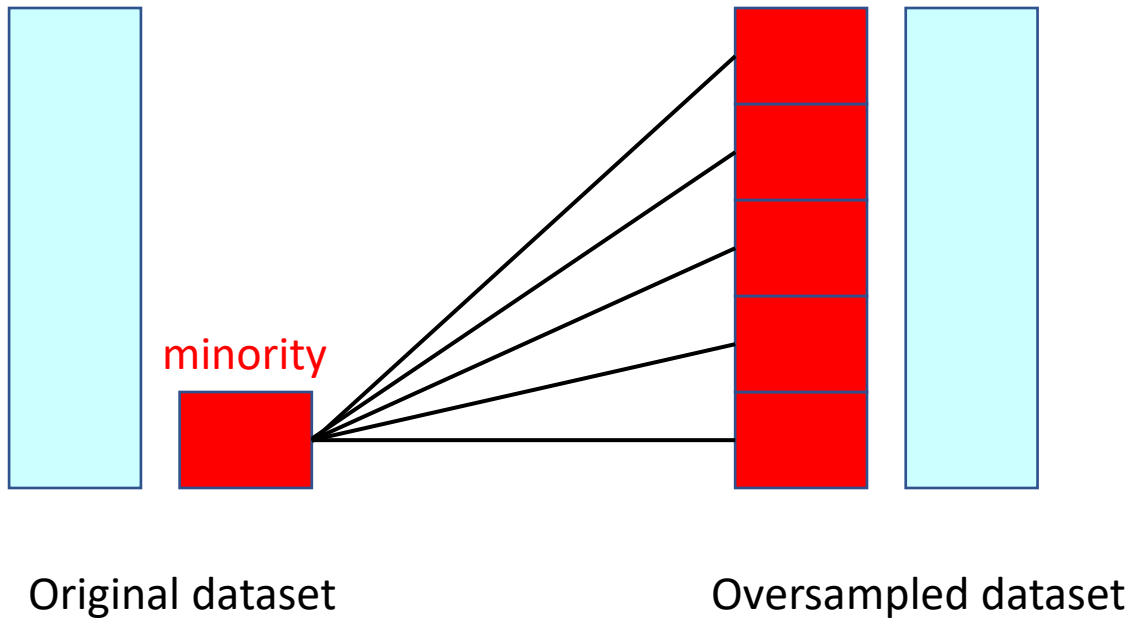
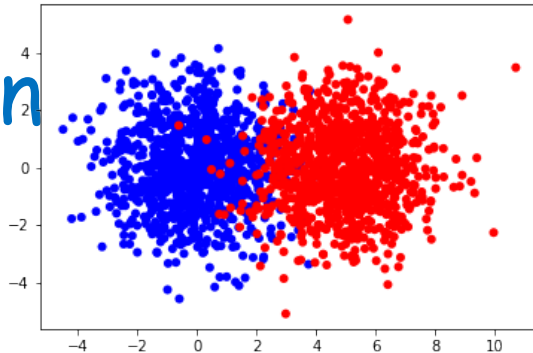
**Algorithm 12** SUBSAMPLETEST( $f^{\text{BINARY}}, \hat{x}$ )

---

```
1: return  $f^{\text{BINARY}}(\hat{x})$ 
```

---

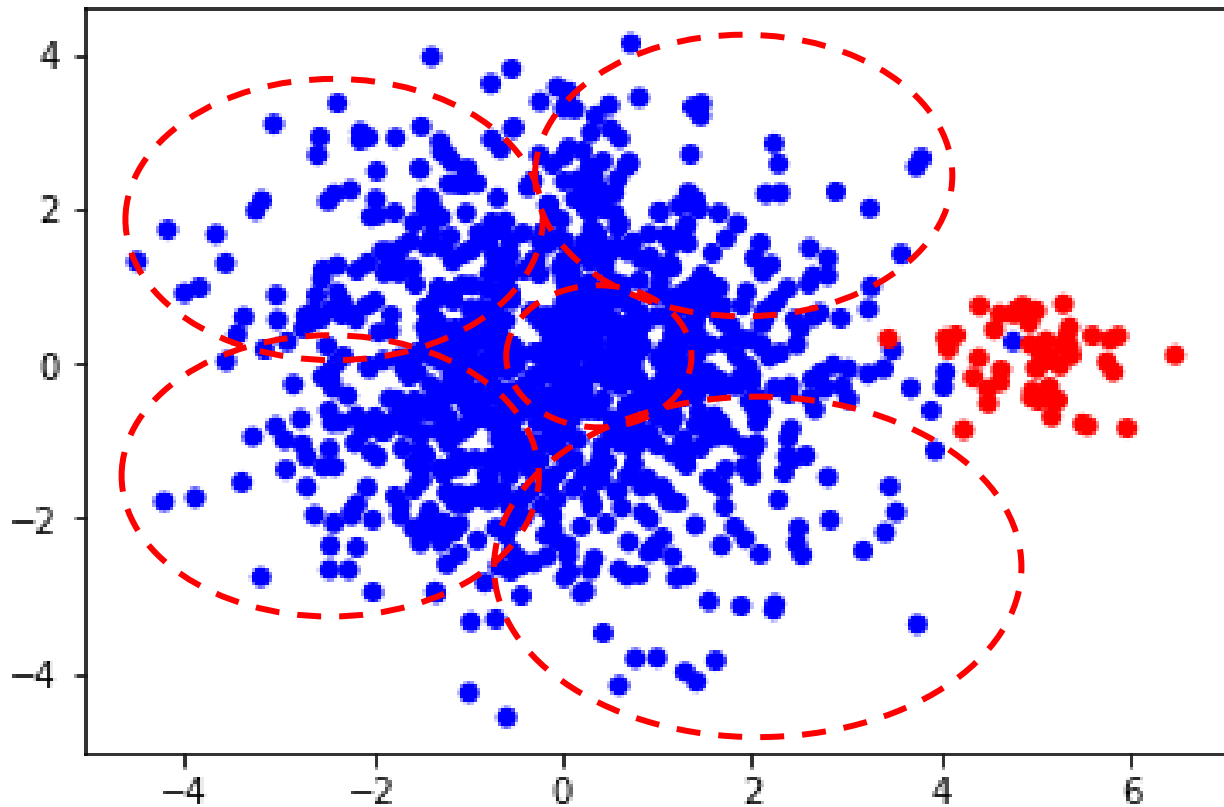
# Sampling (oversampling)



Let's try this out



# Decompose



Let's try this out