

# Rare Category Detection

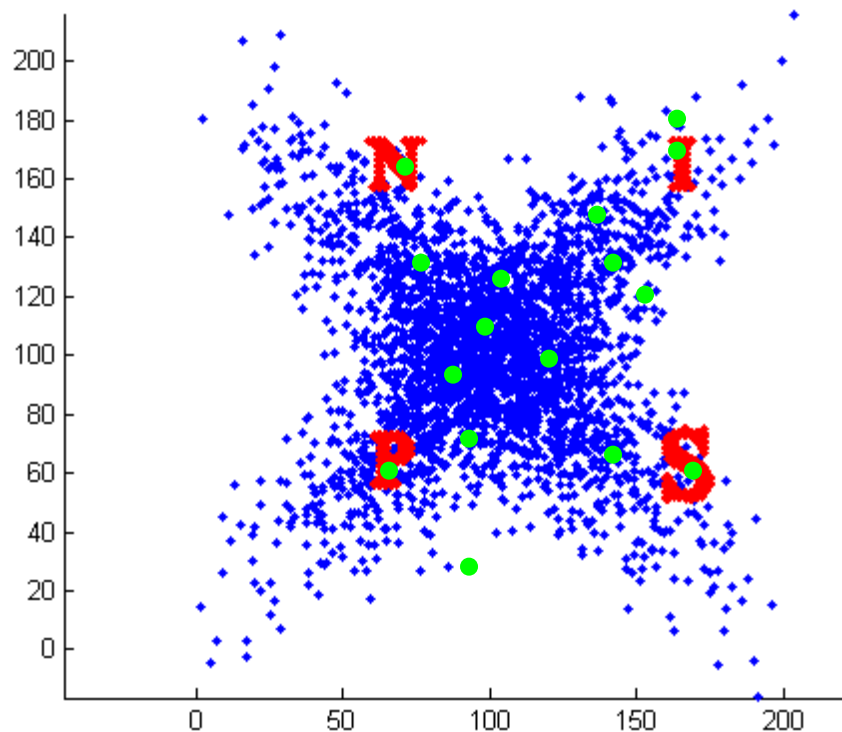


Jingrui He

Machine Learning Department  
Carnegie Mellon University

Joint work with Jaime Carbonell

# What's Rare Category Detection

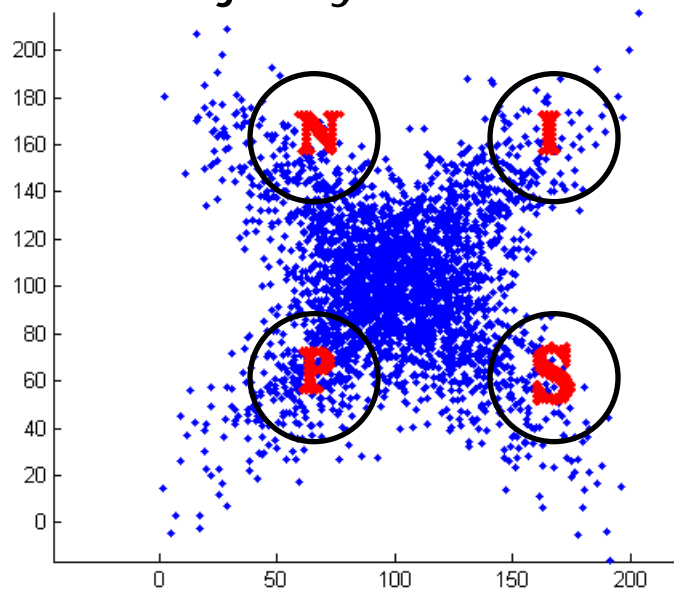


- Start de-novo
- Very skewed classes
  - Majority classes
  - *Minority classes*
- Labeling oracle
- Goal
  - Discover minority classes with *a few* label requests

# Comparison with Outlier Detection

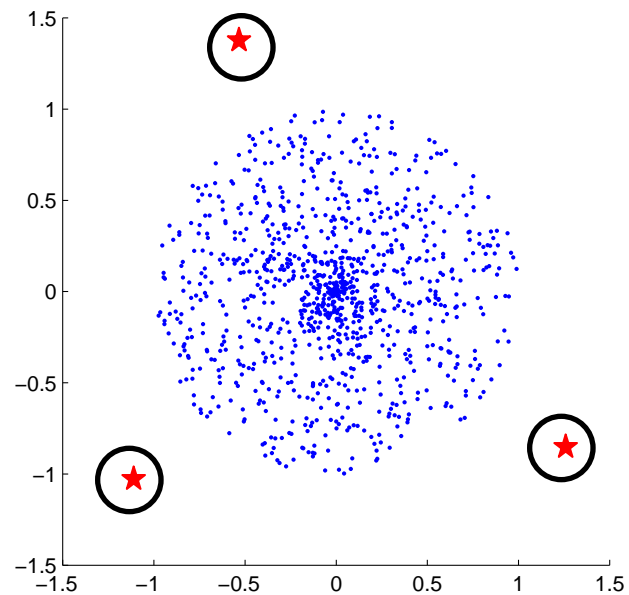
## □ Rare classes

- A group of points
- Clustered
- Non-separable from the majority classes



## □ Outliers

- A single point
- Scattered
- Separable



# Comparison with Active Learning

---

## □ Rare category detection

- Initial condition: *NO* labeled examples
- Goal: *discover* the minority classes with the least label requests

## □ Active learning

- Initial condition: labeled examples from *each* class
- Goal: *improve* the performance of the current classifier with the least label requests

App

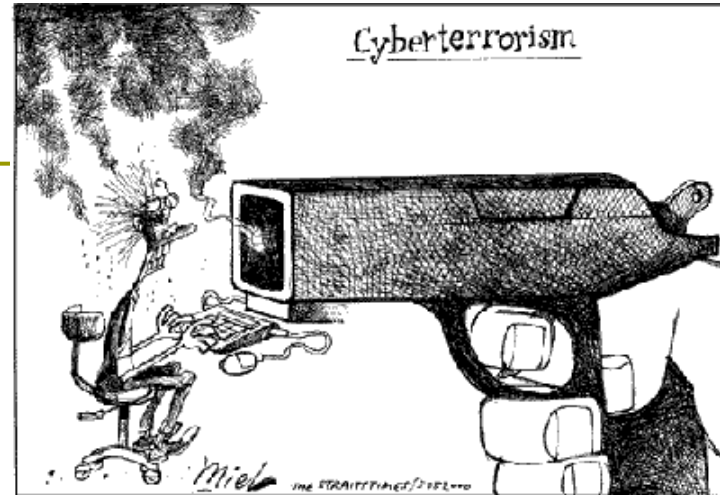
## Fraud detection

© Original Artist  
Reproduction rights obtainable from  
www.CartoonStock.com



"We're saving money this holiday season by heating our home with swiped credit cards."

## Network intrusion detection



## Astronomy



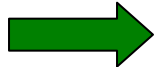
## Spam image detection



# The Big Picture



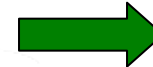
Unbalanced  
Unlabeled  
Data Set



*Rare  
Category  
Detection*



Learning in  
Unbalanced  
Settings



Classifier



Feature  
Extraction



Spatial

Relational

Temporal



Raw  
Data



# Outline

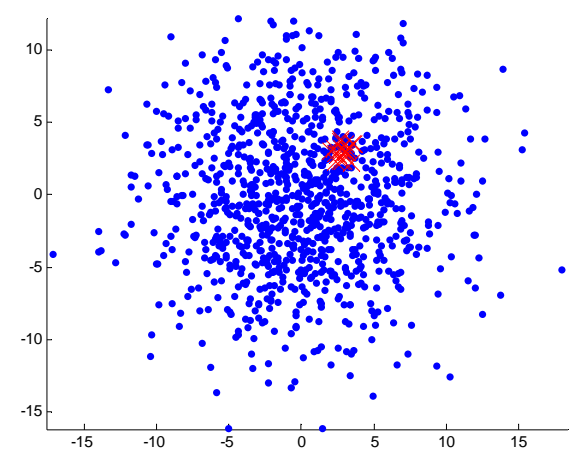
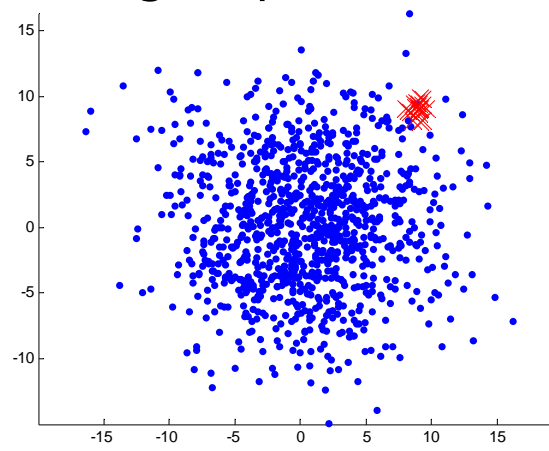
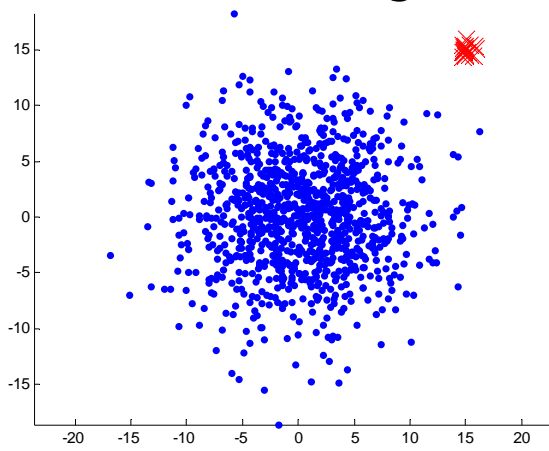
---

- Problem definition
- Related work
- Rare category detection for spatial data
  - Prior-dependent rare category detection
  - Prior-free rare category detection
- Conclusion

# Related Work

- Pelleg & Moore 2004
  - Mixture model
  - Different selection criteria
- Fine & Mansour 2006
  - Generic consistency algorithm
  - Upper bounds and lower bounds
- Papadimitriou et al 2003
  - LOCI algorithm for groups of outliers

*Separable or  
Near-separable*





# Outline

---

- Problem definition
- Related work
- Rare category detection for spatial data
  - Prior-dependent rare category detection
  - Prior-free rare category detection
- Conclusion

# Notations

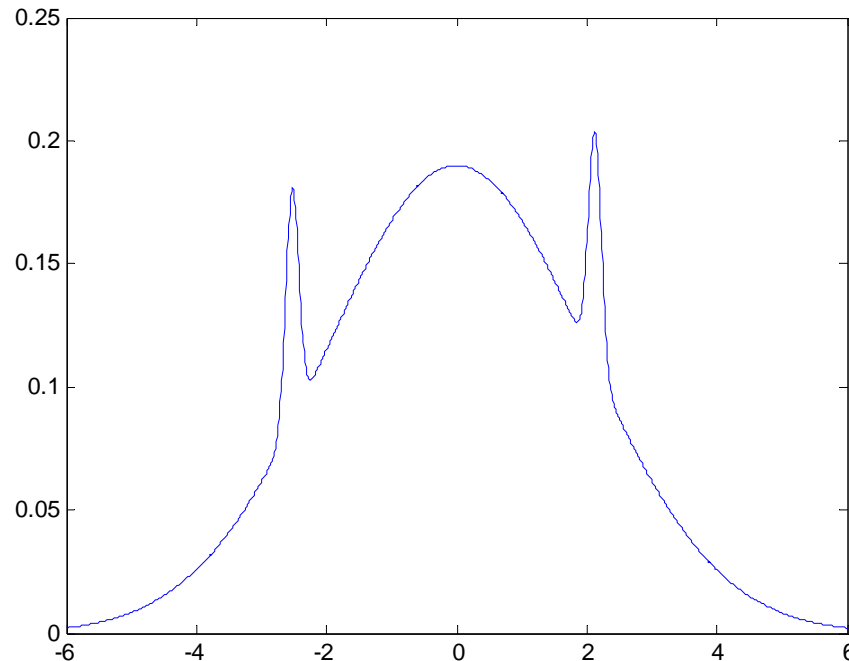
---

- Unlabeled examples:  $S = \{x_1, \dots, x_n\}$ ,  $x_i \in \mathbb{R}^d$
- $m$  Classes:  $y_i \in \{1, \dots, m\}$
- $m-1$  rare classes:  $p^2, \dots, p^m$
- One majority class:  $p^1 \gg p^c$ ,  $2 \leq c \leq m$
  
- Goal: find at least *ONE* example from *each* rare class by requesting *a few* labels

# Assumptions

---

- The distribution of the majority class is sufficiently smooth
- Examples from the minority classes form compact clusters in the feature space

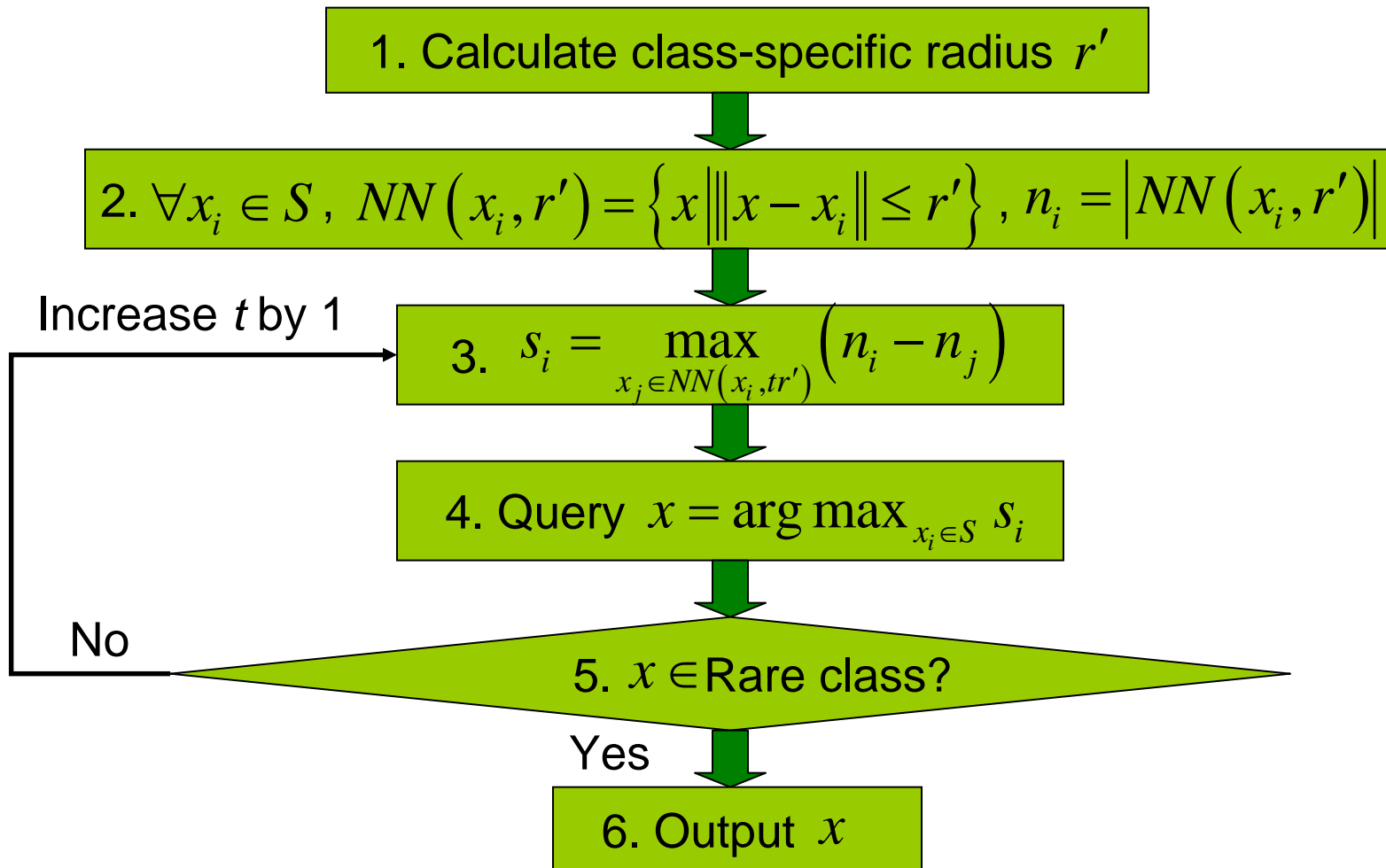


# Overview of the Algorithms

---

- Nearest-neighbor-based methods
  - Methodology: local density differential sampling
  - Intuition: select examples according to the *change in local density*

# Two Classes: NNDB



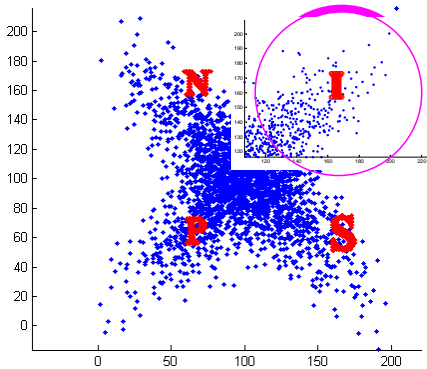
# NNDB: Calculate Class-Specific Radius

---

- Number of examples from the minority class:  $p^2 \rightarrow K = np^2$
- $\forall x_i \in S$ , calculate the distance  $r_i^K$  between  $x_i$  and its  $K^{\text{th}}$  nearest neighbor
- The class-specific radius:

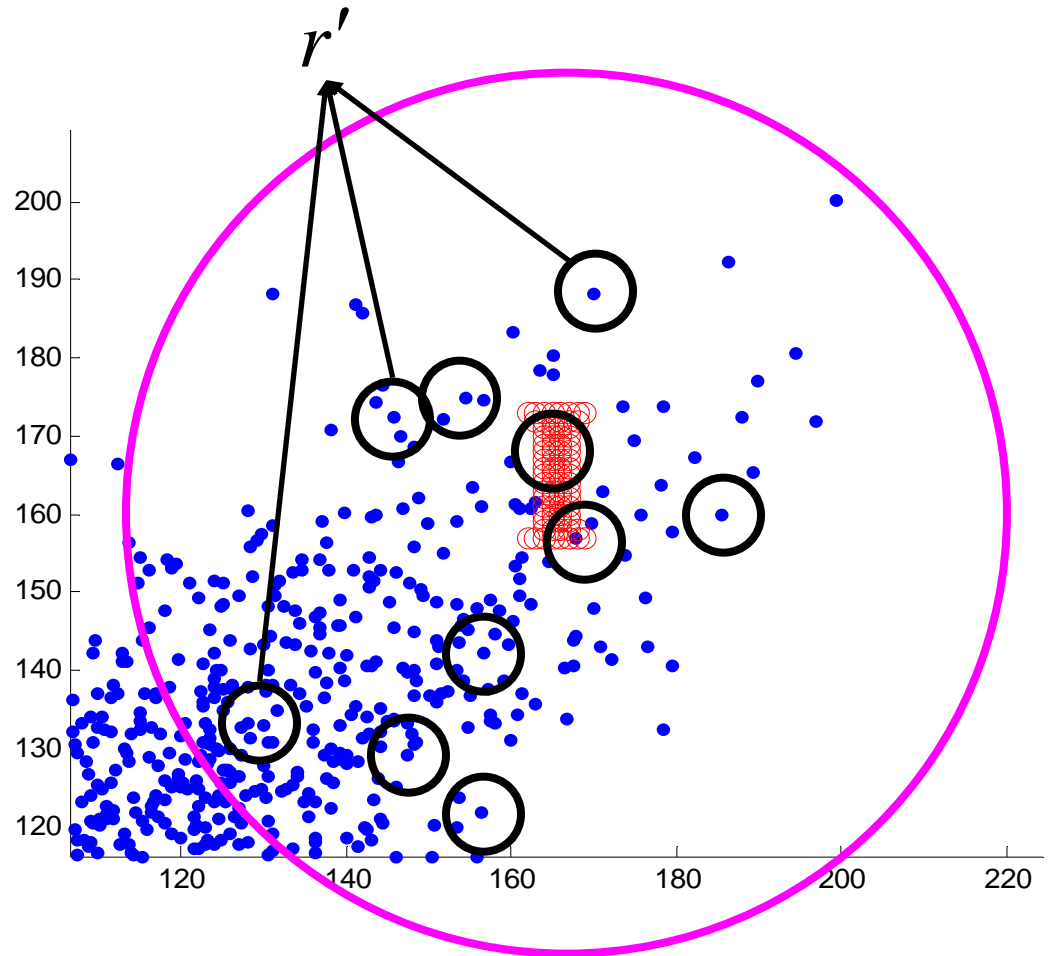
$$r' = \min_{i=1}^n r_i^K$$

# NNDB: Calculate Nearest Neighbors



$$NN(x_i, r') = \{x \mid \|x - x_i\| \leq r'\}$$

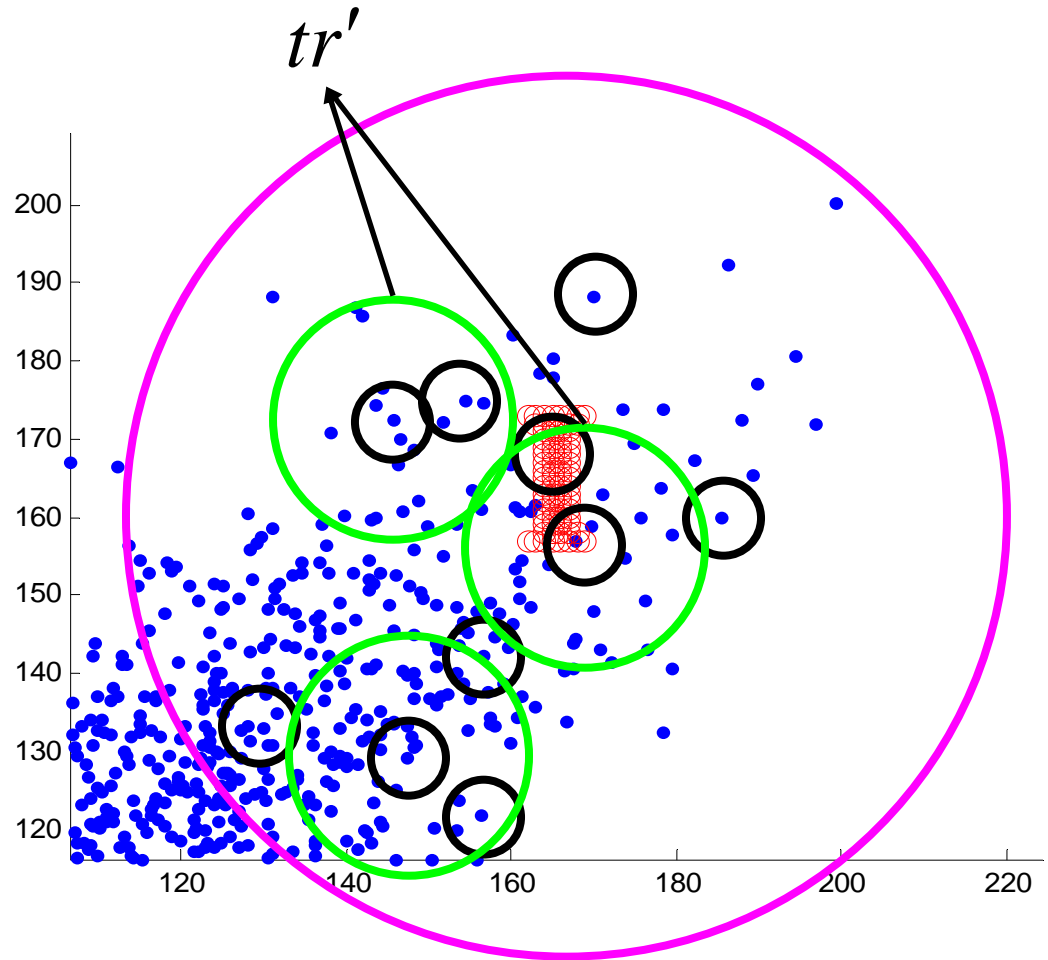
$$n_i = |NN(x_i, r')|$$



# NNDB: Calculate the Scores

$$s_i = \max_{x_j \in NN(x_i, tr')} (n_i - n_j)$$

Query  $x = \arg \max_{x_i \in S} s_i$



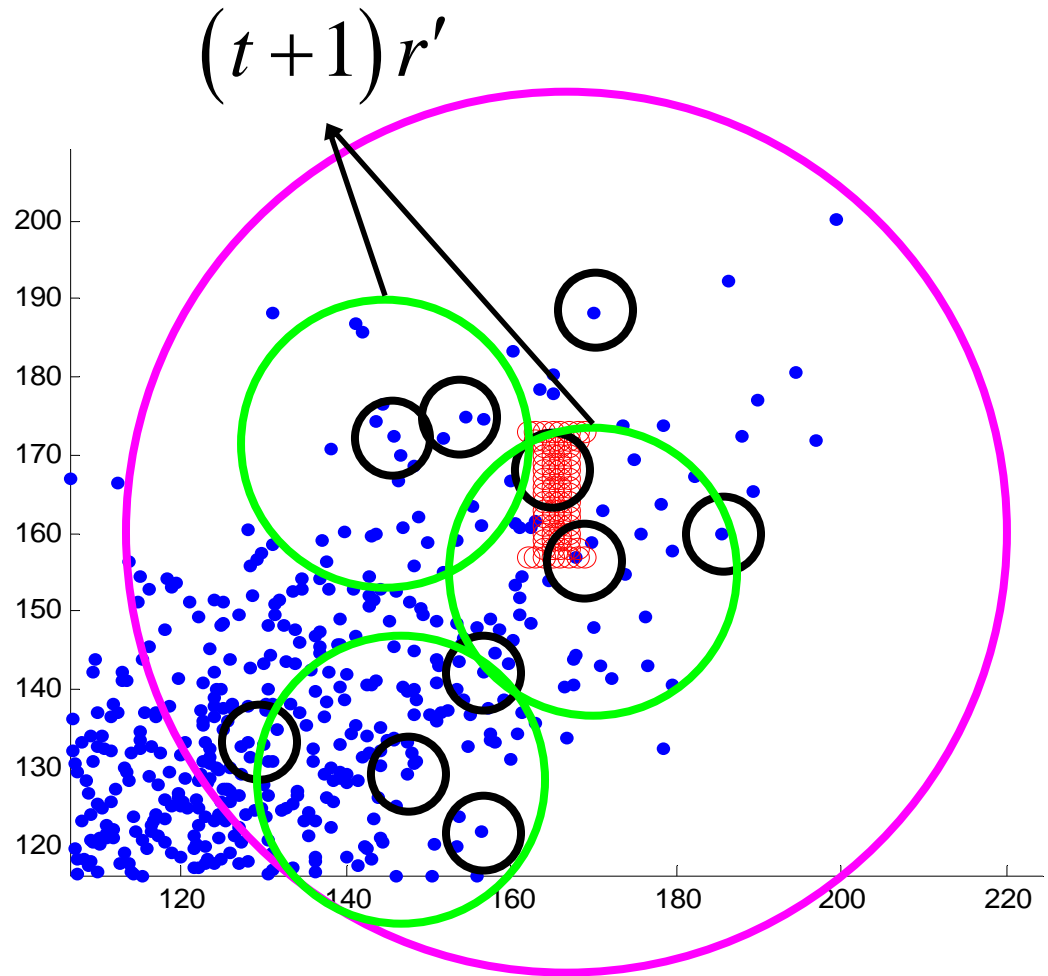


# NNDB: Pick the Next Candidate

Increase  $t$  by 1

$$S_i = \max_{x_j \in NN(x_i, (t+1)r')} (n_i - n_j)$$

Query  $x = \arg \max_{x_i \in S} S_i$



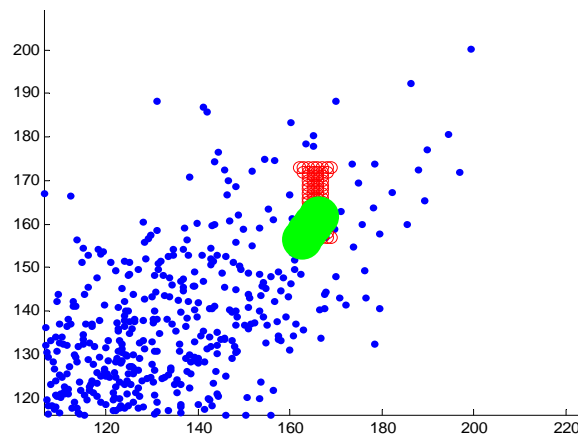
# Why NNDB Works

## □ Theoretically

- **Theorem 1** [He & Carbonell 2007]: under certain conditions, with high probability, after **a few** iteration steps, NNDB queries **at least one** example whose probability of coming from the minority class is **at least 1/3**

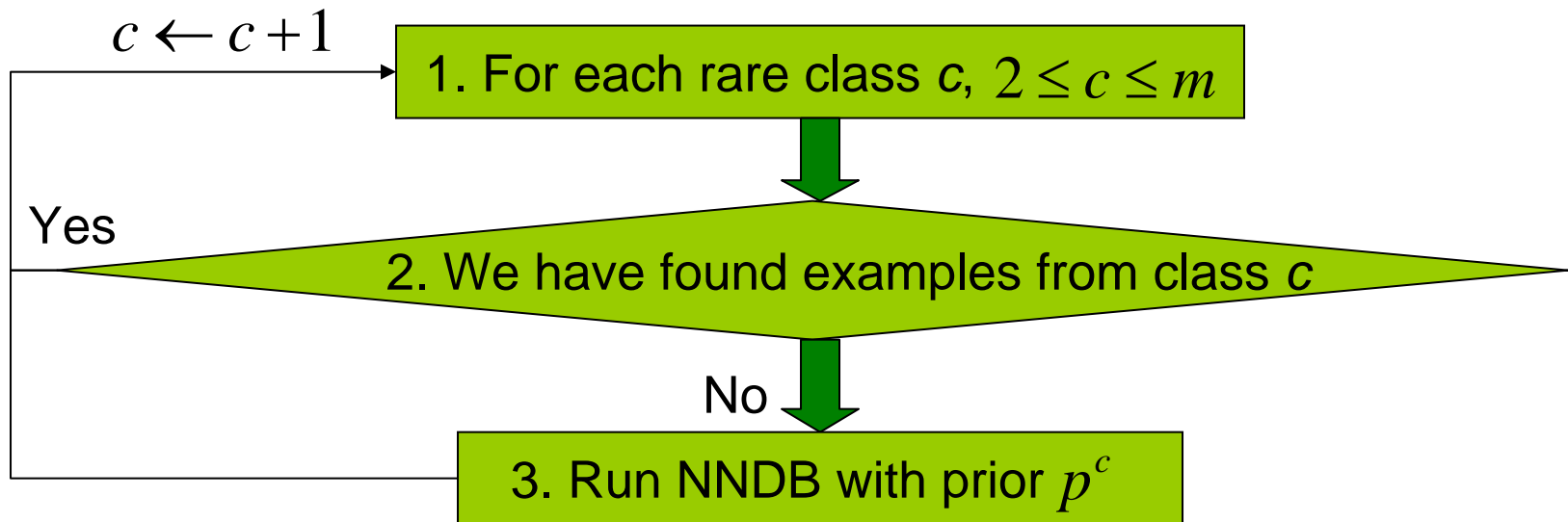
## □ Intuitively

- The score  $s_i$  measures the change in local density



# Multiple Classes: ALICE

- $m-1$  rare classes:  $p^2, \dots, p^m$
- One majority class:  $p^1 \square p^c, 2 \leq c \leq m$



# Why ALICE Works

---

## □ Theoretically

- **Theorem 2** [He & Carbonell 2008]: under certain conditions, with high probability, in **each outer loop** of ALICE, after **a few** iteration steps in NNDB, ALICE queries **at least one** example whose probability of coming from **one** minority class is **at least  $1/3$**

# Implementation Issues

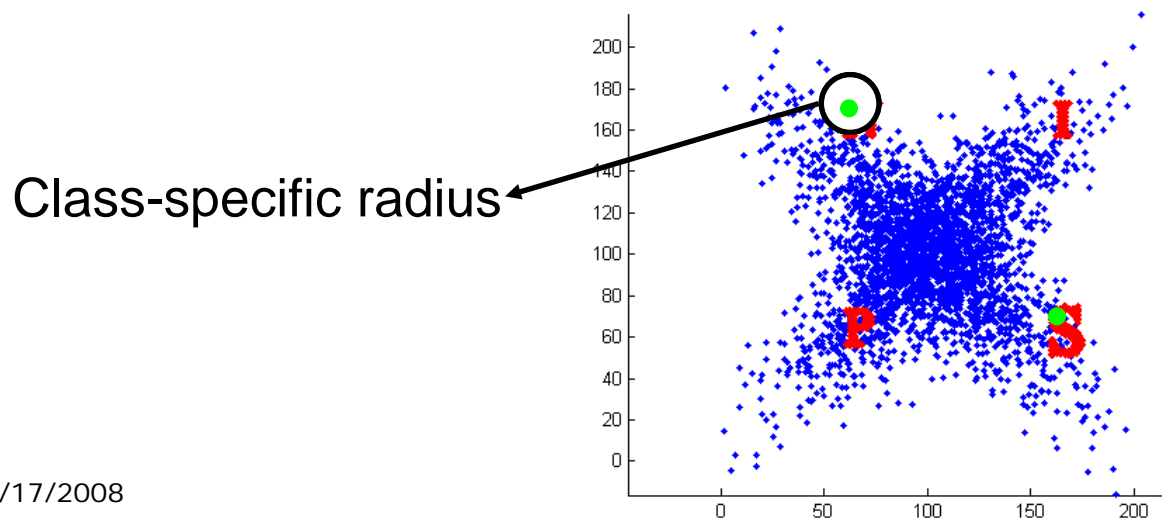
---

## □ ALICE

- Problem: repeatedly sampling from the same rare class

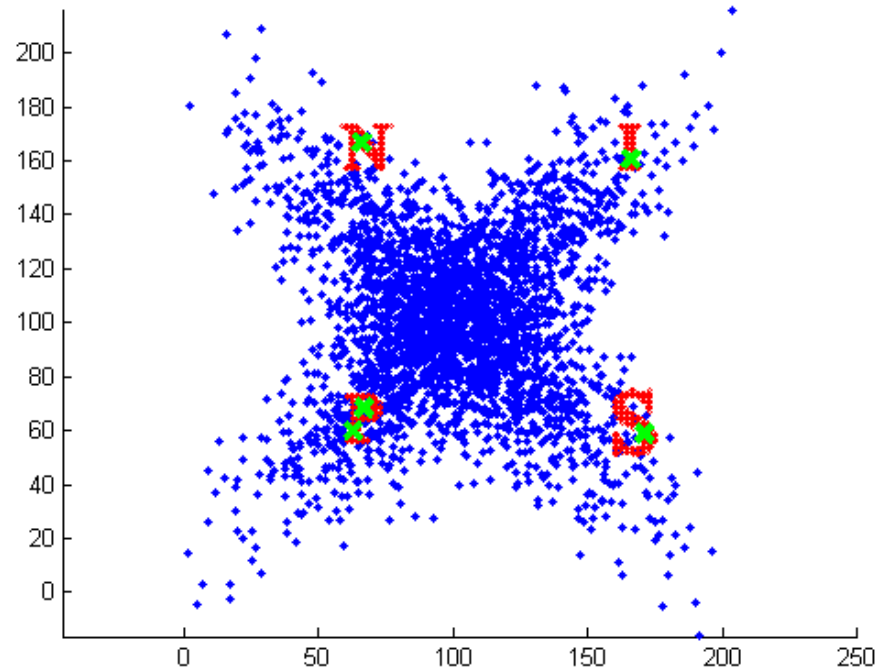
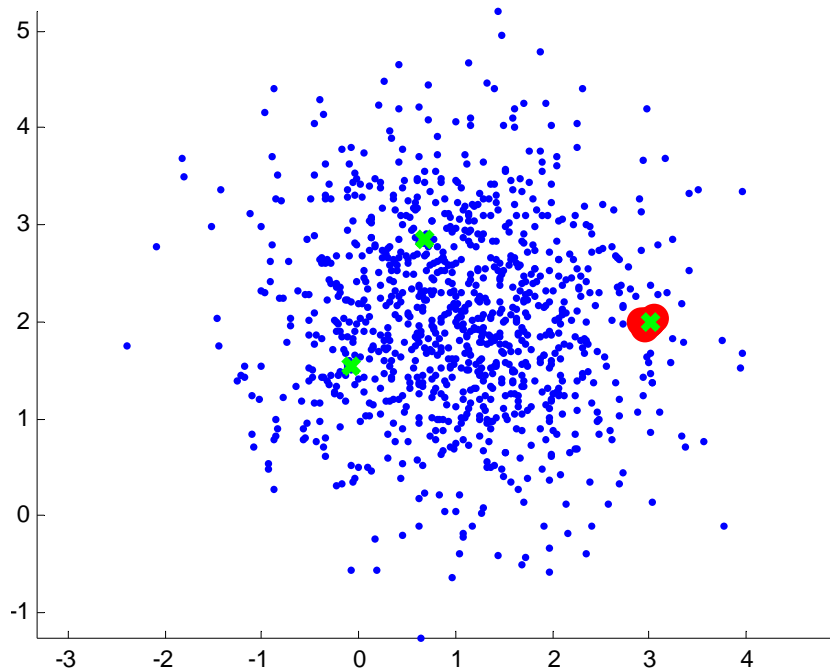
## □ MALICE

- Solution: relevance feedback



# Results on Synthetic Data Sets

---



# Summary of Real Data Sets

---

## □ Abalone

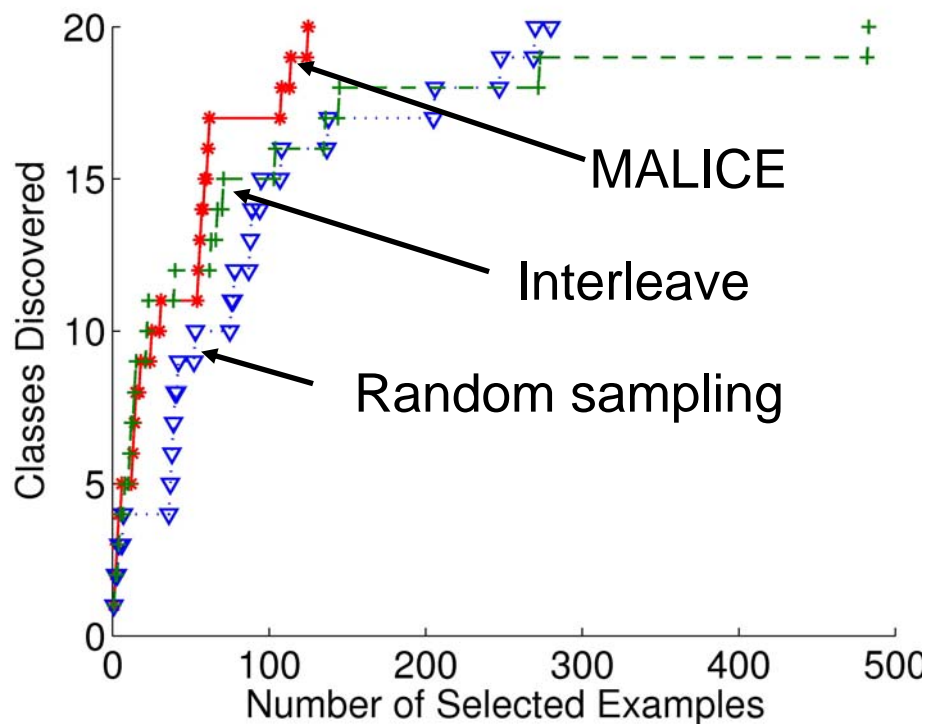
- 4177 examples
- 7-dimensional features
- 20 classes
- Largest class: 16.50%
- Smallest class: 0.34%

## □ Shuttle

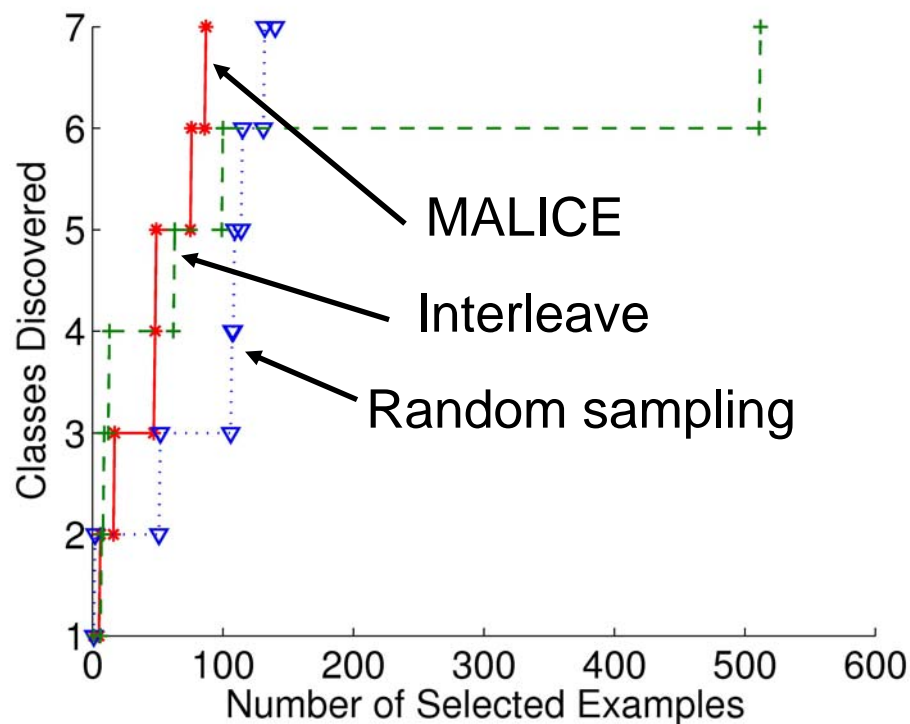
- 4515 examples
- 9-dimensional features
- 7 classes
- Largest class: 75.53%
- Smallest class: 0.13%

# Results on Real Data Sets

## Abalone



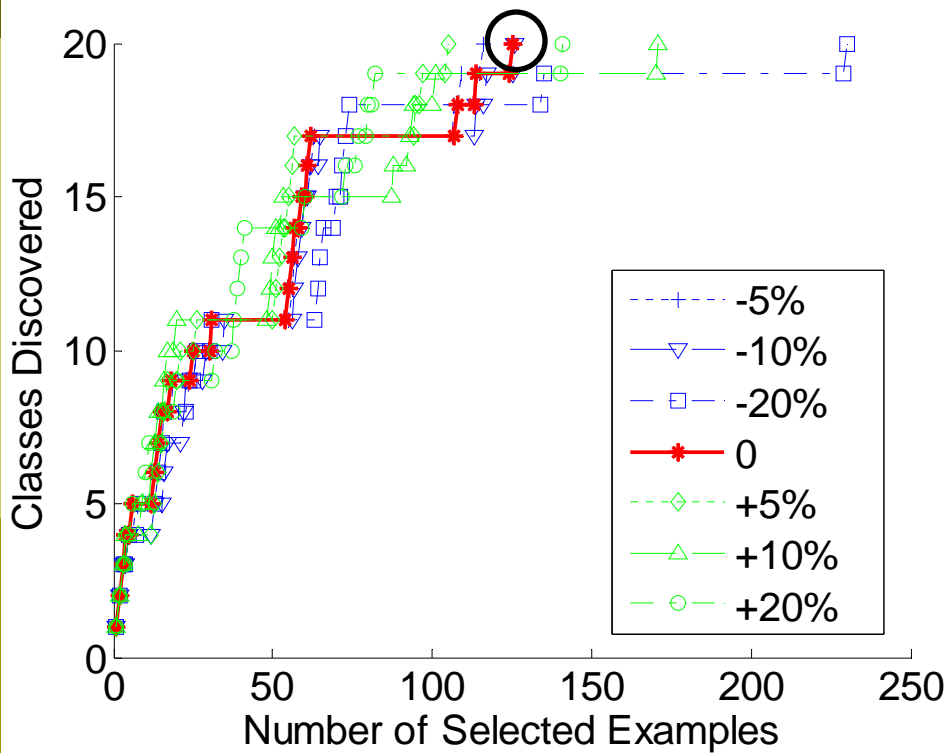
## Shuttle



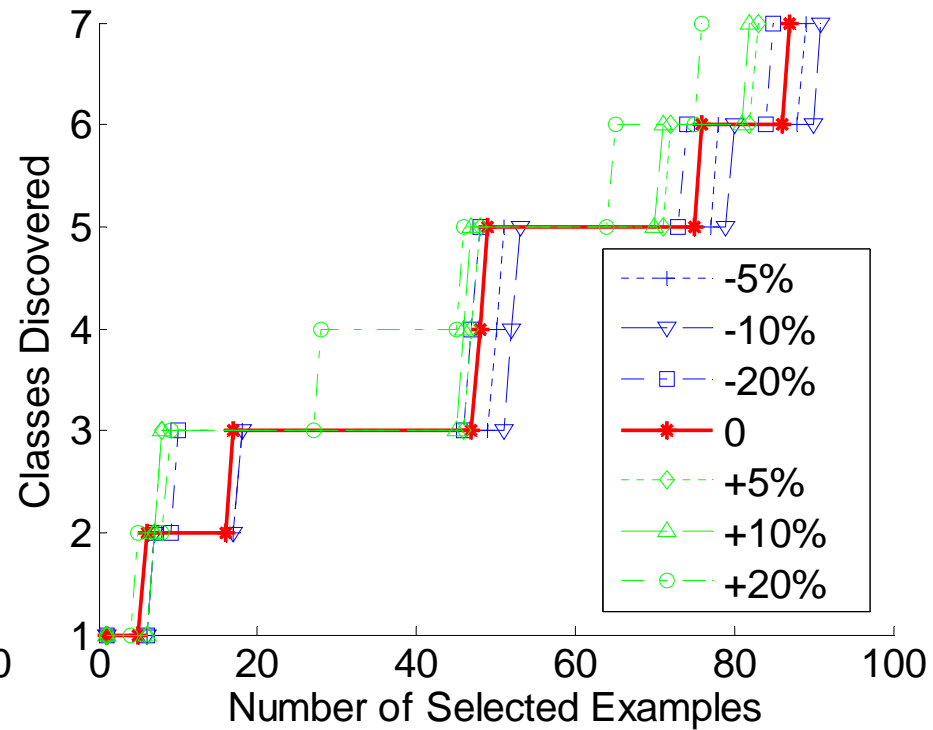


# Imprecise priors

## Abalone



## Shuttle



# Outline

---

- Problem definition
- Related work
- Rare category detection for spatial data
  - Prior-dependent rare category detection
  - Prior-free rare category detection
- Conclusion

# Overview of the Algorithm

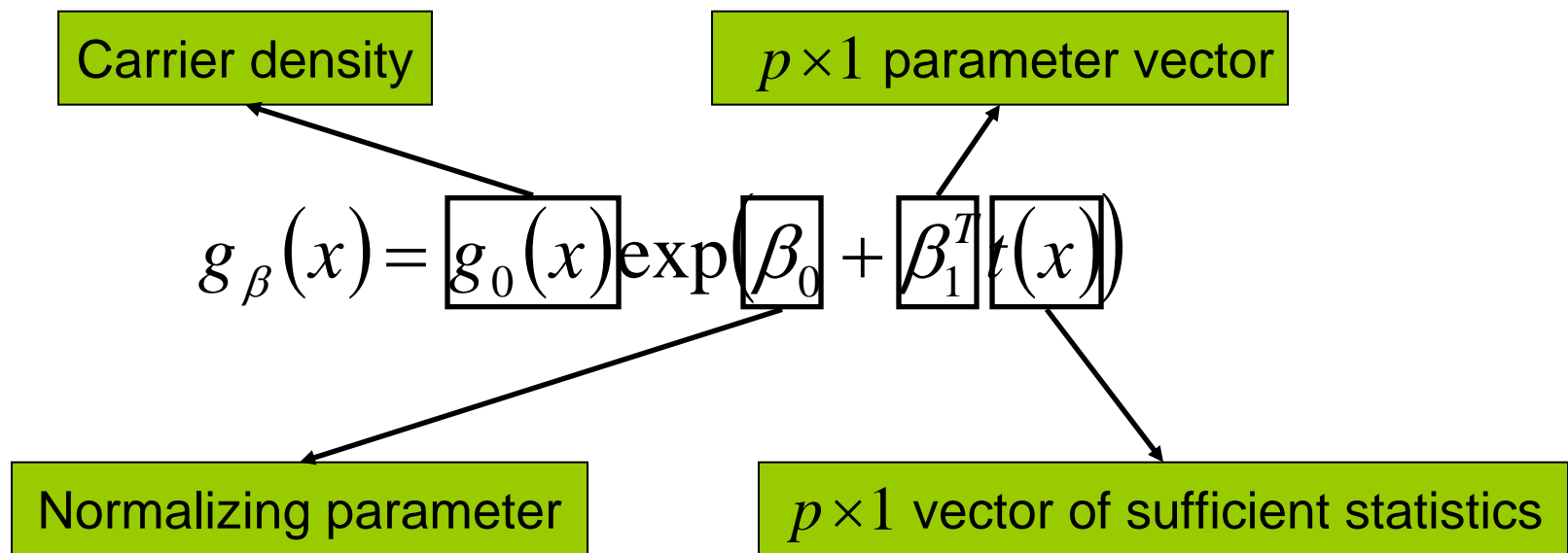
---

## □ Density-based method

- Methodology: specially designed exponential families
- Intuition: select examples according to the change in local density
- Difference from NNDB (ALICE): *NO* prior information needed

# Specially Designed Exponential Families [Efron & Tibshirani 1996]

- Favorable compromise between parametric and nonparametric density estimation
- Estimated density



# SEDER Algorithm

---

- Carrier density: kernel density estimator
- $t(x) = \left[ (x^1)^2, \dots, (x^d)^2 \right]^T$
- To decouple the estimation of different parameters
  - Decompose  $\beta_0 = \sum_{j=1}^d \beta_0^j$
  - Relax the constraint such that

$$\int_{x^j} \frac{1}{\sqrt{2\pi}\sigma^j} \exp\left(-\frac{(x^j - x_i^j)^2}{2(\sigma^j)^2}\right) \exp\left(\beta_{0i}^j + \beta_1^j (x^j)^2\right) dx^j = 1$$

# Parameter Estimation

- **Theorem 3** [To appear]: the maximum likelihood estimate  $\hat{\beta}_1^j$  and  $\hat{\beta}_{0i}^j$  of  $\beta_1^j$  and  $\beta_{0i}^j$  satisfy the following conditions:  $\forall j \in \{1, \dots, d\}$

$$\sum_{k=1}^n (x_k^j)^2 = \sum_{k=1}^n \frac{\sum_{i=1}^n \exp\left(\hat{\beta}_{0i}^j - \frac{(x_k^j - x_i^j)^2}{2(\sigma^j)^2}\right) E_i^j\left((x^j)^2\right)}{\sum_{i=1}^n \exp\left(\hat{\beta}_{0i}^j - \frac{(x_k^j - x_i^j)^2}{2(\sigma^j)^2}\right)}$$

where

$$E_i^j\left((x^j)^2\right) = \int_{x^j} (x^j)^2 \frac{1}{\sqrt{2\pi}\sigma^j} \exp\left(-\frac{(x^j - x_i^j)^2}{2(\sigma^j)^2}\right) \exp\left(\hat{\beta}_{0i}^j + \hat{\beta}_1^j (x^j)^2\right) dx^j$$

# Parameter Estimation cont.

□ Let  $\beta_1^j = \left(1 - \frac{1}{b^j}\right) \frac{1}{2(\sigma^j)^2}$   $\rightarrow b^j$ : positive parameter

□  $\forall j \in \{1, \dots, d\}: \hat{b}^j \approx \frac{-B + \sqrt{B^2 + 4AC}}{2A}$

where  $B = (\sigma^j)^2$ ,  $C = \frac{1}{n} \sum_{k=1}^n (x_k^j)^2$

$$A = \frac{1}{n} \sum_{k=1}^n \frac{\sum_{i=1}^n \exp\left(-\frac{(x_k^j - x_i^j)^2}{2(\sigma^j)^2}\right) (x_i^j)^2}{\sum_{i=1}^n \exp\left(-\frac{(x_k^j - x_i^j)^2}{2(\sigma^j)^2}\right)}$$

$\hat{b}^j \leq 1$   
in most cases

# Scoring Function

---

- The estimated density

$$\tilde{g}_b(x) = \frac{1}{n} \sum_{i=1}^n \prod_{j=1}^d \frac{1}{\sqrt{2\pi b^j} \sigma^j} \exp\left(-\frac{(x^j - b^j x_i^j)^2}{2(\sigma^j)^2 b^j}\right)$$

- Scoring function: norm of the gradient

$$s_k = \sqrt{\sum_{l=1}^d \frac{\left(\sum_{i=1}^n D_i(x_k)(x_k^l - b^l x_i^l)\right)^2}{\left((\sigma^l)^2 b^l\right)^2}}$$

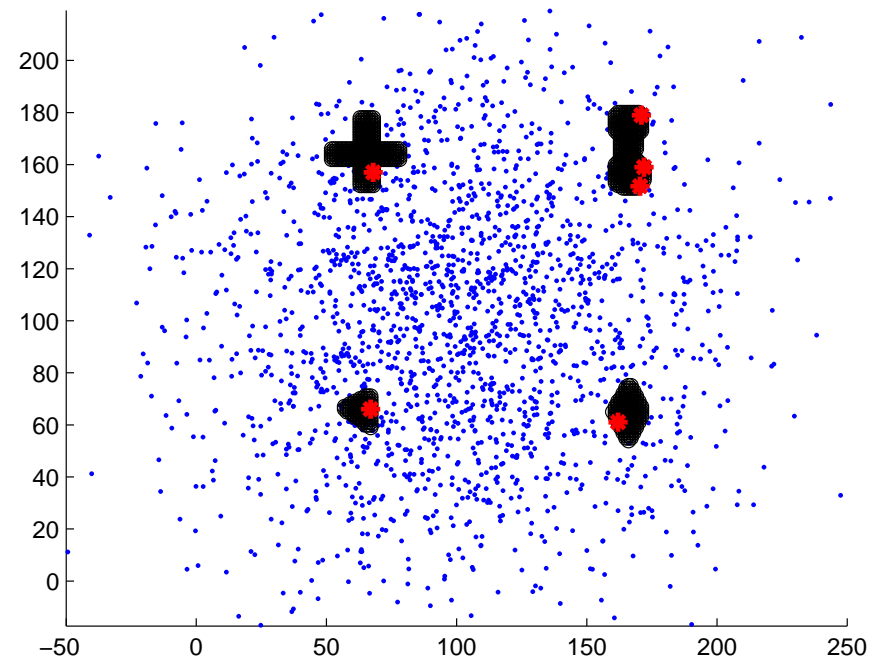
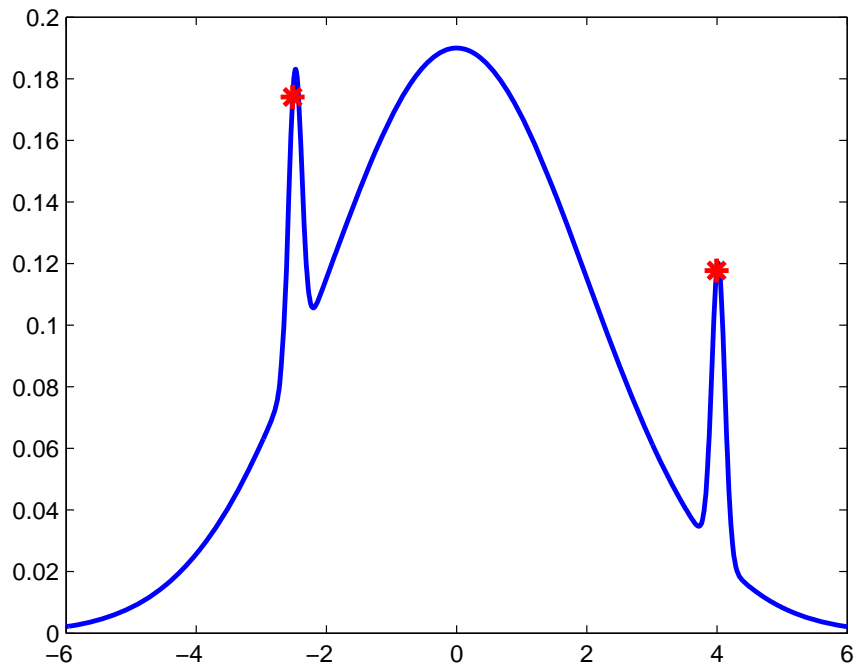
where

$$D_i(x) = \frac{1}{n} \prod_{j=1}^d \frac{1}{\sqrt{2\pi b^j} \sigma^j} \exp\left(-\frac{(x^j - b^j x_i^j)^2}{2(\sigma^j)^2 b^j}\right)$$



# Results on Synthetic Data Sets

---



# Summary of Real Data Sets

---

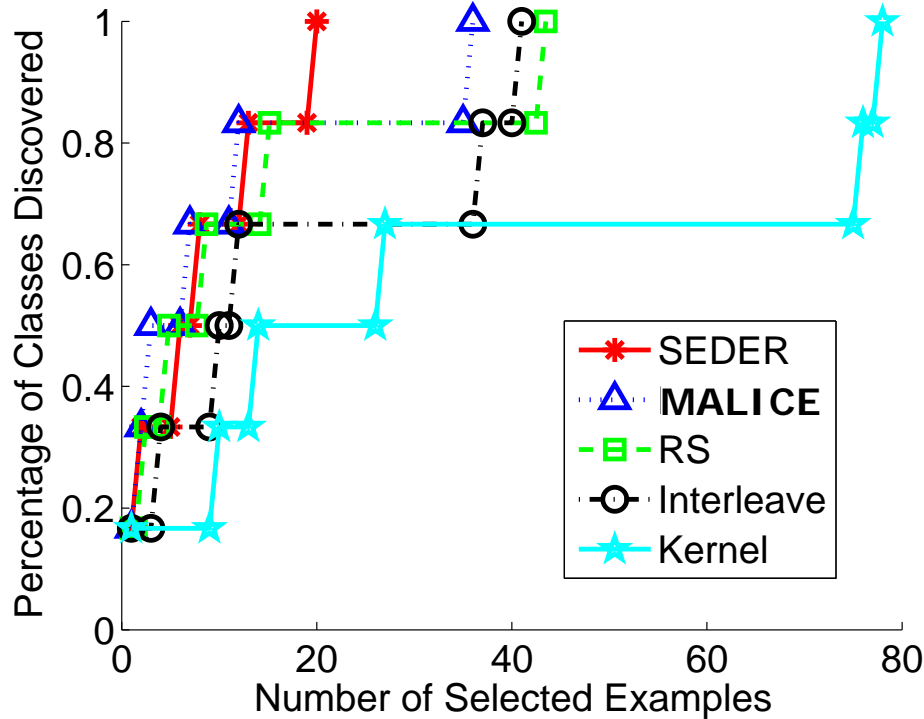
Data Set	$n$	$d$	$m$	Largest Class	Smallest Class
Ecoli	336	7	6	42.56%	2.68%
Glass	214	9	6	35.51%	4.21%
Page Blocks	5473	10	5	89.77%	0.51%
Abalone	4177	7	20	16.50%	0.34%
Shuttle	4515	9	7	75.53%	0.13%

**Moderately Skewed**

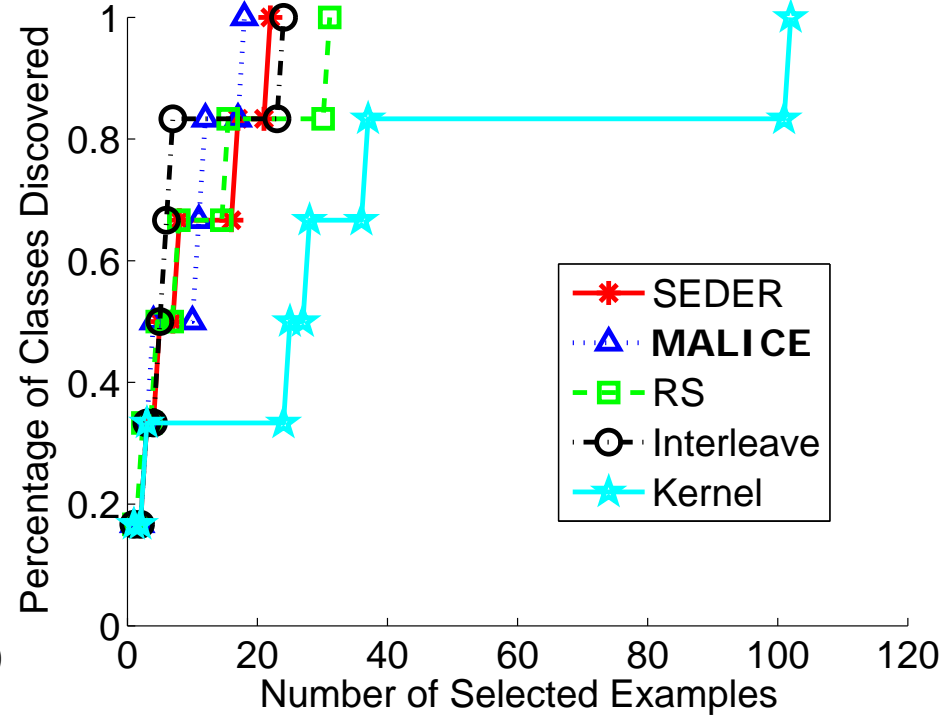
**Extremely Skewed**

# Moderately Skewed Data Sets

## Ecoli

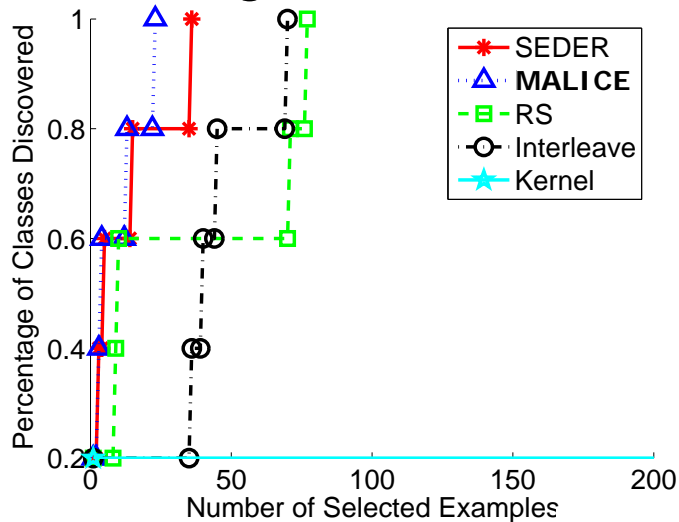


## Glass

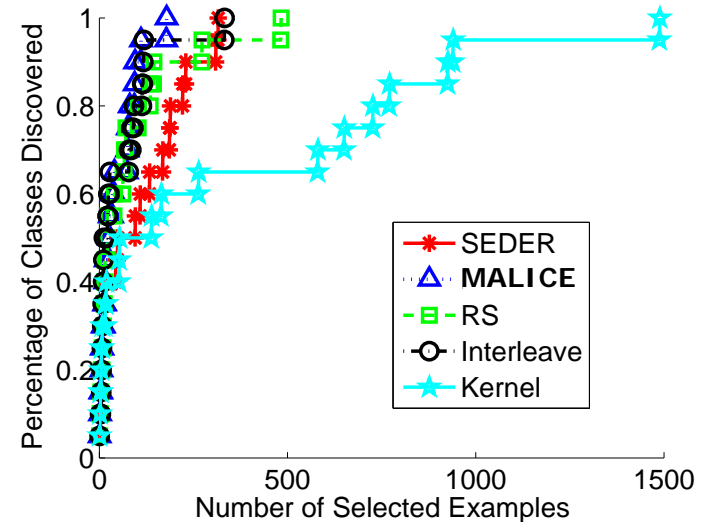


# Extremely Skewed Data Sets

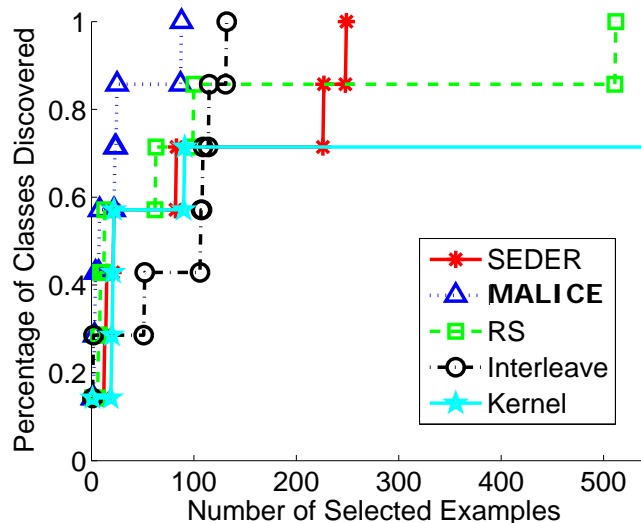
## Page Blocks



## Abalone



## Shuttle



# Conclusion

---

- Rare category detection
  - Open challenge
  - Lack of effective methods
- Nearest-neighbor-based methods
  - Prior-dependent
  - Local density differential sampling
- Density-based method
  - Prior-free
  - Specially designed exponential families

*Thank You!*

