

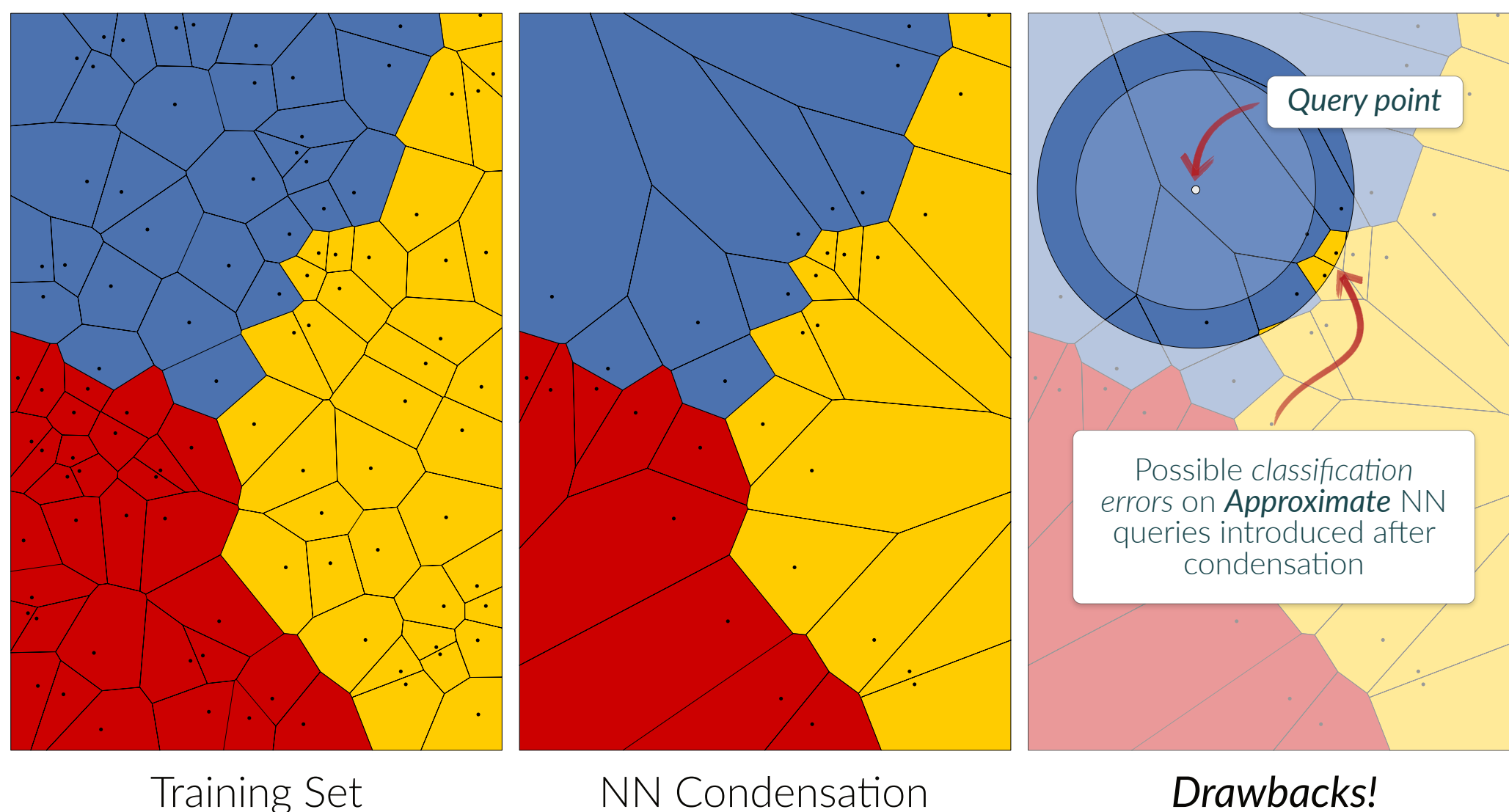
Condensation for the Approximate Nearest-Neighbor Rule

Alejandro Flores-Velazco & David Mount
 afloresv@cs.umd.edu mount@cs.umd.edu

University of Maryland
 Department of Computer Science

Consider the problem of Nearest-Neighbor Condensation

Consider a training set \mathbf{P} of labeled points from a metric space (\mathcal{X}, d) , to be used by the NN rule to classify new query points. The **NN condensation** problem deals with replacing the training set \mathbf{P} with a significantly smaller subset without affecting the classification accuracy under the NN rule.



The goal of *NN condensation* is to find a subset of \mathbf{P} s.t. under the NN rule, every point in \mathbf{P} is correctly classified. Such condensed set is called a **consistent** subset of \mathbf{P} .



The notion of consistency is defined on **exact** NN queries. What if we consider an **approximate** version of this? We propose the α -RSS algorithm to select such subsets.

Preliminaries

An **enemy** of a point $p \in \mathbf{P}$ is any point in \mathbf{P} of a different class. According to the metric d , the **nearest enemy** of p is denoted as $\text{NE}(p)$ and its **NE distance** as $d_{\text{ne}}(p)$.

A parameterized algorithm for NN condensation α -Relaxed Selective Subset

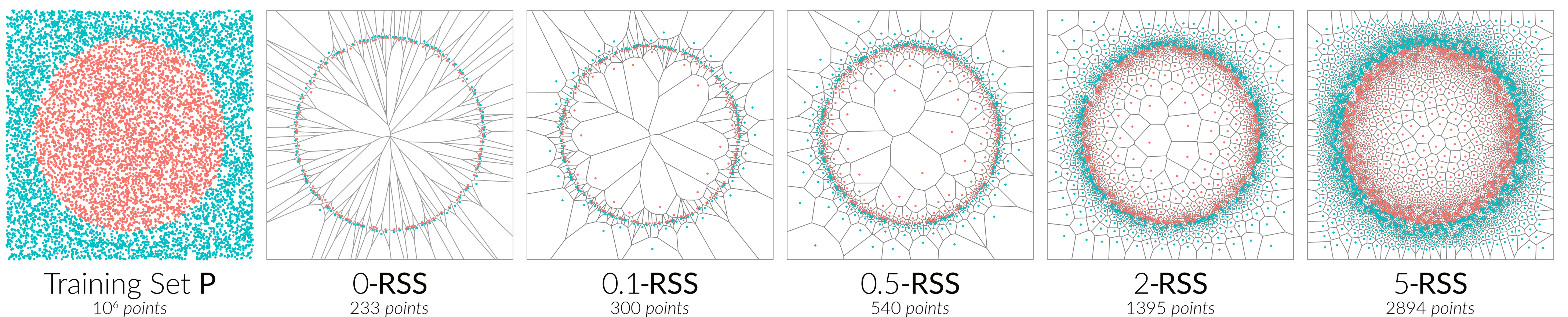
Input: initial training set \mathbf{P} , and value $\alpha \geq 0$
 Output: condensed training set $\alpha\text{-RSS} \subseteq \mathbf{P}$

- 1 Let $\{p_i\}_{i=1}^n$ be the points of \mathbf{P} sorted in increasing order of NE distance $d_{\text{ne}}(p_i)$
- 2 $\alpha\text{-RSS} \leftarrow \emptyset$
- 3 **For each** $p_i \in \mathbf{P}$ where $i = 1 \dots n$ **do**
- 4 **If** $\forall r \in \alpha\text{-RSS}, (1 + \alpha) \cdot d(p_i, r) \geq d_{\text{ne}}(p_i)$ **then**
- 5 $\alpha\text{-RSS} \leftarrow \alpha\text{-RSS} \cup \{p_i\}$
- 6 **Return** $\alpha\text{-RSS}$



The α -RSS algorithm computes a consistent subset of \mathbf{P} in $\mathcal{O}(n^2)$ worst-case time, and it is order independent.

- Order independence means the resulting subset is not determined by the order in which points are considered by the algorithm.
- Every point in \mathbf{P} is correctly classified by α -ANN queries on α -RSS
- 0 -RSS equals RSS and ∞ -RSS equals \mathbf{P} .



Results - Guarantees on the... Classification Accuracy of α -RSS

The **chromatic density** $\delta(q)$ of a query point $q \in \mathcal{X}$ is defined as

$$\delta(q, \mathbf{P}) = \frac{d_{\text{ne}}(q, \mathbf{P})}{d_{\text{nn}}(q, \mathbf{P})} - 1$$

Where $d_{\text{nn}}(q, \mathbf{P})$ is the NN distance of q .

Theorem

Consider two parameters $\varepsilon_1 \geq \varepsilon_2 > 0$ both upper bounded by some constant, and set $\alpha = \Omega(1/(\varepsilon_1 - \varepsilon_2))$. Now, if a query point $q \in \mathcal{X}$ has $\delta(q, \mathbf{P}) > \varepsilon_1$ then $\delta(q, \alpha\text{-RSS}) > \varepsilon_2$.

Theorem

The subset of \mathbf{P} selected by $(2/\varepsilon)$ -RSS is a weak ε -coreset for the chromatic nearest-neighbor of \mathbf{P} on query points $q \in \mathcal{X}$ where $\delta(q, \mathbf{P}) > \varepsilon$.



Sufficient conditions for correct classification after NN condensation using α -RSS.

Results - Upper-bounds for... The size of α -RSS

Let Δ be the **spread** of \mathbf{P} (i.e., the ratio between the *largest* and *smallest* pairwise distances in \mathbf{P}) and κ the **number of NE** points in \mathbf{P} .

Theorem

Consider (\mathcal{X}, d) to be a **doubling space** with **doubling dimension** $\text{ddim}(\mathcal{X})$, then:

$$|\alpha\text{-RSS}| = \mathcal{O}(\kappa \alpha^{\text{ddim}(\mathcal{X})+1} \log \Delta)$$

Theorem

Consider (\mathcal{X}, d) to be the **Euclidean space**, s.t. $\mathbf{P} \subset \mathbb{R}^d$, then:

$$|\alpha\text{-RSS}| = \mathcal{O}(\kappa \alpha^{d-1} \log \Delta)$$