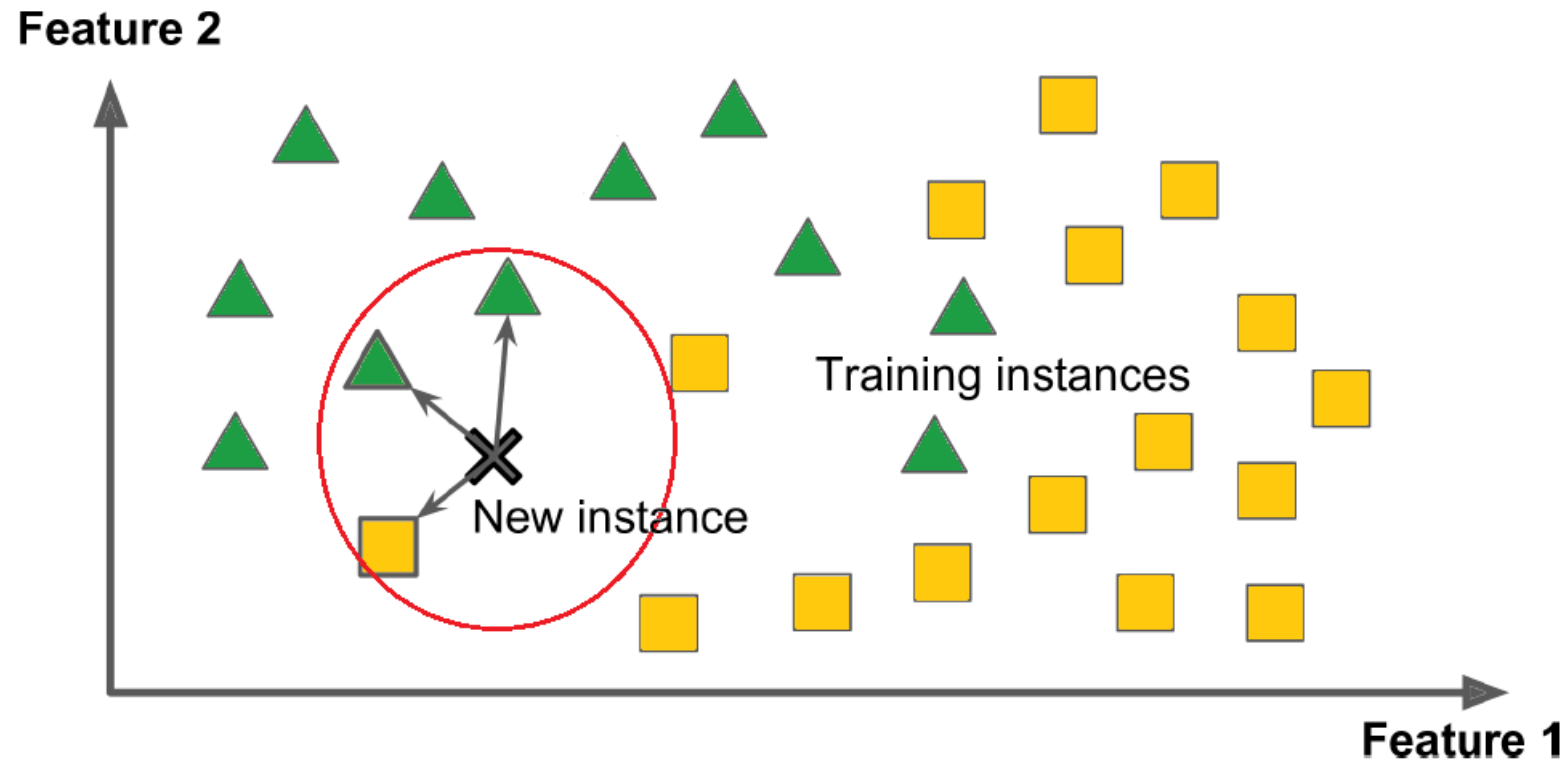


The impact of distance metrics on the performance of KNN algorithm

Angel Emanuel Rodríguez Román

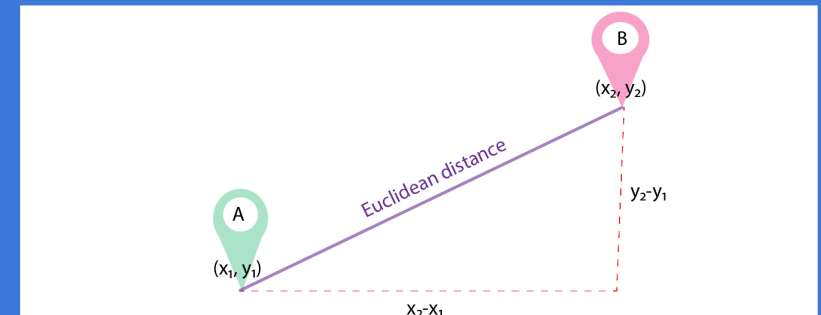
The objective of this work is to evaluate the impact of different distance metrics (such as Euclidean, Manhattan, city block, etc) on the performance of the KNN algorithm with different values for K , using both balanced and imbalanced datasets.

KNN Algorithm





Measure vs Metric



What is a distance metric?

It's a function that defines a way to measure the distance between two points in an space.

It **must** satisfy 3 conditions:

- Non-negativity: $d(x, y) \geq 0$
- Identity: $d(x, y) = 0$ if and only if $x == y$
- Symmetry: $d(x, y) = d(y, x)$
- Triangle Inequality: $d(x, y) + d(y, z) \geq d(x, z)$

Manhattan distance

(Citi block or taxi cab)

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

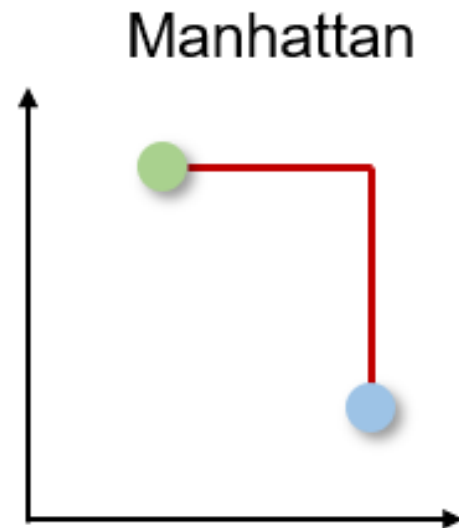
E.G.

$$x = \langle -3, 3 \rangle \quad y = \langle 1, -5 \rangle$$

$$d(x, y) = |-3 - 1| + |3 - (-5)|$$

$$d(x, y) = |-4| + |8|$$

$$d(x, y) = 12$$



Euclidean distance

$$d(x, y) = \left(\sum_1^n |x_i - y_i|^2 \right)^{1/2}$$

E.G.

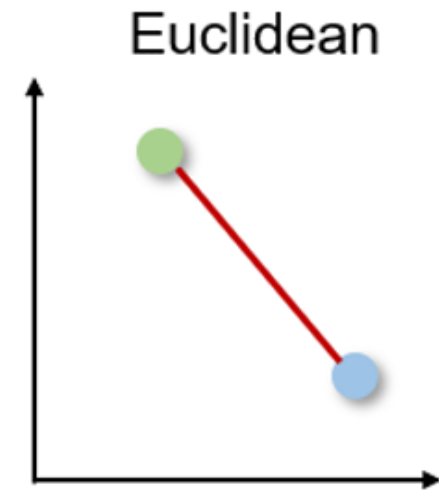
$$x = \langle -3, 3 \rangle \quad y = \langle 1, -5 \rangle$$

$$d(x, y) = [(-3 - 1)^2 + (3 - (-5))^2]^{1/2}$$

$$d(x, y) = [(-4)^2 + (8)^2]^{1/2}$$

$$d(x, y) = (80)^{1/2}$$

$$d(x, y) = 8.94$$



Minkowski

$$d(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

- Manhattan (City block $p=1$)
- Euclidean ($p=2$)
- Chebyshev ($p \rightarrow \infty$)

Chebychev distance (Chessboard)

$$d(x, y) = \sum_{i=1}^n \max |x_i - y_i|$$

E.G.

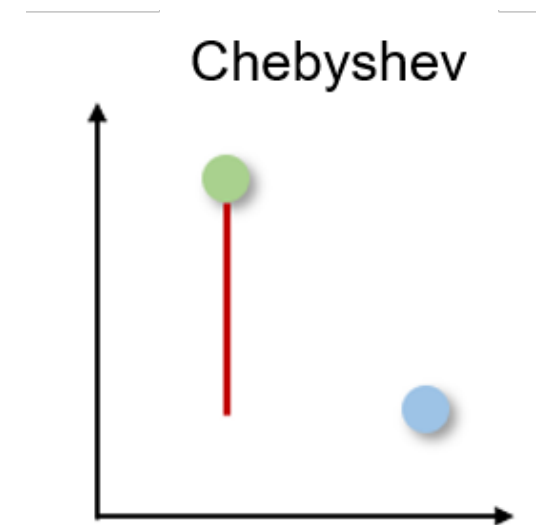
$$x = \langle -3, 3 \rangle$$

$$y = \langle 1, -5 \rangle$$

$$d(x, y) = \max (|-4|, |8|)$$

$$d(x, y) = \max (4, 8)$$

$$d(x, y) = 8$$



Canberra distance

$$d(x, y) = \sum_{i=1}^n \frac{|x_i - y_i|}{|x_i| + |y_i|}$$

The Canberra distance is a weighted version of the Manhattan distance, introduced and refined 1967 by Lance, Williams and Adkins. It is often used for data scattered around an origin, as it is biased for measures around the origin and very sensitive for values close to zero.

Braycurtis distance (Sorensen)

$$d(x, y) = \frac{\sum |x_i - y_i|}{\sum |x_i + y_i|}$$

The Braycurtis distance views the space as grid, similar to the city block distance. If both points are in zero coordinates, the Braycurtis distance is undefined. It's commonly used in botany, ecology and in general, in environmental science fields.

Cosine distance

$$1 - \frac{u \cdot v}{\|u\|_2 \|v\|_2}$$

The cosine distance contains the dot product scaled by the product of the Euclidean distances from the origin. It represents the angular distance of two vectors while ignoring their scale.

This metric is widely used in NLP. In applications such as comparing similarity between words.

Other distance metrics

- Hamming Distance (for **integers** and boolean)
- Dice (boolean)
- Jaccard (boolean)
- Rogerstanimoto (boolean)
- Russellrao (boolean)
- Sokalmichener (boolean)
- Jensenshannon(boolean)
- Mahalanobis $\sqrt{(u - v)V^{-1}(u - v)^T}$

Data

Dataset	Cardinality	Attributes	IR	About	Balanced
Nutt	28	1070	1	This dataset contains attributes related to brain tumors. The goal is to predict between classic glioblastomas and non classic.	Yes
Apendicitis	106	8	4.04	The data represents 7 medical measures taken over 106 patients on which the class label represents if the patient has appendicitis (class label 1) or not (class label 0).	No
Sonar	208	61	1.14	This dataset contains information about rocks and mines (Energy within a particular frequency band).	Yes
Heart	270	14	1.25	This dataset contains valuable information to detect the absence (1) or presence (2) of heart disease.	Yes
Haberman	306	4	2.77	This data set contains cases from a study that was conducted between 1958 and 1970 at the University of Chicago's Billings Hospital on the survival of patients who had undergone surgery for breast cancer. The task is to determine if the patient survived 5 years or longer (positive) or if the patient died within 5 year (negative).	No

Dataset	Cardinality	Attributes	IR	About	Balanced
Ionosphere	351	34	1.78	This dataset provides data from a system that consists of a phased array of 16 high-frequency antennas with a total transmitted power on the order of 6.4 kilowatts. The task is to determine if a given signal is Good (g) or Bad (b).	No
Monk-2	432	7	1.11	The MONK's problems are a collection of three binary artificial classification problems (MONK-1, MONK-2 and MONK-3) over a six-attribute discrete domain. Each problem involves learning a binary function defined over this domain, from a sample of training examples that belong to class 0 or class 1.	Yes
WDBC (Breast Cancer Wisconsin Diagnostic)	569	31	1.68	This database contains 30 features computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image. The task is to determine if a found tumor is benign or malignant (M = malignant, B = benign).	No
Electricity	2400	14	1.47	This data was collected from the Australian New South Wales Electricity Market. The objective is to predict the change of the price (UP or DOWN) .	Yes
Spam	4597	58	1.53	This database contains information about 4597 e-mail messages. The task is to determine whether a given email is spam (class 1) or not (class 2).	No

Experimental Framework

Distances : {“euclidean”, “manhattan”, “chebyshev”, “canberra”, “braycurtis”, “cosine”}

K : {1, 3, 5}

Algorithms: {*Distances* \times *K*}

Validation Method: Leave One Out Cross Validation (LOO)

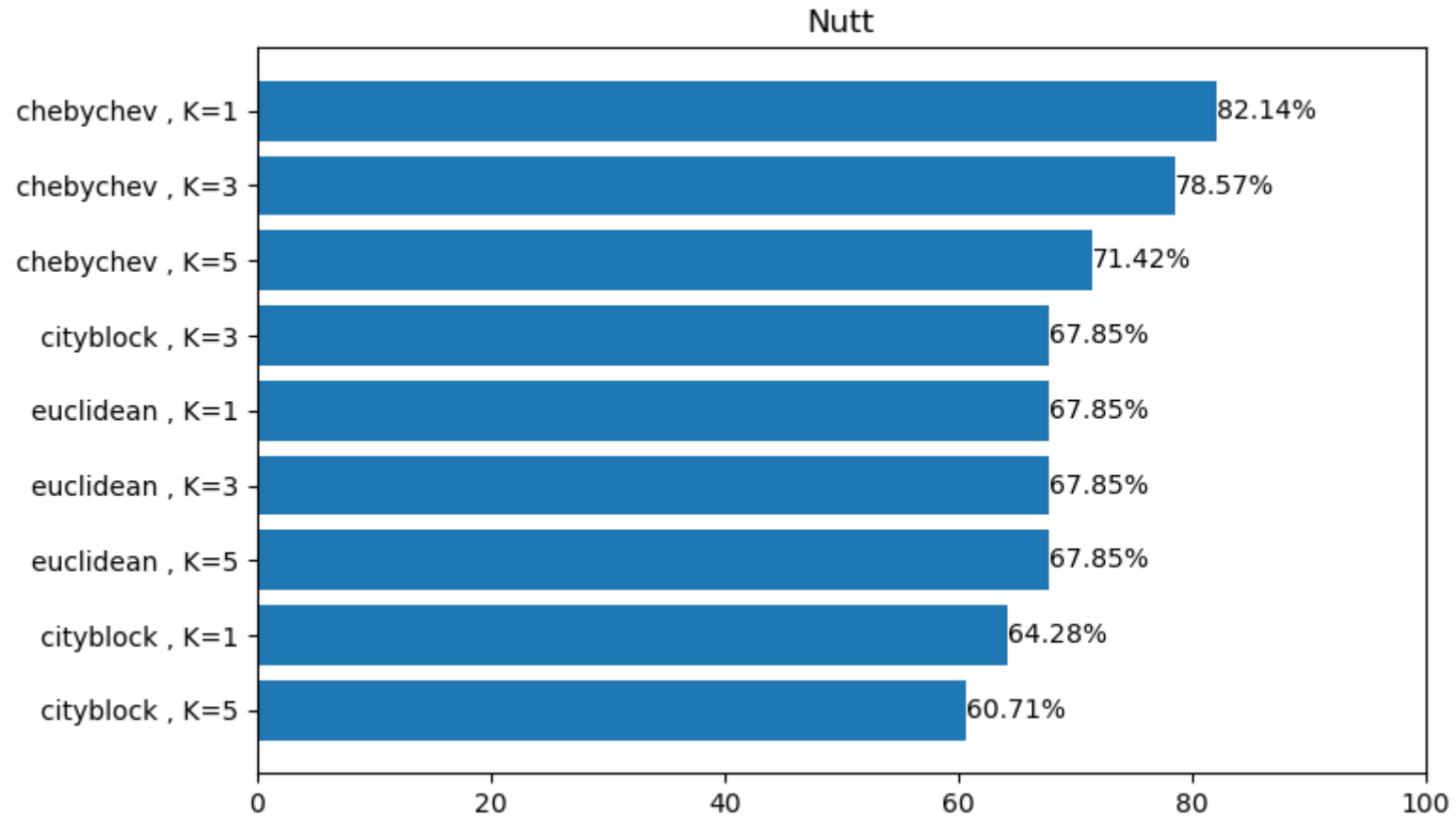
Measure: Balanced Accuracy (BA)

$$BA = \frac{Sensitivity + Specificity}{2}$$

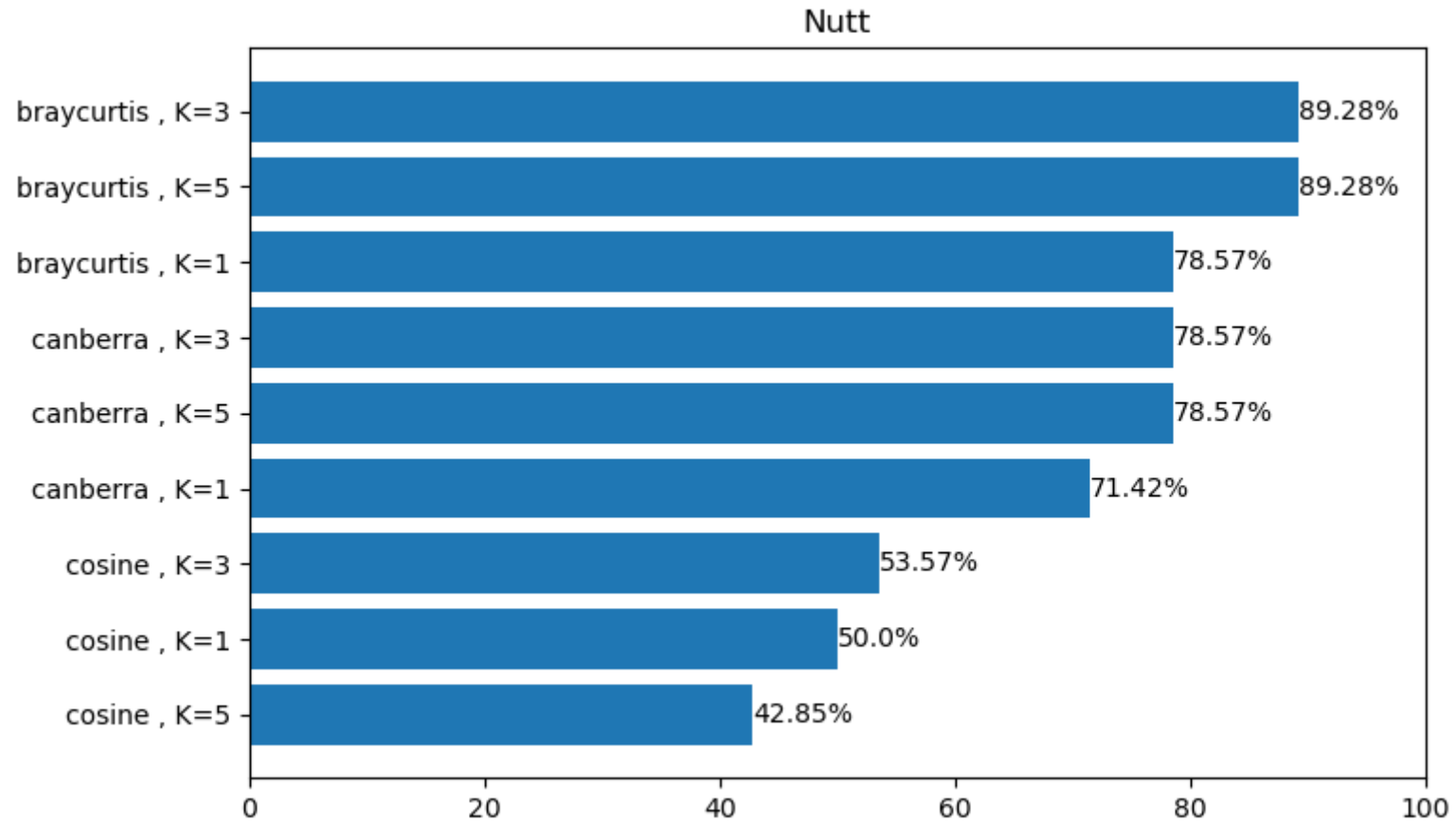


Results

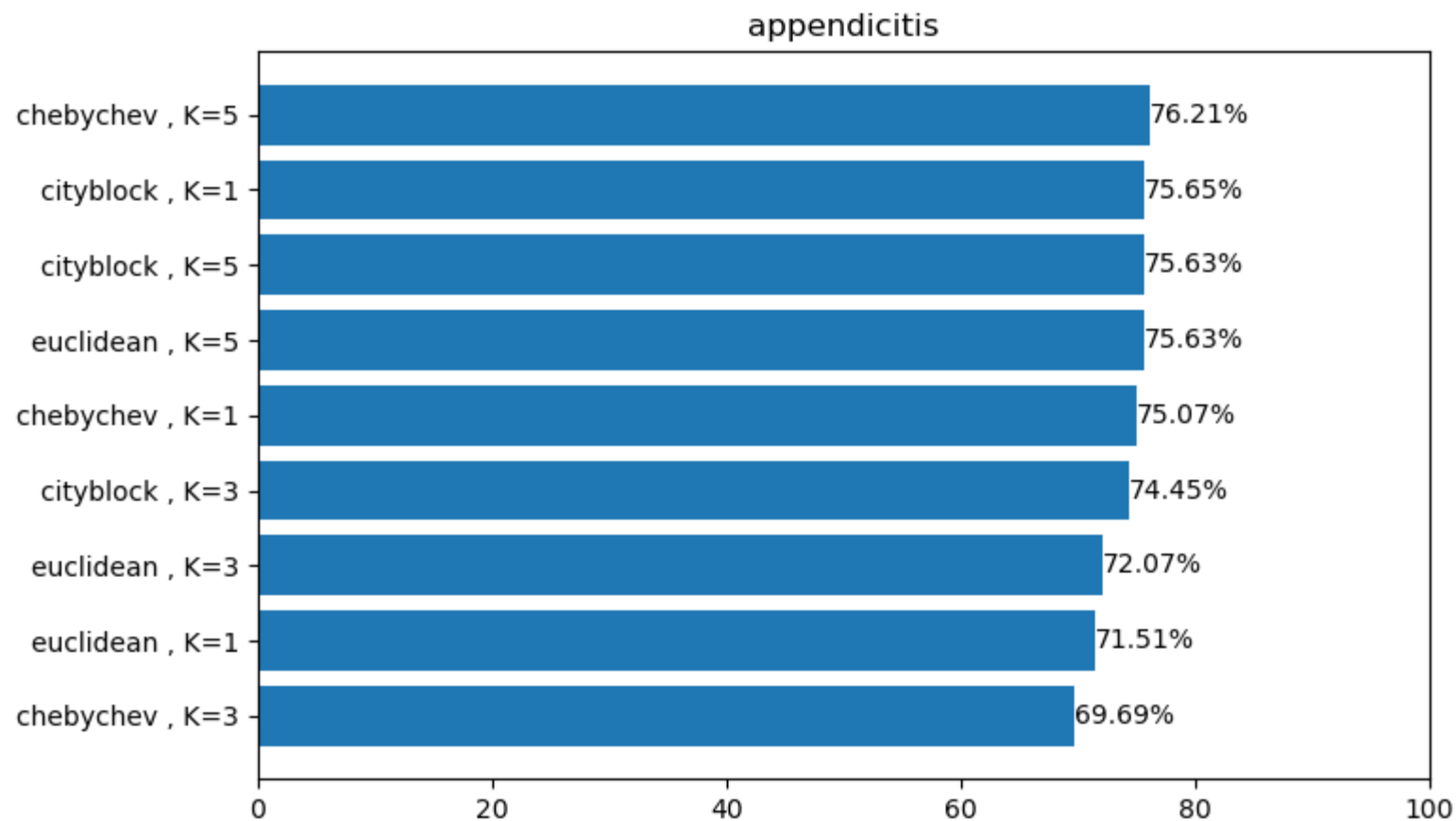
Nutt



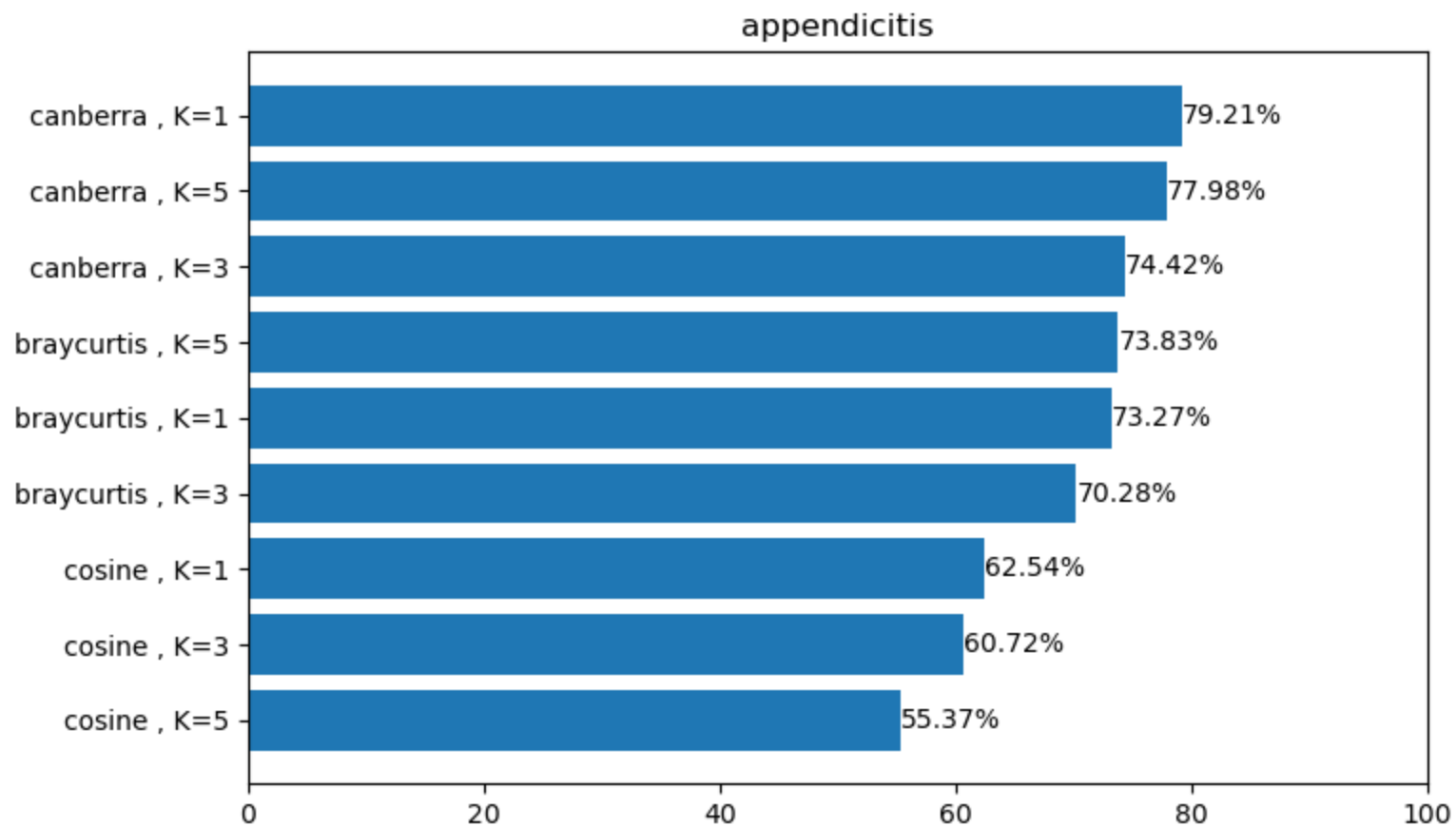
Nutt



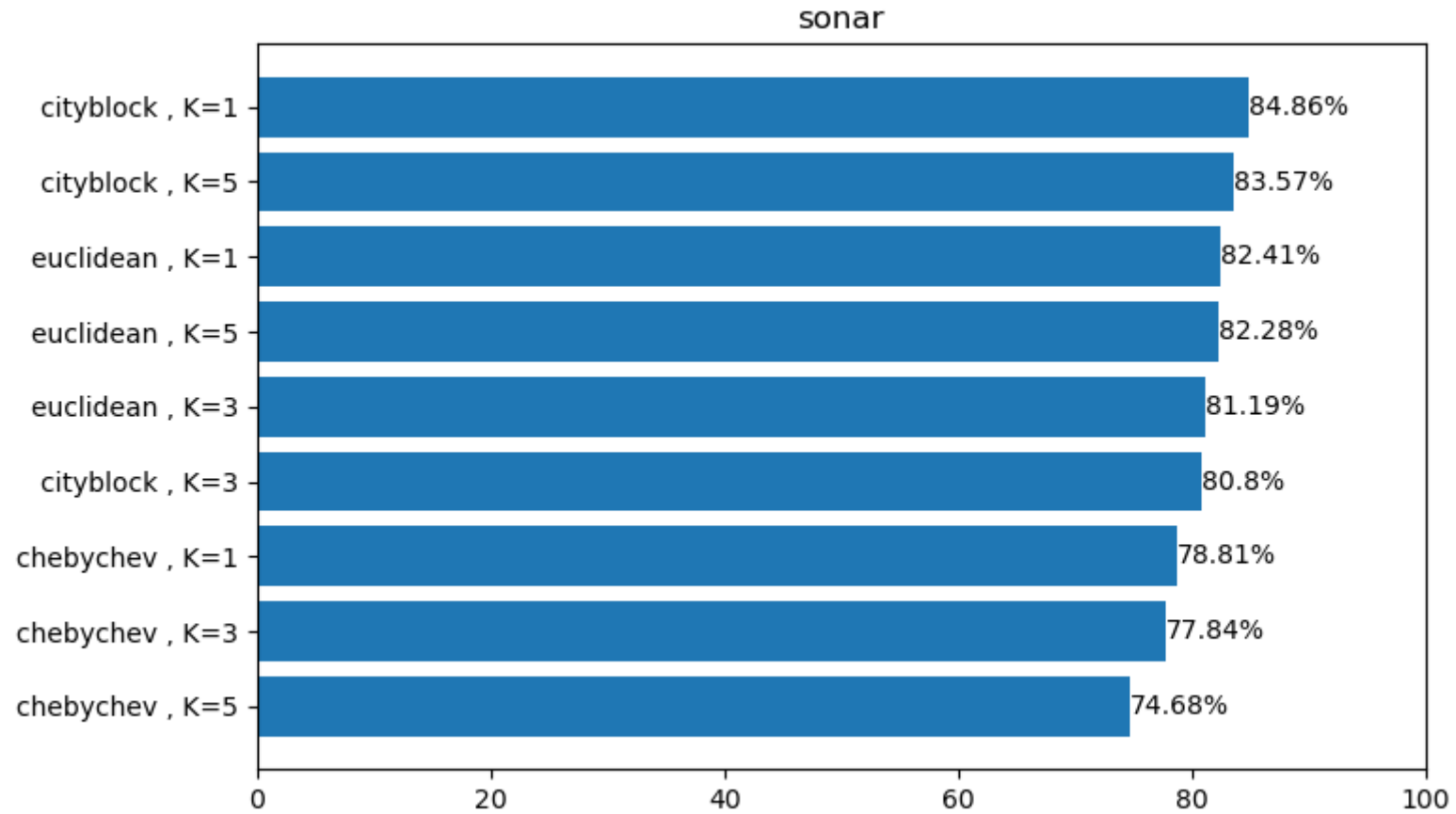
Apendicitis



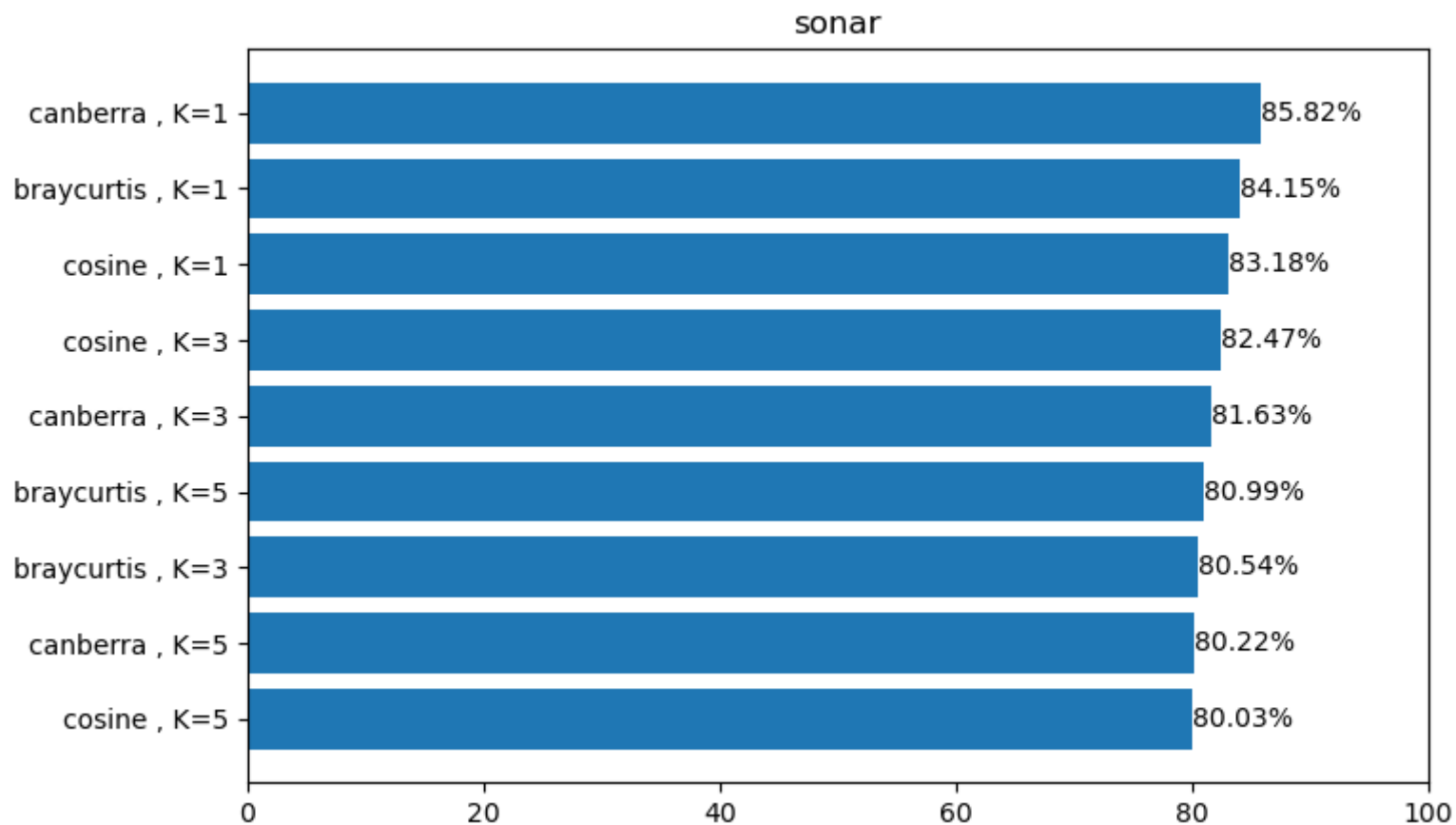
Apendicitis



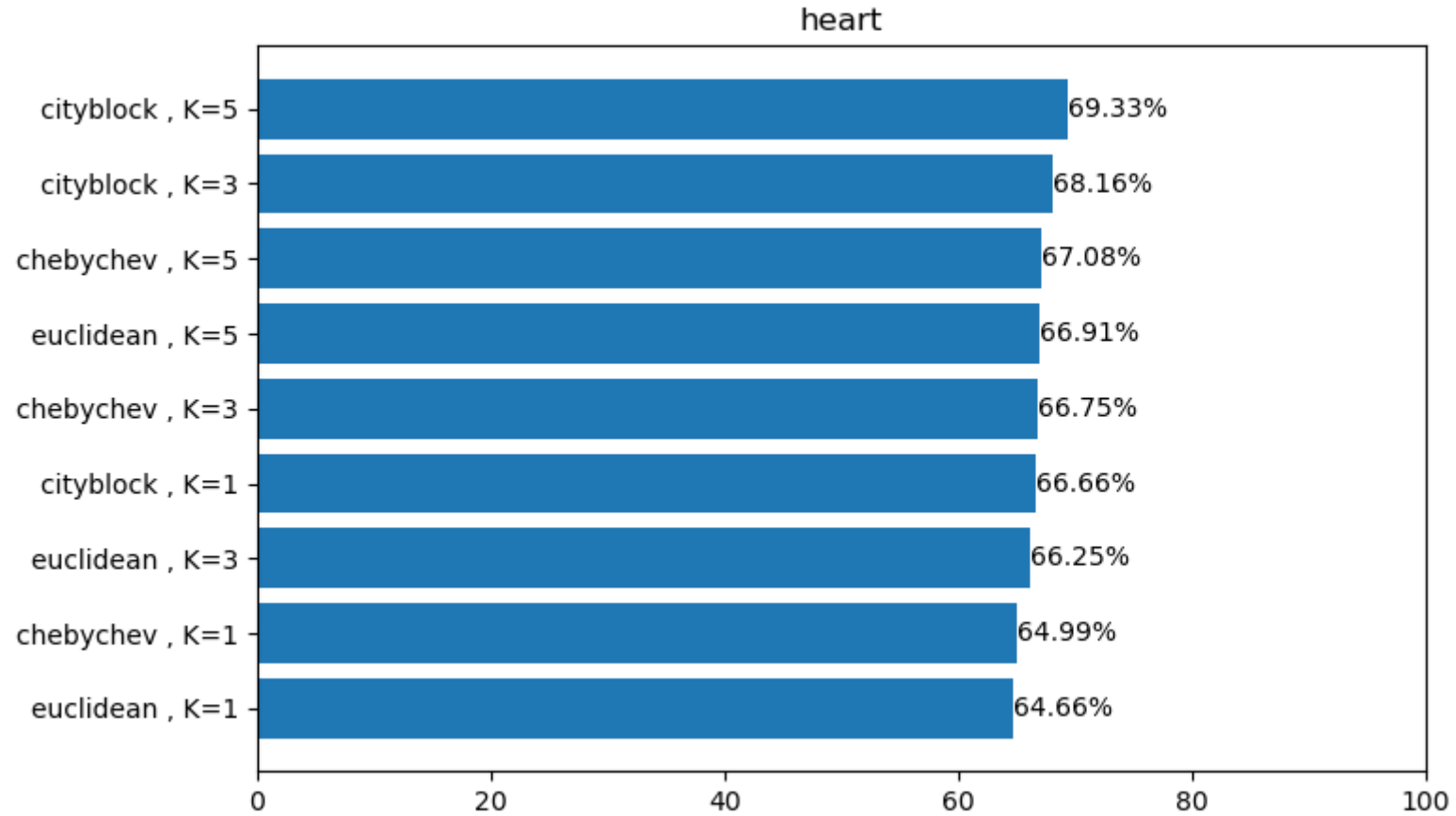
Sonar



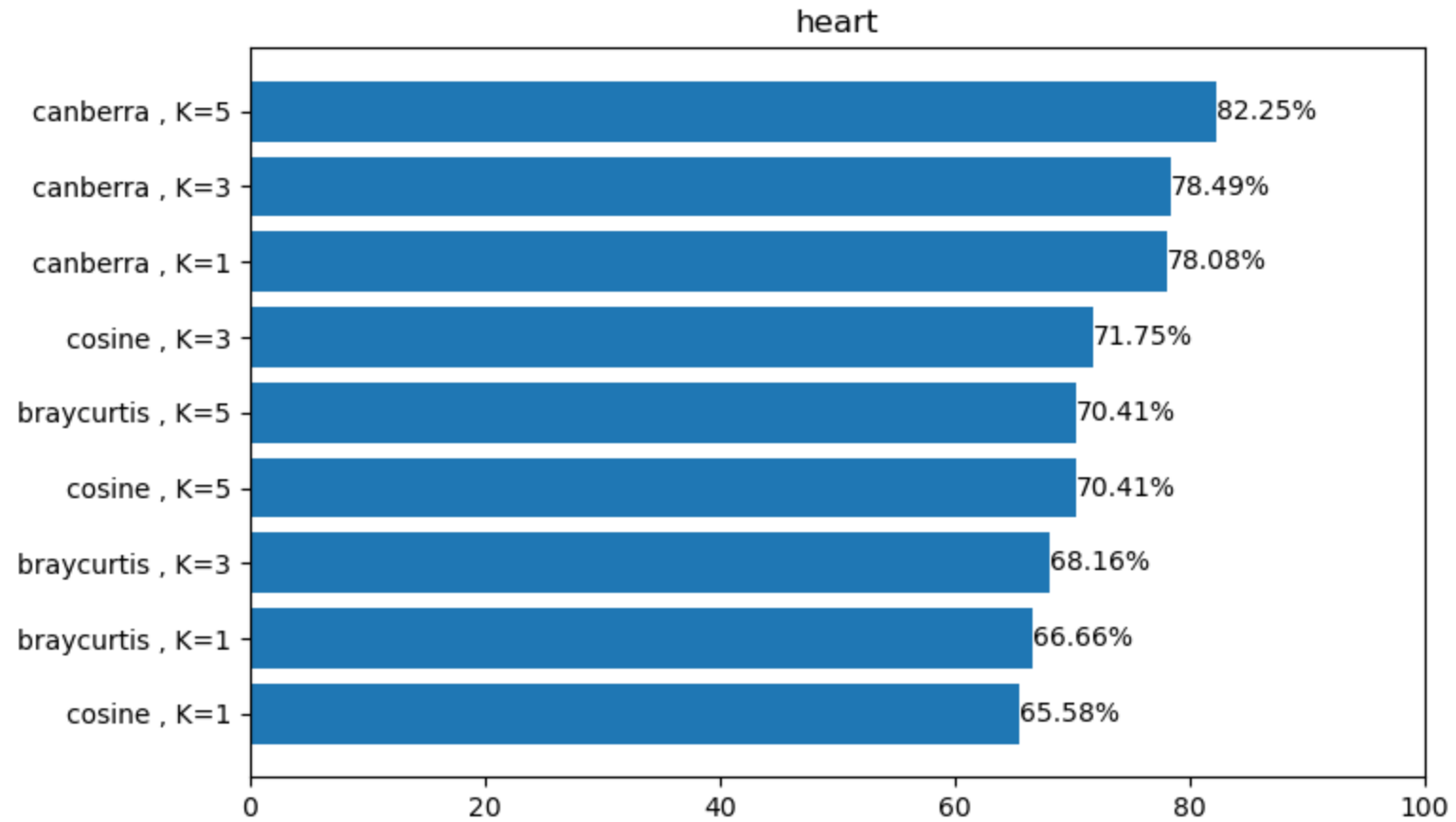
Sonar



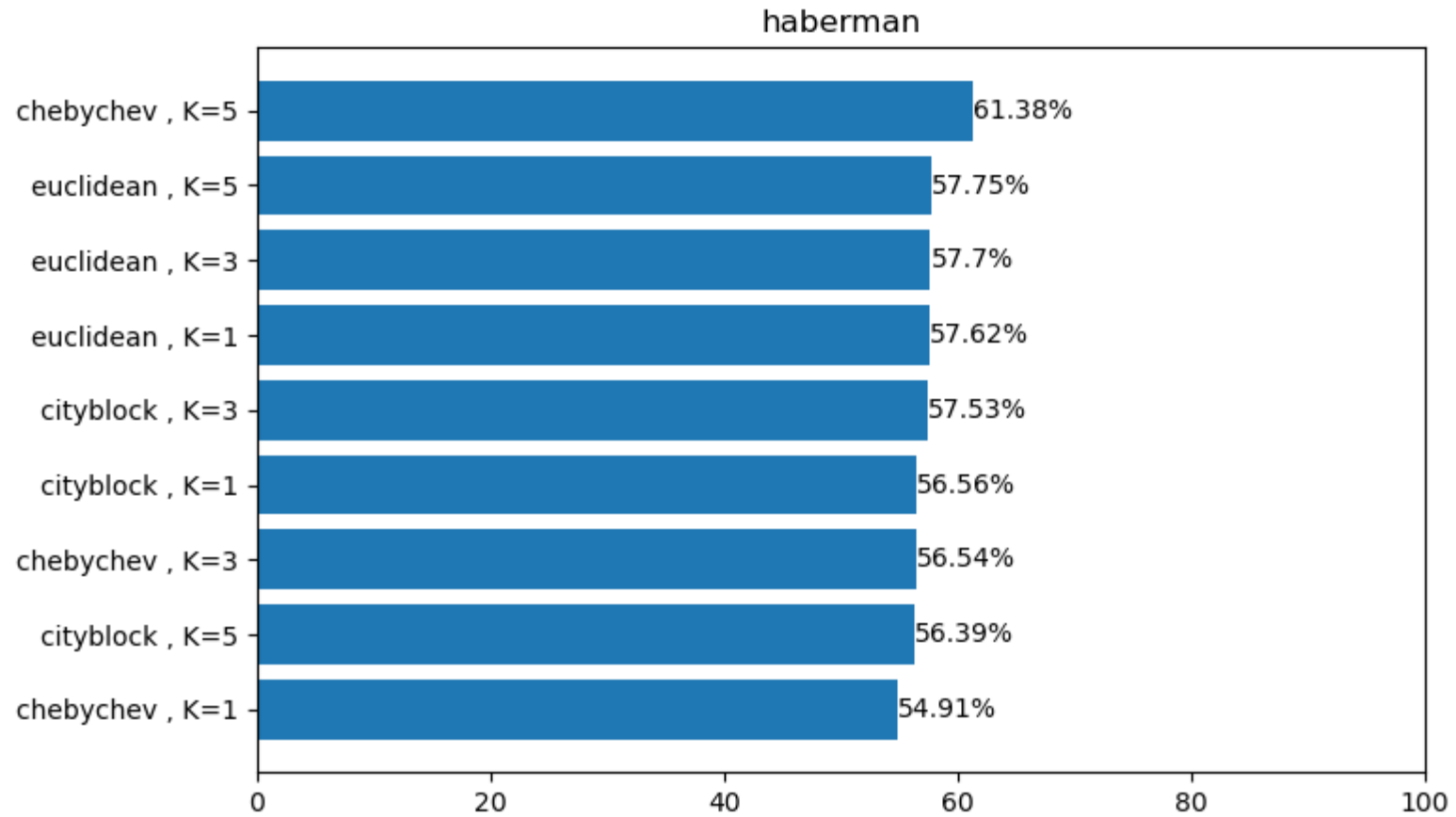
Heart



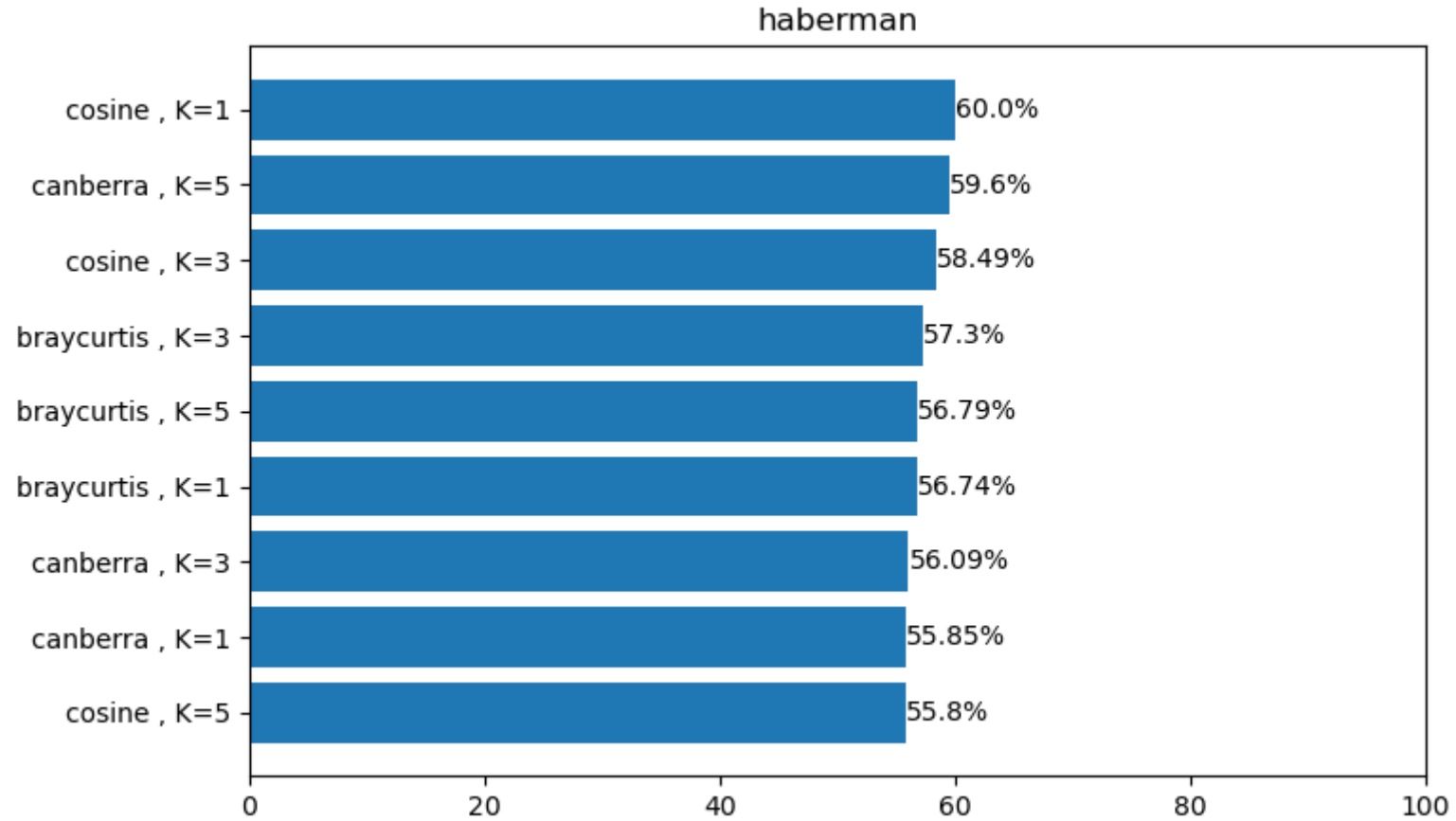
Heart



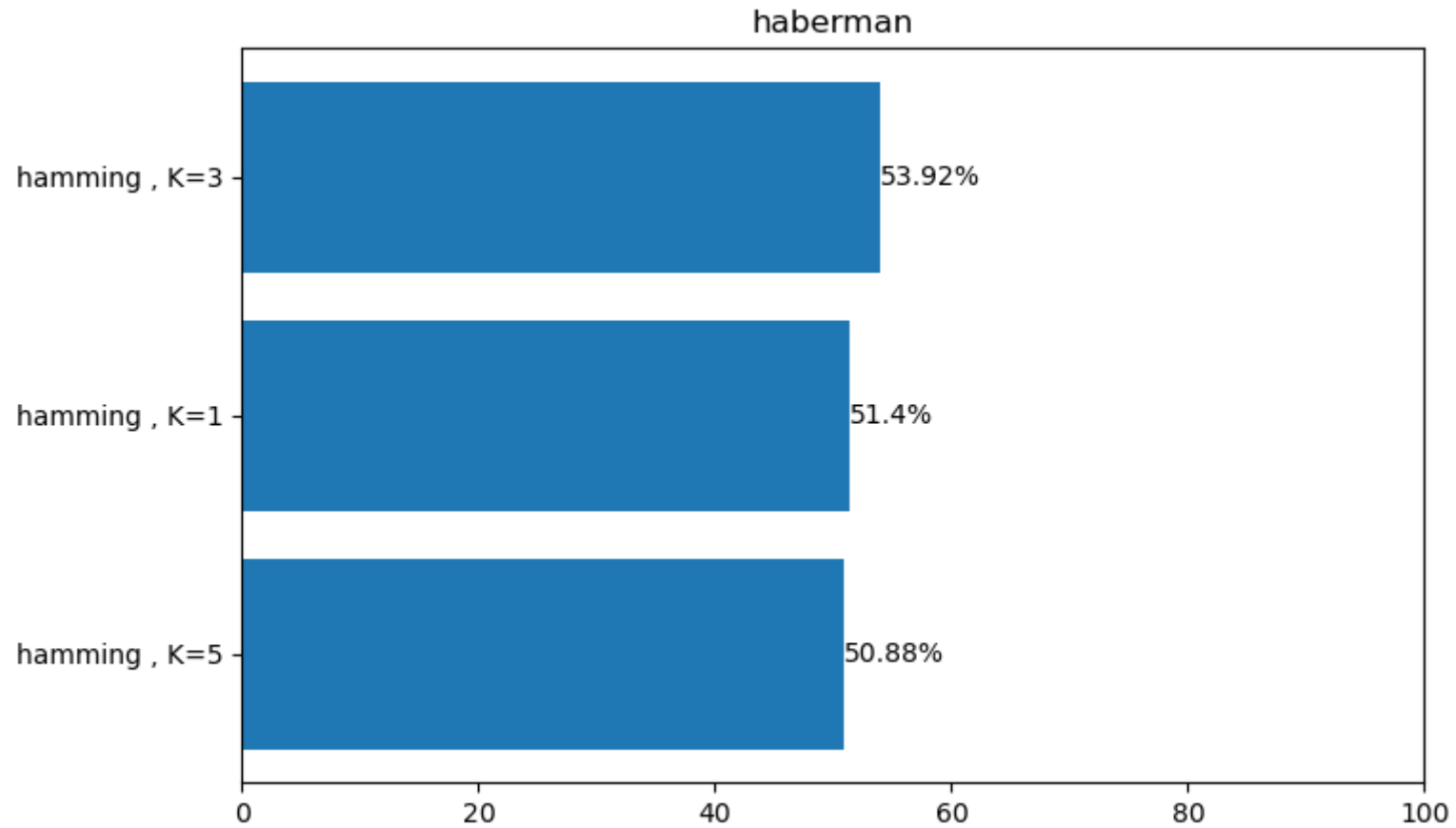
Haberman



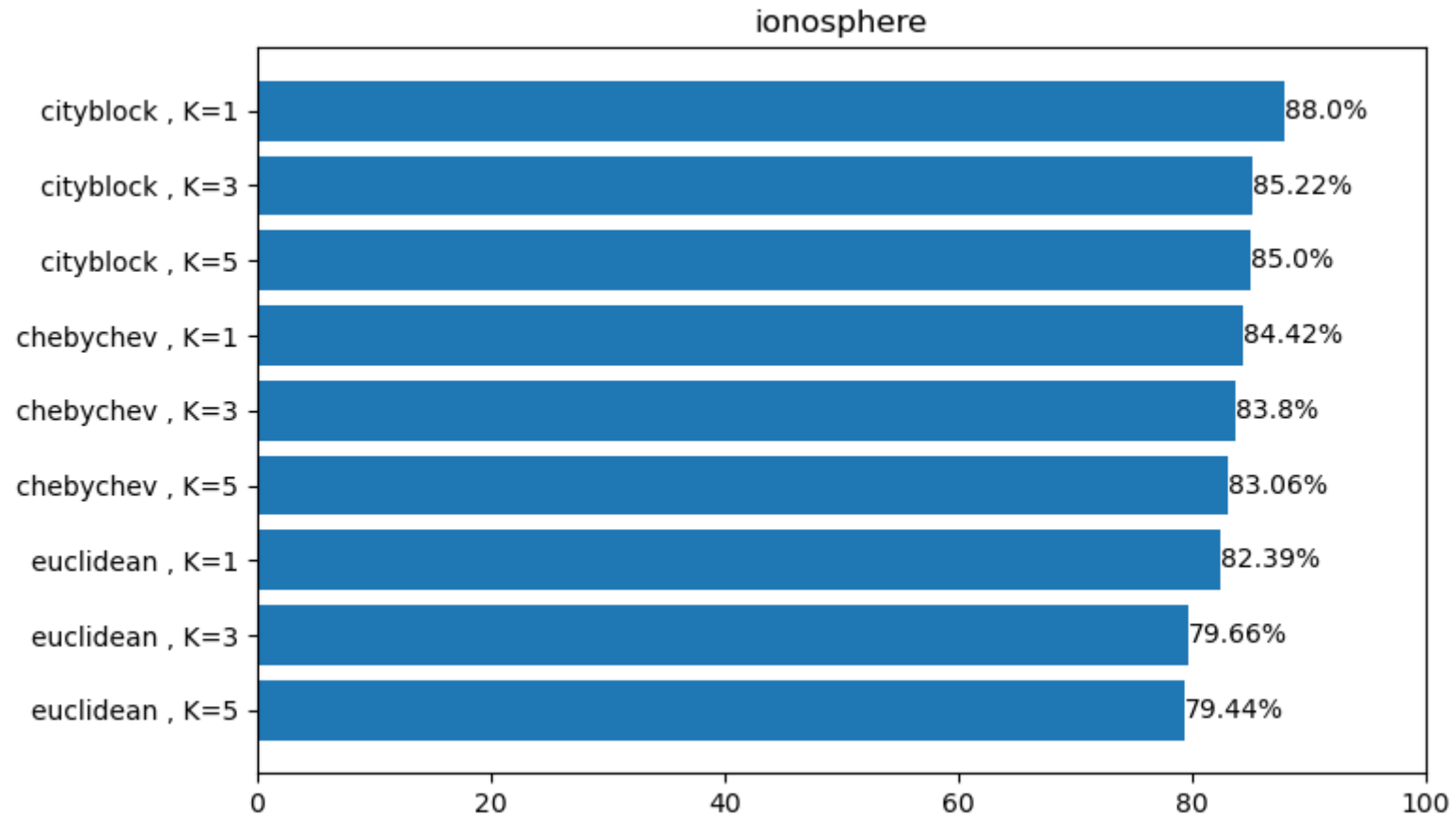
Haberman



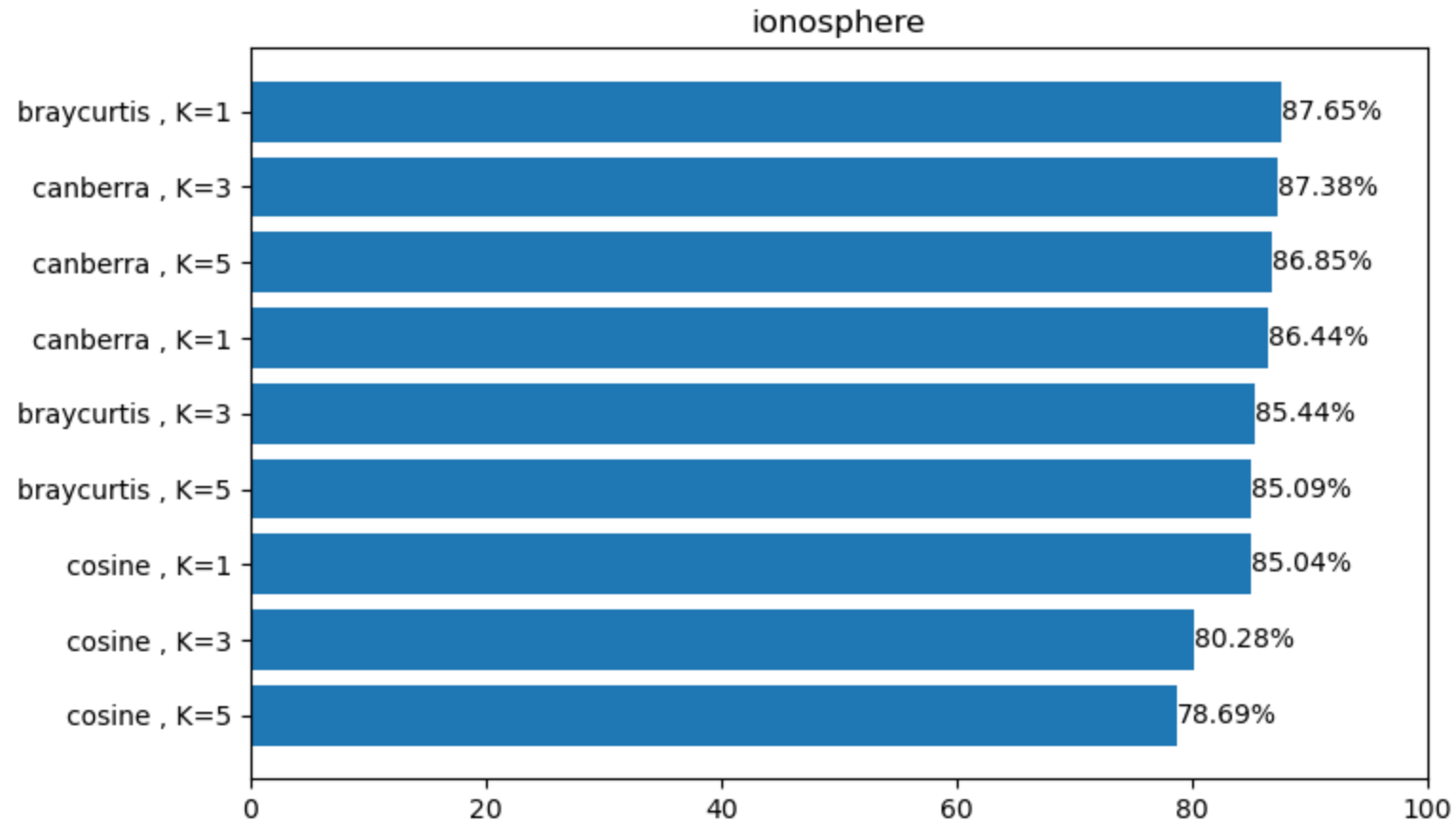
Haberman



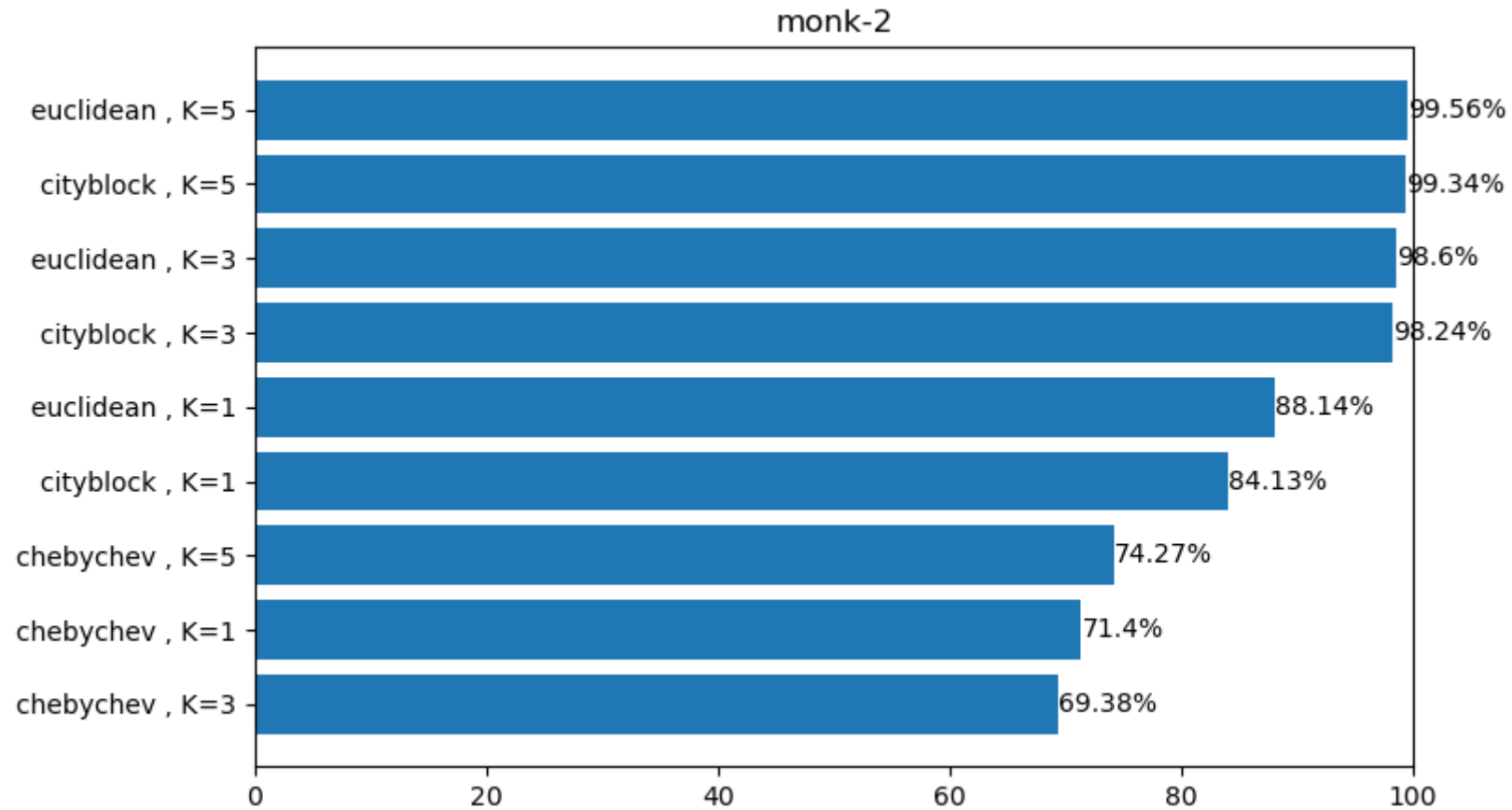
Ionosphere



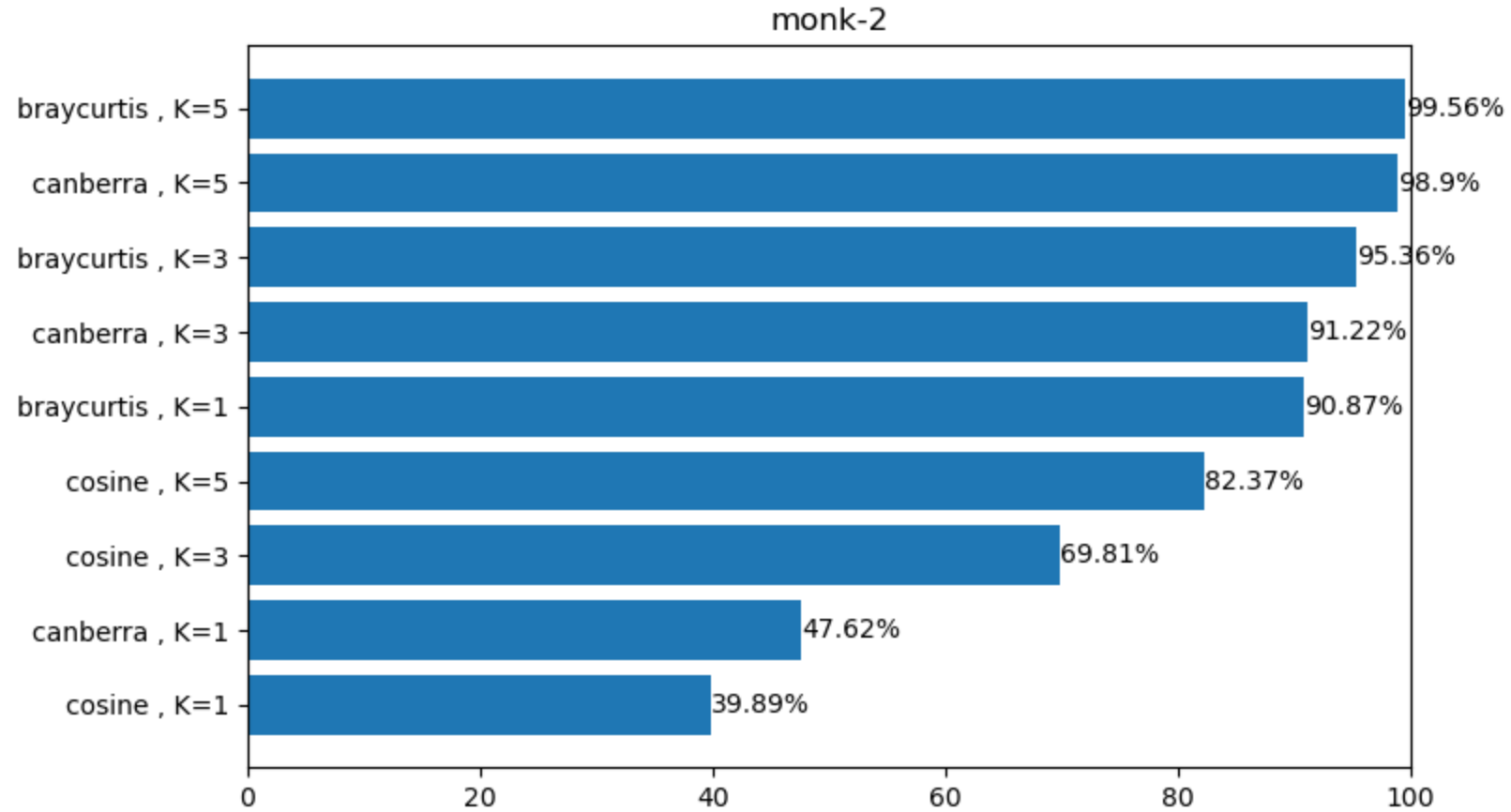
Ionosphere



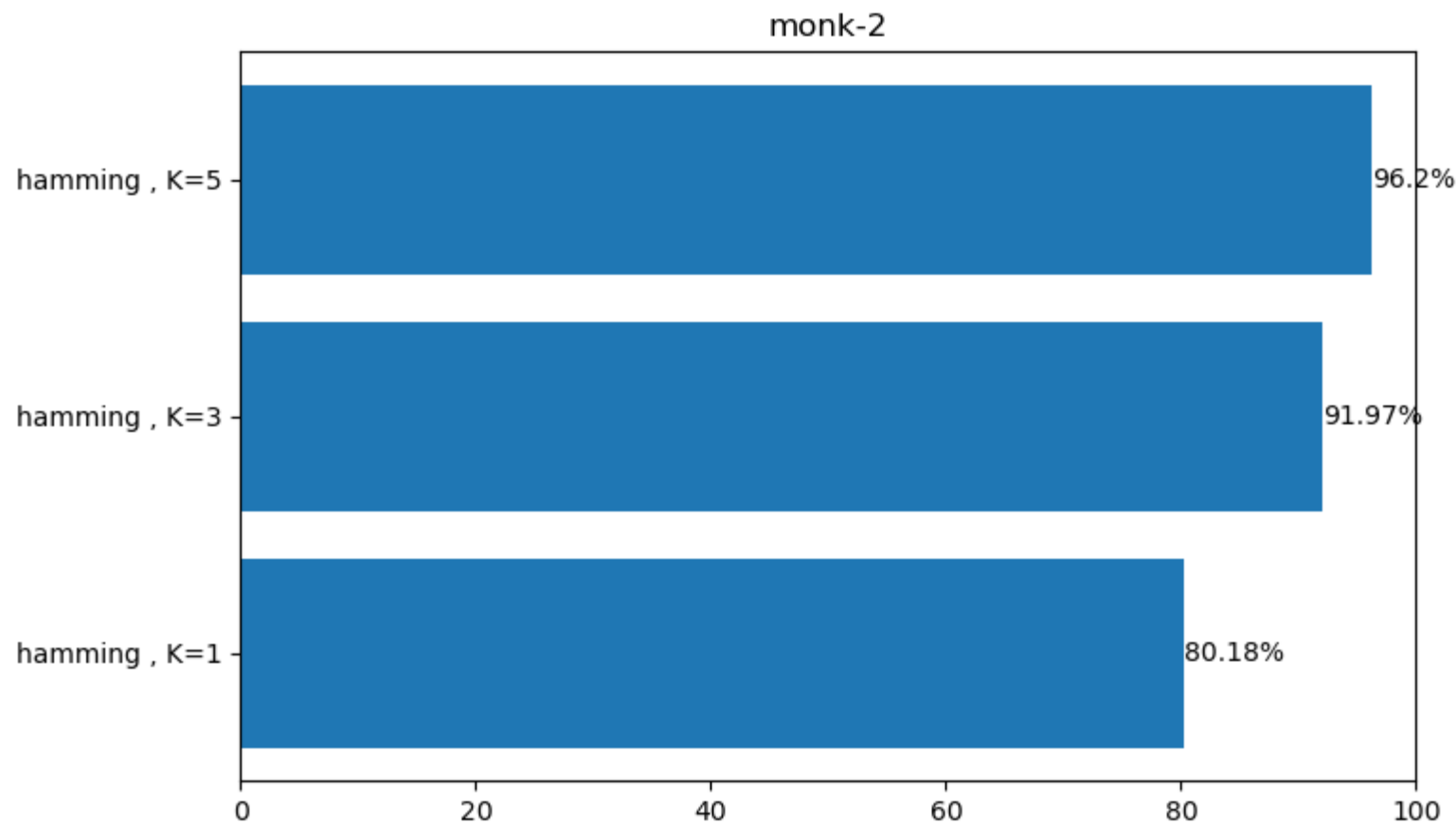
Monk-2



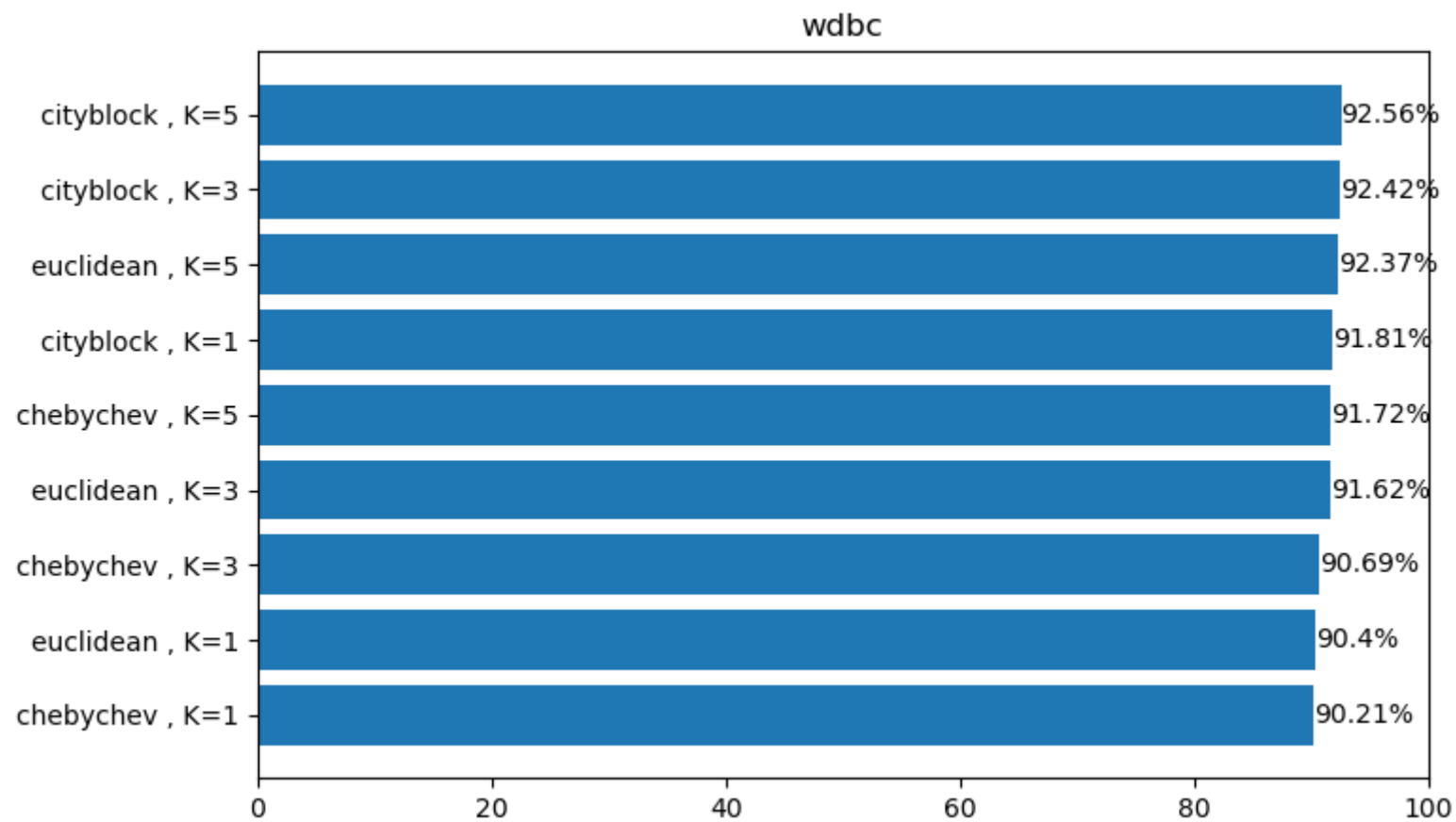
Monk-2



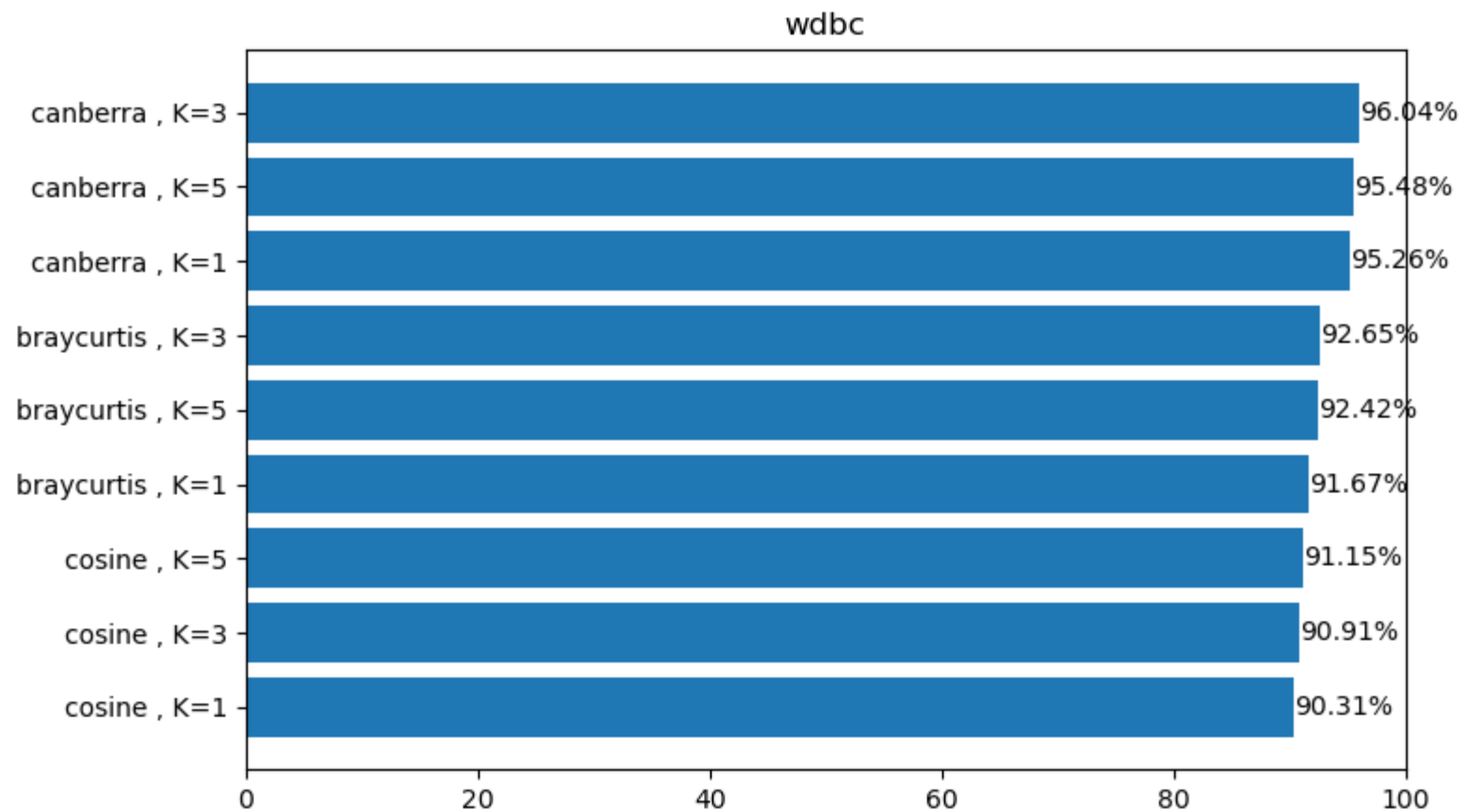
Monk-2



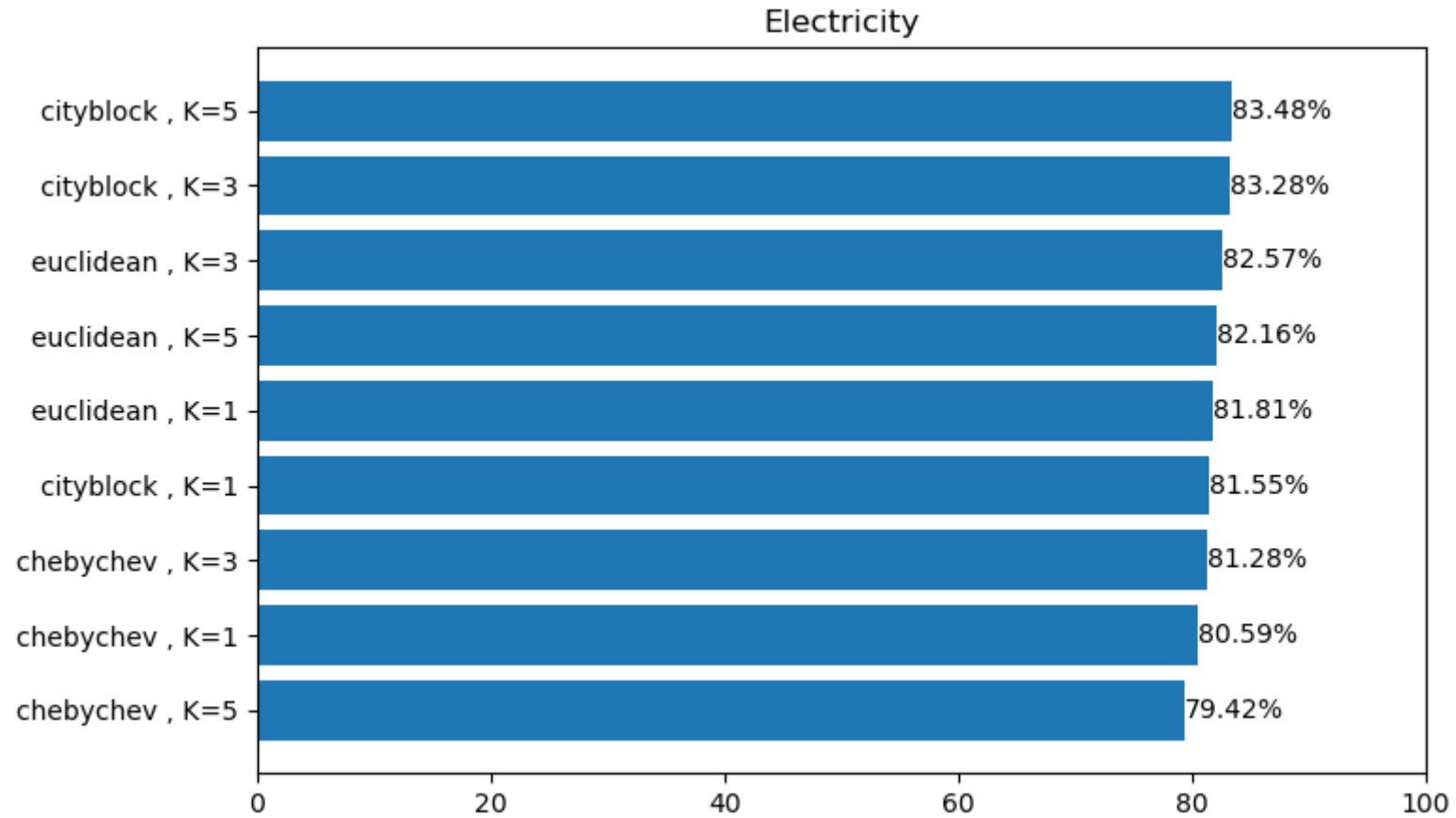
WDBC



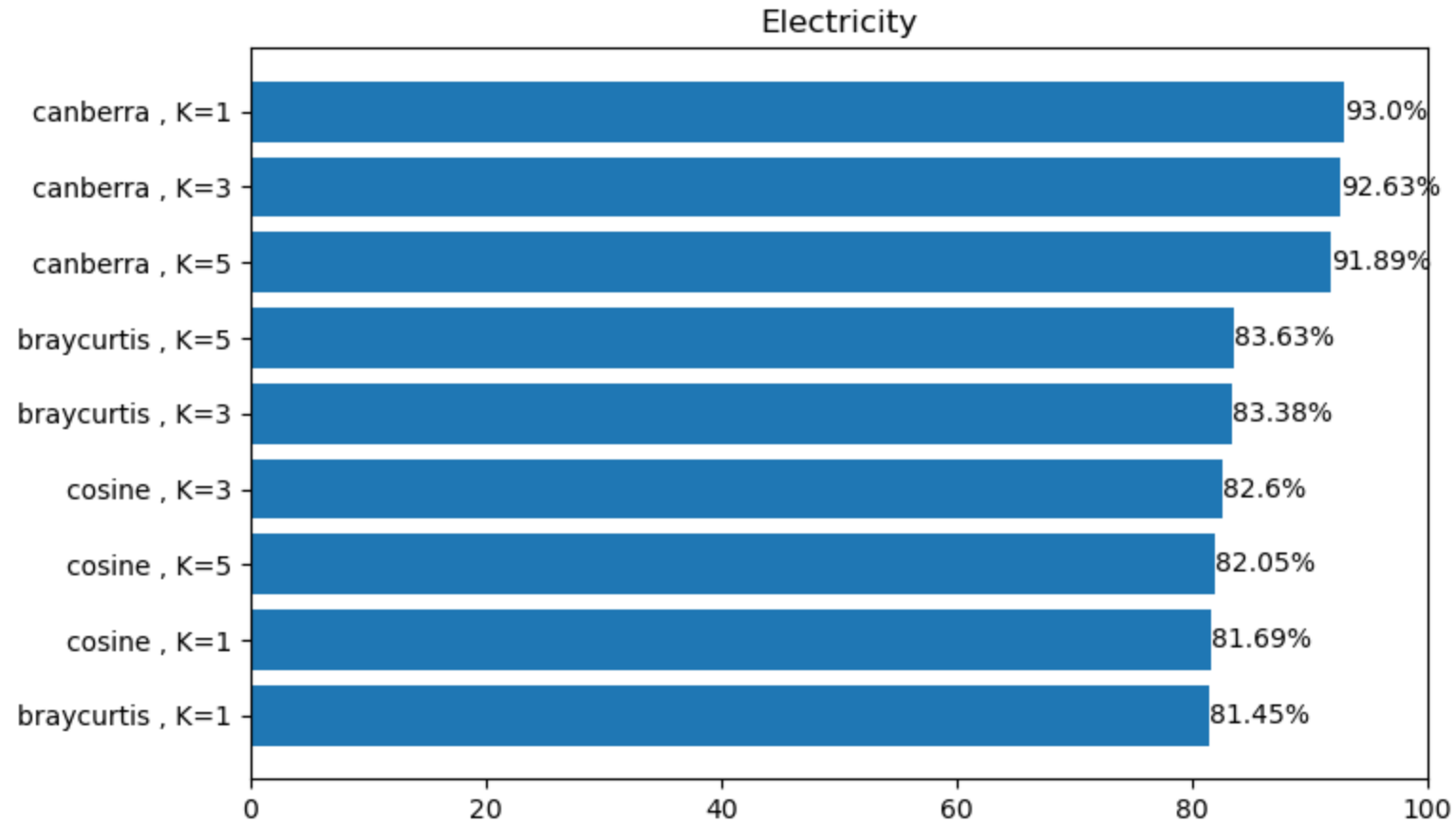
WDBC



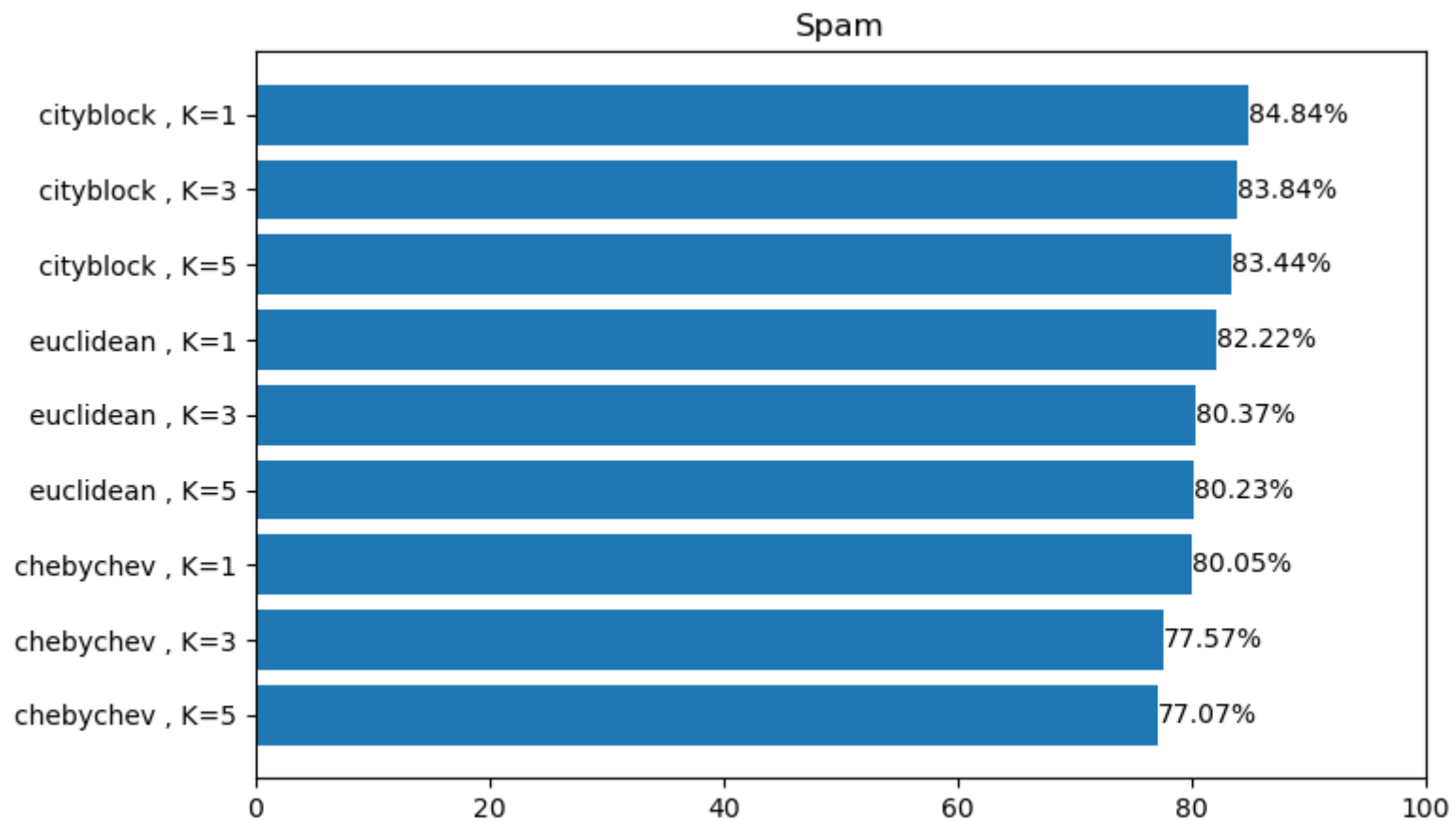
Electricity



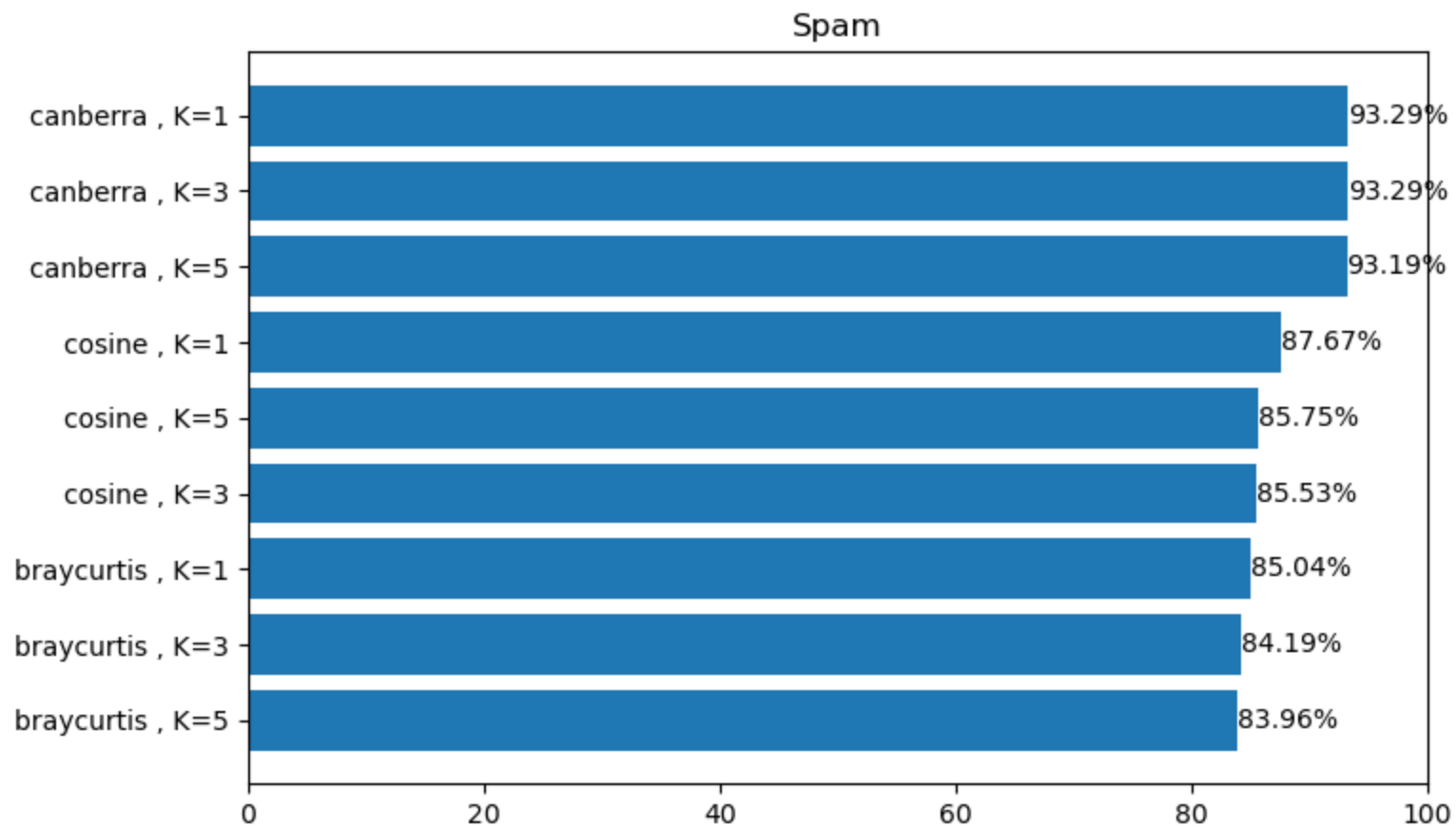
Electricity



Spam



Spam



Top 3 Metrics

Dataset	Metric	K	BA
Nutt	Braycurtis	3	89.28
	Braycurtis	5	89.28
	Chebychev	1	82.14
Apendicitis	Canberra	1	79.21
	Canberra	5	77.98
	Cityblock	1	76.21
Sonar	Canberra	1	85.82
	Cityblock	1	84.86
	Braycurtis	1	84.15
Heart	Canberra	5	82.55
	Canberra	3	78.49
	Canberra	1	78.08
Haberman	Chebychev	5	61.38
	Cosine	1	60
	Canberra	5	59.60

Dataset	Metric	K	BA
Ionosphere	Citiblock	1	88
	Braycurtis	1	87.65
	Canberra	3	87.38
Monk-2	Euclidean	5	99.56
	Braycurtis	5	99.56
	Cityblock	5	99.34
WDBC	Canberra	3	96.04
	Canberra	5	95.48
	Canberra	1	95.26
Electricity	Canberra	1	93
	Canberra	3	92.63
	Canberra	5	91.89
Spam	Canberra	1	93.29
	Canberra	3	93.29
	Canberra	5	93.19

Frequency

Metric	K	Frequency
Canberra	1	6
	5	6
	3	5
Citiblock	1	3
Braycurtis	1	2
	5	2
Braycurtis	3	1
Chebychev	1	1
	5	1
Cityblock	5	1
Cosine	1	1
Euclidean	5	1

Metric	Frequency
Canberra	17
Braycurtis	5
Cityblock	4
Chebychev	2
Cosine	1
Euclidean	1

K	Frequency
1	13
3	11
5	6

Thank you !

Questions?