

A Geometric Theory of Feature Selection and Distance-Based Measures

Kilho Shin and Adrian Pino Angulo

University of Hyogo

Kobe, Japan

yshin@ai.u-hyogo.ac.jp

Abstract

Feature selection measures are often explained by the analogy to a rule to measure the “distance” of sets of features to the “closest” ideal sets of features. An ideal feature set is such that it can determine classes uniquely and correctly. This way of explanation was just an analogy before this paper. In this paper, we show a way to map arbitrary feature sets of datasets into a common metric space, which is indexed by a real number p with $1 \leq p \leq \infty$. Since this determines the distance between an arbitrary pair of feature sets, even if they belong to different datasets, the distance of a feature set to the closest ideal feature set can be used as a feature selection measure. Surprisingly, when $p = 1$, the measure is identical to the Bayesian risk, which is probably the feature selection measure that is used the most widely in the literature. For $1 < p \leq \infty$, the measure is novel and has significantly different properties from the Bayesian risk. We also investigate the correlation between measurements by these measures and classification accuracy through experiments. As a result, we show that our novel measures with $p > 1$ exhibit stronger correlation than the Bayesian risk.

1 Introduction

Feature selection is indeed one of the central focuses of machine learning research. In this paper, we understand the feature selection problem as follows:

Given a dataset, select an appropriately small number of features that faithfully determine the classes of the examples of the dataset.

If we could find a feature subset such that the features correctly determine the classes of the examples, it could be an answer we wanted [Almuallim and Dietterich, 1994]. Such feature subsets are referred to as *reducts* in the rough set theory [Pawlak, 1991] and as *consistent feature sets* in this paper [Liu et al., 1998].

A dataset, however, does not necessarily include a consistent feature subset, and in such cases, we have to be satisfied

with selecting feature subsets that are sufficiently close to being consistent. In other words, we need a measure that can measure the “distance” of a feature subset to the “closest” consistent feature set. However, we only know two necessary conditions for such a measure, namely, *determinacy* and *monotonicity*: Determinacy corresponds to the identity of indiscernibles of the axioms of metrics and requires that the measurement is zero, if, and only if, the feature set is consistent; Monotonicity, on the other hand, requires that the measurement of an arbitrary feature set is no greater than the measurement of its subset. When a measure has the determinacy and monotonicity properties, we also say that it is *determinant* and *monotonous*. With a determinant and monotonous measure, a feature selection algorithm attempts to find feature sets whose measurements are close to zero.

Generally, the feature selection measures proposed in the literature are categorized into two groups. One consists of measures that simply evaluate dependence between two features [Mengle and Goharian, 2009; Bolón-Canedo et al., 2011; Foithong et al., 2012; Suzuki and Sugiyama, 2013; Wang, 2015]. The framework of *mRMR* [Peng et al., 2005] takes advantage of measures of this group and aims to select a feature set that maximizes the difference of the *relevance* from the *redundancy* of the feature set: The relevance is the collective dependence of the feature set to class labels, while the redundancy is the internal mutual dependence of the features of the feature set. On the other hand, the other consists of determinant and monotonous measures [Liu et al., 1998; Shin and Xu, 2009; Arauzo-Azofra et al., 2008], and 14 of the 17 feature selection measures studied in [Molina et al., 2002] are known to belong to this group [Shin et al., 2011]. In this paper, we are interested in the latter group.

Example 1 *Bayesian risk, also known as inconsistency rate [Liu et al., 1998], is defined as follows.*

$$\mu_{br}^D(S) = \sum_{\mathbf{x} \in \Omega_S} \left(p(S = \mathbf{x}) - \max_{c \in \Omega_C} p(S = \mathbf{x}, C = c) \right).$$

S is a feature subset of a dataset D, and C is the random variable to represent classes. Ω_S and Ω_C are the sample spaces of S and C, and p is the empirical probability distribution of D. This measure is, indeed, determinant and monotonous.

Example 2 *Zhao et al. [Zhao and Liu, 2007] proposed Interact, a feature selection algorithm that leverages the Bayes*

risk μ_{br} to evaluate the closeness of feature sets to the state of being consistent. Therefore, *Interact* selects feature sets with small μ_{br} measurements.

Example 3 *Shin et al. [Shin and Xu, 2009] also proposed the conditional entropy defined by*

$$\mu_{ce}^D(S) = \sum_{\mathbf{x} \in \Omega_S} \sum_{c \in \Omega_C} \mathbf{p}(S = \mathbf{x}, C = c) \log \frac{\mathbf{p}(S = \mathbf{x}, C = c)}{\mathbf{p}(S = \mathbf{x})}$$

as a determinant and monotonous measure. On the other hand, the well-known feature selection algorithm that selects the top n features X with respect to the mutual information $I(X; C)$ actually selects those feature sets that minimize

$$\mu^D(\{X_1, \dots, X_n\}) = \sum_{i=1}^n \mu_{ce}^D(\{X_i\}).$$

By using μ^D instead of μ_{ce}^D , this algorithm is fast but is not very accurate, since μ^D is not determinant.

When a measure represents the distance of a feature set to the closest consistent feature set in some metric space, we call it *distance-based*. About distance-based measures, we only knew two necessary conditions, that is, determinacy and monotonicity, and hence, the distance-based measure was only an analogy to explain the desirable natures of feature selection measures. This paper changes this.

In this paper, we present a way to identify a feature subset of a dataset uniquely with a point in a common metric space. The metric space is parameterized by a real number $p \in [1, \infty]$, and hence, we have infinitely many such metric spaces. In the metric space, consistent feature sets form a closed subspace, and hence, we can determine the distance of an arbitrary feature set to the closest consistent feature set. This is nothing other than a distance-based measure. A different value of p determines a different metric function. For $p = 1$, we show that the measure is identical to the well-known Bayesian risk. On the other hand, for $p > 1$, the measure is novel and unknown in the literature. Also, through experiments, we show that measurements by our novel distance-based measures correlate with classification accuracy.

2 Formulating feature selection measures

In this paper, we do not discriminate between features and random variables. Furthermore, for a dataset D , we use the following notations.

- $\overline{\mathcal{F}}_D$ is the entire set of features of D . $\overline{\mathcal{F}}_D$ is finite and includes the feature C_D that represents classes. Also, we denote $\overline{\mathcal{F}}_D \setminus \{C_D\}$ by \mathcal{F}_D .
- $\Omega_{\overline{\mathcal{F}}_D} = \prod_{X \in \overline{\mathcal{F}}_D} \Omega_X$ is the entire sample space, where Ω_X is the sample space of an individual feature X . For a feature set S , Ω_S means the Cartesian product $\prod_{X \in S} \Omega_X$, which is the sample space of S .
- $\mathbf{p}_D : \Omega_{\overline{\mathcal{F}}_D} \rightarrow \mathbb{Q}$ is the empirical probability distribution.
- We also denote D by the triplet $(\overline{\mathcal{F}}_D, \Omega_{\overline{\mathcal{F}}_D}, \mathbf{p}_D)$.

- For $S \subseteq \mathcal{F}_D$, $D|_S$ is the dataset derived from D by eliminating the features in $\mathcal{F}_D \setminus S$.
- For $S \subseteq \mathcal{F}_D$, $D^{\wedge S}$ is the dataset derived from D by replacing the features of S with a single feature $\wedge S$ such that $\Omega_{\wedge S} = \Omega_S$. The values for S of an example in D are replaced with the vector consisting of the values.

Example 4 *For the dataset D determined by (a) below, we have $\overline{\mathcal{F}}_D = \{B, R, G, C_D\}$ and $\Omega_{\overline{\mathcal{F}}_D} = \{A, B, O, AB\} \times \{M, N, C, A\} \times \{M, F\} \times \{p, n\}$. When $S = \{R, G\}$, (b) and (c) determine $D|_S$ and $D^{\wedge S}$.*

(a) D				(b) $D _S$			(c) $D^{\wedge S}$		
B	R	G	C_D	R	G	C_D	B	$\wedge S$	C_D
A	M	M	p	M	M	p	A	(M, M)	p
A	M	M	n	M	M	n	A	(M, M)	n
B	N	M	p	N	M	p	B	(N, M)	p
O	C	F	n	C	F	n	O	(C, F)	n
AB	A	M	n	A	M	n	AB	(A, M)	n

On the other hand, a feature selection measure μ is formulated as a family of μ^D indexed by datasets D , and each μ^D is a real-valued function defined over the power set $\mathfrak{P}(\mathcal{F}_D)$ of \mathcal{F}_D . In this paper, we assume that a feature selection measure supports the following three requirements, which all of the 19 measures surveyed in [Shin et al., 2011] support as well.

Requirement 1 (Naming invariance)

Measurements by the measure are invariant to renaming of features and feature values. This also implies that the measure assumes that every feature is categorical.

Requirement 2 (Localization invariance)

If $T \subseteq S$, $\mu^D(T) = \mu^{D|_S}(T)$ holds.

Requirement 3 (Feature aggregation invariance)

If $T \supseteq S$, $\mu^D(T) = \mu^{D^{\wedge S}}((T \setminus S) \cup \{\wedge S\})$ holds.

Example 5 *Bayesian risk μ_{br} has all of these invariance properties: $\mu_{br}^D(\mathcal{F}_D) = \mu_{br}^E(\mathcal{F}_E) = \frac{1}{5}$ indicates the naming invariance; For $S = \{R, G\}$, $\mu_{br}^D(\{G\})$ and $\mu_{br}^{D|_S}(\{G\})$ are identical to $\frac{2}{5}$, and $\mu_{br}^{D^{\wedge S}}(\{B, \wedge S\}) = \frac{1}{5}$ holds. These indicate the localization and feature aggregation invariance properties.*

				E			
				F_1	F_2	F_3	C_D
1	1	1	$+1$				
1	1	1	-1				
2	2	1	$+1$				
3	3	2	-1				
4	4	1	-1				

In this paper, we will identify (D, S) and (E, T) , pairs of a dataset and a feature subset, if $\mu^D(S) = \mu^E(T)$ holds for any measure μ that supports Requirements 1, 2 and 3. For this, we introduce the notion of *coherent* mappings and then prove Proposition 1.

Definition 1 *A mapping $\varphi : \mathcal{F}_D \rightarrow \mathcal{F}_E$ is coherent between D and E , if, and only if, there exist injections $v_Y : \Omega_{\varphi^{-1}(Y)} \rightarrow \Omega_Y$ for $Y \in \text{Im}(\varphi)$ and an injection $v_{C_E} : \Omega_{C_D} \rightarrow \Omega_{C_E}$ such that*

$$\mathbf{p}_D(\mathbf{x}) = \sum_{\mathbf{y} : \mathbf{y}|_{\text{Im}(\varphi) \cup \{C_E\}} = v_\pi(\mathbf{x})} \mathbf{p}_E(\mathbf{y}),$$

where $v_\pi = \prod_{Y \in \text{Im}(\varphi) \cup \{C_E\}} v_Y : \Omega_{\overline{\mathcal{F}}_D} \rightarrow \Omega_{\text{Im}(\varphi) \cup \{C_E\}}$.

Example 6 We consider the datasets D of Example 4 and E below. When we define φ by $\varphi(B) = F_1$ and $\varphi(R) = \varphi(G) = F_2$, φ is coherent. To show this, we determine v_{F_1} and v_{F_2} as below and v_{C_E} by $v_{C_E}(p) = 1$ and $v_{C_E}(n) = 2$. For example, $\mathbf{p}_D(A, M, M, p) = \frac{1}{5}$ and $\mathbf{p}_E(1, 1, 1, 1) + \mathbf{p}_E(1, 1, 3, 1) = \frac{1}{10} + \frac{1}{10} = \frac{1}{5}$ hold.

E							
F_1	F_2	F_3	C_E	F_1	F_2	F_3	C_E
1	1	1	1	2	3	4	1
1	1	2	2	3	6	2	2
1	1	3	1	3	6	3	2
1	1	3	2	4	7	4	2
2	3	1	1	4	7	4	2

v_{F_1}	\parallel	v_{F_2}
A → 1	\parallel	(M, M) → 1 (C, M) → 5
B → 2	\parallel	(M, F) → 2 (C, F) → 6
O → 3	\parallel	(N, M) → 3 (A, M) → 7
AB → 4	\parallel	(N, F) → 4 (A, F) → 8

Proposition 1 If a feature selection measure μ supports Requirement 1 to 3, $\mu^E(S) = \mu^D(\varphi^{-1}(S))$ holds for an arbitrary coherent mapping $\varphi : \mathcal{F}_D \rightarrow \mathcal{F}_E$ between D and E and an arbitrary $S \subseteq \text{Im}(\varphi)$.

Proof. We let $\varphi = \varphi \circ \bar{\varphi}$ such that $\bar{\varphi} : \mathcal{F}_D \rightarrow \text{Im}(\varphi)$ and $\varphi : \text{Im}(\varphi) \rightarrow \mathcal{F}_E$. Since $\bar{\varphi}$ yields a sequence of feature aggregations that convert D into $E|_{\text{Im}(\varphi)}$ up to renaming of features and feature values, $\mu^E(S) = \mu^{E|_{\text{Im}(\varphi)}}(S) = \mu^D(\varphi^{-1}(S))$ follows from Requirement 1 to 3. ■

Example 7 For the same D , E and φ as Example 6, $\mu_{\text{br}}^E(\{F_1, F_2\}) = \mu_{\text{br}}^D(\{B, R, G\}) = \mu_{\text{br}}^E(\{F_1\}) = \mu_{\text{br}}^D(\{B\}) = \mu_{\text{br}}^E(\{F_2\}) = \mu_{\text{br}}^D(\{R, G\}) = \frac{1}{5}$ holds.

3 Projecting a feature set of a dataset into a subspace $\mathcal{P}_{\mathbb{R}}^{+(\infty)}$ of the ℓ^p space ($p \in [1, \infty]$)

We let $n \in \mathbb{N} \cup \{\infty\}$. We introduce the space $\mathcal{P}_0^{(n)}$ into which feature sets of datasets are projected. We start with preparation. Let \mathbb{K} be either the rational number field \mathbb{Q} or the real number field \mathbb{R} . We define $\mathcal{P}_{\mathbb{K}}^{(\infty)}$ and $\mathcal{P}_{\mathbb{K}}^{+(\infty)}$ as follows, which are sets of functions \mathbf{p} such that $\mathbf{p} : \mathbb{N}^2 \rightarrow \mathbb{K}$.

Definition 2 We define $\mathcal{P}_{\mathbb{K}}^{(\infty)}$ and $\mathcal{P}_{\mathbb{K}}^{+(\infty)}$ by:

$$\mathcal{P}_{\mathbb{K}}^{(\infty)} = \left\{ \mathbf{p} \mid \mathbf{p}(i, j) \geq 0, \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \mathbf{p}(i, j) = 1 \right\};$$

$$\mathcal{P}_{\mathbb{K}}^{+(\infty)} = \left\{ \mathbf{p} \mid \mathbf{p}(i, j) \geq 0, \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \mathbf{p}(i, j) \leq 1 \right\}.$$

We can identify \mathbb{N}^2 with \mathbb{N} (for example, $f : (i, j) \mapsto \frac{(i+j-1)(i+j-2)}{2} + i$ is bijective), and therefore, $\mathcal{P}_{\mathbb{R}}^{(\infty)}$ and $\mathcal{P}_{\mathbb{R}}^{+(\infty)}$ can be viewed as subspaces of ℓ^p for $1 \leq p \leq \infty$. If we focus on datasets that include at most n classes, we can use $\mathcal{P}_{\mathbb{K}}^{(n)}$ and $\mathcal{P}_{\mathbb{K}}^{+(n)}$ instead.

Definition 3 We define $\mathcal{P}_{\mathbb{K}}^{(n)}$ and $\mathcal{P}_{\mathbb{K}}^{+(n)}$ by:

$$\mathcal{P}_{\mathbb{K}}^{(n)} = \left\{ \mathbf{p} \in \mathcal{P}_{\mathbb{K}}^{(\infty)} \mid j > n \Rightarrow \mathbf{p}(i, j) = 0 \right\};$$

$$\mathcal{P}_{\mathbb{K}}^{+(n)} = \left\{ \mathbf{p} \in \mathcal{P}_{\mathbb{K}}^{+(\infty)} \mid j > n \Rightarrow \mathbf{p}(i, j) = 0 \right\};$$

$\mathcal{P}_{\mathbb{R}}^{(n)}$ and $\mathcal{P}_{\mathbb{R}}^{+(n)}$ are metric spaces with the metric derived from the norm $\|\cdot\|_p$ of ℓ^p . $\mathcal{P}_{\mathbb{R}}^{(n)}$ is closed in ℓ^p only when $p = 1$, while $\mathcal{P}_{\mathbb{R}}^{+(n)}$ is the closure of $\mathcal{P}_{\mathbb{R}}^{(n)}$ in ℓ^p for $1 < p \leq \infty$. Therefore, $\mathcal{P}_{\mathbb{R}}^{(n)}$ in ℓ^1 and $\mathcal{P}_{\mathbb{R}}^{+(n)}$ in ℓ^p for $1 < p \leq \infty$ are complete, since ℓ^p is a Banach space. Although $\mathcal{P}_{\mathbb{R}}^{(n)}$ and $\mathcal{P}_{\mathbb{R}}^{+(n)}$ are bounded, none of them is compact.

We define $\mathcal{P}_0^{(n)}$, a subspace of $\mathcal{P}_{\mathbb{R}}^{(n)}$, as follows.

Definition 4 We let $\text{supp}(\mathbf{p})$ be the smallest $\Omega_U \times \Omega_C$ such that $\Omega_U \subseteq \mathbb{N}$, $\Omega_C \subseteq \mathbb{N}$ and $\Omega_U \times \Omega_C \supseteq \text{Supp}(\mathbf{p}) = \{(i, j) \mid \mathbf{p}(i, j) > 0\}$. We define $\mathcal{P}_0^{(\infty)}$ and $\mathcal{P}_0^{(n)}$ for $n \in \mathbb{N}$ by:

$$\mathcal{P}_0^{(\infty)} = \left\{ \mathbf{p} \in \mathcal{P}_{\mathbb{Q}}^{(\infty)} \mid |\text{supp}(\mathbf{p})| < \infty \right\};$$

$$\mathcal{P}_0^{(n)} = \left\{ \mathbf{p} \in \mathcal{P}_0^{(\infty)} \mid j > n \Rightarrow \mathbf{p}(i, j) = 0 \right\}.$$

Example 8 When we define \mathbf{p} by $\mathbf{p}(i, i) = (1-r)r^{i-1}$ for $r \in (0, 1)$, $\mathbf{p} \in \mathcal{P}_{\mathbb{R}}^{(\infty)}$ holds, because $\sum_{i=1}^{\infty} (1-r)r^{i-1} = (1-r) \lim_{n \rightarrow \infty} \frac{1-r^n}{1-r} = 1$. Furthermore, if $r \in \mathbb{Q}$, $\mathbf{p} \in \mathcal{P}_{\mathbb{Q}}^{(\infty)}$ holds. For this \mathbf{p} , we have $\Omega_U = \mathbb{N}$, $\Omega_C = \mathbb{N}$ and $\text{supp}(\mathbf{p}) = \{(i, i) \mid i \in \mathbb{N}\}$. Hence, $\mathbf{p} \notin \mathcal{P}_0^{(\infty)}$.

For $n \in \mathbb{N} \cup \{\infty\}$, $\mathcal{P}_0^{(n)}$ turns out to be dense in $\mathcal{P}_{\mathbb{R}}^{(n)}$. For $\mathbf{p} \in \mathcal{P}_0^{(n)}$, $D_{\mathbf{p}} = (\{U, C\}, \mathbb{N} \times \mathbb{N}, \mathbf{p})$ is not a dataset, because $\mathbb{N} \times \mathbb{N}$ is infinite. Nevertheless, we can think of a coherent mapping $\varphi : \mathcal{F}_D \rightarrow \{U\}$ between a dataset D and $D_{\mathbf{p}}$. Theorem 1 shows how to project a feature set S of a dataset D into $\mathcal{P}_0^{(n)}$.

Theorem 1 We let D be a dataset with n classes and $S \subseteq \mathcal{F}_D$. There exists $\mathbf{p} \in \mathcal{P}_0^{(n)}$ such that $\varphi : S \rightarrow \{U\}$ is coherent between $D|_S$ and $D_{\mathbf{p}}$. Conversely, for arbitrary $\mathbf{p} \in \mathcal{P}_0^{(n)}$, there exists a dataset D with n classes such that $\varphi : \mathcal{F}_D \rightarrow \{U\}$ is coherent between D and $D_{\mathbf{p}}$.

Proof. Let $v_1 : \Omega_S \rightarrow \mathbb{N}$ and $v_2 : \Omega_{C_D} \rightarrow \{1, \dots, n\}$ be injective. It suffices to define $\mathbf{p} \in \mathcal{P}_0^{(n)}$ by $\mathbf{p}(i, j) = \mathbf{p}_{D|_S}(v_1^{-1}(i), v_2^{-1}(j))$ if $(i, j) \in \text{Im}(v_1 \times v_2)$ and $\mathbf{p}(i, j) = 0$ otherwise. To show the converse, we have only to let $D = (\{U, C\}, \text{supp}(\mathbf{p}), \mathbf{p})$. ■

Example 9 To project $D|_S$ of Example 4 into $\mathcal{P}_0^{(2)}$, we determine $\mathbf{p} \in \mathcal{P}_0^{(2)}$ by $\mathbf{p}(1, 1) = \mathbf{p}(1, 2) = \mathbf{p}(3, 1) = \mathbf{p}(6, 2) = \mathbf{p}(7, 4) = \frac{1}{5}$ and $\mathbf{p}(i, j) = 0$ otherwise. When we determine v_U and v_C identical to v_{F_2} and v_{C_E} of Example 6, $\varphi : S \rightarrow \{U\}$ is coherent between $D|_D$ and $D_{\mathbf{p}}$, and hence, $D|_D$ is projected to \mathbf{p} .

Example 10 A single $D|_S$ can be projected to infinitely many points in $\mathcal{P}_0^{(n)}$. For example, we let $D|_S$, \mathbf{p} and v_U be the same as Example 9. For an arbitrary bijection $\pi : \mathbb{N} \rightarrow \mathbb{N}$, we define \mathbf{p}_{π} by $\mathbf{p}_{\pi}(\pi(i), j) = \mathbf{p}(i, j)$. Apparently, $D|_S$ projects to \mathbf{p}_{π} with $\pi \circ v_U$.

Example 11 More than one datasets are projected to a single point in $\mathcal{P}_0^{(n)}$. For example, $D|_S$ of Example 4 and $E|_{\{F_2\}}$ of Example 6 are both projected to any of \mathbf{p}_π of Example 10. To be precise, a single point in $\mathcal{P}_0^{(n)}$ has infinitely many datasets that project to the point.

Based on Proposition 1 and Theorem 1, Theorem 2 asserts that a feature selection measure can be uniquely viewed as a real-valued function defined over $\mathcal{P}_0^{(n)}$.

Theorem 2 For a feature selection measure μ , $\hat{\mu} : \mathcal{P}_0^{(\infty)} \rightarrow \mathbb{R}$ uniquely exists and $\hat{\mu}(\mathbf{p}) = \mu^D(S)$ holds for arbitrary D , S and \mathbf{p} such that $\varphi : S \rightarrow \{U\}$ is coherent between $D|_S$ and $D_{\mathbf{p}}$.

Proof. We let $D_{\mathbf{p}}^*$ be $(\{U, C\}, \text{supp}(\mathbf{p}), \mathbf{p})$ for $\mathbf{p} \in \mathcal{P}_0^{(\infty)}$. We determine $\hat{\mu}(\mathbf{p})$ by $\mu^{D_{\mathbf{p}}^*}(\{U\})$. For a dataset D such that $\varphi : S \rightarrow \{U\}$ is coherent between $D|_S$ and $D_{\mathbf{p}}$ with $v_U : \Omega_S \rightarrow \mathbb{N}$ and $v_C : \Omega_{C_D} \rightarrow \mathbb{N}$, we let $D_{\mathbf{p}}^\# = (\{U, C\}, \text{Im}(v_U) \times \text{Im}(v_C), \mathbf{p})$. $\varphi : S \rightarrow \{U\}$ is coherent between $D|_S$ and $D_{\mathbf{p}}^\#$, and the inclusion $\text{supp}(\mathbf{p}) \subseteq \text{Im}(v_U) \times \text{Im}(v_C)$ yields coherence between $D_{\mathbf{p}}^*$ and $D_{\mathbf{p}}^\#$. Therefore, $\mu^D(S) = \mu^{D_{\mathbf{p}}^\#}(\{U\}) = \mu^{D_{\mathbf{p}}^*}(\{U\}) = \hat{\mu}(\mathbf{p})$. ■

Example 12 Datasets that are projected to the same point in $\mathcal{P}_0^{(n)}$ have the same measurement for an arbitrary feature selection measure μ . For example, as seen in Example 11, $D|_S$ of Example 4 and $E|_{\{F_2\}}$ of Example 6 are both projected to the same point $\mathcal{P}_0^{(2)}$. On the other hand, $\mu_{\text{br}}^D(S) = \mu_{\text{br}}^E(\{F_2\}) = \frac{1}{5}$ holds as seen in Example 7. Therefore, by letting $\hat{\mu}(\mathbf{p}) = \mu^D(S)$ for $D|_S$ that projects to \mathbf{p} , $\hat{\mu}$ is well defined.

The problem of the projection determined by Theorem 1 consists in the fact that a single pair (D, S) is mapped to an infinite number of different points in $\mathcal{P}_0^{(n)}$ (Example 10). In the next section, we will solve this problem.

4 Defining a quotient space $\mathcal{Q}_{\mathbb{R}}^{+(\infty)}$ of $\mathcal{P}_{\mathbb{R}}^{+(\infty)}$

Our idea to solve the problem is to introduce an equivalence relation that unifies points of $\mathcal{P}_0^{(n)}$ that are the images of the same single pair (D, S) and then use the resulting quotient space. The relation \sim defined below is an equivalence relation. We assume $n \in \mathbb{N} \cup \{\infty\}$.

Definition 5 Let \mathbf{p} and \mathbf{q} be in $\mathcal{P}_{\mathbb{R}}^{+(\infty)}$. We define $\mathbf{p} \sim \mathbf{q}$, if, and only if, there exist bijections $v_U : \mathbb{N} \rightarrow \mathbb{N}$ and $v_C : \mathbb{N} \rightarrow \mathbb{N}$ such that $\mathbf{p}(i, j) = \mathbf{q}(v_U(i), v_C(j))$.

Example 13 $\mathbf{p}_\pi \sim \mathbf{p}$ holds for \mathbf{p}_π and π of Example 10.

Definition 6 We let $\mathcal{Q}_{\mathbb{R}}^{+(\infty)} = \mathcal{P}_{\mathbb{R}}^{+(\infty)} / \sim$. Moreover, for the canonical projection $\pi : \mathcal{P}_{\mathbb{R}}^{+(\infty)} \rightarrow \mathcal{Q}_{\mathbb{R}}^{+(\infty)}$, we let $\mathcal{Q}_{\mathbb{R}}^{+(n)} = \pi(\mathcal{P}_{\mathbb{R}}^{+(n)})$, $\mathcal{Q}_{\mathbb{R}}^{(n)} = \pi(\mathcal{P}_{\mathbb{R}}^{(n)})$ and $\mathcal{Q}_0^{(n)} = \pi(\mathcal{P}_0^{(n)})$.

Proposition 2 is easy to see but plays a crucial rule to prove Theorem 3 and 4.

Proposition 2 For \mathbf{p} and \mathbf{q} in $\mathcal{P}_0^{(\infty)}$, the following are equivalent.

1. $\mathbf{p} \sim \mathbf{q}$.
2. There exists (D, S) such that $\varphi : S \rightarrow \{U\}$ is coherent between $D|_S$ and $D_{\mathbf{p}}$ and between $D|_S$ and $D_{\mathbf{q}}$.
3. $\varphi : S \rightarrow \{U\}$ is coherent between $D|_S$ and $D_{\mathbf{p}}$, if, and only if, it is coherent between $D|_S$ and $D_{\mathbf{q}}$.

Theorem 3 and 4 follow from Theorem 1, Theorem 2 and Proposition 2. In the following, $[\mathbf{p}]$ denotes $\pi(\mathbf{p})$, and $[D, S]$ denotes the image of (D, S) in $\mathcal{Q}_0^{(\infty)}$.

Theorem 3 We let D be a dataset with n classes and $S \subseteq \mathcal{F}_D$. There uniquely exists $[\mathbf{p}] \in \mathcal{Q}_0^{(n)}$ such that $\varphi : S \rightarrow \{U\}$ is coherent between $D|_S$ and $D_{\mathbf{p}}$. Conversely, for arbitrary $[\mathbf{p}] \in \mathcal{Q}_0^{(n)}$, there exists a dataset D with n classes such that $\varphi : \mathcal{F}_D \rightarrow \{U\}$ is coherent between D and $D_{\mathbf{p}}$.

Theorem 4 For a feature selection measure μ , there uniquely exists $[\hat{\mu}] : \mathcal{Q}_0^{(\infty)} \rightarrow \mathbb{R}$ such that $[\hat{\mu}]([D, S]) = \mu^D(S)$ holds for any dataset D and $S \subseteq \mathcal{F}_D$.

5 Introducing a metric d_p into $\mathcal{Q}_{\mathbb{R}}^{+(\infty)}$

Although a quotient of a metric space is not always a metric space, we can derive a metric into $\mathcal{Q}_{\mathbb{R}}^{+(\infty)}$ from ℓ^p as follows.

Definition 7 For $[\mathbf{p}]$ and $[\mathbf{q}]$ in $\mathcal{Q}_{\mathbb{R}}^{+(\infty)}$, we define $d_p([\mathbf{p}], [\mathbf{q}]) = \inf \{ \|\mathbf{p}' - \mathbf{q}'\|_p \mid \mathbf{p} \sim \mathbf{p}', \mathbf{q} \sim \mathbf{q}' \}$.

Theorem 5 For $1 \leq p \leq \infty$, d_p is a metric over $\mathcal{Q}_{\mathbb{R}}^{+(\infty)}$.

Proof. We only prove the triangle inequality $d_p([\mathbf{p}], [\mathbf{r}]) \leq d_p([\mathbf{p}], [\mathbf{q}]) + d_p([\mathbf{q}], [\mathbf{r}])$ taking advantage of Lemma 1.

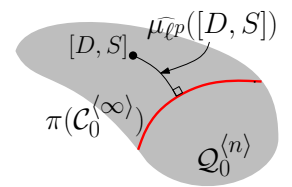
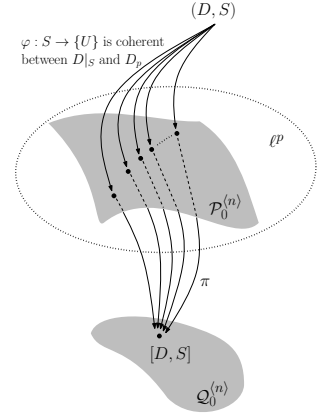
Lemma 1 For $\mathbf{p} \sim \mathbf{p}'$ and \mathbf{q} in $\mathcal{P}_{\mathbb{R}}^{+(\infty)}$, there exists $\mathbf{q}' \sim \mathbf{q}$ such that $\|\mathbf{p} - \mathbf{q}\|_p = \|\mathbf{p}' - \mathbf{q}'\|_p$.

For $\varepsilon > 0$, we assume $\|\mathbf{p} - \mathbf{q}\|_p - d_p([\mathbf{p}], [\mathbf{q}]) < \frac{\varepsilon}{2}$ and $\|\mathbf{q}' - \mathbf{r}\|_p - d_p([\mathbf{q}], [\mathbf{r}]) < \frac{\varepsilon}{2}$ for $\mathbf{q} \sim \mathbf{q}'$. By Lemma 1, we have $\|\mathbf{q} - \mathbf{r}'\|_p = \|\mathbf{q}' - \mathbf{r}\|_p$ and $d_p([\mathbf{p}], [\mathbf{r}]) \leq \|\mathbf{p} - \mathbf{r}'\|_p \leq \|\mathbf{p} - \mathbf{q}\|_p + \|\mathbf{q} - \mathbf{r}'\|_p < d_p([\mathbf{p}], [\mathbf{q}]) + d_p([\mathbf{q}], [\mathbf{r}]) + \varepsilon$. ■

6 Deriving a measure $\widehat{\mu}_{\ell^p}$ from d_p

We see that the minimum of d_p distance from $[D, S]$ to a consistent feature set determines a feature selection measure. We first define consistent points of $\mathcal{P}_{\mathbb{R}}^{+(\infty)}$.

Definition 8 $\mathbf{p} \in \mathcal{P}_{\mathbb{R}}^{+(\infty)}$ is consistent, if, and only if, there exists $\bar{j} : \mathbb{N} \rightarrow \mathbb{N}$ such that $\mathbf{p}(i, j) = 0$ if $j \neq \bar{j}(i)$.



$\mathcal{C}_{\mathbb{R}}^{+(\infty)}$ denotes the entire set of consistent \mathbf{p} , and we let $\mathcal{C}_{\mathbb{R}}^{+(n)} = \mathcal{C}_{\mathbb{R}}^{+(\infty)} \cap \mathcal{P}_{\mathbb{R}}^{+(n)}$, $\mathcal{C}_{\mathbb{R}}^{(n)} = \mathcal{C}_{\mathbb{R}}^{+(\infty)} \cap \mathcal{P}_{\mathbb{R}}^{(n)}$ and $\mathcal{C}_0^{(n)} = \mathcal{C}_{\mathbb{R}}^{+(\infty)} \cap \mathcal{P}_0^{(n)}$ for $n \in \mathbb{N} \cup \{\infty\}$.

Definition 9 For $\bar{j} : \mathbb{N} \rightarrow \mathbb{N}$ such that $\bar{j}(i) \in \operatorname{argmax}\{\mathbf{p}(i, j) \mid j \in \mathbb{N}\}$ and $\mathbf{p} \in \mathcal{P}_0^{(\infty)}$, we define:

For $1 \leq p < \infty$,

$$\widehat{\mu}_{\ell^p}(\mathbf{p}) = \sqrt[p]{\sum_{i \in \mathbb{N}} \left(\sum_{j \in \mathbb{N}} \mathbf{p}(i, j)^p - \mathbf{p}(i, \bar{j}(i))^p \right)};$$

For $p = \infty$,

$$\widehat{\mu}_{\ell^\infty}(\mathbf{p}) = \max \{ \mathbf{p}(i, j) \mid (i, j) \in \mathbb{N}^2, (i, j) \neq (i, \bar{j}(i)) \}.$$

Surprisingly, $\widehat{\mu}_{\ell^1}$ is identical to the Bayesian risk.

Proposition 3 If $\varphi : S \rightarrow U$ is coherent between $D|_S$ and D_p , $\mu_{\text{br}}^D(S) = \widehat{\mu}_{\ell^1}(\mathbf{p})$ holds.

Proof. We assume that $v_U : \Omega_S \rightarrow \mathbb{N}$ and $v_C : \Omega_C \rightarrow \mathbb{N}$ yield the coherence between $D|_S$ and D_p .

$$\begin{aligned} \mu_{\text{br}}^D(S) &= \sum_{\mathbf{x} \in \Omega_S} \left(\mathbf{p}_{D|_S}(S = \mathbf{x}) - \max_{c \in \Omega_C} \mathbf{p}_{D|_S}(S = \mathbf{x}, C = c) \right) \\ &= \sum_{\mathbf{x} \in \Omega_S} \left(\sum_{c \in \Omega_C} \mathbf{p}(v_U(\mathbf{x}), v_C(c)) - \max_{c \in \Omega_C} \mathbf{p}(v_U(\mathbf{x}), v_C(c)) \right) \\ &= \widehat{\mu}_{\ell^1}(\mathbf{p}). \blacksquare \end{aligned}$$

Theorem 6 For $\mathbf{p} \in \mathcal{P}_0^{(\infty)}$, the following holds.

1. $\widehat{\mu}_{\ell^1}(\mathbf{p}) = \frac{1}{2} \min\{d_1([\mathbf{p}], [\mathbf{q}]) \mid \mathbf{q} \in \mathcal{C}_0^{(\infty)}\}$.
2. For $1 < p \leq \infty$, $\widehat{\mu}_{\ell^p}(\mathbf{p}) = \inf\{d_p([\mathbf{p}], [\mathbf{q}]) \mid \mathbf{q} \in \mathcal{C}_0^{(\infty)}\} = \min\{d_p([\mathbf{p}], [\mathbf{q}]) \mid \mathbf{q} \in \mathcal{C}_{\mathbb{R}}^{+(\infty)} \wedge |\operatorname{supp}(\mathbf{q})| < \infty\}$.

Proof. Here, we prove only 1. For $\mathbf{q} \in \mathcal{C}_0^{(\infty)}$, we assume $j : \mathbb{N} \rightarrow \mathbb{N}$ such that $\mathbf{q}(i, j) = 0$ if $j \neq j(i)$. Also, we let $\bar{j} : \mathbb{N} \rightarrow \mathbb{N}$ satisfy $\bar{j}(i) \in \operatorname{argmax}\{\mathbf{p}(i, j) \mid j \in \mathbb{N}\}$. We first fix $i \in \mathbb{N}$.

$$\begin{aligned} &\sum_{j \in \mathbb{N}} |\mathbf{p}(i, j) - \mathbf{q}(i, j)| \\ &= |\mathbf{p}(i, j(i)) - \mathbf{q}(i, j(i))| + \sum_{j \in \mathbb{N} \setminus \{j(i)\}} \mathbf{p}(i, j) \\ &\geq |\mathbf{p}(i, j(i)) - \mathbf{q}(i, j(i))| + \sum_{j \in \mathbb{N}} \mathbf{p}(i, j) - \max_{j \in \mathbb{N}} \mathbf{p}(i, j). \end{aligned}$$

Then, we sum up the both sides across all $i \in \mathbb{N}$.

$$\begin{aligned} \|\mathbf{p} - \mathbf{q}\|_1 &\geq \sum_{i \in \mathbb{N}} |\mathbf{p}(i, j(i)) - \mathbf{q}(i, j(i))| \\ &\quad + \sum_{i \in \mathbb{N}} \left(\sum_{j \in \mathbb{N}} \mathbf{p}(i, j) - \max\{\mathbf{p}(i, j) \mid j \in \mathbb{N}\} \right) \\ &\geq \sum_{i \in \mathbb{N}} \mathbf{q}(i, j(i)) - \sum_{i \in \mathbb{N}} \mathbf{p}(i, j(i)) + \widehat{\mu}_{\ell^1}(\mathbf{p}) \geq 2\widehat{\mu}_{\ell^1}(\mathbf{p}). \end{aligned}$$

On the other hand, for $\mathbf{q} \in \mathcal{C}_0^{(\infty)}$ such that $\mathbf{q}(i, j) = \sum_{j \in \mathbb{N}} \mathbf{p}(i, j)$ if $j = \bar{j}(i)$ and 0 otherwise, we have

$$\|\mathbf{p} - \mathbf{q}\|_1 = 2 \left(1 - \sum_{i \in \mathbb{N}} \max\{\mathbf{p}(i, j) \mid j \in \mathbb{N}\} \right).$$

As $\mathbf{q}' \in \mathcal{C}_{\mathbb{R}}^{(\infty)}$ if $\mathbf{q}' \sim \mathbf{q}$, Lemma 1 implies the assertion. \blacksquare

Determinacy of $\widehat{\mu}_{\ell^p}$ immediately follows from Theorem 6. Also, it is easy to show monotonicity of $\widehat{\mu}_{\ell^p}$. Furthermore, $\widehat{\mu}_{\ell^p}$ turns out continuous under d_q for arbitrary $1 \leq p, q \leq \infty$.

7 Contrasting d_1 and d_p with $p > 1$

Proposition 3 asserts that $\widehat{\mu}_{\ell^1}$ is identical to the Bayesian risk, which is well known in the literature. On the other hand, $\widehat{\mu}_{\ell^p}$ with $p \in (1, \infty]$ is a novel measure unknown in the literature. Also, as a metric function, d_1 and d_p with $p \in (1, \infty]$ are significantly different as stated below without proof.

- $\mathcal{Q}_{\mathbb{R}}^{(\infty)}$ is complete under d_1 , while $\mathcal{Q}_{\mathbb{R}}^{+(\infty)}$ is complete under d_p for $p \in (1, \infty]$.
- $\mathcal{Q}_{\mathbb{R}}^{(n)}$ is not compact under d_1 , while $\mathcal{Q}_{\mathbb{R}}^{+(n)}$ is compact under d_p for $p \in (1, \infty]$.
- $\mathcal{Q}_0^{(n)}$ is dense in $\mathcal{Q}_{\mathbb{R}}^{(n)}$ under d_1 , while it is dense in $\mathcal{Q}_{\mathbb{R}}^{+(n)}$ under d_p for $p \in (1, \infty]$.

Moreover, in Section 8, we will see that $\widehat{\mu}_{\ell^1}$ and $\widehat{\mu}_{\ell^p}$ for $p \in (1, \infty]$ have different properties in terms of correlation between their measurements and classification accuracy. Although we cannot explain the reason for this difference, we imagine that it is related with the aforementioned difference between d_1 and d_p with $p \in (1, \infty]$ as a metric function.

8 Correlation between measurements by $\widehat{\mu}_{\ell^p}$ and classification accuracy

We evidently expect that there exists strong correlation between measurements by a good measure and classification accuracy. From this point of view, we ran experiments with $\widehat{\mu}_{\ell^p}$ for $p = 1, 2, \dots, 5$. In particular, we focus on comparison between $p = 1$ and $p > 1$: $\widehat{\mu}_{\ell^1}$ is identical to the Bayesian risk and is the feature selection measure that is used the most widely in the literature; On the other hand, the other measures are novel and introduced in this paper for the first time.

Table 1: Datasets

NAME	#FEAT.	#EXAM.	#CLASSES
ARRHYTHMIA	279	452	13
AUDIOLOGY	69	226	24
MFEAT-FACTOR	216	2000	10
MFEAT-FOURIER	76	2000	10
MFEAT-KARHUNEN	64	2000	10
MFEAT-PIXEL	240	2000	10
MFEAT-ZERNIKE	47	2000	10
MUSK	166	476	2
OPTIDIGITS	64	5620	10
SONAR	60	208	2
SPAMBASE	57	4601	2
SPECTROMETER	100	531	48

8.1 Datasets

Table 1 shows the datasets that we use in our experiments as well as their important attributes, namely, the numbers of features, examples and classes. All of the datasets are obtained from the UCI repository of machine learning databases [Blake and Merz, 1998].

8.2 Methods

The following are the steps of our experiments.

Sampling feature sets

For each dataset of Table 1, 60 feature sets are selected at random. The sizes of the selected feature sets also vary at random. In total, we obtain $60 \times 12 = 720$ pairs of a feature set and a dataset.

Localizing datasets

For each pair of a feature set and a dataset, we localize the dataset by eliminating all of the features that do not belong to the relevant feature set. Consequently, we obtain 720 localized datasets.

Measuring accuracy of classification

We apply three classifiers, namely, Naïve Bayes, C4.5 and SVM, to each of the localized datasets in the manner of 10-fold cross validation and record the averaged AUC-ROC scores, which we use as classification accuracy scores.

8.3 Results

Figure 1, 2 and 3 display scatter plots of the results of our experiments with Naïve Bayes, C4.5 and SVM, respectively. The x -axis of each chart represents measurements by a measure out of (a) $\widehat{\mu}_{\ell^1}$, (b) $\widehat{\mu}_{\ell^2}$, (c) $\widehat{\mu}_{\ell^3}$, (d) $\widehat{\mu}_{\ell^4}$ and (e) $\widehat{\mu}_{\ell^5}$, while the y -axis does the averaged AUC-ROC scores. From the charts, we have the impression that $\widehat{\mu}_{\ell^p}$ with $p > 1$ shows stronger negative correlation with classification accuracy than $\widehat{\mu}_{\ell^1}$. We look into this more closely in the next subsection.

8.4 Analysis on correlation

To compare $\widehat{\mu}_{\ell^1}$ and $\widehat{\mu}_{\ell^p}$ with $p > 1$ in terms of correlation with classification accuracy, we introduce a function $P_t(x)$ as an index. We let N_x be the number of plots whose $\widehat{\mu}_{\ell^p}$ distances fall within $[x, x + 0.01)$ and $N_{t,x}$ be the number of plots whose AUC-ROC scores exceed t and $\widehat{\mu}_{\ell^p}$ distances fall

within $[x, x + 0.01)$. Then, we define $P_t(x)$ by $P_t(x) = \frac{N_{t,x}}{N_x}$. Intuitively, $P_t(x)$ approximates the probability of the case that the classification accuracy exceeds t when the measurement by the relevant measure is x .

Figure 4, 5 and 6 show the curves of $P_t(x)$ for Naïve Bayes, C4.5 and SVM. The value of t varies in $\{0.95, 0.90, 0.85, 0.80\}$. Since $P_t(x) \geq P_{t'}(x)$ holds for $t < t'$, the curve for t is located above the curve for t' . We observe two properties from the charts.

- The curves for $\widehat{\mu}_{\ell^p}$ for $p > 1$ are akin to one another in shape, while the curves for $\widehat{\mu}_{\ell^1}$ appear significantly different from the other measures.
- The curves for $\widehat{\mu}_{\ell^p}$ with $p > 1$ exhibits clearer negative correlation between measurements by the measure and classification accuracy than the curves for $\widehat{\mu}_{\ell^1}$.

In fact, the table below presents the correlation coefficients between the measurements x and the values of $P_t(x)$. The presented values also support the observation stated above.

$t =$	0.95	0.90	0.85	0.80
NAÏVE BAYES				
$\widehat{\mu}_{\ell^1}$	-0.42	-0.51	-0.60	-0.52
$\widehat{\mu}_{\ell^2}$	-0.58	-0.79	-0.97	-0.85
$\widehat{\mu}_{\ell^3}$	-0.56	-0.69	-0.76	-0.79
$\widehat{\mu}_{\ell^4}$	-0.54	-0.71	-0.89	-0.89
$\widehat{\mu}_{\ell^5}$	-0.54	-0.73	-0.91	-0.89
C4.5				
$\widehat{\mu}_{\ell^1}$	-0.17	-0.41	-0.34	-0.38
$\widehat{\mu}_{\ell^2}$	-0.41	-0.47	-0.79	-0.83
$\widehat{\mu}_{\ell^3}$	-0.52	-0.29	-0.72	-0.69
$\widehat{\mu}_{\ell^4}$	-0.52	-0.38	-0.77	-0.81
$\widehat{\mu}_{\ell^5}$	-0.52	-0.44	-0.82	-0.84
SVM				
$\widehat{\mu}_{\ell^1}$	-0.52	-0.62	-0.55	-0.23
$\widehat{\mu}_{\ell^2}$	-0.52	-0.54	-0.65	-0.80
$\widehat{\mu}_{\ell^3}$	-0.55	-0.52	-0.63	-0.75
$\widehat{\mu}_{\ell^4}$	-0.52	-0.52	-0.59	-0.72
$\widehat{\mu}_{\ell^5}$	-0.52	-0.52	-0.56	-0.71

9 Conclusion

We have introduced a family of feature selection measures $\{\widehat{\mu}_{\ell^p} \mid p \in [1, \infty]\}$ that are derived from the distance functions of some metric spaces. Although $\widehat{\mu}_{\ell^1}$ turns out to be identical to the Bayesian risk, $\widehat{\mu}_{\ell^p}$ for $p \in (1, \infty)$ are novel. Through experiments, we have shown that $\widehat{\mu}_{\ell^p}$ with $p > 1$ exhibits clearer negative correlation between its measurements and classification accuracy than $\widehat{\mu}_{\ell^1}$ does.

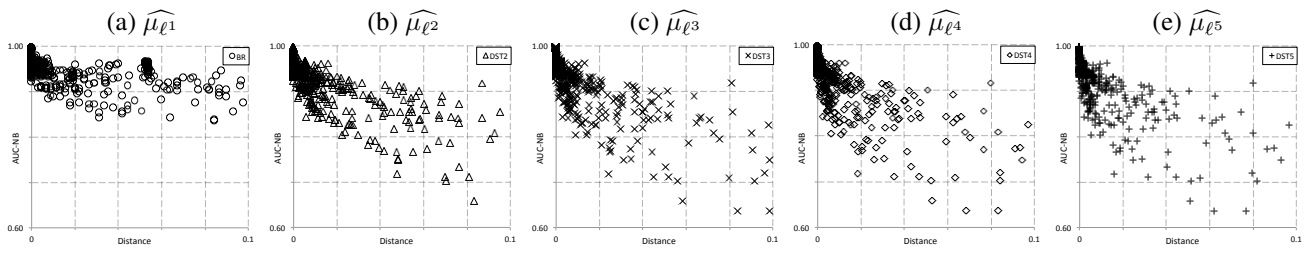


Figure 1: Scatter plots of the experimental results (Naïve Bayes)

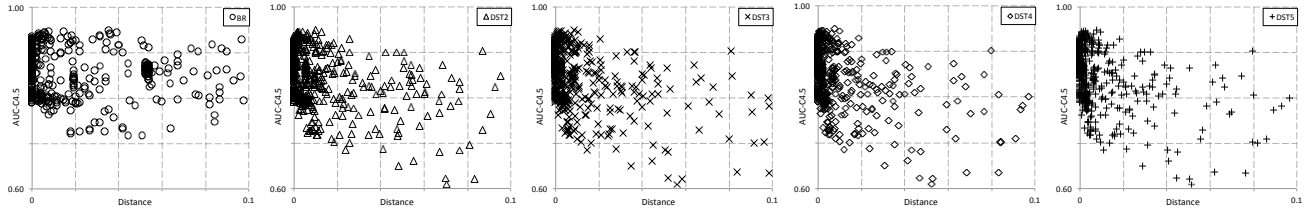


Figure 2: Scatter plots of the experimental results (C4.5)

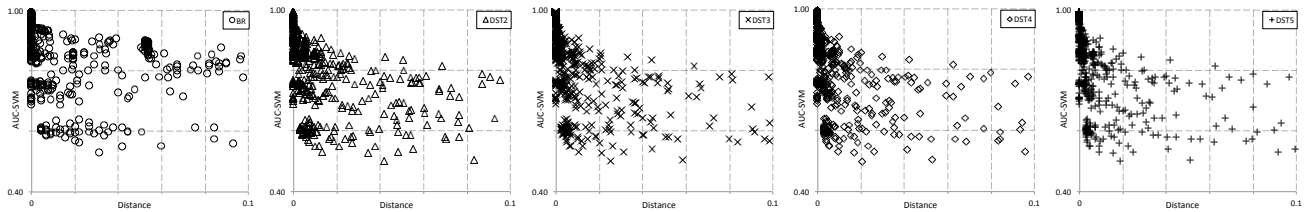


Figure 3: Scatter plots of the experimental results (SVM)

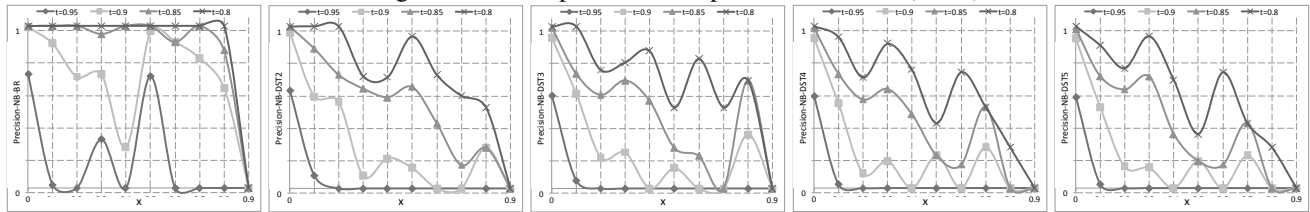


Figure 4: The curves of $P_t(x)$ for $t \in \{0.95, 0.90, 0.85, 0.80\}$ (Naïve Bayes)

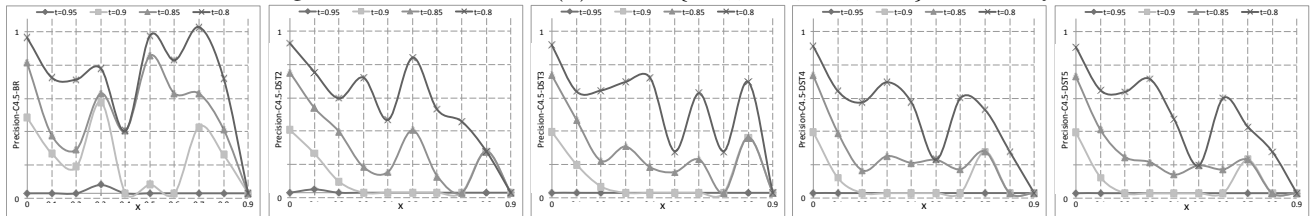


Figure 5: The curves of $P_t(x)$ for $t \in \{0.95, 0.90, 0.85, 0.80\}$ (C4.5)

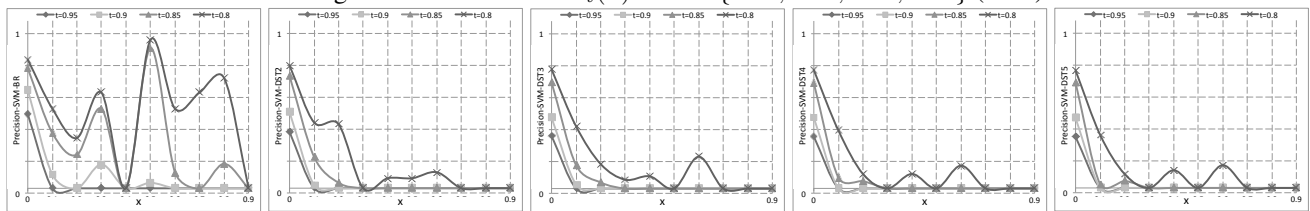


Figure 6: The curves of $P_t(x)$ for $t \in \{0.95, 0.90, 0.85, 0.80\}$ (SVM)

References

- [Almuallim and Dietterich, 1994] H. Almuallim and T. G. Dietterich. Learning boolean concepts in the presence of many irrelevant features. *Artificial Intelligence*, 69(1 - 2), 1994.
- [Arauzo-Azofra *et al.*, 2008] A. Arauzo-Azofra, J. M. Benítez, and J. L. Castro. Consistency measures for feature selection. *Journal of Intelligent Information Systems*, 30(3):273–292, November 2008.
- [Blake and Merz, 1998] C. S. Blake and C. J. Merz. UCI repository of machine learning databases. Technical report, University of California, Irvine, 1998.
- [Bolón-Canedo *et al.*, 2011] V. Bolón-Canedo, S. Seth, and N. Sánchez-Marín. Statistical dependence measure for feature selection in microarray datasets. In *ESANN 2011*, pages 23–28, 2011.
- [Foithong *et al.*, 2012] S. Foithong, O. Pinnigern, and B. Attachoo. Feature subset selection wrapper based on mutual information and rough sets. *Expert Systems with Applications*, 39(1):574–584, 2012.
- [Liu *et al.*, 1998] H. Liu, H. Motoda, and M. Dash. A monotonic measure for optimal feature selection. In *Proceedings of European Conference on Machine Learning*, 1998.
- [Mengle and Goharian, 2009] S.S.R. Mengle and N. Goharian. Ambiguity measure feature-selection algorithm. *Journal of the American Society for Information Science and Technology*, 60(5):1037–1050, 2009.
- [Molina *et al.*, 2002] L. Molina, L. Belanche, and A. Nebot. Feature selection algorithms: A survey and experimental evaluation. In *Proceedings of IEEE International Conference on Data Mining*, pages 306–313, 2002.
- [Pawlak, 1991] Z. Pawlak. *Rough Sets, Theoretical aspects of reasoning about data*. Kluwer Academic Publishers, 1991.
- [Peng *et al.*, 2005] H. Peng, F. Long, and C. Ding. Feature selection based on mutual information: Criteria of max-dependency, max-relevance and min-redundancy. *IEEE Transaction on Pattern Analysis and Machine Intelligence*, 27(8), August 2005.
- [Shin and Xu, 2009] K. Shin and X.M. Xu. Consistency-based feature selection. In *13th International Conference on Knowledge-Based and Intelligent Information & Engineering System*, 2009.
- [Shin *et al.*, 2011] K. Shin, D. Fernandes, and S. Miyazaki. Consistency measures for feature selection: A formal definition, relative sensitivity comparison, and a fast algorithm. In *22nd International Joint Conference on Artificial Intelligence*, pages 1491–1497, 2011.
- [Suzuki and Sugiyama, 2013] T. Suzuki and M. Sugiyama. Sufficient dimension reduction via squared-loss mutual information estimation. *Neural Computation*, 25(3):725–758, 2013.
- [Wang, 2015] L. Wang. Feature selection algorithm based on conditional dynamic mutual information. *International Journal on Smart Sensing and Intelligent Systems*, 8(1), 2015.
- [Zhao and Liu, 2007] Z. Zhao and H. Liu. Searching for interacting features. In *Proceedings of International Joint Conference on Artificial Intelligence*, pages 1156 – 1161, 2007.