

# Efficient Case Based Feature Construction

Ingo Mierswa and Michael Wurst

Artificial Intelligence Unit, Department of Computer Science,  
University of Dortmund, Germany  
{mierswa, wurst}@ls8.cs.uni-dortmund.de

**Abstract.** Feature construction is essential for solving many complex learning problems. Unfortunately, the construction of features usually implies searching a very large space of possibilities and is often computationally demanding. In this work, we propose a case based approach to feature construction. Learning tasks are stored together with a corresponding set of constructed features in a case base and can be retrieved to speed up feature construction for new tasks. The essential part of our method is a new representation model for learning tasks and a corresponding distance measure. Learning tasks are compared using relevance weights on a common set of base features only. Therefore, the case base can be built and queried very efficiently. In this respect, our approach is unique and enables us to apply case based feature construction not only on a large scale, but also in distributed learning scenarios in which communication costs play an important role. We derive a distance measure for heterogeneous learning tasks by stating a set of necessary conditions. Although the conditions are quite basic, they constraint the set of applicable methods to a surprisingly small number.

## 1 Introduction

Many inductive learning problems cannot be solved accurately by using the original feature space. This is due to the fact that standard learning algorithms cannot represent complex relationships as induced for example by trigonometric functions. For example, if only base features  $X_1$  and  $X_2$  are given but the target function depends highly on  $X_c = \sin(X_1 \cdot X_2)$ , the construction of the feature  $X_c$  would ease learning – or is necessary to enable any reasonable predictions at all [1, 2, 3]. Unfortunately, feature construction is a computationally very demanding task often requiring to search a very large space of possibilities [4, 5]. In this work we consider a scenario in which several learners face the problem of feature construction on different learning problems. The idea is to transfer constructed features between similar learning tasks to speed up the generation in such cases in which a successful feature has already been generated by another feature constructor. Such approaches are usually referred to as Meta Learning [6].

Meta Learning was applied to a large variety of problems and on different conceptual levels. The importance of the representation bias, which is closely related to feature construction, was recognized since the early days of Meta

Learning research [7, 8]. The key to many Meta Learning methods is the definition of similarity between different learning tasks [9, 10]. In this work we propose a Meta Learning scheme that compares two learning tasks using only relevance weights assigned to a set of base features by the individual learners.

This is motivated by a set of constraints found in many distributed Meta Learning scenarios. Firstly, the retrieval of similar learning tasks and relevant features usually has to be very efficient, especially for interactive applications. This also means that methods should enable a best effort strategy, such that the user can stop the retrieval process at any point and get the current best result. Secondly, the system should scale well with an increasing number of learning tasks. Also, it has to deal with a large variety of heterogeneous learning tasks, as we cannot make any strict assumptions on the individual problems. Finally, as many Meta Learning systems are distributed, communication cost should be as low as possible. As a consequence, methods that are based on exchanging examples or many feature vectors are not applicable.

## 2 Basic Concepts

Before we state the conditions which must be met by any method comparing learning tasks using feature weights only, we first introduce some basic definitions. Let  $T$  be the set of all learning tasks, a single task is denoted by  $t_i$ . Let  $X_i$  be a vector of numerical random variables for task  $t_i$  and  $Y_i$  another random variable, the target variable. These obey a fixed but unknown probability distribution  $Pr(X_i, Y_i)$ . The components of  $X_i$  are called *features*  $X_{ik}$ . The objective of every *learning task*  $t_i$  is to find a function  $h_i(X_i)$  which predicts the value of  $Y_i$ . We assume that each set of features  $X_i$  is partitioned in a set of *base features*  $X_B$  which are common for all learning tasks  $t_i \in T$  and a set of *constructed features*  $X_i \setminus X_B$ .

We now introduce a very simple model of feature relevance and interaction. The feature  $X_{ik}$  is assumed to be irrelevant for a learning task  $t_i$  if it does not improve the classification accuracy:

**Definition 1.** A feature  $X_{ik}$  is called *IRRELEVANT* for a learning task  $t_i$  iff  $X_{ik}$  is not correlated to the target feature  $Y_i$ , i. e. if  $Pr(Y_i|X_{ik}) = Pr(Y_i)$ .

The set of all irrelevant features for a learning task  $t_i$  is denoted by  $IF_i$ .

Two features  $X_{ik}$  and  $X_{il}$  are alternative for a learning task  $t_i$ , denoted by  $X_{ik} \sim X_{il}$  if they can be replaced by each other without affecting the classification accuracy. For linear learning schemes this leads to the linear correlation of two features:

**Definition 2.** Two features  $X_{ik}$  and  $X_{il}$  are called *ALTERNATIVE* for a learning task  $t_i$  (written as  $X_{ik} \sim X_{il}$ ) iff  $X_{il} = a + b \cdot X_{ik}$  with  $b > 0$ .

This is a very limited definition of alternative features. However, we will show that most weighting algorithms are already ruled out by conditions based on this simple definition.

### 3 Comparing Learning Tasks Efficiently

The objective of our work is to speed up feature construction and improve prediction accuracy by building a case base containing pairs of learning tasks and corresponding sets of constructed features. We assume that a learning task  $t_i$  is completely represented by a feature weight vector  $w_i$ . The vector  $w_i$  is calculated from the base features  $X_B$  only. This representation of learning tasks is motivated by the idea that a given learning scheme approximate similar constructed features by a set of base features in a similar way, e. g. if the constructed feature “ $\sin(X_{ik} \cdot X_{il})$ ” is highly relevant the features  $X_{ik}$  and  $X_{il}$  are relevant as well.

Our approach works as follows: for a given learning task  $t_i$  we first calculate the relevance of all base features  $X_B$ . We then use a distance function  $d(t_i, t_j)$  to find the  $k$  most similar learning tasks. Finally, we create a set of constructed features as union of the constructed features associated with these tasks.

This set is then evaluated on the learning task  $t_i$ . If the performance gain is sufficiently high (above a given fixed threshold) we store task  $t_i$  in the case base as additional case. Otherwise, the constructed features are only used as initialization for a classical feature construction that is performed locally. If this leads to a sufficiently high increase in performance, the task  $t_i$  is also stored to the case base along with the locally generated features.

While feature weighting and feature construction are well studied tasks, the core of our algorithm is the calculation of  $d$  using only the relevance values of the base features  $X_B$ . In a first step, we define a set of conditions which must be met by feature weighting schemes. In a second step, a set of conditions for learning task distance is defined which makes use of the weighting conditions.

**Weighting Conditions.** Let  $w$  be a WEIGHTING FUNCTION  $w : X_B \rightarrow \mathbb{R}$ . Then the following must hold:

(W1)  $w(X_{ik}) = 0$  if  $X_{ik} \in X_B$  is irrelevant

(W2)  $F_i \subseteq X_B$  is a set of alternative features. Then

$$\forall S \subset F_i, S \neq \emptyset : \sum_{X_{ik} \in S} w(X_{ik}) = \sum_{X_{ik} \in F_i} w(X_{ik}) = \hat{w}$$

(W3)  $w(X_{ik}) = w(X_{il})$  if  $X_{ik} \sim X_{il}$

(W4) Let  $AF$  be a set of features where

$$\forall X_{ik} \in AF : (X_{ik} \in IF_i \vee \exists X_{il} \in X_B : X_{ik} \sim X_{il}).$$

Then

$$\forall X_{il} \in X_B : \exists X_{ik} \in AF : X_{il} \sim X_{ik} \wedge w'(X_{il}) = w(X_{il})$$

where  $w'$  is a weighting function for  $X'_B = X_B \cup AF$ .

These conditions state that irrelevant features have weight 0 and that the sum of weights of alternative features must be constant independently of the

actual number of alternative features used. Together with the last conditions this guarantees that a set of alternative features is not more important than a single feature of this set. Obviously, this is a desired property of a weighting function used for the comparison of learning tasks. In the following we assume that for a modified space of base features  $X'_B$  the function  $w'$  denotes the weighting function for  $X'_B$  according to the definition in (W4).

Additionally, we can define a set of conditions which must be met by distance measures for learning tasks which are based on feature weights only:

**Distance Conditions.** A DISTANCE MEASURE  $d$  for learning tasks is a mapping  $d: T \times T \rightarrow \mathbb{R}^+$  which should fulfill at least the following conditions:

- (D1)  $d(t_1, t_2) = 0 \Leftrightarrow t_1 = t_2$
- (D2)  $d(t_1, t_2) = d(t_2, t_1)$
- (D3)  $d(t_1, t_3) \leq d(t_1, t_2) + d(t_2, t_3)$
- (D4)  $d(t_1, t_2) = d(t'_1, t'_2)$  if  $X'_B = X_B \cup IF$  and  $IF \subseteq IF_1 \cap IF_2$
- (D5)  $d(t_1, t_2) = d(t'_1, t'_2)$  if  $X'_B = X_B \cup AF$  and  $\forall X_k \in AF: \exists X_l \in X_B: X_k \sim X_l$

(D1)–(D3) represent the conditions for a metric. These conditions are required for efficient case retrieval and indexing. (D4) states that irrelevant features should not have an influence on the distance. Finally, (D5) states that adding alternative features should not have an influence on distance.

## 4 Negative Results

In this section we will show that many feature weighting approaches do not fulfill the conditions (W1)–(W4). Furthermore, one of the most popular distance measures, the euclidian distance, cannot be used as a learning task distance measure introduced above.

**Lemma 1.** Any feature selection method does not fulfill the conditions (W1)–(W4).

*Proof.* For a feature selection method, weights are always binary, i. e.  $w(X_{ik}) \in \{0, 1\}$ . We assume a learning task  $t_i$  with no alternative features and  $X'_B = X_B \cup \{X_{ik}\}$  with  $\exists X_{il} \in X_B: X_{il} \sim X_{ik}$ , then either  $w'(X_{il}) = w'(X_{ik}) = w(X_{il}) = 1$ , leading to a contradiction with (W2), or  $w'(X_{il}) \neq w'(X_{ik})$  leading to a contradiction with (W3).  $\square$

**Lemma 2.** Any feature weighting method for which  $w(X_{ik})$  is calculated independently of  $X_B \setminus X_{ik}$  does not fulfill the conditions (W1)–(W4).

*Proof.* We assume a learning task  $t_i$  with no alternative features and  $X'_B = X_B \cup \{X_{ik}\}$  with  $\exists X_{il} \in X_B: X_{il} \sim X_{ik}$ . If  $w$  is independent of  $X_B \setminus X_{ik}$  adding  $X_{ik}$  would not change the weight  $w'(X_{il})$  in the new feature space  $X'_B$ . From (W3) follows that  $w'(X_{ik}) = w'(X_{il}) = w(X_{il})$  which is a violation of (W2).  $\square$

Lemma 2 essentially covers all feature weighting methods that treat features independently such as information gain [11] or Relief [12]. The next theorem states that the euclidian distance cannot be used as a distance measure based on feature weights.

**Theorem 3.** *Euclidean distance does not fulfill the conditions (D1)–(D5).*

*Proof.* We give a counterexample. We assume that a weighting function  $w$  is given which fulfills the conditions (W1)–(W4). Further assume that learning tasks  $t_i, t_j$  are given with no alternative features. We add an alternative feature  $X_{ik}$  to  $X_B$  and get  $X'_B = X_B \cup \{X_{ik}\}$  with  $\exists X_{il} \in X_B : X_{il} \sim X_{ik}$ . We infer from conditions (W2) and (W3) that

$$w'(X_{ik}) = w'(X_{il}) = \frac{w(X_{il})}{2} \quad \text{and} \quad w'(X_{jk}) = w'(X_{jl}) = \frac{w(X_{jl})}{2}$$

and from condition (W4) that

$$\forall p \neq k : w'(X_{ip}) = w(X_{ip}) \quad \text{and} \quad \forall p \neq k : w'(X_{jp}) = w(X_{jp}).$$

In this case the following holds for the euclidian distance

$$\begin{aligned} d(t'_i, t'_j) &= \sqrt{S + 2(w'(X_{ik}) - w'(X_{jk}))^2} = \sqrt{S + 2\left(\frac{w(X_{ik})}{2} - \frac{w(X_{jk})}{2}\right)^2} \\ &= \sqrt{S + \frac{1}{2}(w(X_{ik}) - w(X_{jk}))^2} \neq \sqrt{S + (w(X_{ik}) - w(X_{jk}))^2} = d(t_i, t_j) \end{aligned}$$

with

$$S = \sum_{p=1, p \neq k}^{|X_B|} (w'(X_{ip}) - w'(X_{jp}))^2 = \sum_{p=1, p \neq k}^{|X_B|} (w(X_{ip}) - w(X_{jp}))^2. \quad \square$$

## 5 Positive Results

In this section we will prove that a combination of feature weights delivered by a linear Support Vector Machine (SVM) with the Manhattan distance obeys the proposed conditions. Support Vector Machines are based on the work of Vapnik in statistical learning theory [13]. They aim to minimize the regularized risk  $R_{reg}[f]$  of a learned function  $f$  which is the weighted sum of the empirical risk  $R_{emp}[f]$  and a complexity term  $\|w\|^2$ :

$$R_{reg}[f] = R_{emp}[f] + \lambda \|w\|^2.$$

The result is a linear decision function  $y = \text{sgn}(w \cdot x + b)$  with a minimal length of  $w$ . The vector  $w$  is the normal vector of an optimal hyperplane with a maximal margin to both classes. One of the strengths of SVMs is the use of kernel functions to extend the feature space and allow linear decision boundaries after efficient

nonlinear transformations of the input [14]. Since our goal is the construction of (nonlinear) features during preprocessing we can just use the most simple kernel function which is the dot product. In this case the components of the vector  $w$  can be interpreted as weights for all features.

**Theorem 4.** *The feature weight calculation of SVMs with linear kernel function meets the conditions (W1)–(W4).*

*Proof.* Since these conditions can be proved for a single learning task  $t_i$  we write  $X_k$  and  $w_k$  as a shortcut for  $X_{ik}$  and  $w(X_{ik})$ .

(W1) *Sketch* We assume that the SVM finds an optimal hyperplane. The algorithm tries to minimize both the length of  $w$  and the empirical error. This naturally corresponds to a maximum margin hyperplane where the weights of irrelevant features are 0 if enough data points are given.

(W2) SVMs find the optimal hyperplane by minimizing the weight vector  $w$ . Using the optimal classification hyperplane with weight vector  $w$  can be written as  $y = \text{sgn}(w_1x_1 + \dots + w_ix_i + \dots + w_mx_m + b)$ . We will show that this vector cannot be changed by adding the same feature more than one time. We assume that all alternative features can be transformed into identical features by normalizing the data. Adding  $k - 1$  alternative features will result in

$$y = \text{sgn} \left( \dots + \underbrace{(w_i^1 + \dots + w_i^k)}_{\text{alternative features}} x_i + \dots + b \right).$$

However, the optimal hyperplane will remain the same and does not depend on the number of alternative attributes. This means that the other values  $w_j$  will not be changed. This leads to  $w_i = \sum_{l=1}^k w_i^l$  which proves condition (W2).

(W3) The SVM optimization minimizes the length of the weight vector  $w$ . This can be written as

$$w_1^2 + \dots + w_i^2 + \dots + w_m^2 \stackrel{!}{=} \min.$$

We replace  $w_i$  using condition (W2):

$$w_1^2 + \dots + \left( \hat{w} - \sum_{j \neq i} w_j \right)^2 + \dots + w_m^2 \stackrel{!}{=} \min.$$

In order to find the minimum we have to partially differentiate the last equation for all weights  $w_k$ :

$$\begin{aligned} \frac{\partial}{\partial w_k} \left( \dots + \left( \hat{w} - \sum_{j \neq i} w_j \right)^2 + w_k^2 + \dots \right) &= 0 \\ \Leftrightarrow 2w_k - 2 \left( \hat{w} - \sum_{j \neq i} w_j \right) &= 0 \quad \Leftrightarrow \quad w_k + \sum_{j \neq i} w_j = \hat{w} \end{aligned}$$

The sum on the left side contains another  $w_k$ . This leads to a system of linear equations of the form  $\dots + 0 \cdot w_i + \dots + 2 \cdot w_k + \dots = \hat{w}$ . Solving this system of equations leads to  $w_p = w_q$  (condition (W3)).

(W4) *Sketch* We again assume that a SVM finds an optimal hyperplane given enough data points. Since condition (W1) holds adding an irrelevant feature would not change the hyperplane and thus the weighting vector  $w$  for the base features will remain. The proofs of conditions (W2) and (W3) state that the optimal hyperplane is not affected by alternative features as well.  $\square$

In order to calculate the distance of learning tasks based only on a set of base feature weights we still need a distance measure that met the conditions (D1)–(D5).

**Theorem 5.** *Manhattan distance does fulfill the conditions (D1)–(D5).*

*Proof.* The conditions (D1)–(D3) are fulfilled due to basic properties of the manhattan distance. Therefore, we only give proofs for conditions (D4) and (D5).

(D4) We follow from the definition of the manhattan distance that

$$\begin{aligned} d(t'_i, t'_j) &= \sum_{X_{ip}, X_{jp} \in X_B} |w'_i(X_{ip}) - w'_j(X_{jp})| + \underbrace{\sum_{X_{iq}, X_{jq} \in IF} |w'_i(X_{iq}) - w'_j(X_{jq})|}_0 \\ &= d(t_i, t_j) \end{aligned}$$

from (W4).

(D5) *Sketch* We show the case for adding  $k$  features with  $\forall X_{ik} : X_{ik} \sim X_{il}$  for a fixed  $X_{il} \in X_B$ :

$$\begin{aligned} d(t'_i, t'_i) &= \sum_{p=1, p \neq k}^{|X_B|} |w'_i(X_{ip}) - w'_j(X_{jp})| + (k+1) \cdot |w'_i(X_{ik}) - w'_j(X_{jk})| \\ &= \sum_{p=1, p \neq k}^{|X_B|} |w_i(X_{ip}) - w_j(X_{jp})| + |w_i(X_{ik}) - w_j(X_{jk})| = d(t_i, t_j) \end{aligned}$$

from (W4) and (W2).  $\square$

Therefore, we conclude that SVM feature weights in combination with manhattan distance fulfill the necessary constraints for an efficient learning task distance measure based on feature weights.

## 6 Conclusion and Outlook

We presented a Meta Learning approach to feature construction that compares tasks using relevance weights on a common set of base features only. After stating some very basic conditions for such a distance measure, we have shown that a SVM as base feature weighting algorithm and the manhattan distance fulfill

these conditions, while several other popular feature weighting methods and distance measures do not. In [15] we have presented experimental results indicating that our method can speed up feature construction considerably. Our approach is therefore highly relevant for practical problems involving feature construction. Some limitations of the work presented here are the following. Firstly, our definition for alternative or exchangeable features is rather simple and should be generalized to a weaker concept as e. g. highly correlated features. Also, complex interactions between features are not covered by our conditions. However, it is very interesting that the conditions stated in this work are already sufficient to rule out large sets of feature weighting methods and distance measures. Finally, the assumption of estimating the distance of constructed features by the distance of base features is well motivated, though it would be interesting to analyze this relationship analytically to get a better estimation in which cases our approach can be successfully applied.

## References

1. Blum, A.L., Langley, P.: Selection of relevant features and examples in machine learning. *Artificial Intelligence* (1997) 245–271
2. Dash, M., Liu, H.: Feature selection for classification. *International Journal of Intelligent Data Analysis* **1** (1997) 131–156
3. Koller, D., Sahami, M.: Toward optimal feature selection. In: *Proc. of the ICML*. (1996) 129–134
4. Mierswa, I., Morik, K.: Automatic feature extraction for classifying audio data. *Machine Learning Journal* **58** (2005) 127–149
5. Wolpert, D., Macready, W.: No free lunch theorems for optimisation. *IEEE Trans. on Evolutionary Computation* **1** (1997) 67–82
6. Vilalta, R., Drissi, Y.: A perspective view and survey of meta-learning. *Artificial Intelligence Review* **18** (2002) 77–95
7. Baxter, J.: Learning internal representations. In: *Proc. of the eighth annual conference on Computational learning theory '95*, ACM Press (1995) 311–320
8. Baxter, J.: A model of inductive bias learning. *Journal of Artificial Intelligence Research* **12** (2000) 149–198
9. Ben-David, S., Schuller, R.: Exploiting task relatedness for multiple task learning. In: *Proc. of the Sixteenth Annual Conference on Learning Theory 2003*. (2003)
10. Thrun, S., O'Sullivan, J.: Discovering structure in multiple learning tasks: The TC algorithm. In Saitta, L., ed.: *Proc. of the ICML, San Mateo, CA, Morgan Kaufmann* (1996)
11. Quinlan, R.: Induction of decision trees. *Machine Learning* **1** (1986) 81–106
12. Kira, K., Rendell, I.A.: The feature selection problem: Traditional methods and a new algorithm. In: *10th National Conference on Artificial Intelligence*, MIT Press (1992) 129–134
13. Vapnik, V.N.: *The Nature of Statistical Learning Theory*. Springer, New York (1995)
14. Schölkopf, B., Smola, A.J.: *Learning with Kernels – Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press (2002)
15. Mierswa, I., Wurst, M.: Efficient feature construction by meta learning – guiding the search in meta hypothesis space. In: *Proc. of the ICML Workshop on Meta Learning*. (2005)